



# Optimering över finita diskreta mått

och dess tillämpning på bayesianska neurala nätverk

## Optimization of Finite Discrete Measures

*Kandidatarbete inom civilingenjörsutbildningen vid Chalmers Tekniska Högskola*

Filip Berglund

Karl Egeland

Samuel Martinsson

Hampus Olsson

Joel Wall

Johan Ödesjö



# Optimering över finita diskreta mått

och dess tillämpning på bayesianska neurala nätverk

*Kandidatarbete i matematik inom civilingenjörsprogrammet Maskinteknik vid Chalmers*  
Hampus Olsson Karl Egeland

*Kandidatarbete i matematik inom civilingenjörsprogrammet Teknisk matematik vid Chalmers*

Filip Berglund Samuel Martinsson  
Joel Wall Johan Ödesjö

Handledare: Sergei Zuyev

Institutionen för Matematiska vetenskaper  
CHALMERS TEKNISKA HÖGSKOLA  
GÖTEBORGS UNIVERSITET  
Göteborg, Sverige 2023



# Förord

Detta är ett kandidatarbete vid institutionen för Matematiska Vetenskaper vid Chalmers Tekniska Högskola under vårterminen 2023. Vi vill inledningsvis tacka handledare Sergei Zuyev för kontinuerligt stöd och vägledning under arbetet. Vi vill även tacka Hans Malmström för hans respons på ett tidigt utkast av rapporten och Marco Schirone för hans hjälp med källor och källhänvisning.

Arbetet har skrivits av sex författare och vad som skrivits av vem hittas i tabell (1). Det finns utöver detta en loggbok där det varje vecka antecknats vad alla medlemmar gjort under veckan och hur mycket tid de lagt. Denna loggbok är uppdelad i två delar, en del med hur mycket tid medlemmarna lagt med en kort förklaring av vad de gjort samt ett längre dokument där det står mer noggrant vad medlemmarna gjort. Samtliga författare är överens om att alla bidragit till projektet och att alla har varit delaktiga i utvecklandet av Python-koden. Vilka delar som varje författare kan tillskrivas framgår från tidigare nämnda loggbok.

Under arbetets gång har dessutom versionshanteringsystemet Git använts och alla implementationer i Python finns att tillgå i ett delat GitHub-projekt. Detta kan hämtas från <https://github.com/Miralleyan/Kandidatarbete>.

Tabell 1: Bidragsrapport - författare för avsnitt

Del i rapport	Författare	Kommentar
Populärvetenskaplig presentation Abstract Sammandrag	Karl, Samuel Johan Hampus	Huvudförfattare: Karl
<b>Inledning</b> Syfte Avgränsningar	Karl, Samuel Johan Karl	
<b>Bakgrund</b> Bayesiansk statistik Neurala nätverk Bayesianska neurala nätverk Olika former av regression PyTorch Automatisk differentiering Mått Differentiering av mått Kärndensitetsuppskattning	Samuel, Johan Samuel, Johan, Joel Samuel Samuel Filip, Joel Johan, Samuel Joel, Filip Filip Johan	
<b>Metod</b> Parametriska modeller Icke-parametriska modeller Målfunktioner	Johan Hampus, Filip, Joel Joel, Filip, Karl, Johan Hampus, Joel, Samuel	Texten innan sektion 3.1 börjar 3.1.2: Hampus 3.2.1: Filip, 3.2.3: Karl och Johan 3.3.1-3.3.2: Samuel, 3.3.3: Joel, 3.3.4: Hampus
<b>Resultat</b> Jämförelse mellan parametriska och icke-parametriska metoder Resultat från tester	Samuel, Johan Hampus	Huvudförfattare: Samuel
<b>Diskussion</b> Jämförelse Val för KDE Tidskomplexitet för större modeller med OFDM Vidare forskning Samhälleliga och etiska aspekter Slutsats	Samuel Johan Joel Samuel Karl Joel	
<b>Appendix</b> Resultattabeller Degenererad summa av väntevärden Linjärkombination av stokastiska variabler Yttre strafffunktion Lagrangemultiplikator-metoden Icke-parametriska diskreta sannolikhetsmått med softmax Poängfunktioner och poängregler	Johan, Samuel Filip Filip Hampus Samuel Joel, Johan Filip	

## Populärvetenskaplig presentation

Under de senaste åren har det blivit allt mer vanligt att använda sig av applikationer som bygger på artificiell intelligens. För att skapa artificiell intelligens används olika former av maskininlärning. Maskininlärning kan beskrivas som metoder för att få en dator att lära sig ett antal regler från en mängd data, utan explicit instruktion från användaren. Med hjälp av dessa regler ska exempelvis en uppgift kunna lösas eller förutsägelser ges utifrån ny indata.

Ett område där maskininlärningsprocesser har implementerats och som växt mycket de senaste åren är i rekrytering till företag vilket är en etisk aspekt som tas upp i texten. Här ges förutsägelser på hur väl en viss individ hade presterat i rollen den söker, utifrån data som beskriver hur väl individer med liknande egenskaper har presterat i en liknande roll tidigare. Det blir i detta fall tydligt hur verktyg som bygger på maskininlärning påverkar människor och vilka samhällsliga konsekvenser de kan ge. Det är i dessa tillämpningar, som kan påverka människor och samhället i stort, viktigt att verktygen kan ge goda förutsägelser för ny indata.

I maskininlärningsprocesserna användes ofta så kallade artificiella neurala nätverk. Traditionella artificiella neurala nätverk består av noder eller neuroner som sammankopplas med hjälp av länkar mellan de olika noderna. Detta är inspirerat av hur den mänskliga hjärnan fungerar. Noderna struktureras i ett antal lager där antalet noder i varje lager kan variera. Antalet lager kan också variera men i det enklaste fallet finns ett lager för indata, ett lager för utdata samt ett mellanliggande lager som kopplar indata till utdata. Varje länk går då från varje nod i ett lager till varje nod i ett närliggande lager. Länkarna har också en associerad vikt vilken beskriver hur stark en viss sammankoppling är mellan två noder. Dessa vikter kan sedan modifieras och det är just det som händer när det neurala nätverket lär sig regler för hur en viss indata bör ge en viss utdata.

Ett problem som vanligen uppstår när det neurala nätverket lär sig regler utifrån träningsdata är överanpassade nätverk. Överanpassning (eng. overfitting) innebär att det neurala nätverket är allt för väl anpassat på den data som datorn lärt sig regler utifrån. När det neurala nätverket sedan ska ge förutsägelser på ny indata blir förutsägelseerna bristfälliga. Detta kan antingen bero på att den data som nätverket tränat på ej är representativ för alla möjliga värden på indata eller att mängden träningsdata helt enkelt är för liten. Det är i denna typ av fall som så kallade bayesianska neurala nätverk kan vara användbara.

Medan traditionella neurala nätverk har länkar bestående av numeriska vikter har istället bayesianska neurala nätverk länkar där ett tal slumpas varje gång enligt en given regel. Varje länk är alltså inte ett enstaka värde utan istället en slumpvariabel. Modeller som bygger på bayesianska neurala nätverk kallas stokastiska i motsats till de traditionella, deterministiska modellerna. Eftersom varje länk i nätverket har en sannolikhetsfördelning finns det en inbyggd variation i nätverket. Detta kan göra att nätverket inte blir överanpassat till träningsdata och genom att studera länkarnas sannolikhetsfördelningar ges en bild av hur säkert nätverket är, samt vilka andra möjliga värden modellens resultat kan anta.

I detta projekt utforskas bayesianska neurala nätverk, främst genom att undersöka metoder för att hitta den lämpligaste sannolikhetsfördelningen för en viss datamängd. Metoder som syftar till att hitta de lämpligaste lösningarna till ett specifikt problem kallas optimeringsalgoritmer. Vi studerar optimeringsalgoritmer ämnade att hitta en sannolikhetsfördelning som kan ha gett upphov till observerad data. Specifikt testas en metod som tidigare i stor utsträckning inte använts för just neurala nätverk. Metoden testas på olika typer av data och jämförs sedan med andra metoder med liknande syfte.

## Abstract

Artificial intelligence and neural networks is currently a highly debated topic. In this paper, we will investigate some methods for optimization on measures which could be used to implement so called Bayesian neural networks. These can potentially be used to avoid overfitting and to model the uncertainty in the output of neural networks.

Using Python and PyTorch we have implemented three models designed for Bayesian regression: Steepest descent algorithm for fixed total mass problem, polynomial regression and linear regression. These are three distinct ways to model data, and we will investigate how effective and precise they are given the same data.

Linear combination of stochastic variables describes a measure as a sum of stochastic variables multiplied with an arbitrary function. By assuming gaussian distributions for the stochastic variables some assumptions can be made, which make calculations simpler. During testing, this model showed high precision but was the slowest to converge to an answer.

The polynomial neural network is similar to the linear combination in that the functions are terms in a polynomial. However, the method to solve it is different as it uses a neural network and never calculates the parameters of the stochastic variables explicitly. It returns functions that for input  $x$  returns the mean values or standard deviation in that  $x$ . The precision for this method is somewhat poorer, but it converges faster.

Steepest descent algorithm for fixed total mass problem has never been used together with neural networks and automatic differentiation. For this method, we make no assumptions on the distribution of the measure and instead represent it as a finite amount of discrete atoms. The mass of the atoms are adjusted to minimize a loss function, but the total mass is preserved. With this function, any function can be described and stochastic variables can be either dependent or independent. This method appears to need more data to train on for the results to be within acceptable precision, but is the fastest to converge.

Different loss functions are tried, however, we conclude that the logarithmic loss function is the most accurate for us. This loss function is used in all tests. For future papers, it could be valuable to investigate the method's ability as a part of larger networks, as this paper only cover simpler problems with few parameters. We conclude that further reaserch is needed before the steepest decent algorithm for fixed total mass is practical for use in large stochastic models like Bayesian neural networks.

## Sammandrag

Artificiell intelligens och neurala nätverk är just nu ett mycket omdebatterat ämne. I denna rapport undersöks metoder för optimering av mått, vilka kan användas i implementationen av bayesianska neurala nätverk. Dessa nätverk kan potentiellt förhindra överanpassning till träningsdata, och möjliggör kvantifiering av osäkerheten för utdata.

Med hjälp av Python och PyTorch har tre modeller designade för bayesiansk regression implementerats: optimering över diskreta mått, polynomisk modell och linjärkombination av normalfördelade stokastiska variabler. Dessa har distinkta tillvägagångsätt att modellera data. I synnerhet undersöks hur effektivt och exakt de kan anpassas till given data.

Linjärkombination av stokastiska variabler beskriver ett mått som en summa av stokastiska variabler multiplicerat med en godtycklig funktion. Genom att anta att stokastiska variablerna är normalfördelade kan stora förenklingar göras vilket gör algoritmen relativt okomplicerad. Vid testning påvisar denna metod god förmåga att modellera data med hög precision, men konvergerar långsammast.

Det polynomiska neurala nätverket efterliknar linjärkombinationen vid fallet att funktionerna är termer i ett polynom. Lösningsgången är dock annorlunda då ett neuralt nätverk används och att stokastiska variablernas parametrar inte beräknas explicit. Istället matar nätverket ut en sammansatt funktion för medelvärde respektive standardavvikelse. Denna metod har något sämre precision men snabbare konvergenstid.

Optimeringsalgoritmen över finita diskreta mått är tidigare utforskad i samband med neurala nätverk och automatisk differentiering. Här görs inga antagelser om måttets form utan det presenteras helt av ett finit antal diskreta atomer. Atomernas massa justeras för att minimera en målfunktion men måttet bibehåller dess totala massa. Med denna metod kan valfri funktion beskrivas och stokastiska variablerna kan antas vara beroende eller oberoende. Testning tyder på att denna metod kräver fler datapunkter att träna på för att erhålla en acceptabel precision men är den snabbaste att konvergera av de tre.

Olika målfunktioner undersöks, men det fastställs att log-likelihoodfunktionen är den bästa för stokastiska variabler. Denna används sedan för alla tester. Det kvarstår att undersöka metodernas effektivitet i större nätverk då denna rapport endast behandlar enklare problem med ett få antal parametrar. Slutsatsen är att vidare forskning för optimeringsalgoritmen över finita diskreta mått krävs innan metoden är tillämpbar i stora bayesianska neurala nätverk.



# Innehållsförteckning

<b>1</b>	<b>Inledning</b>	<b>1</b>
1.1	Syfte . . . . .	1
1.2	Avgränsningar . . . . .	1
<b>2</b>	<b>Bakgrund</b>	<b>2</b>
2.1	Bayesiansk statistik . . . . .	2
2.2	Neurala nätverk . . . . .	2
2.3	Bayesianska neurala nätverk . . . . .	3
2.4	Olika former av regression . . . . .	3
2.4.1	Linjär regression . . . . .	3
2.4.2	Polynomiell regression . . . . .	4
2.4.3	Stokastisk polynomiell regression . . . . .	4
2.5	Matematisk optimering . . . . .	4
2.6	PyTorch . . . . .	4
2.7	Automatisk differentiering . . . . .	5
2.8	Mått . . . . .	6
2.9	Differentiering av mått . . . . .	7
2.10	Kärndensitetsuppskattning . . . . .	7
<b>3</b>	<b>Modeller och metoder</b>	<b>8</b>
3.1	Parametrisk modell . . . . .	8
3.1.1	Linjärkombination av stokastiska variabler . . . . .	8
3.1.2	Polynomisk modell med normalfördelade koefficienter . . . . .	9
3.2	Icke-parametrisk modell . . . . .	9
3.2.1	Optimeringsalgoritm över finita diskreta mått (OFDM) . . . . .	9
3.2.2	Steg för implementering . . . . .	10
3.2.3	Implementation i Python . . . . .	11
3.3	Målfunktioner . . . . .	11
3.3.1	Maximum likelihood-metoden . . . . .	11
3.3.2	Chi-kvadrat-metoden . . . . .	12
3.3.3	ESSR-metoden . . . . .	13
3.3.4	Val av målfunktion . . . . .	13
<b>4</b>	<b>Resultat</b>	<b>14</b>
4.1	Jämförelse mellan parametriska och icke-parametriska metoder . . . . .	14
4.1.1	Prestanda . . . . .	14
4.1.2	Konvergenstid för optimeringsalgoritmen över finita diskreta mått . . . . .	15
4.1.3	Konvergenstid för parametriska metoder . . . . .	15
4.2	Resultat från tester . . . . .	15
<b>5</b>	<b>Diskussion</b>	<b>17</b>
5.1	Jämförelse . . . . .	17
5.2	Val för KDE . . . . .	17
5.3	Tidskomplexitet för större modeller med OFDM . . . . .	18
5.4	Vidare forskning . . . . .	18
5.5	Samhälleliga och etiska aspekter . . . . .	19
5.6	Slutsats . . . . .	20
<b>A</b>	<b>Appendix 1 - Resultattabeller</b>	<b>i</b>
A.1	Resultat av tester . . . . .	i

<b>B</b>	<b>Appendix 2 – Teori</b>	<b>ii</b>
B.1	Degenererad summa av väntevärden . . . . .	ii
B.2	Linjärkombination av stokastiska variabler . . . . .	ii
B.3	Yttre strafffunktion . . . . .	iii
B.4	Lagrangemultiplikator-metoden . . . . .	iv
B.5	Icke-parametriska diskreta sannolikhetsmått med softmax . . . . .	iv
B.6	Poängfunktioner och poängregler . . . . .	iv
	B.6.1 Exempel på poängregler . . . . .	v
<b>C</b>	<b>Appendix 3 – Källkod</b>	<b>vi</b>
C.1	Måttklassen . . . . .	vi
C.2	Optimeringsklassen . . . . .	viii
C.3	Checkklassen . . . . .	xiv
C.4	Generering av data . . . . .	xvi
C.5	Tester för OFDM . . . . .	xvii
	C.5.1 Fallet $\alpha_i$ . . . . .	xvii
	C.5.2 Fallet $\alpha_i \cdot x$ . . . . .	xviii
	C.5.3 Fallet $\alpha_i + \beta_i \cdot x$ . . . . .	xix
	C.5.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$ . . . . .	xx
C.6	Tester för linjärkombination av stokastiska variabler . . . . .	xxii
	C.6.1 Fallet $\alpha_i$ . . . . .	xxii
	C.6.2 Fallet $\alpha_i \cdot x$ . . . . .	xxiii
	C.6.3 Fallet $\alpha_i + \beta_i \cdot x$ . . . . .	xxiv
	C.6.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$ . . . . .	xxv
C.7	Tester för polynomiskt neuralt nätverk . . . . .	xxvi
	C.7.1 Fallet $\alpha_i$ . . . . .	xxvi
	C.7.2 Fallet $\alpha_i \cdot x$ . . . . .	xxviii
	C.7.3 Fallet $\alpha_i + \beta_i \cdot x$ . . . . .	xxx
	C.7.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$ . . . . .	xxxii
C.8	Implementation av linjärkombination av stokastiska variabler . . . . .	xxxiv

# 1 Inledning

Neurala nätverk har visat sig användbara inom maskininlärning, specifikt då de applicerats inom bild- och språkanalys. Ett stort problem med neurala nätverk är det fenomen som kallas överanpassning. Överanpassning innebär att en modell, i detta fall ett neuralt nätverk, anpassar sig för väl till träningsdata och därför inte kan beskriva ny data som den inte tidigare observerat. Det är intressant att försöka studera hur detta problem kan undvikas.

Bayesianska neurala nätverk (BNN) är en typ av neurala nätverk som utnyttjar stokastiska parametrar. De kvantifierar osäkerheten i en modells förutsägelser, vilket kan göra dem robustare vid tillämpning. Då parametrarna är stokastiska variabler kan även problemet med överanpassning undvikas, då modellen alltid kommer variera kring medelvärden.

Ett problem med BNN är att de är mer beräkningstunga än traditionella neurala nätverk (TNN). Detta beror på att parametrarna i ett BNN är sannolikhetsmått och träningen av dessa kräver mer beräkningskraft då varje sannolikhetsmått innehåller mer information än parametrar i ett TNN. Därav kan de bli opraktiska eftersom det tar lång tid att träna och utvärdera på en större datamängd. Det är därför av intresse att studera algoritmer som effektiviserar optimeringen över sannolikhetsmått.

En algoritm för optimering över mått som visat sig effektiv är *Steepest descent algorithm for fixed total mass* [21], vilken i denna rapport benämns ”optimeringsalgoritm över finita diskreta mått” (OFDM). Algoritmen påstås i rapporten vara en sann ekvivalent till ”sjunkgradient” (eng. gradient descent), som är en etablerad algoritm inom matematisk optimering [3]. Detta motiverar att den teoretiskt sett bör prestera väl vid optimering av funktioner över sannolikhetsmått.

Denna algoritm har studerats i syfte att implementeras som en del av ett BNN tidigare [12] och då påstods att algoritmen ej var praktisk för denna tillämpning, då den var svår att skala för högre dimensioner. Författarna beskriver hur det inte är uppenbart hur metoden ska kunna integreras med moderna tekniker för neurala nätverk. OFDM utvärderas dock ej explicit i rapporten. Vi anser att det är intressant att utvärdera metoden vidare i syfte att implementeras i neurala nätverk. Vi ämnar att jämföra OFDM med andra algoritmer och visa att algoritmen kan vara effektivare än dessa alternativ.

Den tidigare implementationen av algoritmen skedde i programmeringsspråket R. En implementering med Python-biblioteket PyTorch är ett framsteg då vi med detta kan beräkna gradienten för funktioner automatiskt. En sådan implementation skulle därför kunna vara beräkningsmässigt snabbare än tidigare och därför kunna appliceras på större nätverk.

## 1.1 Syfte

Syftet med arbetet är att studera olika typer av optimeringsalgoritmer för BNN. I synnerhet ska en algoritm [21] som tidigare inte utnyttjat automatisk differentiering implementeras med hjälp av Python-biblioteket PyTorch. Denna metod ska sedan jämföras med andra metoder, specifikt ”linjärkombination av stokastiska variabler” samt ”polynomisk modell med normalfördelade koefficienter”.

## 1.2 Avgränsningar

Projektet kommer enbart att implementeras i programmeringsspråket Python där det redan existerande biblioteket PyTorch utnyttjas. Detta för att upprätthålla en tillräcklig hög abstraktionsnivå med användning av redan implementerade funktioner för maskininlärning. En specifik funktion som utnyttjas är att gradienten automatiskt beräknas för tensorer i PyTorch.

Det finns många möjliga vägar för projektet, men på grund av begränsad tid har fokuserat rapporten på en jämförelse. På grund av en given tidsram undersöks och behandlas endast de algoritmer som anses mest relevanta. Någon omfattande jämförelse med traditionella neurala nätverk kommer heller inte genomföras.

## 2 Bakgrund

I detta avsnitt etableras nödvändig teori bakom bland annat bayesiansk statistik, neurala nätverk samt mått.

### 2.1 Bayesiansk statistik

Bayesiansk statistik grundar sig i Bayes lag, vilken kan definieras enligt följande:

Låt  $A_1, \dots, A_n$  vara  $n$  disjunkta händelser, där  $\cup_{i=1}^n A_i = \Omega$  utgör rummet av alla möjliga händelser. Bayes lag säger då att sannolikheten,  $P$ , för händelse  $A_i$  givet en händelse  $B$  ges av

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{P(B)}. \quad (1)$$

Bayes lag används i en metod kallad bayesiansk inferens, vilket är en typ av statistisk inferens, som används för att bestämma egenskaper hos en okänd sannolikhetsfördelning.

Bayesiansk inferens utgår från att vi har en hypotes för fördelningen av en parameter. Denna hypotes kallas en a priori-fördelning och vi söker den riktiga fördelningen för parametern. För att beräkna den fördelningen behöver vi uppdatera vår a priori-fördelning med hjälp av data som vi observerar. Vi ställer upp en statistisk modell för hur sannolik vår observerade data är givet ett värde på parametern. Denna modell kallas för en likelihoodfunktion. En a posteriori-fördelning, d.v.s. vår nya hypotes om hur fördelning av parametern ser ut givet observerad data, kan då beräknas via

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}, \quad (2)$$

där  $H$  är hypotesen för parametern och  $E$  är data som observerats.  $P(E|H)$  anger likelihoodfunktionen,  $P(H)$  är a priori-fördelningen och  $P(E)$  anger den så kallade marginella fördelningen, vilket är fördelningen för vår data oberoende av hypotesen. Denna metod kan användas för att uppdatera en initial hypotes gällande en parameters fördelning efter att data har observerats, vilket är användbart i situationer där det är svårt eller till och med omöjligt att bestämma en fördelning på förhand.

### 2.2 Neurala nätverk

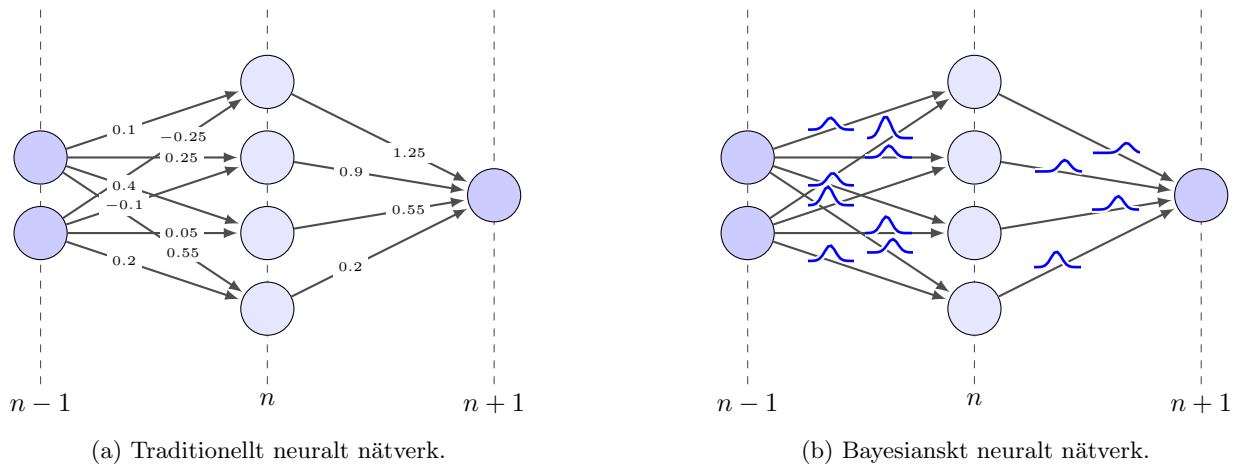
Neurala nätverk (NN) är strukturer som används inom maskininlärning för att skapa modeller som exempelvis klassificerar data eller implementerar regression till olika funktioner. Som namnet antyder imiterar dessa nätverk biologiska neurala nätverk, som exempelvis en mänsklig hjärna, och är på så sätt uppbyggda av noder och kopplingar mellan dessa noder. Dessa noder, eller neuroner, betraktas som funktioner som tar in en uppsättning variabler som indata och returnerar därefter ett värde.

Dessa noder delas in i lager vilka innehåller de noder som tar utdatan från föregående lager och vars utdata skickas till nästa lager av noder. Detta illustreras i figur 1a. Varje nod består av en linjär funktion plus en konstant, som sedan ges till en så kallad aktiveringsfunktion. Istället för en transformation för varje nod beskrivs ofta istället hela lager där ett lager beskrivs enligt

$$f(\mathbf{x}) = g(\mathbf{Ax} + \mathbf{b}), \quad f: \mathbb{R}^n \mapsto \mathbb{R}^m. \quad (3)$$

$\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  är utdatan från det föregående lagret,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  innehåller nodernas vikter,  $\mathbf{b} \in \mathbb{R}^m$  deras så kallade bias, och  $g: \mathbb{R}^m \rightarrow \mathbb{R}^m$  är aktiveringsfunktionen[10]. Här representerar varje element i utdatan värdet av en nod i detta lager.

Att flera lager av dessa funktioner kan approximera okända funktioner visas av de så kallade universella approximationssatserna, vilka kommer i lite olika former beroende på vilka aktiveringsfunktioner som



Figur 1: Schematiska bilder över olika neurala nätverk.

används, hur många noder ett lager får innehålla, samt hur många lager som används [17] [18]. Detta är den teoretiska bakgrunden till hur NN har lyckats göra allt från identifiering av objekt i bilder till röstigenkänning med mera.

### 2.3 Bayesianska neurala nätverk

Genom att sätta ihop teorin i de två sektionerna ovanför, 2.1, och 2.2, kan en ny typ av neurala nätverk skapas. Dessa kallas bayesianska neurala nätverk (BNN). Ett bayesianskt neuralt nätverk använder sig av stokastiska variabler som parametrar istället för att använda entydiga reella tal, se figur 1b. Dessa stokastiska variabler uppdateras i träningskedet av nätverket genom att strategiskt använda sig av bayesiansk inferens. Det bayesianska nätverket är tänkt att kunna träna en mer robust modell, som undviker problemet där en modell blir övertränad på en viss datamängd och därav presterar dåligt på ny indata. BNN påstås även kunna ge en bättre bild av osäkerheten i en modells prediktion, då prediktionen är beroende av stokastiska variabler [28].

### 2.4 Olika former av regression

I neurala nätverk används regression av olika slag under träningskedet för att uppdatera dess parametrar. I följande kapitel presenteras teorin för relevanta varianter av regression.

#### 2.4.1 Linjär regression

Låt  $y$  vara en stokastisk variabel given av

$$y = ax + b + \epsilon, \tag{4}$$

där  $a, b \in \mathbb{R}$  är konstanter och  $\epsilon \sim N(0, \sigma^2)$  är en normalfördelad stokastisk variabel. Detta innebär att avvikelser från  $y = ax + b$  är normalfördelade, det vill säga

$$y \sim N(a + bx, \sigma^2).$$

Detta kallas för en linjär regressionsmodell. Linjär regression innefattar anpassningen, eller regressionen, av denna modell  $y$  till data  $x$ , vilken görs genom maximering av sannolikheten att vår data  $x$  skulle genereras från modellen, vilket görs genom att skatta parametrarna  $a$  och  $b$  i modellen. Det existerar teoretiska metoder för att göra detta exakt [2], men regressionen kan även utföras via matematisk optimering, med en funktion för sannolikheten av vår data givet modellen som målfunktion.

### 2.4.2 Polynomiell regression

Vi kan även studera modeller som beskrivs av generella polynom. Vi antar då att  $y$  är en stokastisk variabel beskriven av

$$y = \sum_{k=0}^K a_k x^k + \epsilon, \quad (5)$$

där  $a_k \in \mathbb{R}$ ,  $\forall k \in \{0, 1, \dots, K\}$  och  $\epsilon \sim N(0, \sigma^2)$ . Detta kallas för en polynomiell regressionsmodell, där polynomet i fråga har grad  $K$ . Denna generalisering av regressionsproblem från linjära till polynomiella gjordes först av Gergonne [15].

### 2.4.3 Stokastisk polynomiell regression

Parametrarna i dessa regressionsmodeller behöver inte vara bestämda värden, utan kan även antas vara stokastiska variabler [8]. Låt  $\alpha_k$ ,  $k \in \{0, 1, \dots, K\}$ , för något  $K \in \mathbb{N}$ , vara stokastiska variabler med någon fördelning, det vill säga  $\alpha_k \sim \mu_k$ , där  $\mu_k$  är en sannolikhetsfördelning. Då kan en stokastisk regressionsmodell ställas upp

$$y = \sum_{k=0}^K \alpha_k x^k. \quad (6)$$

En sådan modell kan bättre beskriva samband där de beskrivande parametrarna inte kan antas vara bestämda värden, utan bör antas vara stokastiska. Det är enbart denna typ av stokastiska regressionsmodeller som kommer studeras i denna rapport. En stokastisk regressionsmodell kan ses som det enklaste exemplet av ett BNN, där vi enbart studerar en nod med en funktion beskriven av regressionsmodellen.

## 2.5 Matematisk optimering

Det centrala problemet inom matematisk optimering handlar om att minimera (eller maximera) en viss funktion, kallad målfunktion, av variabler. Variablerna kan även ha likhets- eller olikhetskrav, som begränsar vilka värden de kan anta. Detta kan formuleras matematiskt som

$$\begin{aligned} & \underset{x}{\text{minimera}} && f(x) \\ & \text{så att} && x \geq g(x), \\ & && x = h(x), \\ & && x \in X. \end{aligned} \quad (7)$$

Här är  $f$  målfunktionen som ska minimeras,  $x$  variablerna,  $X$  definitionsmängd för variablerna,  $g$  och  $h$  begränsande funktioner för variablerna. Ett optimeringsproblem på denna formen kan sedan lösas med en optimeringsalgoritm. Exempelvis kan sjunkgradient användas för problem utan begränsande funktioner [3].

## 2.6 PyTorch

Projektet skrivs i Python och Python-biblioteket PyTorch används för att realisera neurala nätverk och för möjligheten att automatiskt differentiera tensorer. Informationen i följande avsnitt är tagen från PyTorch officiella dokumentation [14].

En tensor i PyTorch är ett objekt som representerar en multidimensionell matris. En instans av objektet kan vara endimensionell (vektor), tvådimensionell (matris) eller av högre dimension. PyTorch har även diverse inbyggda funktioner för vanliga transformationer av data. Ett exempel är en linjär

transformation (adderat med en konstant) för att modellera ett lager i ett neuralt nätverk. Detta implementeras med hjälp av multiplikation av indata med en matris innehållande vikter, och sedan vektoraddition för tillägg av en konstant.

Linjära funktioner tillsammans med så kallade aktiveringsfunktioner är de byggstenar som utgör majoriteten av neurala nätverk. En aktiveringsfunktion är en deriverbar olinjär funktion placerad mellan linjära funktioner. Utan dessa skulle neurala nätverk vara linjära funktioner och därmed högst begränsade i vad de kan användas till. Två exempel på vanliga aktiveringsfunktioner är ReLU (eng. rectified linear unit) vilken definieras enligt  $f(x) = \max(0, x)$  och den så kallade sigmoid funktionen definierad som  $f(x) = 1/\exp(-x)$ .

I träningen av ett neuralt nätverk är det nödvändigt att beräkna derivatan av felet (ett mått på skillnaden mellan prediktion och faktiskt värde på träningsdata) med avseende på modellens parametrar. Detta görs i PyTorch med hjälp av funktionen `backward()`. Algoritmen bakom denna kallas automatisk differentiering och kommer förklaras i nästa sektion, 2.7. Denna funktions syfte är att genom små förändringar av parametervärdena kunna gå i den riktning (ändring av parametervärden) som minimerar felet snabbast.

## 2.7 Automatisk differentiering

När metoden `backward()` åberopas på en tensor i PyTorch beräknas gradienten för denna tensor. Det görs med hjälp av metoden automatisk differentiering [5] som varken är symbolisk eller numerisk, utan använder sig istället av duala tal [6]. För att beskriva duala tal introducerar vi det så kallade nilpotenta elementet  $\varepsilon$ , vilket är ett nollskilt tal som uppfyller att  $\varepsilon^2 = 0$ . Duala tal är då tal på formen  $a + b\varepsilon$ , där  $a, b \in \mathbb{R}$ . Låt oss betrakta Taylorutvecklingen av  $f(a + b\varepsilon)$  kring punkten  $a$ ,

$$f(a + b\varepsilon) = \sum_{k=0}^{\infty} \frac{f^{(k)}(a)b^k\varepsilon^k}{k!}. \quad (8)$$

Genom definitionen av  $\varepsilon$ , följer det att

$$f(a + b\varepsilon) = f(a) + f'(a)b\varepsilon. \quad (9)$$

Både värdet av funktionen och dess derivata för ett reellt värde  $a$  kan därför beräknas samtidigt genom att utvärdera funktionen i det duala talet  $a + b\varepsilon$ , då vi observerar att koefficienten för  $f'(a)$  kommer vara  $b\varepsilon$ . Detta fungerar även för sammansatta funktioner, eftersom

$$f(g(a + b\varepsilon)) = f(g(a) + g'(a)b\varepsilon) = f(g(a)) + f'(g(a))g'(a)b\varepsilon. \quad (10)$$

Automatisk differentiering nyttjar kedjeregeln och beräknar de partiella derivatorna i varje delsteg av funktionsevalueringen. Evalueringen av en funktion delas därför upp i delsteg, där varje delsteg representeras av en elementär funktion. Exempelvis skulle funktionen

$$f(x_1, x_2) = \cos(x_1 + x_2) - \log(x_2) \quad (11)$$

kunna delas upp i delfunktionerna

$$\begin{aligned} v_0(x_1, x_2) &= x_1, & v_1(x_1, x_2) &= x_2, \\ v_2(x_1, x_2) &= x_1 + x_2 = v_0 + v_1, & v_3(x_1, x_2) &= \log(x_2) = \log(v_1), \\ v_4(x_1, x_2) &= \cos(x_1 + x_2) = \cos(v_2), & v_5(x_1, x_2) &= f(x_1, x_2) = v_4 - v_3. \end{aligned} \quad (12)$$

På så sätt kan derivator för dessa delfunktioner beräknas och sedan multipliceras med redan beräknade derivator för att erhålla derivatan för varje variabel.

I vilken ordning derivatorna beräknas kan variera. Antingen kan derivatorna beräknas utgående från variablerna, vilket kallas framåt läge (eng. forward mode), eller kan derivatorna beräknas med utgångspunkt från sista beräkningen i funktionsevalueringen, så kallat omvänt läge (eng. reverse mode).

## 2.8 Mått

I träningskedet av BNN kommer optimering över sannolikhetsfördelningar behöva implementeras, vilket kräver en grundläggande etablering av teorin kring mått. De förenklingar som görs för att kunna implementera optimeringen innebär att de mått som kommer arbetas med är finita, diskreta mått.

Ett mått  $\mu$  är en funktion  $\mu : \mathcal{A} \rightarrow \mathbb{R}$ , där  $\mathcal{A}$  är en  $\sigma$ -algebra i domänen  $X$ [13]. Måttet  $\mu$  har egenskaperna

$$\mu(\emptyset) = 0, \quad (13a)$$

$$\mu(A \cup B) = \mu(A) + \mu(B), \quad A \cap B = \emptyset, \quad A, B \in \mathcal{A}, \quad (13b)$$

$$\mu\left(\bigcap_{i=1}^{\infty} A_i\right) = \lim_{i \rightarrow \infty} \mu(A_i), \quad \text{för följder } A_1 \supseteq A_2 \supseteq \dots \quad (13c)$$

Mått definieras som linjära och uppfyller därmed

$$(\mu + \eta)(A) = \mu(A) + \eta(A), \quad \mu, \eta \in \mathbb{M}, \quad A \in \mathcal{A} \quad (14)$$

$$(t\mu)(A) = t\mu(A), \quad \mu \in \mathbb{M}, \quad A \in \mathcal{A}, \quad t \in \mathbb{R}. \quad (15)$$

Ett mått kan delas upp i en positiv del och en negativ del. Detta kallas för Jordan-uppdelning och skrivs  $\mu = \mu^+ - \mu^-$ , där  $\mu^+, \mu^- \geq 0$  och ortogonala[20]. Ortogonalitet betyder att måtten uppfyller

$$\begin{cases} \text{Om } \mu^+(A) > 0 \Rightarrow \mu^-(A) = 0, \quad \forall A, \\ \text{Om } \mu^-(A) > 0 \Rightarrow \mu^+(A) = 0, \quad \forall A. \end{cases} \quad (16)$$

Ett mått  $\mu$  som uppfyller

$$\mu(A) \geq 0, \quad \forall A \in \mathcal{A}, \quad (17)$$

kallas för ett positivt mått. Med Jordan-uppdelning kan därmed mått uttryckas som skillnaden mellan två positiva mått. Totala variationen av ett mått definieras som

$$\|\mu\| = \mu^+(X) + \mu^-(X). \quad (18)$$

Totala variationen utgör en norm för mått över  $X$ . Om totala variationen är finit över  $X$  är även  $\mu(X)$  finit. Mängden av alla mått med finit norm  $\mathbb{M}$  utgör därmed ett Banach rum över  $X$ [20]. Väntevärdet och variansen för ett mått definieras respektive som

$$E[\mu] = \int x\mu(dx), \quad (19)$$

$$\text{Var}[\mu] = \int (x - E[\mu])^2 \mu(dx). \quad (20)$$

Sannolikhetsfördelningar är ett exempel på mått, vilka även har begränsningen att  $\mu(X) = 1$ . Antag att ett diskret finit mått  $\mu$  har norm 1 och enbart har massa i punkter  $x = (x_1, \dots, x_n)$  Variansen av  $\mu$  över  $X$  beräknas enligt:

$$\text{Var}_X(\mu) = \sum_{i=1}^n (x_i - \bar{x})^2 \mu(x_i), \quad (21)$$

där  $\bar{x}$  är väntevärdet av  $\mu$  över  $X$ [20].



## 2.9 Differentiering av mått

I följande underavsnitt är källan till teorin [20] om inget annat anges.

En funktion  $\psi : \mathbb{M} \rightarrow \mathbb{R}$  är Frechét differentierbar om

$$\psi(\mu + \eta) - \psi(\mu) = D\psi(\mu)[\eta] + o(\|\eta\|), \text{ när } \|\eta\| \rightarrow 0, \quad (22)$$

där  $D\psi(\mu)[\eta]$  är en begränsad linjär kontinuerlig funktional av  $\eta$ . Låt  $X = \{1, \dots, n\}$ . Finita mått på  $X$  är  $\mu = (m_1, \dots, m_n)$  vilket ger att  $\mathbb{M} = \mathbb{R}^n$ .  $D\psi(\mu)$  är i det fallet en linjär avbildning vid punkten  $\mu \in \mathbb{R}^n$ , alltså finns det en vektor  $d(x, \mu) = (d_1, \dots, d_n)$ ,  $x \in X$ , så att för en inkrementering  $\eta(x) = (\eta_1, \dots, \eta_n)$  är

$$D\psi(\mu)[\eta] = \sum_{x=1}^n d_x \eta_x = \int_X d(x, \mu) \eta(dx). \quad (23)$$

I det fallet är  $d(x, \eta)$  en gradient.

## 2.10 Kärndensitetsuppskattning

Då diskreta sannolikhetsmått arbetas med behövs ibland metoder för att approximera måttet med en kontinuerlig fördelning för att undvika problem. En möjlighet är kärndensitetsuppskattning (KDE från eng. Kernel Density Estimation) som är en välkänd metod för att approximera täthetsfunktionen  $p$  av ett okänt sannolikhetsmått  $P$  genom stickprov [11]. Via ett stickprov  $X_1, X_2, \dots, X_n \in \mathbb{R}^d$  från  $P$  som är oberoende och likafördelade kan vi approximera  $p$  som

$$\hat{p}_n(x) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right), \quad (24)$$

där  $K : \mathbb{R}^d \mapsto \mathbb{R}$  är en glatt funktion som kallas kärnfunktionen och  $h > 0$  är utjämningskonstanten som bestämmer hur mycket funktionen utjämnas.

Ett annat sätt att skriva ekvation (24) är

$$\hat{p}_m(x) = \frac{1}{mh^d} \sum_{i=1}^m w_i \cdot K\left(\frac{x - Z_i}{h}\right), \quad (25)$$

där  $Z_1, Z_2, \dots, Z_m$  är utfallsrummet av fördelningen  $P$  och  $w_i$  vikten av utfallet  $Z_i$ , det vill säga  $P(Z_i) = w_i$  för  $i = 1, 2, \dots, m$ . För tillräckligt stora  $n$ , det vill säga stora stickprov, är ekvation (24) och (25) ekvivalenta enligt 'De stora talens lag' då frekvensen av  $K\left(\frac{x - Z_i}{h}\right)$  kommer närma sig  $w_i$  [27].

Ett vanligt exempel på en kärnfunktion är

$$K(x) = \frac{\exp(-x^2/2)}{\sqrt{2\pi}}, \quad (26)$$

det vill säga en standardiserad normalfördelning, även kallad den gaussiska kärnfunktionen.

Av funktioner som fungerar som kärnfunktioner är det inte stor skillnad på deras effektivitet och därmed är inte valet av kärnfunktion särskilt viktigt. Däremot påverkar valet av utjämningskonstanten,  $h$ , vår approximation avsevärt [7]. En metod för att få ett värde på  $h$  är Silvermans tumregel,

$$h = 0.9 \cdot \min(\sigma, IQR/1.34) \cdot n^{-\frac{1}{5}} \quad (27)$$

där  $IQR$  är kvartilavståndet av fördelningen  $P$ , det vill säga avståndet mellan första och tredje kvartilen [7]. En nackdel med Silvermans tumregel är att den inte är lika användbar i fallet då fördelningen vars täthetsfunktion vi vill uppskatta inte är Gaussisk [7]. Ett annat sätt att välja utjämningskonstanten är Sheath och Jones metod vilket ger ett  $h$  som i många fall kommer vara bättre då metoden anpassar sig lättare efter situationen [24].

## 3 Modeller och metoder

Skattning av mått kan genomföras med två olika typer av modeller: parametrisk och icke-parametrisk. Med en parametrisk modell görs antagandet att fördelningen följer en tidigare känd fördelning som kan beskrivas fullkomligt av några få parametrar, exempelvis en normalfördelning. En icke-parametrisk modell däremot är generellt sett mer komplext då stora antaganden om hur fördelningen ser ut inte kan göras. Istället behöver värden skattas i en större mängd punkter för att beskriva fördelningen [22]. Nedan kommer teorin bakom ett urval av metoder för att optimera modellerna i båda kategorierna beskrivas som sedan kommer implementeras i Python. Under projektets gång har andra alternativ alternativ än modellerna nedan testats, exempelvis strafffunktioner- och Lagrangemultiplikator-metoden, se appendix B.3-B.4. Dessa används dock inte i jämförelsen då de inte presterade tillräckligt bra under de givna förhållandena.

### 3.1 Parametrisk modell

När en modell antas vara parametrisk kan vi fullständigt beskriva den med några enstaka parametrar. Ett exempel på detta är normalfördelningen som har två parametrar: medelvärde och standardavvikelse. Två metoder som optimerar denna modellen testas.

#### 3.1.1 Linjärkombination av stokastiska variabler

*Följande teori är opublicerad men tillskrivs handledare Sergei Zuyev.*

Låt  $\eta(x)$  vara en funktion definierad enligt

$$\eta(x) = \sum_{k=1}^K \alpha_k h_k(x), \quad (28)$$

där  $\alpha_k \sim \mu_k$  är oberoende stokastiska variabler och  $h_k(x)$  antas vara kontinuerliga funktioner av  $x \forall k$ . I många fall kan det antas att en normalfördelning eller en sammansättning av normalfördelningar kan beskriva vår data. Det antagandet kommer göras här, men det generella fallet samt dess lösning finns beskrivet i appendix B.2. Om vi antar att  $\alpha_k$  är oberoende och normalfördelade, det vill säga  $\alpha_k \sim \mathcal{N}(m_k, \sigma_k^2) \forall k$ , kommer  $\eta(x_i)$  också vara normalfördelad, eftersom det är en linjär kombination av normalfördelade stokastiska variabler med icke-stokastiska koefficienter  $h_k(x)$  för alla  $x$  [19]. Fördelningen följer

$$\eta(x) \sim \mathcal{N}(m(x), \sigma^2(x)), \quad \forall x, \quad (29)$$

där

$$m(x) = \sum_{k=1}^K m_k h_k(x), \quad (30)$$

$$\sigma^2(x) = \sum_{k=1}^K (\sigma_k h_k(x))^2. \quad (31)$$

Detta ger täthetsfunktionen för  $\eta(x)$  som

$$f_{\eta(x)}(y) = \frac{1}{\sqrt{2\pi\sigma^2(x)}} \exp\left(-\frac{(y - m(x))^2}{2\sigma^2(x)}\right), \quad (32)$$

som kan kombineras med log-likelihoodfunktion i avsnitt 3.3.1 för att skatta parametrar.

### 3.1.2 Polynomisk modell med normalfördelade koefficienter

Om den tidigare nämnda linjärkombinationen antas vara ett polynom av grad  $K$  kan vi förenkla ytterligare till stokastisk polynomiell regression, se avsnitt 2.4,

$$\eta(x) = \sum_{k=0}^K \alpha^{(k)} x^k. \quad (33)$$

Istället för att explicit söka varje parametrar för varje stokastisk variabel i detta polynom kan ett kombinerat medelvärde med associerad varians hittas genom att separera dem till två separata polynom. Medelvärdet  $\mu$  är ett polynom av grad  $K$  som det ursprungliga polynomet, där vi ersätter stokastiska variablerna med skalära parametrar  $a_\mu$ ,

$$\mu(x) = \sum_{k=0}^K a_\mu^{(k)} x^k. \quad (34)$$

Eftersom varians är kvadratisk beroende av  $x$  enligt ekvation (21) kommer variansen  $\sigma^2$  kräva ett polynom av graden  $2K$ ,

$$\sigma^2(x) = \sum_{k=0}^K a_\sigma^{(k)} x^{2k}. \quad (35)$$

Ekvation (34) och (35) är båda polynom som kan lösas med generell polynomiell regression. I praktiken implementeras detta i PyTorch med hjälp av två parallella linjära lager utan bias. Ett lager för medelvärde och ett för variansen. En tensor med  $x$  för varje exponent matas in i respektive lager:  $[1, x, x^2 \dots x^K]$ . Lösningen är en kurva för  $\mu$  som visar medelvärdet i varje punkt  $x$  samt en kurva för  $\sigma^2$  som visar totala variansen. Detta gör det enkelt att visualisera resultatet och om så önskas kan stokastiska variablernas parametrar härledas från de skalära parametrarna.

## 3.2 Icke-parametrisk modell

Nackdelen med en parametrisk modell är att antaganden om måttet måste göras. Användaren måste i förväg bestämma vilka parametrar modellen ska anpassa och hur dessa beskriver måttets form. En lösning är att istället diskretisera måttet. Ett mått med ett stöd på ett visst intervall kan modelleras med jämnt fördelade atomer (punkter) från detta intervall vilka har varsin vikt eller massa associerad med sig. I appendix B.5 illustreras hur fallet med diskreta sannolikhetsmått kan implementeras med den så kallade softmax-funktionen. Denna funktion är enkel att integrera i neurala nätverk, men har potentiellt långsam konvergens beroende på var i parameterrummet minimipunkten ligger. Nedan introduceras en alternativ metod som effektiviserar optimering av diskreta mått avsevärt.

### 3.2.1 Optimeringsalgoritm över finita diskreta mått (OFDM)

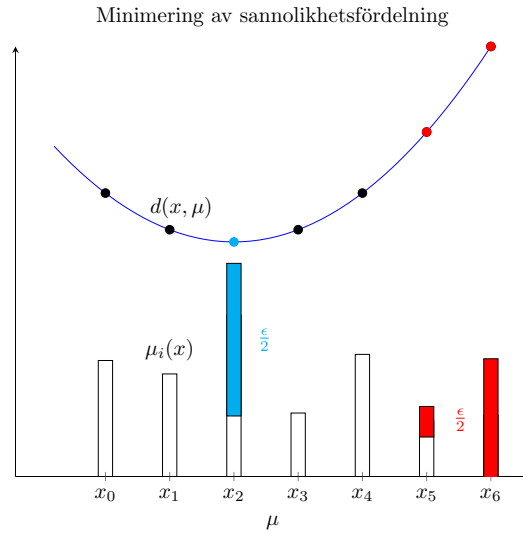
En metod för att minimera en funktion över ett diskret mått är att flytta en liten mängd massa i fördelningen [21, s. 119]. För att måttet ska fortsätta vara en sannolikhetsfördelning, alltså summera till 1, behöver all massa som subtraheras vid en punkt adderas vid en annan punkt. Låt  $\mu$  vara en fördelning och  $F$  en funktion. Vi vill hitta ett  $\eta$  som minimerar funktionen  $F$ , alltså

$$\min_{\eta} F(\mu + \eta) \quad (36a)$$

$$\text{så att } 0 \leq \mu_i + \eta_i \leq 1, \forall i \in \{1, \dots, n\}, \quad (36b)$$

$$\|\mu\| = 1, \quad (36c)$$

$$\sum_{i=1}^n \eta_i = 0. \quad (36d)$$



Figur 2: För att minimera måttet över  $F$  adderas positiv massa (turkos stapel) till fördelningen  $\mu$  där gradienten  $d(x, \mu)$  är som minst (turkos punkt) och negativ massa (röda staplar) där gradienten  $d(x, \mu)$  är som störst (röda punkter).

Eftersom  $\eta$  är konstant med avseende på  $\mu$  kan minimering av (22) lösa minimering av (36a). Vi får då

$$F(\mu + \eta) = F(\mu) + DF(\mu)[\eta] + \|\eta\|, \quad (37)$$

$$F(\mu + \eta) \approx F(\mu) + \sum_{i=1}^n d_i \eta_i, \quad (38)$$

enligt (22) och (23) och när  $\|\eta\|$  är liten. Med hjälp av (22) kan alltså minimeringsproblemet approximeras genom att placera ut massorna av  $\eta_i$  på olika delar av gradienten  $d(x, \mu)$ . I figur 2 placerar vi ut den positiva massan där gradienten är som minst och negativ massa där gradienten är som störst för att minimera måttet över  $F$  (36a).

Det finns två utfall vid fördelning av den negativa massan. Antingen kan all massa placeras vid punkten där gradienten är störst eller bara en del av massan för att fortfarande uppfylla (36b). Efter att all möjlig negativ massa har placerats ut på en punkt itereras metoden på punkten med näst störst gradient med den återstående massan. Detta upprepas tills all massa har placerats.

### 3.2.2 Steg för implementering

Vi kan definera  $\theta^+$  och  $\theta^-$  som den positiva respektive negativa massan som har distribuerats på  $\eta$ . Mer konkret är

$$\theta^+ = \sum_{i=1}^n \eta_i^+, \quad (39)$$

$$\theta^- = \sum_{i=1}^n \eta_i^-. \quad (40)$$

Nedan följer steg för hur massorna distribueras från given gradient  $d(x, \mu)$ .

1. Låt  $\theta^+, \theta^- = \frac{\epsilon}{2}$  och  $I = \emptyset$  och  $k = 0$ .

2. Låt  $i$  vara punkten där gradienten är minst,  $i = \arg \min_{i \in \{1, \dots, n\}} d(x, \mu_i)$ . Placera massan  $\theta^+$  på  $\mu_i$ , vilket ger  $\mu_i := \mu_i + \theta^+$ .
3. Låt  $i$  vara punkten där gradienten är störst,  $i = \arg \max_{i \in \{1, \dots, n\} \setminus I} d(x, \mu_i)$ . Bestäm massa enligt  $\eta_i^- = \min(\theta_k^-, \mu_i)$  och fördela  $\mu_i := \mu_i - \theta^-$ . Uppdatera  $\theta_{k+1}^- = \theta_k^- - \eta_i^-$ .

**Om  $\theta_{k+1}^- = 0$ : Stopp!**

**Annars:** Uppdatera  $k := k + 1$ , lägg till  $I := I \cup \{i\}$  och **upprepa steg 3**.

Eftersom  $\theta^+ = \theta^-$  har storleken av  $\mu$  inte ändrats utan massan har enbart förflyttats för att minimera måttet över  $F$ . Dessa steg upprepas tills gradienten är konstant och minimal i alla atomer med nollskild massa. Måttet uppfyller då nödvändiga villkor för minimeringsproblemet.

### 3.2.3 Implementation i Python

Den huvudsakliga klassen `Measure`, med tillhörande metoder, implementerar ett diskret mått. Klassens metoder är dokumenterade och en instans av klassen beskrivs i huvudsak av dess vikter (parametern `self.weights`), samt en tillhörande plats (parametern `self.locations`) vilka tillsammans beskriver atomer. Koden hittas i appendix C.1. Funktionerna i denna klass uppfyller vanliga funktioner kring räkning med mått, så som att se om ett givet mått är ett sannolikhetsmått, vad ett mått har för total massa, att hitta stödet för ett mått, hur ett mått visualiseras, med mera.

En instans av klassen `Optimizer` skapas sedan med hjälp av en instans (eller lista av instanser) av klassen `Measure` samt en textsträng som specificerar vilken målfunktion som bör användas. Dess metod `minimize` används sedan för att optimera måtten. Koden hittas i appendix C.2. Kärnan av `minimize` funktionen är algoritmen som beskrivs ovan. Mellan varje iteration verifierar funktionen om resultatet uppfyller kraven för att avsluta som presenteras i sektion 3.2.2 och avslutar i så fall. Annars kollar den om resultatet faktiskt blev bättre, om det inte är fallet minskas steglängden och den senaste förändringen återställs; den kontrollerar även om steglängden har blivit mindre än ett givet tröskelvärde och i så fall anses optimeringen vara färdig. Om detta kriteriet ej uppnås avslutar algoritmen automatiskt efter ett givet antal iterationer.

Dessutom finns klassen `Check` som används för att undersöka hur likt ett anpassat mått är den ursprungliga fördelningens data. En instans av `Check` tar in en instans av `Optimizer`, modellen som används och data från den ursprungliga sannolikhetsfördelningen. Koden finns i appendix C.3.

## 3.3 Målfunktioner

Vid optimering av mått behövs en metod för att avgöra om den nya hypotesen är bättre än den förgående. Detta kan uppnås med en målfunktion som avbildar måttet på en tallinje vars värde indikerar hur bra måttet representerar given data. Ett värde närmare 0 indikerar att måttet är bättre anpassat till given data. Det finns flera relevanta målfunktioner att applicera på mått, deras relevans motiveras av poängfunktioner, se appendix B.6. Nedan presenteras de tre som undersöks.

### 3.3.1 Maximum likelihood-metoden

En vanlig målfunktion med en stark grund i matematisk statistik är maximum likelihood [2]. Låt  $x_1, \dots, x_n$  vara ett stickprov från en stokastisk variabel  $X$  med täthetsfunktion  $f_\theta$ . Vi definierar den så kallade likelihoodfunktionen

$$L(\theta) = \prod_{i=1}^n f_\theta(x_i). \quad (41)$$

Det parametervärde  $\hat{\theta}$  som maximerar likelihoodfunktionen kallas maximum likelihood-skattningen (ML-skattningen). Då det i många fall kan vara lättare att hantera en summa, kan vi definiera den logaritmiska sannolikheten:

$$l(\theta) = \log \left( \prod_{i=1}^n f_{\theta}(x_i) \right) = \sum_{i=1}^n \log(f_{\theta}(x_i)) \quad (42)$$

Även denna funktion kommer maximeras av ML-skattningen. Båda dessa funktioner kan ses som ett mått på sannolikheten att ett stickprov producerades av en viss fördelning med parameter  $\theta$ . Då det ofta är minimeringsproblem som studeras inom matematisk optimering, betraktas den additiva inversen av likelihoodfunktionen. Framöver benämns den som NLL-funktionen (negativ log-likelihood).

### 3.3.2 Chi-kvadrat-metoden

En annan målfunktion som är av intresse att testa är Chi-kvadrat ( $\chi^2$ ). Den så kallade Chi-kvadrat-metoden grundar sig i Chi-kvadrat-testet. Nedan följer en beskrivning av testet taget från [26].

Låt  $X$  vara en diskret stokastisk variabel, som kan anta värdena  $x_i$ ,  $i = 1, 2, \dots, n$ , och en täthetsfunktion enligt  $P(X = x_i) = p_i$ . Vi kan definiera

$$C = \sum_{i=1}^n \frac{(f_i - kp_i)^2}{kp_i} = \sum_{i=1}^n \frac{f_i^2}{kp_i} - k, \quad (43)$$

där  $f_i$  är frekvensen av värdet  $x_i$  från ett stickprov på  $k$  värden från fördelningen. Då  $k$  är stort approximerar  $C$  Chi-kvadrat-fördelningen med  $n - 1$  frihetsgrader. Värdet kan därför användas som mått på hur nära stickprovet är den förväntade fördelningen. Om vi nu låter de relativa frekvenserna  $\frac{f_i}{k}$  bestämmas av vikterna från ett förutspått sannolikhetsmått  $\hat{Y} = (\hat{y}_1, \dots, \hat{y}_n)$  och  $Y = (y_1, \dots, y_n)$  vara fördelningen från observerad data. Då kan vi skriva  $C$  som  $k \left( \sum_{i=1}^n \frac{\hat{y}_i^2}{y_i} - 1 \right)$ . För att  $\hat{Y}$  ska bli så nära  $Y$  som möjligt vill vi då minimera  $\sum_{i=1}^n \frac{\hat{y}_i^2}{y_i}$ . Vi kan då definiera målfunktionen som kallas Chi-kvadrat förlust (eller bara Chi-kvadrat):

$$L = \sum_{i=1}^n \frac{\hat{y}_i^2}{y_i} \quad (44)$$

Denna funktion föreslogs först i *On the criterion that a given system of deviations...* [23] och analyserades vidare som en del av ett TNN i *Chi-square Loss for Softmax: an Echo of Neural Network Structure* [26].

Något som påpekas i den senare artikeln är att det kan förekomma frekvenser  $y_i$  i observerad data som ej är nollskilda. Detta blir problematiskt då dessa förekommer i nämnaren i Chi-kvadrat. För att lösa detta kan så kallad label smoothing användas. Denna teknik omfördelar vikterna så att alla vikter är nollskilda, utan att påverka den generella formen av sannolikhetsmättet. En exakt metod som kan användas för detta är att justera vikterna enligt

$$y_i^{LS} = y_i(1 - \alpha) + \frac{\alpha}{n}, \quad (45)$$

där  $\alpha$  är ett litet positivt tal ( $\alpha \in \mathbb{R}^+ : \alpha \approx 0$ ). Dessa justerade vikter kan sedan användas i funktionen.

### 3.3.3 ESSR-metoden

Väntevärdet av summan av kvadrerade rester eller ESSR (eng. Expected Sum of Square Residuals) innebär beräkning av

$$\arg \min_{\theta} E_{\theta} \left[ \sum_{i=1}^n (\hat{y}_i(\theta) - y_i)^2 \right],$$

där  $y_i$  är värden tagna från den riktiga fördelning och  $\hat{y}_i$  är en uppskattning av  $y_i$  gjord av modellen.  $\theta$  är fördelningen av parametrar för det diskreta måttet. I fallet linjär regression skulle värden från fördelningen  $\theta$  kunna vara par av lutning och skärningspunkt på y-axel för en linje som skall anpassas till en datamängd. ESSR-metoden kan betraktas som den probabilistiska versionen av minstakvadrat-metoden, vilken är standardmetoden för regression där parametrarna är vanliga skalärer.

### 3.3.4 Val av målfunktion

ESSR-metoden har egenskapen att lösningen degenererar till ML-skattningen av fördelningen, för mer information läs appendix B.1. Detta innebär att all vikt (eller vid parametrerade fördelningar, så mycket vikt som är möjligt) kommer att läggas på vad modellen tror är det mest sannolika värdet. Därmed reflekterar ej den resulterande fördelningen det faktum att vi fortfarande har viss osäkerhet kring parametervärdet. Eftersom syftet med det bayesianska tillvägagångsättet är att beskriva osäkerheten i resultatet är ESSR-metoden inte ett lämpligt val. Detta är ej fallet för Chi-kvadrat-metoden och log-likelihood. Fördelningen bibehåller en osäkerhet i form av en mängd olika möjliga världar. Efter vidare testning av dessa två dras slutsatsen att log-likelihood är det mer lämpliga valet då Chi-kvadrat vid flera tillfällen fastnade i lokala minimum. Eftersom ett diskret mått har ett begränsat antal vikter kommer  $f_{\theta}(x_i)$  inte existera för alla  $x_i$  vid användning av log-likelihood. En lösning är att avrunda till närmaste  $x$  som existerar i det diskreta måttets utfallsrum. Detta kommer dock med sina nackdelar. Avrundningar minskar inte bara precisionen utan kan även skapa problem vid optimering då gradienten blir diskontinuerlig. Därför approximeras först måttets täthetsfunktion med hjälp av kärndensitetsuppskattning för att få mer exakta värden på  $f_{\theta}(x_i)$ .

## 4 Resultat

I detta avsnitt presenteras metodiken för tester som görs samt dess resultat. I syfte att förenkla återskapandet av resultaten visas här versioner av Python och andra relevanta bibliotek som använts:

- Python: Version 3.10.6
- PyTorch: Version 1.13.1
- Numpy: Version 1.24.1

### 4.1 Jämförelse mellan parametriska och icke-parametriska metoder

För att utvärdera OFDM och jämföra den med andra modeller som beskrivits i rapporten, ska fyra olika tester genomföras. Låt  $\alpha_i, \beta_i, \gamma_i$  vara tre normalfördelade stokastiska variabler med parametrar  $\mu_i$  och  $\sigma_i$ ,  $i = 1, 2, \dots, 50$ , som också är tagna från två normalfördelningar. Testerna består av anpassning till följande modeller:

1. En normalfördelad stokastisk variabel ( $\alpha_i$ )
2. En linjär modell med en stokastisk variabel som koefficient ( $\alpha_i \cdot x$ )
3. En linjär modell med två stokastiska variabler som koefficienter ( $\alpha_i + \beta_i \cdot x$ )
4. En kvadratisk modell med tre stokastiska variabler som koefficienter ( $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$ )

De tre metoderna som ska jämföras är de två parametriska metoderna Linjärkombination av stokastiska variabler och Polynomisk modell med normalfördelade koefficienter, samt den icke-parametriska OFDM. För varje test används metoden med 50 olika värden av parametrarna  $\mu_i$  och  $\sigma_i$  och för varje sådant värde används tre olika storlekar på datamängder, 100, 500 respektive 1000 datapunkter vardera. Detta innebär 150 anpassningar för varje test och för varje metod. Alla dessa datamängder genereras via koden i appendix C.4.

I de olika testerna mäts tre olika metriker. **Prestanda** skattar sannolikheten att metoden kan anpassa en modell väl till godtycklig träningsdata. Alltså ett värde 0-1 där högre är bättre. Dessutom beräknas ett 95% konfidensintervall för prestandan. Se avsnitt 4.1.1 för en mer genomgående förklaring. Sedan mäts även det genomsnittliga antalet epoker som varje metod behöver för att konvergera, samt genomsnittlig körtid för metodernas konvergens. Dessa metriker kallas **konvergenstid (epoker)** och **konvergenstid (s)**. Resultatet av testerna kan ses i tabellerna i appendix A.1 och sammanfattas i avsnitt 4.2.

#### 4.1.1 Prestanda

Implementationen av beräkningen av prestanda skiljer sig mellan OFDM och de parametriska metoderna, men idén är den samma. Först beräknas för varje invärde  $x$  ett stickprov av modell-värden  $\hat{y}$  baserat på den anpassade modellen. Genom att bortse från de 2.5% lägsta, respektive högsta värdena i stickprovet erhålles ett approximativt konfidensintervall på 95%. Detta konfidensintervall kommer närma sig det teoretiska då antalet datapunkter växer. Sedan testas om den faktiska datapunkten  $y$  ligger i detta intervall och om inte räknas det som en "miss" (se figur 3). När detta gjorts för varje datapunkt  $x$ , beräknas ett teoretiskt konfidensintervall på 95% för hur många datapunkter  $y$  som bör missa. Detta görs genom att se antalet missar som en binomialfördelad stokastisk variabel, där antalet försök är lika med antalet datapunkter och sannolikheten för ett positivt utfall (en miss i detta fall) är 5% (se figur 4).

Metriken som beräknas i slutändan är andelen anpassningar där antalet missar ligger i konfidensintervallet. Denna parameter antas i sig vara binomialfördelad. Värdet kommer vara en skattning av det verkliga värdet på sannolikheten för att en modell producerad av metoden är välanpassad till



godtycklig träningsdata. En osäkerhet för denna metrik erhålles därför också, vilken ges av ett konfidensintervall på 95%. I beräkningarna antar vi i de flesta fall att metriken är normalfördelad, då detta är en bra approximation av binomialfördelningen när värdena på prestandan inte är nära 0 eller 1. Detta på grund av att normalfördelningen har svansartill skillnad från binomialfördelningen. I de fall prestandan var mycket nära eller lika med 0 eller 1 används därför ett 95% konfidensintervall för binomialfördelningen istället. För att se hur den faktiska implementationen gjordes för OFDM, se klassen `Check` i C.3, och för de parametriska metoderna, se funktionen `misses` i C.6 och C.7.

#### 4.1.2 Konvergenstid för optimeringsalgoritmen över finita diskreta mått

Som det står ovanför undersöks konvergenstiden på två sätt, mätt i sekunder och mätt i antal epoker. Dessa mätningar sker på lite olika sätt för de olika metoderna och i detta avsnitt detaljeras hur detta görs för OFDM. När metoden `minimize` kör kan koden sluta av tre anledningar, den har kört max antal epoker som tillåts, den har nått ett optimum eller om steglängden är för liten. I fallen av att den nått ett optimum returnerar metoden den dåvarande tiden och epoken den är i. Om den däremot stannar av steglängden eller epokerna kommer den undersöka om målfunktionen inte har ändrat värde under de senaste fem eller fler epokerna. Om den inte ändrat under fem eller fler epoker returneras den första tiden och epoken där målfunktionen inte ändrades, annars returneras den nuvarande tiden och epoken.

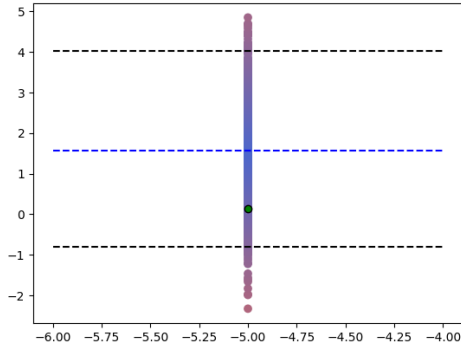
Värt att nämna är att metoden körs med ett max-antal epoker, vilket betyder att den inte kan köra för evigt. Den kan därför stanna på grund av att den når max-antalet innan den konvergerar och då räknas max-antalet som konvergenstiden. Max-antalet valdes till 4 000 epoker.

#### 4.1.3 Konvergenstid för parametriska metoder

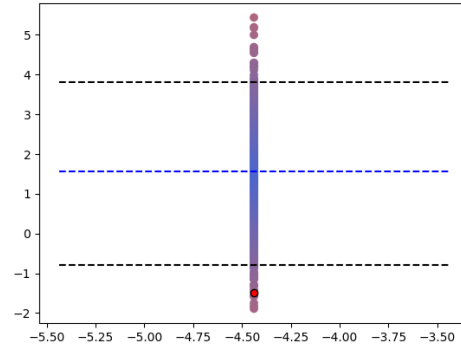
För de parametriska metoderna beräknas konvergenstid genom att i varje steg beräkna skillnaden i värde på målfunktionen mellan föregående och nuvarande steg. Denna skillnad jämförs sedan med en tolerans och om skillnaden är mindre än denna anses konvergens ha uppnåtts. Toleranserna väljs till  $10^{-4}$  för parametriskt neuralt nätverk och  $10^{-15}$  för metoden med linjärkombinationer av stokastiska variabler. Toleransen valdes till högre för parametriskt neuralt nätverk då denna metod observerades fluktuera runt en konvergenspunkt ibland, därför ansågs det att den bör anses konvergera för högre skillnader.

## 4.2 Resultat från tester

Det fullständiga resultatet finns i appendix A.1. Linjärkombination av normalfördelade stokastiska variabler får en ideal prestanda (1.0) på näst intill alla tester. Konvergenstiden verkar öka linjärt med antal datapunkter medan antal parametrar knappt påverkar tiden. Polynomisk modell med normalfördelade koefficienter påvisar sämre prestanda för fler datapunkter. Prestanda minskar även med antal datapunkter. OFDM har överlägset lägst prestanda i de flesta testerna men prestanda ökar med antalet datapunkter. För 1000 datapunkter uppnås en maximal prestanda av 0.92. Både OFDM och polynomisk modell har en betydligt lägre konvergenstid för större antal datapunkter än linjärkombination av stokastiska variabler.

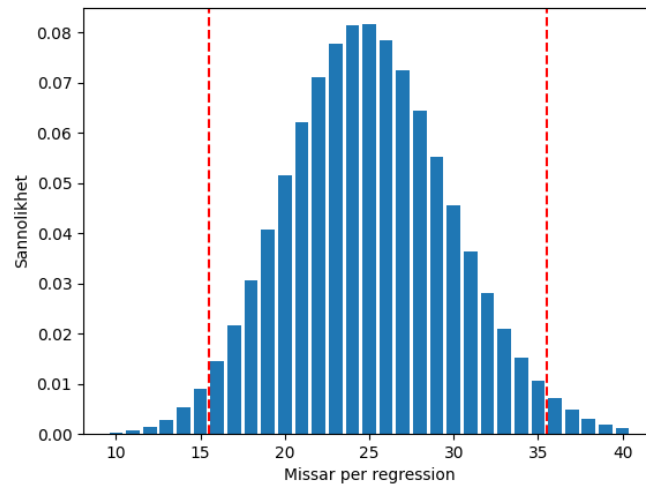


(a) Datapunkt innanför konfidensintervall



(b) Datapunkt utanför konfidensintervallet (räknas som en miss)

Figur 3: Stickprov av värden från normalfördelningar, med approximativt 95% konfidensintervall markerat med svart streckad linje och medelvärde markerat med blå streckad linje. I vänster figur är datapunkten innanför konfidensintervallet (grön punkt) och i höger figur är datapunkten utanför intervallet (röd punkt).



Figur 4: Täthetsfunktion för en binomialfördelad stokastisk variabel med sannolikhet för positivt utfall  $p = 0.05$  och antal försök  $n = 500$

## 5 Diskussion

I detta avsnitt görs jämförelser mellan OFDM-algoritmen och parametriska metoder och förslag på vidare forskning studeras. Sedan diskuteras kärdensitetsuppskattning och tidskomplexitet för större modeller samt samhälleliga och etiska aspekter utifrån arbetets syfte.

### 5.1 Jämförelse

I detta avsnitt hänvisas det till tabellerna i appendix A.1 om inget annat anges. Vi kan observera från resultaten att OFDM överlag har sämre prestanda än de parametriska metoderna för 100 och 500 datapunkter. Särskilt i fallet med modellen  $a$ , då den har prestanda 0.00 för alla mängder datapunkter. Detta kan bero på olika orsaker. En hypotes är att det beror på låg varians i vår data. I övriga modeller förstärks variansen då den multipliceras med  $x$ . För modellen  $a$  sker inte detta och denna mindre varians av vår data kan vara orsaken till låg varians i modellerna OFDM anpassar.

OFDM verkar ha relativt bra prestanda för högre antal datapunkter och presterar nära 0.95-1.00 då 1 000 datapunkter används. Den polynomiska modellen presterar bra för 100 och 500 datapunkter, men verkar prestera sämre för de flesta modellerna vid 1 000 datapunkter. Metoden med linjärkombination av stokastiska variabler verkar prestera inom det förväntade intervallet i alla tester.

Gällande prestandan för OFDM är det dock värt att nämna att den kan nå maxantalet epoker utan att konvergera. Prestandan skulle därför kunna öka om vi höjer detta maxantal, då den skulle kunna konvergera till bra anpassade modeller i fler fall. Detta påverkar såklart även konvergenstiden för OFDM, då denna metrik också skulle öka om maxantalet epoker ökas och många av regressionerna i nuläget inte konvergerar.

Om vi ser till konvergenstiden i sekunder är den som högst för linjärkombinations-metoden. Detta implicerar att denna är mer beräkningstung, då den verkar ha lägre antal epoker till konvergens jämfört med övriga. Detta kan innebära att den är mindre lämplig att applicera i ett neuralt nätverk, då träning av ett sådant kan ta längre tid än vad som är praktiskt genomförbart. Vi kan dock även observera att konvergenstiden tycks vara relativt oförändrad vid en ökning av antalet parametrar i modellen som skattas. Den skulle därför kunna vara att föredra vid modeller med många parametrar jämfört med OFDM och polynomiska modellen, vars beräkningstid tycks öka med antalet parametrar. För de modeller som studeras här verkar dock OFDM och polynomiska modellen vara att föredra rent tidsmässigt. Men som tidigare nämnts, påverkas den genomsnittliga konvergenstiden av maxantalet epoker om metoderna inte hinner konvergera.

Med avseende på antal epoker till konvergens är linjärkombinations-metoden mest effektiv. Den verkar som sagt inte heller påverkas i sin konvergenstid, även med avseende på antal epoker, av antalet parametrar. Detta skulle kunna förändras vid ett högre antal parametrar, men är intressant att observera. Polynomiska modellen verkar vara den metod som konvergerar efter högst antal epoker av metoderna i jämförelsen. Den verkar dock ha låg beräkningstid per epok, vilket gör antalet epoker mindre intressant i praktiken.

### 5.2 Val för KDE

När KDE i målfunktionen beräknas används varianten från ekvation (25) istället för ekvation (24). Anledningen till valet är att målfunktionen kräver en referens till  $w_i$  för att `backward()` ska kunna beräkna gradienten med avseende på vikterna. Om ekvation (24) används måste ett stickprov genereras vilket bryter kopplingen med vikten  $w_i$ . Genom att använda sannolikhetsmättet direkt i uppskattningen undviks detta problem. Ekvationen uttrycks dock något annorlunda i koden då fördelningen  $X$  för oss följer en modell som generellt kan skrivas  $y = \eta(x, \theta)$  där  $\theta = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m$  är en fördelning med vikten  $w_i$  för ett utfall  $\theta_i$ . Detta ger oss

$$\hat{p}_m(y) = \frac{1}{mh} \sum_{i=1}^m w_i \cdot K\left(\frac{y - \eta(x, \theta_i)}{h}\right). \quad (46)$$

I alla fall där KDE har använts i denna rapport har Silvermans tumregel utnyttjats för att välja lämplig bandbredd  $h$ . Detta fungerar bra då alla mått som använts har varit unimodala normalfördelningar. Försök gjordes även med Sheath and Jones som är en mer generell metod och fungerar i nästan alla exempel, men då erhöles problem vid implementation via biblioteket KDEpy. Paketets funktioner gav enbart värden på  $\hat{p}_m(y)$  för jämnt fördelad data  $y$ , d.v.s. lika stort avstånd mellan alla datapunkter. I fall då  $\hat{p}_m(y)$  sökes för data som kommer från en sannolikhetsfördelning är den inte jämnt fördelad. I och med detta problem och att en egen implementering av Sheath and Jones hade varit alltför tidskrävande gjordes valet att ändå använda Silvermans tumregel.

### 5.3 Tidskomplexitet för större modeller med OFDM

För OFDM metoden betraktas alla möjliga kombinationer av vikter från de olika måtten i beräkningen av felet. Vid regression med upp till tre koefficienter är detta ej ett stort problem men detta komplexitetsproblem lär bli signifikant när antalet mått ökar; antalet möjliga kombinationer av vikter från de olika måtten växer nämligen exponentiellt med antalet mått för nuvarande implementation.

Då neurala nätverk har många koefficienter inblandade innebär detta att motsvarande BNN skulle involvera många fördelningar. Därav är möjliga metoder för att hantera dessa problem en central del i de framtida efterforskningar som behöver göras om denna metod ska bli lämplig för användning i större BNN.

### 5.4 Vidare forskning

Fler kvantitativa tester behöver troligtvis utföras för att se hur OFDM presterar vid regression och för att avgöra dess lämplighet för träning av BNN. Dels undersöks inte algoritmen som en del av träningsprocessen för ett BNN i denna rapport, utan enbart småskalig optimering med ett få antal parametrar studeras. Det hade varit intressant att implementera ett BNN med denna algoritm och se hur nätverket presterar i en relevant tillämpning. Ytterligare vore det relevant att undersöka redan implementerade paket för BNN, exempelvis python biblioteket Pyro [9].

Denna jämförelse har heller inte tagit hänsyn till inte att metoden är beroende av antalet vikter som måtten består av. Det är av intresse att testa hur metoden presterar i de observerade metriker för olika mängder av vikter och då undersöka om det finns ett optimalt antal vikter. Något som också kan variera prestandan är intervallet för måtten som används. I utvecklingen av metoden kunde vi observera att prestandan varierade med längden och placeringen av intervallet relativt parametrarna i modellen som genererade data. Detta fenomen har dock ej kvantifierats i denna rapport och det skulle vara intressant att undersöka vidare. Detta kan troligtvis även samspela med antalet vikter. I testerna anpassades intervallet och antalet vikter för att förbättra anpassningen till testdata. För exakt utseende av denna anpassning, se appendix C.5.1-C.5.4.

Det hade även varit intressant att återskapa testerna med andra datamängder, för att se om detta påverkar prestandan. Möjligheten finns att dessa datamängder leder till resultat som än missvisande för metodernas allmänna prestanda. Att beräkna andra metriker för metoderna som kvantifierar deras prestanda kan också vara intressant. Exempelvis kan ett lågt värde i metriken som vi kallar prestanda bero på antingen en överanpassning, med för låg varians, eller en dålig anpassning. En annan metrik skulle möjligtvis kunna visa detta och då ge en bättre bild av metodernas prestanda.

Att även testa med data som genererats från andra fördelningar än normalfördelningar hade varit av intresse. Eftersom OFDM är icke-parametrisk och inte anpassad till specifikt normalfördelningar, som de parametriska i denna rapport, borde den prestera bättre för icke-normalfördelad data än de

andra metoderna. Även tester av algoritmen för data som ej genererats specifikt för testerna hade varit intressant. Ett exempel skulle vara att nyttja OFDM för att anpassa en modell till någon form av ekonomisk data, så som aktievärden eller fastighetspriser.

## 5.5 Samhälleliga och etiska aspekter

Arbetets undersökning av optimeringen över mått görs med syftet att kunna göra träningen av bayesianska nätverken effektivare. Det är därför viktigt att också studera vilka samhälleliga konsekvenser ett bayesianskt neuralt nätverk kan ge upphov till när de tillämpas.

Neurala nätverk och deras tillämpningar kan i många fall lägga grunden för kraftfulla tekniska verktyg och idag finns det redan implementerade neurala nätverk som används av industri och näringsliv. Europeiska kommissionen har delat in områden, där artificiell intelligens kan tänkas implementeras, i olika riskklasser för att på så sätt kunna bestämma om vissa typer av tillämpningar bör regleras samt vilka krav på transparens som bör uppfyllas. Exempelvis är anställning samt personalhantering områden klassas som högriskområde där implementationer av maskininlärningsprocesser bör bedömas noggrant innan det släpps ut på marknaden [25].

När neurala nätverk implementeras är det dock av stor vikt att överväga konsekvenser som dessa tillämpningar kan ge upphov till, både på samhället i stort såväl som på individnivå. För att undvika att etiska problem uppstår i arbetet med implementeringen av bayesianska neurala nätverk har etiska aspekter som är generella för alla neurala nätverk beaktas under arbetets gång. Många av dessa finns beskrivna på EU-kommisionens hemsida [4].

De etiska aspekter som är generella för neurala nätverk är exempelvis partiskhet i träningsstadiet av de neurala nätverken, säkerhetsmässiga aspekter, integritet och transparens hos nätverken samt ansvar och gottgörelse vid negativ påverkan.

Neurala nätverk tränas med hjälp av stora mängder data och liknas ibland med en svart låda då det kan vara svårt att veta exakt vad som sker i träningsstadiet. Hur väl dessa nätverk sedan kan förutspå resultat på ny indata är helt beroende av vilken data nätverket tränats på. Det tränade nätverket kan i sig bli partiskt om den data som det tränats med är partisk på något vis. Risken för en partiskhet i det resulterande nätverket är viktig att vara medveten om i valet av träningsdata så att denna exempelvis inte innehåller mer data från vissa etniska grupper än andra samt att den är mångfaldig och representativ.

Konsekvensen av ett partiskt nätverk kan vara att orättvisa eller diskriminerande förutsägelser uppstår. Eftersom neurala nätverk och deras tillämpningar skulle kunna användas i många delar av samhället, som exempelvis domstol eller polisiär verksamhet, och därför påverka både individer och samhället i stort är det viktigt att säkerställa att förutsägelser är opartiska. Ett annat exempel på dessa tillämpningar är bankverksamhet där statistiska modeller ibland används för att uppskatta kreditvärdighet [25].

I tillämpningar såsom självkörande bilar eller medicinsk utrustning, där säkerheten är vital, är det viktigt att överväga hur pålitliga och säkra de neurala nätverken är samt att de inte utgör en risk för människor.

Neurala nätverk används för att hämta information från data. Om den data som ska analyseras är känslig, exempelvis finansiell eller hälsomässig information är det viktigt att den skyddas från obehöriga samt att GDPR [1] och andra förordningar följs.

Skulle en tillämpning av ett neuralt nätverk ändå få negativa konsekvenser är det viktigt att säkerställa vilka som bär ansvar för deras skapande och användning. Detta är viktigt för att kunna ge ersättning eller återupprättelse till någon som påverkats negativt.

Den största relevanta skillnaden mellan bayesianska neurala nätverk och traditionella neurala nätverk i detta sammanhang är att BNN låter oss få en bättre förståelse för osäkerheten i förutsägelseerna. Detta låter oss potentiellt även få bättre insikt i hur nätverket fungerar och därav bättre bedöma dess lämplighet för olika tillämpningar.

## 5.6 Slutsats

Sammanfattningsvis lyckades vi implementera OFDM-algoritmen i PyTorch. Jämfört med regression med normalfördelade koefficienter har vi funnit att OFDM ger sämre förutsägelser, särskilt då modellen har tränats på mindre datamängder. Vi har observerat att modeller vilka använder OFDM går snabbare att anpassa jämfört med modeller baserade på normalfördelade mått, specifikt metoden med en linjärkombination av stokastiska variabler. Vi tror att OFDM underskattar den underliggande fördelningens varians, vilket skulle förklara varför den ofta har fler antal datapunkter utanför konfidensintervallet än förväntat, men att detta problem verkar minska i modeller där fler mått ingår. Detta är dock endast en hypotes och resultatet kan ej bekräfta denna hypotes.

Det återstår att se om den förutspådda variansen konvergerar mot det verkliga värdet eller kommer fortsätta öka vid ytterligare expansion av antal mått som ingår i modellen. Ytterligare efterforskningar krävs för att förstå mekanismerna bakom detta fenomen.

Problemet med tidskomplexiteten för beräkning av felet behöver även lösas innan OFDM-algoritmen praktiskt kan integreras i större neurala nätverk.

Överlag anser vi att fler tester behöver göras för att förstå hur OFDM anpassar modeller för olika typer av data och för olika val av parametrar i modellen, exempelvis definitionsintervallen för måtten, antal vikter per mått och initial steglängd. Med implementationen som gjorts i och med denna rapport finns goda möjligheter att genomföra ytterligare tester för att få en bättre insyn i metodens effektivitet och lämplighet i praktisk tillämpning.

## Referenser

- [1] "Allmän dataskyddsförordning", GDPR Info. (april 2016), URL: <https://gdprinfo.eu/sv>.
- [2] M. Andersson och P. Olofsson, *Probability, Statistics and Stochastic Processes*, 2 utg. John Wiley & Sons, 2012, ISBN: 9780470889749.
- [3] N. Andréasson, A. Evgrafov, M. Patriksson m. fl., *An Introduction to Continuous Optimization*, 3. utg. Lund, Sverige: Studentlitteratur, 2016, s. 374–375, ISBN: 9789144115290.
- [4] "Artificiell intelligens: Möjligheter och risker: Nyheter: Europaparlamentet", EU-parlamentet. (maj 2022), URL: <https://www.europarl.europa.eu/news/sv/headlines/society/20200918ST087404/artificiell-intelligens-mojligheter-och-risker>.
- [5] *Automatic differentiation package - torch.autograd*, 2022. URL: <https://pytorch.org/docs/1.13/autograd.html>.
- [6] A. G. Baydin, B. A. Pearlmutter och A. A. Radul, "Automatic differentiation in machine learning: a survey", *CoRR*, årg. abs/1502.05767, 2015. DOI: [10.48550/arXiv.1502.05767](https://doi.org/10.48550/arXiv.1502.05767).
- [7] S. W. Bernard, *Density Estimation for Statistics and Data Analysis*. Bristol: J.W. Arrowsmith Ltd, 1986, ISBN: 9780412246203.
- [8] S. Bharali och J. Hazarika, "Regression models with stochastic regressors: An expository note", *International Journal of Agricultural and Statistical Sciences*, årg. 15, nr 2, s. 873–880, 2019. URL: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85077907409&partnerID=40&md5=1dd7db7534130bfda6b170156d83e876>.
- [9] E. Bingham, J. P. Chen, M. Jankowiak m. fl., "Pyro: Deep Universal Probabilistic Programming", *Journal of Machine Learning Research*, 2018. DOI: [10.48550/arXiv.1810.09538](https://doi.org/10.48550/arXiv.1810.09538).
- [10] C. M. Bishop, "Neural networks and their applications", *Review of Scientific Instruments*, årg. 65, nr 6, s. 1803–1832, juni 1994. DOI: [10.1063/1.1144830](https://doi.org/10.1063/1.1144830).
- [11] Y.-C. Chen, "A tutorial on kernel density estimation and recent advances", *Biostatistics & Epidemiology*, årg. 1, s. 161–187, 2017. DOI: [10.1080/24709360.2017.1396742](https://doi.org/10.1080/24709360.2017.1396742).
- [12] C. Chu, J. Blanchet och P. Glynn, *Probability Functional Descent: A Unifying Perspective on GANs, Variational Inference, and Reinforcement Learning*, jan. 2019. DOI: [10.48550/arXiv.1901.10691](https://doi.org/10.48550/arXiv.1901.10691).
- [13] G. B. Folland, *Real Analysis: Modern Techniques and Their Applications*, 2 utg. John Wiley & Sons, 1999, s. 19–41, ISBN: 0471317160.
- [14] T. L. Foundation, *Pytorch*, <https://pytorch.org>. (hämtad 2023-01-23).
- [15] J. Gergonne, "The application of the method of least squares to the interpolation of sequences", *Historia Mathematica*, årg. 1, nr 4, s. 439–447, 1974. DOI: [10.1016/0315-0860\(74\)90034-2](https://doi.org/10.1016/0315-0860(74)90034-2).
- [16] T. Gneiting och A. E. Raftery, "Strictly Proper Scoring Rules, Prediction, and Estimation", *Journal of the American Statistical Association*, årg. 102, nr 477, s. 359–378, 2007. DOI: [10.1198/016214506000001437](https://doi.org/10.1198/016214506000001437).
- [17] B. Hanin och M. Sellke, *Approximating Continuous Functions by ReLU Nets of Minimal Width*, 2018. DOI: [10.48550/arXiv.1710.11278](https://doi.org/10.48550/arXiv.1710.11278).
- [18] K. Hornik, M. Stinchcombe och H. White, "Multilayer feedforward networks are universal approximators", *Neural Networks*, årg. 2, nr 5, s. 359–366, 1989. DOI: [10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8).
- [19] D. S. Lemons, *An Introduction to Stochastic Processes in Physics*. The John Hopkins University Press, 2002, s. 34, ISBN: 0801868661.
- [20] I. Molchanov och S. Zuyev, "Tangent Sets in the Space of Measures: With Applications to Variational Analysis", *Journal of Mathematical Analysis and Applications*, årg. 249, nr 2, s. 539–552, 2000. DOI: [10.1006/jmaa.2000.6906](https://doi.org/10.1006/jmaa.2000.6906).

- [21] I. Molchanov och S. Zuyev, "Steepest descent algorithms in a space of measures", *Statistics and computing*, årg. 12, nr 2, s. 115–123, 2002. DOI: [10.1023/A:1014878317736](https://doi.org/10.1023/A:1014878317736).
- [22] "Parametric and non-parametric distributions", EPIX analytics. (), URL: <https://modelassist.epixanalytics.com/display/EA/Parametric+and+non-parametric+distributions>.
- [23] K. Pearson och F.R.S., "X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling", *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, årg. 50, nr 302, s. 157–175, 1900. DOI: [10.1080/14786440009463897](https://doi.org/10.1080/14786440009463897).
- [24] S. J. Sheather och M. C. Jones, "A Reliable Data-Based Bandwidth Selection Method for Kernel Density Estimation", *Journal of the Royal Statistical Society. Series B (Methodological)*, årg. 53, nr 3, s. 683–690, 1991. URL: <http://www.jstor.org/stable/2345597> (hämtad 2023-04-11).
- [25] "Spetskompetens inom och förtroende för AI", EU-kommissionen. (2021), URL: [https://commission.europa.eu/strategy-and-policy/priorities-2019-2024/europe-fit-digital-age/excellence-and-trust-artificial-intelligence\\_sv](https://commission.europa.eu/strategy-and-policy/priorities-2019-2024/europe-fit-digital-age/excellence-and-trust-artificial-intelligence_sv).
- [26] Z. Wang och M. Wang, "Chi-square Loss for Softmax: an Echo of Neural Network Structure", *CoRR*, årg. abs/2108.13822, 2021. arXiv: [2108.13822](https://arxiv.org/abs/2108.13822). URL: <https://arxiv.org/abs/2108.13822>.
- [27] E. Weisstein, *Weak Law of Large Numbers*, MathWorld A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/WeakLawofLargeNumbers.html>.
- [28] Y. Wong, *Why You Should Use Bayesian Neural Network*, okt. 2021. URL: <https://towardsdatascience.com/why-you-should-use-bayesian-neural-network-aaf76732c150>.



# A Appendix 1 - Resultattabeller

## A.1 Resultat av tester

I tabellerna nedan återfinns resultaten från jämförelsen beskriven i avsnitt 4.1. Tabellerna anger prestanda och 95% konfidensintervall för dessa värden, samt konvergenstid i sekunder och epoker, för samtliga metoder och regressionsmodeller. Konvergenstiderna är medelvärden för 50 tester.

Optimeringsalgorithm över finita diskreta mått (OFDM)

Antal datapunkter	Modell	Prestanda och 95% CI	Konvergenstid (s)	Konvergenstid (Epoker)
100	$\alpha_i$	0.00 [0.00 0.06]	0.487	539.36
	$\alpha_i \cdot x$	0.52 [0.38 0.66]	3.195	2543.62
	$\alpha_i + \beta_i \cdot x$	0.48 [0.34 0.62]	2.219	1546.98
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.84 [0.74 0.94]	2.945	1205.00
500	$\alpha_i$	0.00 [0.00 0.06]	3.251	3394.98
	$\alpha_i \cdot x$	0.58 [0.44 0.72]	4.908	3837.14
	$\alpha_i + \beta_i \cdot x$	0.66 [0.53 0.79]	5.624	3572.5
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.86 [0.76 0.96]	13.030	3467.70
1000	$\alpha_i$	0.00 [0.00 0.06]	3.731	3817.72
	$\alpha_i \cdot x$	0.70 [0.57 0.83]	4.863	3645.60
	$\alpha_i + \beta_i \cdot x$	0.90 [0.82 0.98]	6.927	3961.00
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.92 [0.84 1.00]	18.652	3850.54

Polynomiell regression

Antal datapunkter	Modell	Prestanda och 95% CI	Konvergenstid (s)	Konvergenstid (Epoker)
100	$\alpha_i$	1.00 [0.94 1.00]	4.358	5483.66
	$\alpha_i \cdot x$	0.98 [0.94 1.00]	8.323	16030.50
	$\alpha_i + \beta_i \cdot x$	1.00 [0.94 1.00]	9.500	11333.84
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.98 [0.94 1.00]	12.487	14667.54
500	$\alpha_i$	0.98 [0.94 1.00]	5.530	6515.72
	$\alpha_i \cdot x$	0.88 [0.79 0.97]	8.095	10815.06
	$\alpha_i + \beta_i \cdot x$	0.92 [0.84 1.00]	10.834	12415.40
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.84 [0.74 0.94]	14.421	16834.62
1000	$\alpha_i$	1.00 [0.94 1.00]	6.260	6947.02
	$\alpha_i \cdot x$	0.72 [0.60 0.84]	9.407	29693.18
	$\alpha_i + \beta_i \cdot x$	0.74 [0.62 0.86]	10.559	12516.06
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	0.70 [0.57 0.83]	16.696	20543.82

Linjärkombination av normalfördelade stokastiska variabler

Antal datapunkter	Modell	Prestanda och 95% CI	Konvergenstid (s)	Konvergenstid (Epoker)
100	$\alpha_i$	1.00 [0.94 1.00]	4.376	308.14
	$\alpha_i \cdot x$	1.00 [0.94 1.00]	4.443	303.36
	$\alpha_i + \beta_i \cdot x$	1.00 [0.94 1.00]	4.317	300.12
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	1.00 [0.94 1.00]	4.765	336.6
500	$\alpha_i$	1.00 [0.94 1.00]	25.156	311.86
	$\alpha_i \cdot x$	1.00 [0.94 1.00]	25.487	308.20
	$\alpha_i + \beta_i \cdot x$	1.00 [0.94 1.00]	41.128	521.64
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	1.00 [0.94 1.00]	27.529	344.16
1000	$\alpha_i$	1.00 [0.94 1.00]	53.383	311.78
	$\alpha_i \cdot x$	1.00 [0.94 1.00]	54.570	313.70
	$\alpha_i + \beta_i \cdot x$	0.98 [0.94 1.00]	51.743	302.84
	$\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$	1.00 [0.94 1.00]	54.917	323.52

## B Appendix 2 – Teori

### B.1 Degenererad summa av väntevärden

Denna härledning utarbetades tillsammans med handledaren. Antag att vi försöker hitta en fördelning  $\theta$  så att  $f(x, \theta)$  är nära  $y$ , alltså  $|f(x, \theta) - y|$  är liten, där  $\mu$  är den sanna fördelningen. Vi försöker alltså lösa

$$\min \sum_i E(f(x_i, \theta_i) - y_i)^2. \quad (47)$$

För ett generellt problem kan vi skriva  $y_i = f(x_i, \theta_i)$  där  $\theta_i \sim \mu$ . Problemet kan formuleras om till

$$\min \sum_i E_\theta (f(x_i, \theta) - y_i)^2, \quad (48)$$

$$\min \sum_i \int (f(x_i, \theta) - y_i)^2 d\mu, \quad (49)$$

$$\min \left\langle \underbrace{\sum_i (f(x_i, \cdot) - y_i)^2}_{g(\theta)}, \mu \right\rangle, \quad (50)$$

där (50) är en linjär funktion med avseende på  $\mu$  och  $g(\theta)$  är gradienten. Vi får alltså nödvändiga och tillräckliga villkor för minimeringen som

$$\sum_i (f(x_i, \theta) - y_i)^2 = \begin{cases} = C, \mu, \\ \geq C, \forall \theta. \end{cases} \quad (51)$$

Lösningen till problemet är att distribuera  $\mu$  över

$$\theta^* = \arg \min_\theta \sum_i (f(x_i, \theta) - y_i)^2. \quad (52)$$

För regression får vi

$$\min \sum_i (\alpha + \beta x_i - y_i)^2, \quad (53)$$

och  $\mu$  är en degenererad lösningen till regressionsproblemet i punkten  $(a, b)$ .

### B.2 Linjärkombination av stokastiska variabler

Följande teori är opublicerad men tillskrivs handledare Sergei Zuyev.

Låt  $\eta(x)$  vara en funktion definierad enligt

$$\eta(x) = \sum_{k=1}^K \alpha^{(k)} h_k(x), \quad (54)$$

där  $\alpha^{(k)} \sim \mu_k$  är oberoende stokastiska variabler och  $h_k(x)$  är funktioner av  $x$ . Om vi antar att  $\mu_k(dx) = f_k(x)dx$ , där  $f_k(x)$  är täthetsfunktioner och att  $h_k$  är kontinuerliga  $\forall k$ , då har  $\eta(x)$  kumulativa täthetsfunktion

$$F_{\eta(x)}(y) = P(\eta(x) \leq y). \quad (55)$$

Med omskrivning blir den kumulativa täthetsfunktionen

$$F_{\eta(x)}(y) = P \left( \alpha^{(1)} h_1(x) \leq y - \sum_{k=2}^K \alpha^{(k)} h_k(x) \right), \quad (56)$$

$$F_{\eta(x)}(y) = E \left[ F_1 \left( \frac{y - \sum_{k=2}^K \alpha^{(k)} h_k(x)}{h_1(x)} \right) \middle| \alpha^{(2)}, \dots, \alpha^{(K)} \right], \quad (57)$$

$$F_{\eta(x)}(y) = \int \dots \int_{\mathbb{R}^{K-1}} F_1 \left( \frac{y - \sum_{k=2}^K \alpha^{(k)} h_k(x)}{h_1(x)} \right) f_2(\alpha^{(2)}) \dots f_K(\alpha^{(K)}) d\alpha^{(2)} \dots d\alpha^{(K)}. \quad (58)$$

Täthetsfunktionen kan beskrivas som derivatan av den kumulativa täthetsfunktionen

$$f_{\eta(x)}(y) = \frac{d}{dy} F_{\eta(x)}(y), \quad (59)$$

$$\Rightarrow f_{\eta(x)}(y) = \int \dots \int_{\mathbb{R}^{K-1}} f_1 \left( \frac{y - \sum_{k=2}^K \alpha^{(k)} h_k(x)}{h_1(x)} \right) f_2(\alpha^{(2)}) \dots f_K(\alpha^{(K)}) d\alpha^{(2)} \dots d\alpha^{(K)}. \quad (60)$$

Sannolikheten för att observera en viss data  $(x_1, y_1), \dots, (x_N, y_N)$  givet en viss skattning  $\bar{y}$  ges av

$$\text{Lik}(\bar{y}) = \prod_{i=1}^N f_{\eta(x_i)}(y_i), \quad (61)$$

vilket beräknas enligt (60). För en bra skattning vill vi maximera sannolikheten att observera given data, alltså erhålls lösningen som

$$\max \text{Lik}(\bar{y}) = \max \prod_{i=1}^N f_{\eta(x_i)}(y_i). \quad (62)$$

### B.3 Yttre strafffunktion

Metoden bygger på att bestraffa målfunktionen när ett kriterium ej uppfylls med hjälp av en strafffunktion. Enligt [3] kan den önskade mängden av  $x$  formuleras som en kombination av likheter och olikheter enligt:

$$S := \{x \in \bar{\mathbb{R}}^n \mid g_i(x) \leq 0, \quad i = 1, \dots, m \\ h_j(x) = 0, \quad j = 1, \dots, l\} \quad (63)$$

Vilket sedan kan approximeras den yttre strafffunktionen [3]:

$$\chi_S(x) \approx \nu \hat{\chi}(x) := \nu \left( \sum_{i=1}^m (\max\{0, g_i(x)\})^2 + \sum_{j=1}^l (h_j(x))^2 \right) \quad (64)$$

Med denna metod kan begränsningen  $0 \leq p \leq 1$  formuleras enligt följande:

$$\nu(\max(0, p-1)^2 + \max(0, -p)^2), \quad (65)$$

där  $\nu$  är en skalfaktor för straffvärdet. Genom att addera strafffunktionen till målfunktionen erhålls följande ekvation:

$$L = p(1-p) + \nu(\max(0, p-1)^2 + \max(0, -p)^2) \quad (66)$$

som kan styra optimeringen till att konvergera inom det önskade området eftersom  $\max()$  termerna kommer vara nollskilda då  $0 \leq p \leq 1$  ej uppfylls. Denna metod kan framgångsrikt implementeras för det enkla fallet med en parameter men problem uppstår då antalet parametrar ökar. Konvergensen är mycket känslig för straffaktorn  $\nu$  och resultatet har en stor felmarginal. Detta fel skulle kunna minskas genom att dynamiskt justera steglängden då parametern närmar sig det korrekta värdet.

## B.4 Lagrangemultiplikator-metoden

En metod som utvärderades i början av projektet var metoden med Lagrangemultiplikatorer som är vanlig inom matematisk optimering. Metoden konverterar ett begränsat problem till ett problem utan begränsningar, så att metoder för obegränsad optimering, så som gradient descent, kan användas. Metoden utgår från att först bilda den så kallade Lagrangefunktionen

$$\mathcal{L}(x, \lambda) = f(x) + \sum_k \lambda_k g_k(x), \quad (67)$$

där  $f(x)$  är funktionen som vi vill minimera,  $g_k(x)$  är en begränsning,  $\lambda_k$  är en så kallad Lagrange-multiplikator,  $k = 1, \dots, n$ , där  $n$  är antalet begränsningar för problemet. Det kan visas att de stationära punkterna till Lagrangefunktionen kommer sammanfalla med de lokala extrempunkterna för  $f(x)$ .

Anledningen till att vi inte fortsatte arbetet med denna metoden var att den globala extrempunkten inte alltid är en stationärpunkt, exempelvis då extrempunkten ligger på randen av det begränsade området.

## B.5 Icke-parametriska diskreta sannolikhetsmått med softmax

Softmax funktionen är definierad enligt

$$(\text{Softmax}(\mathbf{x}))_i = \frac{e^{x_i}}{\sum_j e^{x_j}},$$

där  $\mathbf{x}$  är vektorn med element  $x_1, x_2, \dots$ . Det är alltså en funktion som tar en vektor och ger tillbaka en vektor med lika många element, där varje element har transformerats enligt ovan. Detta skapar en vektor med element i intervallet  $[0, 1]$  vilka har summan 1.

Ett exempel där detta kan appliceras är vid optimering av någonting som beror av en Bernoullifördelning. Bernoullifördelning är den fördelningen som ger värdet 1 med sannolikhet  $p$  och värdet 0 med sannolikhet  $q = 1 - p$ . Här måste båda världena vara mellan noll och ett, och summera till ett. Parameterrummet kan med andra ord betraktas som ett diskret sannolikhetsmått med två möjliga utfall. Om vi vill använda oss av vanlig sjunkgradient kan vi låta våra parametrar vara obegränsade reella tal och sedan applicera softmax på dem då vi vill finna  $p$  och  $q$ .

Det problem som uppstår är att det tar extremt många iterationer för optimerings-algoritmen att komma nära fallet där till exempel  $q$  är noll. Ett exempel där vi stöter på detta problem är minimeringen av Bernoullifördelningens varians; denna varians uttrycks algebraiskt  $pq$ , och kräver då att  $p$  eller  $q$  blir noll för att minimeras. För att uppnå detta behöver skillnaden mellan parametrarna (före applicering av softmax) vara oändlig. Därav kan det ta väldigt många iterationer för att uppnå en bra approximation av de optimala parametervärdena.

I detta fall kan vi få en bättre prestanda genom att sätta en högre steglängd (eng. Learning rate), men detta kan vara problematiskt vid mer komplicerade problem då målfunktionen kräver mer än att tvinga en parameter till noll för att finna minimum.

## B.6 Poängfunktioner och poängregler

*I följande underavsnitt är källan till teorin [16] om inget annat anges.*

Vi har ett utfallsrum  $\Omega$ , en  $\sigma$ -algebra  $\mathcal{A}$  av delmängder på  $\Omega$  och  $\mathcal{P}$  en konvex klass av mått på  $(\Omega, \mathcal{P})$ . En funktion  $S$  är en poängregel (eng. scoring rule) om  $S : \mathcal{P} \times \Omega \rightarrow \mathbb{R}$ , där  $\mathbb{R} = [-\infty, \infty]$  är den utökade reella tallinjen. Om  $\omega$  är ett utfall i  $\Omega$  ger  $S(P, \omega)$  belöningen för vår förutsägelse  $P$  när  $\omega$

materialiserar. Den förväntade belöningen av förutsägelsen  $P$  när  $Q$  är den sanna fördelningen skriver vi som

$$S(P, Q) := E_Q[S(P, \cdot)] = \int S(P, \omega) Q(d\omega). \quad (68)$$

En poängregel är en funktion som sätter ett mått på hur väl en sannolikhetsfördelning förutsäger utfall. Exempelvis studeras de i syfte av att ge en rättvis evaluering av förutsägelser inom meteorologi.

**Definition B.1.** En poängregel sägs vara *äkta* (eng. proper) om

$$S(Q, Q) \geq S(P, Q), \quad \forall P, Q \in \mathcal{P}. \quad (69)$$

Ytterligare sägs en poängregel vara *strikt äkta* om och endast om likhet i (69) uppfylls då  $P = Q$ .

En äkta poängregel ger en ärlig belöning. Oavsett vad den sanna fördelningen  $Q$  är finns det ingen annan fördelning  $P$  som är bättre, och för en strikt äkta poängregel är det endast den sanna fördelningen som maximerar  $S$ .

**Definition B.2.** En poängregel  $S : \mathcal{P} \times \Omega \rightarrow \bar{\mathbb{R}}$  sägs vara *reguljär* (eng. regular) med avseende på  $\mathcal{P}$  om  $S(P, Q)$  är reell värd för alla  $P, Q \in \mathcal{P}$  förutom möjligen att  $S(P, Q) = -\infty$  om  $P \neq Q$ .

**Definition B.3.** Om en poängregel  $S$  är äkta och reguljär kallar vi

$$d(P, Q) = S(Q, Q) - S(P, Q), \quad P, Q \in \mathcal{P}. \quad (70)$$

den associerade *divergensfunktionen* (eng. divergence function).

Divergensfunktionen är icke-negativ och då poängregeln är strikt äkta är divergensfunktionen strikt positiv förutom då  $P = Q$ . Vid optimering används avstånd för att mäta hur nära målet skattningen är den verkliga datan. Divergensfunktionen delar likheter med avstånd vilket motiverar teorin för poängregler. Då en divergensfunktion associerar med en äkta poängregel styrker det användningen denna som poängregel.

### B.6.1 Exempel på poängregler

Ett exempel är Brier score. Om  $\mathbf{p} = \{p_1, \dots, p_n\}$  är en sannolikhetsfördelning är poängregeln

$$S(\mathbf{p}, i) = - \sum_{j=1}^n (\delta_{ij} - p_j)^2 = 2p_i - \sum_{j=1}^n p_j^2 - 1, \quad i \in 1 \dots n, \quad (71)$$

där  $\delta_{ij} = 1$  då  $i = j$  och  $\delta_{ij} = 0$  då  $i \neq j$ . Divergensfunktionen för Brier score är

$$d(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^n (p_j - q_j)^2. \quad (72)$$

Ett annat exempel är logaritmisk poäng regeln. Poängregeln blir

$$S(\mathbf{p}, i) = \log p_i, \quad (73)$$

och divergensfunktionen blir

$$d(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^n q_j \log(q_j/p_j). \quad (74)$$

Notera att om förutsägelsen påstår att  $p_i = 0$  men utfall  $i$  ändå uppstår blir  $S(\mathbf{p}, i) = -\infty$ .

## C Appendix 3 – Källkod

För att se mer exakt hur filerna är uppbyggda och vilka bibliotek som importeras i de olika koderna se vårt GitHub-projekt: <https://github.com/Miralleyan/Kandidatarbete>

### C.1 Mättklassen

```
1 class Measure:
2     def __init__(self, locations: torch.tensor, weights: torch.tensor, device='cpu', optim_locations = False):
3         self.locations = locations
4         self.weights = torch.nn.parameter.Parameter(weights)
5         self.device = device
6
7     def __str__(self) -> str:
8         """
9         Returns the locations and weights of the measure as a string.
10        :returns: str
11        """
12        out = "\033[4mLocations:\033[0m".ljust(28) + "\033[4mWeights:\033[0m \n"
13        for i in range(len(self.weights)):
14            out += f'{self.locations[i].item():<20.9f}{self.weights[i].item():<.9f}\n'
15        return out
16
17    def __repr__(self) -> str:
18        """
19        Returns the locations and weights of the measure as a string.
20        :returns: str
21        """
22        return self.__str__()
23
24    def is_probability(self, tol=1e-6):
25        """
26        Returns True if the measure is a probability measure.
27        :param tol: Tolerance for summing to 1
28        """
29        if torch.any(self.weights < 0):
30            return False
31        if torch.abs(self.weights.sum() - 1) > tol:
32            return False
33        return True
34
35    def total_mass(self) -> float:
36        """
37        Returns the sum of all weights in the measure:  $\sum_{i=1}^n w_i$ 
38        :returns: float
39        """
40        return sum(self.weights).item()
41
42    def total_variation(self) -> float:
43        """
44        Returns the sum of the absolute value of all weights in the measure:  $\sum_{i=1}^n |w_i|$ 
45        :returns: float
46        """
47        return sum(abs(self.weights)).item()
48
49    def support(self, tol=5e-3):
50        """
51        :param tol: proportion of total variation that can be un-accounted for by the support.
52        :returns: all index where the weights are non-zero
53        """
54        sorted_idx = torch.argsort(self.weights.abs())
55        accum_weight = torch.cumsum(self.weights[sorted_idx].abs(), dim=0)
56        cutoff = tol * self.total_variation()
57        return sorted_idx[cutoff < accum_weight]
```

```

58
59 def positive_part(self):
60     """
61     Returns the positive part of the Lebesgue decomposition of the measure
62     """
63     return Measure(self.locations, torch.max(self.weights, torch.zeros(len(self.weights))))
64
65 def negative_part(self):
66     """
67     Returns the negative part of the Lebesgue decomposition of the measure
68     """
69     return Measure(self.locations, -torch.min(self.weights, torch.zeros(len(self.weights))))
70
71 def sample(self, size):
72     """
73     Returns a sample of indeces from the locations of the measure
74     given by the distribution of the measures weights
75     :param size: Number of elements to sample
76     :returns: sample of indeces for random numbers based on measure
77     """
78     if torch.any(self.weights < 0):
79         assert ValueError("You can't have negative weights in a probability measure!")
80
81     sample_idx = torch.multinomial(self.weights, size, replacement=True)
82     sample = self.locations[sample_idx] + torch.randn(size)*(self.locations[1]-self.locations[0])/2
83     return sample
84
85 def copy(self):
86     return Measure(self.locations, self.weights)
87
88 def zero_grad(self):
89     self.weights.grad = None
90
91 def visualize(self):
92     """
93     Visualization of the weights
94     """
95     plt.bar(self.locations.detach(), self.weights.detach(), width=0.1, label="Measure")
96     plt.axhline(y=0, c="grey", linewidth=0.5)
97     plt.draw()

```

## C.2 Optimeringsklassen

```
1 class Optimizer:
2
3     def __init__(self, measures, loss : str, lr : float = 0.1):
4         # Create list of measures
5         if type(measures) == Measure:
6             self.measures = [measures]
7         elif type(measures) != list:
8             Exception('Error: measures has to be of type Measure or list')
9         else:
10            self.measures = measures
11        # Create list of lr's
12        if type(lr) == float or type(lr) == int:
13            self.lr = [lr]*len(self.measures)
14        elif type(lr) != list:
15            Exception('Error: lr has to be of type float or list')
16        else:
17            self.lr = lr
18
19        loss_dict = {'essr':self.essr, 'nll':self.nll, 'KDEnll':self.KDEnll, 'chi_squared':self.chi_squared}
20        self.loss = loss_dict[loss]
21        self.state = {'measure':self.measures, 'lr':self.lr}
22        self.is_optim = False
23
24    def put_mass(self, meas_index, mass, location_index):
25        """
26        In current form, this method puts mass at a specified location, s.t. the location still
27        has mass less at most 1 and returns how much mass is left to distribute.
28        :param mass: Mass left to take
29        :param location_index: Index of location to take mass from
30        :returns: mass left to add to measure after adding at specified location
31        """
32        with torch.no_grad():
33            self.measures[meas_index].weights[location_index] += mass
34
35    def take_mass(self, meas_index, mass, location_index) -> float:
36        """
37        In current form, this method takes mass from a specified location, s.t. the location still
38        has non-negative mass and returns how much mass is left to take.
39        :param mass: Mass left to take
40        :param location_index: Index of location to take mass from
41        :returns: mass left to remove from measure after removing from specified location
42        """
43        with torch.no_grad():
44            if mass > self.measures[meas_index].weights[location_index].item():
45                mass_removed = self.measures[meas_index].weights[location_index].item()
46                self.measures[meas_index].weights[location_index] = 0.
47            else:
48                self.measures[meas_index].weights[location_index] -= mass
49                mass_removed = mass
50        return mass_removed
51
52    def stop_criterion(self, tol_supp=1e-6, tol_const=1e-2, adaptive = False):
53        """
54        Checks if the difference between the maximum and minimum gradient is
55        within a certain range.
56        :param tol_supp: lower bound for wieghts considered
57        :param tol_const: stop value, when the maximum difference of gradients
58        is smaller than this value the minimization should seize
59        """
60        if adaptive:
61            return min([measure.weights.grad[measure.support(tol_supp)].max()
62                        - measure.weights.grad.min() < tol_const*(measure.weights.grad.max() -
63                          measure.weights.grad.min()) for measure in self.measures])
64        else:
```



```

65         return min([measure.weights.grad[measure.support(tol_supp)].max()
66                    - measure.weights.grad.min() < tol_const for measure in self.measures])
67
68     def step(self, meas_index):
69         """
70         Steepest decent with fixed total mass
71         """
72
73         # Sort gradient
74         grad_sorted = torch.argsort(self.measures[meas_index].weights.grad)
75
76         # Distribute positive mass
77         mass_pos = self.lr[meas_index]
78         self.put_mass(meas_index, mass_pos, grad_sorted[0].item())
79
80         # Distribute negative mass
81         mass_neg = self.lr[meas_index]
82         for i in torch.flip(grad_sorted, dims=[0]):
83             mass_neg -= self.take_mass(meas_index, mass_neg, i.item())
84             if mass_neg <= 0.:
85                 break
86
87     def update_lr(self, index, fraction=0.7):
88         """
89         Updates learning rate for the optimizer
90
91         :param fraction: multiply lr with this value
92         """
93         self.lr[index]*=fraction
94
95     def state_dict(self):
96         """
97         Updates the state dictionary for the optimizer
98         """
99         print("\n".join("\t{:}": {}".format(k, v) for k, v in self.state.items()))
100        return self.state
101
102     def load_state_dict(self, state_dict):
103         """
104         Overloads the current state dictionary
105
106         :param state_dict: state dictionary to load
107         """
108         self.state = state_dict
109
110     def lr_decrease_criterion(self, loss_fn, measure, old_measure):
111         """
112         Checks if learning rate should be decreased
113
114         :param loss_fn: loss function
115         :param measure: measure to compare current measure to
116         """
117         return loss_fn(old_measure) < loss_fn(measure)
118
119     def minimize(self, data, model, h = 0, alpha = 0.001, max_epochs=2000, smallest_lr=1e-6, verbose=False,
120                tol_supp=1e-6, tol_const=1e-2, print_freq=100, adaptive=False, test=False):
121         """
122         :param data: list of tensors of data points. If x and y, then data should be on the form
123                    [x_tensor,y_tensor]. If only one series of data points, just this tensor is needed
124         :param model: model written as a function.
125         Should take a data input (x) and a set of parameters (params).
126         :param max_epochs: Max number of iterations
127         :param smallest_lr: Minimizer wil stop when lr is below this value
128         :param verbose: Print information about each epoch
129         :param print_freq: How frequently the minimizer should print information

```

```

130     :param tol_supp:
131     :param tol_const:
132     :param adaptive:
133     :return: Optimized measures
134     """
135     lr = self.lr
136     loss_fn = self.loss
137
138     if type(data) == torch.tensor:
139         data = [data, data]
140
141     perms, prep = self.prep_step(data, model, h, alpha)
142
143     if test:
144         tid=[]
145         LossNotChanged=0
146         t1=time.time()
147
148     for epoch in range(max_epochs):
149         # Backup current measures and reset grad
150         old_measures = copy.deepcopy(self.measures)
151         for m in self.measures:
152             m.zero_grad()
153
154
155         # Compute loss and grad
156         loss_old = loss_fn(perms, *prep)
157         loss_old.backward()
158
159         # Stop criterion
160         if self.stop_criterion(tol_supp, tol_const, adaptive):
161             print(f'\nOptimum is attained. Loss: {loss_old}. Epochs: {epoch} epochs.')
162             self.is_optim = True
163             if test:
164                 t2=time.time()
165                 return self.measures,t2-t1,epoch
166             else:
167                 return self.measures
168
169         if min(lr) < smallest_lr:
170             print(f'The step size is too small: {lr}')
171             if test:
172                 t2=time.time()
173                 if LossNotChanged<5:
174                     return self.measures,t2-t1,epoch
175                 else:
176                     return self.measures, tid[0][0],tid[0][1]
177             else:
178                 return self.measures
179
180         # Step
181         maxima = []
182         for meas_index in range(len(self.measures)):
183             sup_index = self.measures[meas_index].support()
184             grads = copy.deepcopy(self.measures[meas_index].weights.grad)
185             #print(grads)
186             maxima.append(torch.max(grads[sup_index]))
187         max_index = maxima.index(sorted(maxima)[-1])
188         self.step(max_index)
189
190         loss_new = loss_fn(perms, *prep)
191         loss_new.backward()
192
193         # bad step
194         if loss_old < loss_new:

```

```

195         # Revert to the backup measure and decrease lr
196         self.measures = copy.deepcopy(old_measures)
197         self.update_lr(max_index, fraction=0.1)
198
199         if verbose:
200             print(f'Epoch: {epoch:<10} Lr was reduced to: {lr}')
201     elif loss_old == loss_new and verbose:
202         if test and LossNotChanged<5:
203             LossNotChanged+=1
204             t2=time.time()
205             tid.append([t2-t1,epoch])
206             print(f'Epoch: {epoch:<10} Loss did not change ({loss_new}'))
207
208     # successful step
209     else:
210         if test and LossNotChanged<5:
211             LossNotChanged=0
212             tid=[]
213         if epoch % print_freq == 0:
214             if verbose:
215                 print(f'Epoch: {epoch:<10} Loss: {loss_new:<10.9f} LR: {lr}')
216             else:
217                 print('.')
218
219
220     print('Max epochs reached')
221     if test:
222         t2=time.time()
223         if LossNotChanged<5:
224             return self.measures,t2-t1,epoch
225         else:
226             return self.measures, tid[0][0],tid[0][1]
227     else:
228         return self.measures
229
230 def visualize(self):
231     """
232     Visualizes the measures of the optimizer. Currently requires that
233     a gradient has been stored in the weights of the measures.
234     """
235     cols = int(torch.ceil(torch.sqrt(torch.tensor(len(self.measures)))).item())
236     rows = int(torch.ceil(torch.tensor(len(self.measures)/cols)).item())
237     fig, axs = plt.subplots(rows,cols)
238     axs = axs.flatten()
239     fig.suptitle('Optimizer Visualization')
240     grads = [measure.weights.grad for measure in self.measures]
241     for i, measure in enumerate(self.measures):
242         M, m = grads[i].max(), grads[i].min()
243         wm = measure.weights.max()
244         scaled_weights = measure.weights * (M - m) / wm * 0.25 + m
245         support = measure.support()
246         with torch.no_grad():
247             # Support locations
248             axs[i%cols+(i//cols)*cols].plot(measure.locations[support], torch.zeros(len(measure.weights))[support]
249             + m, '.', c='red', label=' Measure Support')
250
251             # Gradient
252             axs[i%cols+(i//cols)*cols].plot(measure.locations, grads[i], c='green', label=' Gradient')
253             # Measure weights where there is support
254             axs[i%cols+(i//cols)*cols].vlines(measure.locations[support], torch.zeros(len(measure.weights))[support]
255             + 1.3 * m, scaled_weights[support], colors='blue', label=' Measure Weights')
256             axs[i%cols+(i//cols)*cols].axhline(y=m, c="orange", linewidth=0.5)
257             axs[i%cols+(i//cols)*cols].legend(loc='upper right')
258             axs[i%cols+(i//cols)*cols].set_ylim([m, M])
259
260 # Loss functions

```

```

260 def essr(self, perms, errors):
261     """
262     Calculates the expected sum of square residuals loss function
263
264     :param perms: list of the possible permutations of one location from each measure
265     :param errors: tensor of the sum of errors, compared with true data, for each permutation of locations
266     """
267     probs = torch.cat([self.measures[i].weights[perms[:, i]].unsqueeze(1) for i in range(len(self.measures))],
268                       1).prod(1)
269     return errors.dot(probs)
270
271 def nll(self, perms, loc_index):
272     """
273     Calculates the negative log-likelihood loss function
274
275     :param perms: list of the possible permutations of one location from each measure
276     :param loc_index: list of location permutation closest with the least absolute error for each data point
277     """
278     probs = torch.cat([self.measures[i].weights[perms[:, i]].unsqueeze(1) for i in range(len(self.measures))],
279                       1).prod(1)
280     return -sum(torch.log(probs[loc_index]))
281
282 def KDEnll(self, perms, kde_mat, h):
283     """
284     Calculates the negative log-likelihood loss function, with a KDE
285
286     :param perms: list of the possible permutations of one location from each measure
287     :param kde_mat: matrix of kernels for each location permutation
288     :param h: bandwidth for the KDE
289     """
290     probs = torch.cat([self.measures[i].weights[perms[:, i]].unsqueeze(1) for i in range(len(self.measures))],
291                       1).prod(1)
292     return -(torch.matmul(kde_mat, probs) / h).log().sum()
293
294 def chi_squared(self, perms, bins_freq):
295     """
296     Calculates the chi-squared loss function
297
298     :param perms: list of the possible permutations of one location from each measure
299     :param bins_freq: frequencies of data points closest to each location permutation
300     """
301     probs = torch.cat([self.measures[i].weights[perms[:, i]].unsqueeze(1) for i in range(len(self.measures))],
302                       1).prod(1)
303     return (probs**2/bins_freq).sum()
304
305 # Preparational step for loss functions
306 def prep_step(self, data, model, h = 0, alpha = 0.001):
307     """
308     Does a preparatory step and returns the prepared data needed for the optimizers loss function
309
310     :param data: data to fit the model to
311     :param model: model that should be fit to data, should be a function taking the data (x) and
312                  locations as parameters
313     :param h: bandwidth parameter for KDE
314     :param alpha: label smoothing parameter for chi-squared loss function
315     """
316     # permutations of all measures
317     perms = torch.tensor([item for item in itertools.product(*[range(measure.weights.size(dim=0))
318                                                                for measure in self.measures])])
319     locs = torch.cat([self.measures[i].locations[perms[:, i]].unsqueeze(1) for i in range(len(self.measures))], 1)
320     prep = []
321     if self.loss == self.essr:
322         prep.append(torch.tensor([(model(data[0], locs[i]) - data[1]).pow(2).sum() for i in range(len(perms))]))
323     elif self.loss == self.nll:
324         loc_idx = []

```

```

325     for i in range(len(data[0])):
326         ab = torch.abs(model(data[0][i], [locs[:,i] for i in range(locs.size(dim=1))]) - data[1][i])
327         loc_idx.append(torch.argmax(torch.tensor(ab)))
328     prep.append(torch.tensor(loc_idx))
329 elif self.loss == self.KDEnll:
330     if h == 0:
331         sigma=torch.std(data[0])
332         A=min(sigma,(torch.quantile(data[0],0.75)-torch.quantile(data[0],0.25))/1.34)
333         h=0.9*A*len(data[0])**(-1/5)
334     kde_mat = 1/np.sqrt(2*np.pi)*np.exp(-((data[1].view(-1,1) - model(data[0].view(-1,1), locs.transpose(0,1)))
335                                     / h)**2/2)
336     prep.append(kde_mat)
337     prep.append(h)
338 elif self.loss == self.chi_squared:
339     loc_idx = []
340     for i in range(len(data[0])):
341         ab = torch.abs(model(data[0][i], [locs[:,i] for i in range(locs.size(dim=1))]) - data[1][i])
342         loc_idx.append(torch.argmax(ab))
343     bins = torch.tensor(loc_idx)
344     bins_freq = torch.bincount(bins, minlength=np.cumprod([self.measures[i].weights.size(dim=0)
345                                                         for i in range(len(self.measures))])[-1])/len(data[0])**2
346     bins_freq = bins_freq*(1-alpha)+alpha / len(bins_freq)
347     prep.append(bins_freq)
348 return perms, prep

```

## C.3 Checkklassen

```
1 class Check():
2     def __init__(self, opt: Optimizer, model,
3         x: torch.tensor,y: torch.tensor, alpha=0.05,
4         normal=False, Return=False,sample_size=2000):
5         """
6         A class that will check how close a fitted measure is to the original by creating a confidence
7         intervall at each x-value and checking to see if the corresponding y-value is inside the
8         confidence interval.
9
10        :param opt: An instance of the class Optimizer
11        :param model: A function with input: a list of x-values and a list of lists containing weights from measures
12        :param x: The x-values used in the model
13        :param y: The y-values from the original distribution that we use to fitt the measure
14        :param alpha: The amount of confidence we want for our donfidence intervals, standard is 0.05
15                    which corresponds to a 95% CI
16        :param normal: Set to true if you know that the distribution is normal
17        :param Return: Set to true if you wish the class to return the amount of misses and the bounds for the 95% CI
18        """
19        self.opt=opt
20        self.model=model
21        self.data=[x,y]
22        self.N=len(x)
23        self.alpha=alpha
24        self.normal=normal
25        self.Return=Return
26        self.sample_size=sample_size
27
28
29    def check(self):
30        '''
31        Calculates the amount of the original data (y) that is outside
32        the boundaries of a 95% confidence intervall (if no value is given to
33        the variable alpha) and then calculates the probability of
34        that amount of misses
35
36        '''
37        bounds=[]
38        for x in self.data[0]:
39            input=[]
40            for meas in self.opt.measures:
41                input.append(meas.sample(self.sample_size))
42                bounds.append(self.CI(self.model(x,input)))
43        miss=self.misses(self.data[1],bounds)
44
45
46        # b_lower = [b[0].item() for b in bounds]
47        # b_upper = [b[1].item() for b in bounds]
48        # plt.scatter(self.data[0], self.data[1], sizes=[20]*len(self.data[0]))
49        # plt.plot(self.data[0], b_lower, 'r--')
50        # plt.plot(self.data[0], b_upper, 'r--')
51        # plt.show()
52
53        lci=scipy.stats.binom.ppf(self.alpha/2,self.N,self.alpha)
54        hci=scipy.stats.binom.ppf(1-self.alpha/2,self.N,self.alpha)
55        if (miss >= lci) and (miss <= hci):
56            print(f'{miss} is inside the confidence interval ({lci}, {hci}):')
57            print(f'No contradiction with the fitted model at {100*(1-self.alpha)}% confidence level')
58        else:
59            print(f'{miss} is outside the confidence interval ({lci}, {hci}):')
60            print(f'Number of misses is significantly at {100*(1-self.alpha)}% confidence level different
61                from expected for the fitted model!')
62        if self.Return==True:
63            return lci,hci, miss
64
```

```

65 def CI(self, data:list[float]):
66     '''
67     Calculates the bounds of an approximate 95% confidence intervall
68     for the given data in output
69     :param data: List of values that the confidence interval is calculated from
70     '''
71     if self.normal:
72         mean=torch.mean(data)
73         std=torch.std(data)
74         q=scipy.stats.norm.ppf(self.alpha/2)
75         cil=mean+q*std
76         cih=mean-q*std
77         bounds=[cil,cih]
78     else:
79         edge=int(self.alpha/2*self.sample_size)
80         idx_sorted_cropped=torch.argsort(data)[edge:self.sample_size-edge]
81         bounds=data[idx_sorted_cropped[[0,-1]]]
82     return bounds
83
84
85 def misses(self,y:list[float],bounds:list[list[float,float]]):
86     '''
87     Calculates the amount of values in y that are not within the
88     corresponding boundary in bounds
89
90     :param y: The data which we will se if it is in the corresponding confidence interval
91     :param bounds: list containing the pairs of bounds for each x
92     '''
93     miss=0
94     for i in range(len(y)):
95         if y[i]>bounds[i][1] or y[i]<bounds[i][0]:
96             miss+=1
97     return miss

```

## C.4 Generering av data

Kod för att generera data som användes för testerna beskrivna i avsnitt 4.1, vars resultat återfinns i tabellerna i avsnitt A.1.

```
1 import numpy as np
2 np.random.seed(0)
3 tests = 50
4
5 filename = 'data.npy'
6 means = np.array([])
7 stds = np.array([])
8
9 for i in range(tests):
10     mean = np.random.normal(0,1,3)
11     std = np.abs(np.random.normal(1,0.1,3))
12     means = np.append(means, mean)
13     stds = np.append(stds, std)
14     for N in [100, 500, 1000]:
15         a = np.random.normal(mean[0], std[0], N)
16         b = np.random.normal(mean[1], std[1], N)
17         c = np.random.normal(mean[2], std[2], N)
18
19         x = np.linspace(-5, 5, N)
20         y = a
21         y_lin = a + x * b
22         y_sqr = a + x * b + x**2 * c
23         y_nonNorm = np.concatenate((a[0:N//2], b[0:N//2]))
24         y_ax = a * x
25
26         data = {
27             'y': y,
28             'y_lin': y_lin,
29             'y_sqr': y_sqr,
30             'y_nonNorm': y_nonNorm,
31             'y_ax': y_ax
32         }
33
34         for k, v in data.items():
35             file, ext = filename.split('.', 1)
36             name = file + '_' + str(N) + '_' + k + '_' + str(i) + '.' + ext
37
38             with open(name, 'wb') as f:
39                 np.save(f, np.array(v))
40
41 params = np.array([means, stds])
42 # Saves means and stds in npy file in format: [[means],[stds]]
43 np.save('params.npy', params)
```



## C.5 Tester för OFDM

### C.5.1 Fallet $\alpha_i$

```
1
2 import torch
3 import pytorch_measure as pm
4 import numpy as np
5 from tqdm import tqdm
6 import json
7 import math
8
9 # Regression model
10 def regression_model(x,list):
11     return list[0]
12
13 #Load data of parameters from the file where you ran data_generator.py
14 param=np.load(f'../Finalized/test_data/params.npy')
15
16 for length in [100,500,1000]:
17     success=[]
18     tid=[]
19     epoch=[]
20     measures=[]
21     for i in tqdm(range(50)):
22         data=np.load(f'../Finalized/test_data/data_{length}_y_{i}.npy')
23         y=torch.from_numpy(data)
24         x = torch.linspace(-5, 5, length)
25
26         M=length #Amount of datapoints
27
28         s=2 #Decides width of the measures
29         aU=math.ceil(param[0][3*i]+s*param[1][3*i])
30         aL=math.floor(param[0][3*i]-s*param[1][3*i])
31         N=2*(aU-aL)+1
32
33         #Creates the measure
34         measure = pm.Measure(torch.linspace(aL, aU, N), torch.ones(N).double() / N)
35
36         # Instance of optimizer
37         opt = pm.Optimizer([measure],"KDEnll" ,lr=0.1)
38
39         # Call to minimizer
40         new_mes,time,iteration=opt.minimize([x,y], regression_model,verbose=False,adaptive=False,
41                                         max_epochs=4000,test=True)
42
43         # Visualize measures and gradient
44         new_mes[0].visualize()
45         #plt.show()
46
47
48         #Run tests and save the results
49         check=pm.Check(opt,regression_model,x,y,normal=False,Return=True)
50         l,u,miss=check.check()
51
52         success.append(l<miss and miss<u)
53         tid.append(time)
54         epoch.append(iteration)
55         measures.append([new_mes[0].locations.tolist(),new_mes[0].weights.tolist()])
56
57
58 data=[measures,sum(tid)/len(tid),sum(epoch)/(len(epoch)),sum(success)/len(success)]
59 with open(f"resultat/Sergey1_{M}.json", "w") as outfile:
60     outfile.write(json.dumps(data))
```

## C.5.2 Fallet $\alpha_i \cdot x$

```
1 import torch
2 import pytorch_measure as pm
3 import numpy as np
4 from tqdm import tqdm
5 import json
6 import math
7
8 # Linear regression model
9 def regression_model(x,list):
10     return list[0]*x
11
12 #Load data of parameters from the file where you ran data_generator.py
13 param=np.load(f'../Finalized/test_data/params.npy')
14
15 for length in [100,500,1000]:
16     success=[]
17     tid=[]
18     epoch=[]
19     measures=[]
20     for i in tqdm(range(50)):
21         data=np.load(f'../Finalized/test_data/data_{length}_y_ax_{i}.npy')
22         y=torch.from_numpy(data)
23         x = torch.linspace(-5, 5, length)
24
25         M=length #Amount of datapoints
26
27         s=2 #Decides width of the measures
28         aU=math.ceil(param[0][3*i]+s*param[1][3*i])
29         aL=math.floor(param[0][3*i]-s*param[1][3*i])
30         N=2*(aU-aL)+10
31
32         #Creates the measure
33         measure = pm.Measure(torch.linspace(aL, aU, N), torch.ones(N).double() / N)
34
35         # Instance of optimizer
36         opt = pm.Optimizer([measure], "KDEnll" ,lr=0.1)
37
38         # Call to minimizer
39         new_mes,time,iteration=opt.minimize([x,y], regression_model,verbose=False,adaptive=False,
40                                           max_epochs=4000,test=True)
41
42         # Visualize measures and gradient
43         new_mes[0].visualize()
44         #plt.show()
45
46         #Run tests and save the results
47         check=pm.Check(opt,regression_model,x,y,normal=False,Return=True)
48         l,u,miss=check.check()
49
50         success.append(l<miss and miss<u)
51         tid.append(time)
52         epoch.append(iteration)
53         measures.append([new_mes[0].locations.tolist(),new_mes[0].weights.tolist()])
54
55     data=[measures,sum(tid)/len(tid),sum(epoch)/(len(epoch)),sum(success)/len(success)]
56     with open(f"resultat/Sergey1a_{M}.json", "w") as outfile:
57         outfile.write(json.dumps(data))
```

### C.5.3 Fallet $\alpha_i + \beta_i \cdot x$

```
1 import torch
2 import pytorch_measure as pm
3 import numpy as np
4 from tqdm import tqdm
5 import json
6 import math
7
8 # Linear regression model
9 def regression_model(x,list):
10     return list[0]*x+list[1]
11
12 #Load data of parameters from the file where you ran data_generator.py
13 param=np.load(f'../Finalized/test_data/params.npy')
14
15 for length in [100,500,1000]:
16     success=[]
17     tid=[]
18     epoch=[]
19     measures=[]
20     for i in tqdm(range(50)):
21         data=np.load(f'../Finalized/test_data/data_{length}_y_lin_{i}.npy')
22         y=torch.from_numpy(data)
23         x = torch.linspace(-5, 5, length)
24
25         M=length #Amount of datapoints
26
27         s=2 #Decides width of the measures
28         aU=math.ceil(param[0][3*i+1]+s*param[1][3*i+1])
29         aL=math.floor(param[0][3*i+1]-s*param[1][3*i+1])
30         bU=math.ceil(param[0][3*i]+s*param[1][3*i])
31         bL=math.floor(param[0][3*i]-s*param[1][3*i])
32         Nb=2*(bU-bL)+1
33         Na=2*(aU-aL)+1
34
35         #Creates the measures
36         a = pm.Measure(torch.linspace(aL, aU, Na), torch.ones(Na).double() / Na)
37         b = pm.Measure(torch.linspace(bL, bU, Nb), torch.ones(Nb).double() / Nb)
38         measure= [a,b]
39
40         # Instance of optimizer
41         opt = pm.Optimizer(measure, "KDEnll", lr = 0.1)
42
43         # Call to minimizer
44         new_mes,time,iteration=opt.minimize([x,y],regression_model,max_epochs=4000,verbose = False,
45                                           print_freq=100, smallest_lr=1e-10,test=True)
46
47         # Visualize measures and gradient
48         new_mes[0].visualize()
49         #plt.show()
50         new_mes[1].visualize()
51         #plt.show()
52
53         #Run tests and save the results
54         check=pm.Check(opt,regression_model,x,y,normal=False,Return=True)
55         l,u,miss=check.check()
56
57         success.append(l<=miss and miss<=u)
58         tid.append(time)
59         epoch.append(iteration)
60         for i in range(len(new_mes)):
61             measures.append([new_mes[i].locations.tolist(),new_mes[i].weights.tolist()])
62
63 data=[measures,sum(tid)/len(tid),sum(epoch)/len(epoch),sum(success)/len(success)]
64 with open(f"resultat_samuel/Sergey2M_{M}.json", "w") as outfile:
65     outfile.write(json.dumps(data))
```

### C.5.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$

```
1 import torch
2 import pytorch_measure as pm
3 import numpy as np
4 from tqdm import tqdm
5 import json
6 import math
7
8 # Linear regression model
9 def regression_model(x,list):
10     return list[0]*x**2+list[1]*x+list[2]
11
12 #Load data of parameters from the file where you ran data_generator.py
13 param=np.load('./Finalized/test_data/params.npy')
14
15 for length in [100,500,1000]:
16     success=[]
17     tid=[]
18     epoch=[]
19     measures=[]
20     for i in tqdm(range(50)):
21         data=np.load(f'./Finalized/test_data/data_{length}_y_sqr_{i}.npy')
22         y=torch.from_numpy(data)
23         x = torch.linspace(-5, 5, length)
24
25         #plt.scatter(x,y)
26         #plt.show()
27         M=length #Amount of datapoints
28
29         s=2 #Decides width of the measures
30         aU=math.ceil(param[0][3*i+2]+s*param[1][3*i+2])
31         aL=math.floor(param[0][3*i+2]-s*param[1][3*i+2])
32         bU=math.ceil(param[0][3*i+1]+s*param[1][3*i+1])
33         bL=math.floor(param[0][3*i+1]-s*param[1][3*i+1])
34         cU=math.ceil(param[0][3*i]+s*param[1][3*i])
35         cL=math.floor(param[0][3*i]-s*param[1][3*i])
36         Nc=2*(cU-cL)+1
37         Nb=2*(bU-bL)+1
38         Na=2*(aU-aL)+1
39
40
41         # Measure for slope (a) and intercept (b) of linear model
42         a = pm.Measure(torch.linspace(aU, aL, Na), torch.ones(Na).double() / Na)
43         b = pm.Measure(torch.linspace(bU, bL, Nb), torch.ones(Nb).double() / Nb)
44         c = pm.Measure(torch.linspace(cU, cL, Nc), torch.ones(Nc).double() / Nc)
45         measure = [a,b,c]
46
47         # Instance of optimizer
48         opt = pm.Optimizer(measure, "KDEnll", lr = 0.1)
49
50         # Call to minimizer
51         new_mes,time,iteration=opt.minimize([x,y],regression_model,max_epochs=4000,verbose = False,
52                                         print_freq=100, smallest_lr=1e-10,test=True)
53
54         # Visualize measures and gradient
55         new_mes[0].visualize()
56         #plt.show()
57         new_mes[1].visualize()
58         #plt.show()
59         new_mes[2].visualize()
60         #plt.show()
61
62         #Run tests and save the results
63         check=pm.Check(opt,regression_model,x,y,normal=False,Return=True)
64         l,u,miss=check.check()
```

```
65
66     success.append(l<=miss and miss<=u)
67     tid.append(time)
68     epoch.append(iteration)
69     for i in range(len(new_mes)):
70         measures.append([new_mes[i].locations.tolist(),new_mes[i].weights.tolist()])
71
72 data=[measures,sum(tid)/len(tid),sum(epoch)/len(epoch),sum(success)/len(success)]
73 with open(f"resultat/Sergey3_{M}.json", "w") as outfile:
74     outfile.write(json.dumps(data))
```

## C.6 Tester för linjärkombination av stokastiska variabler

### C.6.1 Fallet $\alpha_i$

```
1 import torch
2 import numpy as np
3 import scipy as sp
4 import linear_combination_optimizer as lco
5 import json
6
7 # Calculate confidence intervals and returns amount of misses
8 def misses(x, y, mu, sigma):
9     miss = 0
10    S = 1000
11    for i in range(x.size(dim=0)):
12        # Generate sample from model
13        sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
14        sample = torch.sort(sample)[0]
15        if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
16            miss += 1
17    # Calculate confidence interval
18    c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
19    c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
20    print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
21    return c1, c2, miss
22
23 # For each data size
24 for length in [100, 500, 1000]:
25     # Lists for storing results
26     success=[]
27     time=[]
28     epoch=[]
29     means=[]
30     std=[]
31     # Run regressions 50 times
32     for i in range(50):
33         x = torch.linspace(-5, 5, length)
34         y = np.load(f'../Finalized/test_data/data_{length}_y_{i}.npy')
35         opt = lco.Optimizer(x, y, order=1)
36         mu, sigma, conv_epoch, conv_time = opt.optimize(epochs = 3000, test = True)
37
38         l, u, miss = misses(x,torch.tensor(y),mu,sigma)
39         success.append(l<=miss and miss<=u)
40         time.append(conv_time)
41         epoch.append(conv_epoch)
42         means.append(mu)
43         std.append(sigma)
44
45     # Save results in a json file
46     results = [means, std, sum(time)/len(time), sum(epoch)/(len(epoch)),
47               float(100*sum(success)/len(success))]
48     with open(f"Lin_comb_results/lin_comb_results_{length}_y.json", "w") as outfile:
49         outfile.write(json.dumps(results))
```

## C.6.2 Fallet $\alpha_i \cdot x$

```
1 import torch
2 import numpy as np
3 import scipy as sp
4 import linear_combination_optimizer as lco
5 import json
6
7 # Confidence intervals
8 def misses(x, y, mu, sigma):
9     miss = 0
10    S = 1000
11    for i in range(x.size(dim=0)):
12        # Generate sample from model
13        sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
14        sample = torch.sort(sample)[0]
15        if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
16            miss += 1
17    # Calculate confidence interval
18    c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
19    c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
20    print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
21    return c1, c2, miss
22
23 # For each data size
24 for length in [100, 500, 1000]:
25     # Lists for storing results
26     success=[]
27     time=[]
28     epoch=[]
29     means=[]
30     std=[]
31     # Run regressions 50 times
32     for i in range(50):
33         x = torch.linspace(-5, 5, length)
34         y = np.load(f'../Finalized/test_data/data_{length}_y_ax_{i}.numpy')
35         opt = lco.Optimizer(x, y, order=1, ax = True)
36         mu, sigma, conv_epoch, conv_time = opt.optimize(epochs = 3000, test = True)
37
38         l, u, miss = misses(x,torch.tensor(y),mu,sigma)
39         success.append(l<=miss and miss<=u)
40         time.append(conv_time)
41         epoch.append(conv_epoch)
42         means.append(mu)
43         std.append(sigma)
44
45     # Save results in a json file
46     results = [means, std, sum(time)/len(time), sum(epoch)/(len(epoch)),
47               float(100*sum(success)/len(success))]
48     with open(f"Lin_comb_results/lin_comb_results_{length}_y_ax.json", "w") as outfile:
49         outfile.write(json.dumps(results))
```

### C.6.3 Fallet $\alpha_i + \beta_i \cdot x$

```
1 import torch
2 import numpy as np
3 import scipy as sp
4 import linear_combination_optimizer as lco
5 import json
6
7 # Calculate confidence intervals and returns amount of misses
8 def misses(x, y, mu, sigma):
9     miss = 0
10    S = 1000
11    for i in range(x.size(dim=0)):
12        sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
13        sample = torch.sort(sample)[0]
14        if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
15            miss += 1
16    # Calculate confidence interval
17    c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
18    c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
19    print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
20    return c1, c2, miss
21
22 # For each data size
23 for length in [100, 500, 1000]:
24     # Lists for storing results
25     success=[]
26     time=[]
27     epoch=[]
28     means=[]
29     std=[]
30     # Run regressions 50 times
31     for i in range(50):
32         x = torch.linspace(-5, 5, length)
33         y = np.load(f'../Finalized/test_data/data_{length}_y_lin_{i}.npy')
34         opt = lco.Optimizer(x, y, order=2)
35         mu, sigma, conv_epoch, conv_time = opt.optimize(epochs = 3000, test = True)
36
37         l, u, miss = misses(x,torch.tensor(y),mu,sigma)
38         success.append(l<=miss and miss<=u)
39         epoch.append(conv_epoch)
40         means.append(mu)
41         std.append(sigma)
42
43     # Save results in a json file
44     results = [means, std, sum(time)/len(time), sum(epoch)/(len(epoch)),
45               float(100*sum(success)/len(success))]
46     with open(f"Lin_comb_results/lin_comb_results_{length}_y_lin.json", "w") as outfile:
47         outfile.write(json.dumps(results))
```



### C.6.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$

```
1 import torch
2 import numpy as np
3 import scipy as sp
4 import linear_combination_optimizer as lco
5 import json
6
7 # Calculate confidence intervals and returns amount of misses
8 def misses(x, y, mu, sigma):
9     miss = 0
10    S = 1000
11    for i in range(x.size(dim=0)):
12        # Generate sample from model
13        sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
14        sample = torch.sort(sample)[0]
15        if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
16            miss += 1
17    # Calculate confidence interval
18    c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
19    c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
20    print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
21    return c1, c2, miss
22
23 # For each data size
24 for length in [100, 500, 1000]:
25     # Lists for storing results
26     success=[]
27     time=[]
28     epoch=[]
29     means=[]
30     std=[]
31     # Run regressions 50 times
32     for i in range(50):
33         x = torch.linspace(-5, 5, length)
34         y = np.load(f'../Finalized/test_data/data_{length}_y_sqr_{i}.numpy')
35         opt = lco.Optimizer(x, y, order=3)
36         mu, sigma, conv_epoch, conv_time = opt.optimize(epochs = 3000, test = True)
37
38         l, u, miss = misses(x,torch.tensor(y),mu,sigma)
39         success.append(l<=miss and miss<=u)
40         time.append(conv_time)
41         epoch.append(conv_epoch)
42         means.append(mu)
43         std.append(sigma)
44
45     # Save results in a json file
46     results = [means, std, sum(time)/len(time), sum(epoch)/(len(epoch)),
47               float(100*sum(success)/len(success))]
48     with open(f"Lin_comb_results/lin_comb_results_{length}_y_sqr.json", "w") as outfile:
49         outfile.write(json.dumps(results))
```

## C.7 Tester för polynomiskt neuralt nätverk

### C.7.1 Fallet $\alpha_i$

```
1 import torch
2 import numpy as np
3 import scipy as sp
4 import time
5 import json
6
7 # Polynomial nn method
8 def eval_powers_of_x(x, n):
9     return x.pow(torch.arange(n))
10 def eval_even_powers_of_x(x, n):
11     return x.pow(2 * torch.arange(n))
12
13 # Calculate confidence intervals and returns amount of misses
14 def misses(x, y, mu, sigma):
15     miss = 0
16     S = 1000
17     for i in range(x.size(dim=0)):
18         # Generate sample from model
19         sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
20         sample = torch.sort(sample)[0]
21         if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
22             miss += 1
23     # Calculate confidence interval
24     c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
25     c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
26     print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
27     return c1, c2, miss
28
29 # Polynomial nn method
30 class NormalPolynomialModel(torch.nn.Module):
31     def __init__(self):
32         super().__init__()
33
34         self.mean = None
35         self.std = None
36
37         self.N_mean = 1 # constant, linear, quadratic, ... terms
38         self.N_var = 1
39
40         self.poly_multipliers_mean = torch.nn.parameter.Parameter(torch.rand(self.N_mean))
41         self.poly_multipliers_var = torch.nn.parameter.Parameter(torch.rand(self.N_var))
42         self.var_shift = torch.nn.parameter.Parameter(torch.tensor(0.1))
43
44         self.mean_layer = torch.nn.Linear(self.N_mean, 1)
45         self.var_layer = torch.nn.Sequential(
46             torch.nn.Linear(self.N_var, 1),
47             torch.nn.ReLU()
48         )
49
50     def forward(self, x):
51         self.mean = self.mean_layer(self.poly_multipliers_mean * eval_powers_of_x(x, self.N_mean))
52         self.var = self.var_layer(self.poly_multipliers_var * eval_even_powers_of_x(x, self.N_var)) +
53             self.var_shift ** 2
54         return self.mean, self.var
55
56 # Loss function for nn
57 def log_k_with_var(mean, var, y):
58     return -0.5*torch.log(var) - (y - mean)**2 / (2 * var)
59
60 # For each data size
61 for length in [100, 500, 1000]:
62     # Lists for storing results
```

```

63 success=[]
64 tid=[]
65 end_epoch=[]
66 means=[]
67 std=[]
68 # Run regressions 50 times
69 for i in range(50):
70     x = torch.linspace(-5, 5, length)
71     y = np.load(f'../Finalized/test_data/data_{length}_y_{i}.npy')
72     y = torch.tensor(y)
73     x_unsq = x.view(-1, 1)
74     y = y.view(-1, 1)
75     model = NormalPolynomialModel()
76     opt = torch.optim.Adam(model.parameters())
77
78     max_epoch = 300000
79     conv_epoch = max_epoch
80     conv_time = float('inf')
81     old_loss = float('inf')
82     t1 = time.time()
83     # Main loop of regression
84     for epoch in range(max_epoch):
85         opt.zero_grad()
86         mean, var = model(x_unsq)
87         log_likelihood = log_k_with_var(mean, var, y).sum()
88         loss = -log_likelihood
89         loss.backward()
90         opt.step()
91         if torch.abs(old_loss-loss) < 1e-4:
92             t2 = time.time()
93             conv_time = t2-t1
94             conv_epoch = epoch
95             break
96         old_loss = loss
97
98         if epoch % 1000 == 0:
99             print(f'{epoch}:: Loss = {loss.item()}')
100
101     m, var = model(x_unsq)
102     s = var.sqrt()
103     m = m.detach().numpy()
104     s = s.detach().numpy()
105
106     l, u, miss = misses(x,torch.tensor(y),m,s)
107     success.append(l<=miss and miss<=u)
108     if conv_time != float('inf'):
109         tid.append(conv_time)
110     end_epoch.append(conv_epoch)
111     means.append(m.tolist())
112     std.append(s.tolist())
113
114 # Save results in a json file
115 results = [means, std, sum(tid)/len(tid), sum(end_epoch)/(len(end_epoch)),
116           float(100*sum(success)/len(success))]
117 with open(f"Poly_results/poly_results_{length}_y.json", "w") as outfile:
118     outfile.write(json.dumps(results))

```

## C.7.2 Fallet $\alpha_i \cdot x$

```
1 import torch
2 import pytorch_measure as pm
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import scipy as sp
6 from scipy import stats
7 import time
8 import linear_combination_optimizer as lco
9 import json
10
11 # Polynomial nn method
12 def eval_powers_of_x(x, n):
13     return x
14 def eval_even_powers_of_x(x, n):
15     return x.pow(2 * torch.arange(n))
16
17 # Calculate confidence intervals and returns amount of misses
18 def misses(x, y, mu, sigma):
19     miss = 0
20     S = 1000
21     for i in range(x.size(dim=0)):
22         # Generate sample from model
23         sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
24         sample = torch.sort(sample)[0]
25         if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
26             miss += 1
27     # Calculate confidence interval
28     c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
29     c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
30     print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
31     return c1, c2, miss
32
33 # Polynomial nn method
34 class NormalPolynomialModel(torch.nn.Module):
35     def __init__(self):
36         super().__init__()
37
38         self.mean = None
39         self.std = None
40
41         self.N_mean = 1 # constant, linear, quadratic, ... terms
42         self.N_var = 1
43
44         self.poly_multipliers_mean = torch.nn.parameter.Parameter(torch.rand(self.N_mean))
45         self.poly_multipliers_var = torch.nn.parameter.Parameter(torch.rand(self.N_var))
46         self.var_shift = torch.nn.parameter.Parameter(torch.tensor(0.1))
47
48         self.mean_layer = torch.nn.Linear(self.N_mean, 1)
49         self.var_layer = torch.nn.Sequential(
50             torch.nn.Linear(self.N_var, 1),
51             torch.nn.ReLU()
52         )
53
54     def forward(self, x):
55         self.mean = self.mean_layer(self.poly_multipliers_mean * x)
56         self.var = self.var_layer(self.poly_multipliers_var * x**2) + self.var_shift ** 2
57         return self.mean, self.var
58
59 # Loss function for nn
60 def log_k_with_var(mean, var, y):
61     return -0.5*torch.log(var) - (y - mean)**2 / (2 * var)
62
63 # For each data size
64 for length in [100, 500, 1000]:
```

```

65 # Lists for storing results
66 success=[]
67 tid=[]
68 end_epoch=[]
69 means=[]
70 std=[]
71 # Run regressions 50 times
72 for i in range(50):
73     x = torch.linspace(-5, 5, length)
74     y = np.load(f'../Finalized/test_data/data_{length}_y_ax_{i}.npy')
75     y = torch.tensor(y)
76     x_unsq = x.view(-1, 1)
77     y = y.view(-1, 1)
78     model = NormalPolynomialModel()
79     opt = torch.optim.Adam(model.parameters())
80
81     max_epoch = 300000
82     conv_epoch = max_epoch
83     conv_time = float('inf')
84     old_loss = float('inf')
85     t1 = time.time()
86     # Main loop of regression
87     for epoch in range(max_epoch):
88         opt.zero_grad()
89         mean, var = model(x_unsq)
90         log_likelyhood = log_k_with_var(mean, var, y).sum()
91         loss = -log_likelyhood
92         loss.backward()
93         opt.step()
94         if torch.abs(old_loss-loss) < 1e-4:
95             t2 = time.time()
96             conv_time = t2-t1
97             conv_epoch = epoch
98             break
99         old_loss = loss
100
101         if epoch % 1000 == 0:
102             print(f'{epoch}:: Loss = {loss.item()}')
103
104     m, var = model(x_unsq)
105     s = var.sqrt()
106     m = m.detach().numpy()
107     s = s.detach().numpy()
108
109     if sum([np.isnan(m[i]) for i in range(len(m))]) == 0 and
110        sum([np.isnan(s[i]) for i in range(len(s))]) == 0:
111         l, u, miss = misses(x, torch.tensor(y), m, s)
112         success.append(l<=miss and miss<=u)
113     if conv_time != float('inf'):
114         tid.append(conv_time)
115     end_epoch.append(conv_epoch)
116     means.append(m.tolist())
117     std.append(s.tolist())
118
119 # Save results in a json file
120 results = [means, std, sum(tid)/len(tid), sum(end_epoch)/(len(end_epoch)),
121           float(100*sum(success)/len(success))]
122 with open(f'Poly_results/poly_results_{length}_y_ax.json", "w") as outfile:
123     outfile.write(json.dumps(results))

```

### C.7.3 Fallet $\alpha_i + \beta_i \cdot x$

```

1 import torch
2 import numpy as np
3 import scipy as sp
4 import time
5 import linear_combination_optimizer as lco
6 import json
7
8 # Polynomial nn method
9 def eval_powers_of_x(x, n):
10     return x.pow(torch.arange(n))
11 def eval_even_powers_of_x(x, n):
12     return x.pow(2 * torch.arange(n))
13
14 # Calculate confidence intervals and returns amount of misses
15 def misses(x, y, mu, sigma):
16     miss = 0
17     S = 1000
18     for i in range(x.size(dim=0)):
19         # Generate sample from model
20         sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
21         sample = torch.sort(sample)[0]
22         if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
23             miss += 1
24     # Calculate confidence interval
25     c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
26     c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
27     print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
28     return c1, c2, miss
29
30 # Polynomial nn method
31 class NormalPolynomialModel(torch.nn.Module):
32     def __init__(self):
33         super().__init__()
34
35         self.mean = None
36         self.std = None
37
38         self.N_mean = 2 # constant, linear, quadratic, ... terms
39         self.N_var = 2
40
41         self.poly_multipliers_mean = torch.nn.parameter.Parameter(torch.rand(self.N_mean))
42         self.poly_multipliers_var = torch.nn.parameter.Parameter(torch.rand(self.N_var))
43         self.var_shift = torch.nn.parameter.Parameter(torch.tensor(0.1))
44
45         self.mean_layer = torch.nn.Linear(self.N_mean, 1)
46         self.var_layer = torch.nn.Sequential(
47             torch.nn.Linear(self.N_var, 1),
48             torch.nn.ReLU()
49         )
50
51     def forward(self, x):
52         self.mean = self.mean_layer(self.poly_multipliers_mean * eval_powers_of_x(x, self.N_mean))
53         self.var = self.var_layer(self.poly_multipliers_var * eval_even_powers_of_x(x, self.N_var)) +
54             self.var_shift ** 2
55         return self.mean, self.var
56
57 # Loss function for nn
58 def log_k_with_var(mean, var, y):
59     return -0.5*torch.log(var) - (y - mean)**2 / (2 * var)
60
61 # For each data size
62 for length in [100, 500, 1000]:
63     success=[]
64     tid=[]

```

```

65 end_epoch=[]
66 means=[]
67 std=[]
68 # Run regressions 50 times
69 for i in range(50):
70     x = torch.linspace(-5, 5, length)
71     y = np.load(f'../Finalized/test_data/data_{length}_y_lin_{i}.npy')
72     y = torch.tensor(y)
73     x_unsq = x.view(-1, 1)
74     y = y.view(-1, 1)
75     model = NormalPolynomialModel()
76     opt = torch.optim.Adam(model.parameters())
77
78     max_epoch = 300000
79     conv_epoch = max_epoch
80     conv_time = float('inf')
81     old_loss = float('inf')
82     t1 = time.time()
83     # Main loop of regression
84     for epoch in range(max_epoch):
85         opt.zero_grad()
86         mean, var = model(x_unsq)
87         log_likelyhood = log_k_with_var(mean, var, y).sum()
88         loss = -log_likelyhood
89         loss.backward()
90         opt.step()
91         if torch.abs(old_loss-loss) < 1e-4:
92             t2 = time.time()
93             conv_time = t2-t1
94             conv_epoch = epoch
95             break
96         old_loss = loss
97
98         if epoch % 1000 == 0:
99             print(f'{epoch}:: Loss = {loss.item()}')
100
101     m, var = model(x_unsq)
102     s = var.sqrt()
103     m = m.detach().numpy()
104     s = s.detach().numpy()
105
106     l, u, miss = misses(x,torch.tensor(y),m,s)
107     success.append(l<=miss and miss<=u)
108     if conv_time != float('inf'):
109         tid.append(conv_time)
110     end_epoch.append(conv_epoch)
111     means.append(m.tolist())
112     std.append(s.tolist())
113
114 # Save results in a json file
115 results = [means, std, sum(tid)/len(tid), sum(end_epoch)/(len(end_epoch)),
116           float(100*sum(success)/len(success))]
117 with open(f'Poly_results/poly_results_{length}_y_lin.json", "w") as outfile:
118     outfile.write(json.dumps(results))

```

### C.7.4 Fallet $\alpha_i + \beta_i \cdot x + \gamma_i \cdot x^2$

```

1 import torch
2 import pytorch_measure as pm
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import scipy as sp
6 from scipy import stats
7 import time
8 import linear_combination_optimizer as lco
9 import json
10
11 # Polynomial nn method
12 def eval_powers_of_x(x, n):
13     return x.pow(torch.arange(n))
14 def eval_even_powers_of_x(x, n):
15     return x.pow(2 * torch.arange(n))
16
17 # Calculate confidence intervals and returns amount of misses
18 def misses(x, y, mu, sigma):
19     miss = 0
20     S = 1000
21     for i in range(x.size(dim=0)):
22         # Generate sample from model
23         sample = torch.normal(mean = float(mu[i]), std = float(sigma[i]), size = (1,S)).squeeze(dim=0)
24         sample = torch.sort(sample)[0]
25         if y[i] < sample[int(np.floor(S*0.025))] or y[i] > sample[int(np.ceil(S*0.975))-1]:
26             miss += 1
27     # Calculate confidence interval
28     c1 = sp.stats.binom.ppf(0.025,y.size(dim=0),0.05)
29     c2 = sp.stats.binom.ppf(0.975,y.size(dim=0),0.05)
30     print(f"CI: [{c1}, {c2}], Misses: {miss}, Within CI: {c1<=miss<=c2}")
31     return c1, c2, miss
32
33 # Polynomial nn method
34 class NormalPolynomialModel(torch.nn.Module):
35     def __init__(self):
36         super().__init__()
37
38         self.mean = None
39         self.std = None
40
41         self.N_mean = 3 # constant, linear, quadratic, ... terms
42         self.N_var = 3
43
44         self.poly_multipliers_mean = torch.nn.parameter.Parameter(torch.rand(self.N_mean))
45         self.poly_multipliers_var = torch.nn.parameter.Parameter(torch.rand(self.N_var))
46         self.var_shift = torch.nn.parameter.Parameter(torch.tensor(0.1))
47
48         self.mean_layer = torch.nn.Linear(self.N_mean, 1)
49         self.var_layer = torch.nn.Sequential(
50             torch.nn.Linear(self.N_var, 1),
51             torch.nn.ReLU()
52         )
53
54     def forward(self, x):
55         self.mean = self.mean_layer(self.poly_multipliers_mean * eval_powers_of_x(x, self.N_mean))
56         self.var = self.var_layer(self.poly_multipliers_var * eval_even_powers_of_x(x, self.N_var)) +
57             self.var_shift ** 2
58         return self.mean, self.var
59
60 # Loss function for nn
61 def log_k_with_var(mean, var, y):
62     return -0.5*torch.log(var) - (y - mean)**2 / (2 * var)
63
64 # For each data size

```



```

65 for length in [100, 500, 1000]:
66     # Lists for storing results
67     success=[]
68     tid=[]
69     end_epoch=[]
70     means=[]
71     std=[]
72     # Run regressions 50 times
73     for i in range(50):
74         x = torch.linspace(-5, 5, length)
75         y = np.load(f'../Finalized/test_data/data_{length}_y_sqr_{i}.npy')
76         y = torch.tensor(y)
77         x_unsq = x.view(-1, 1)
78         y = y.view(-1, 1)
79         model = NormalPolynomialModel()
80         opt = torch.optim.Adam(model.parameters())
81
82         max_epoch = 300000
83         conv_epoch = max_epoch
84         conv_time = float('inf')
85         old_loss = float('inf')
86         t1 = time.time()
87         # Main loop of regression
88         for epoch in range(max_epoch):
89             opt.zero_grad()
90             mean, var = model(x_unsq)
91             log_likelihood = log_k_with_var(mean, var, y).sum()
92             loss = -log_likelihood
93             loss.backward()
94             opt.step()
95             if torch.abs(old_loss-loss) < 1e-4:
96                 t2 = time.time()
97                 conv_time = t2-t1
98                 conv_epoch = epoch
99                 print(f'{epoch}:: Loss = {loss.item()}')
100                break
101                old_loss = loss
102
103                if epoch % 1000 == 0:
104                    print(f'{epoch}:: Loss = {loss.item()}')
105
106                m, var = model(x_unsq)
107                s = var.sqrt()
108                m = m.detach().numpy()
109                s = s.detach().numpy()
110
111                l, u, miss = misses(x,torch.tensor(y),m,s)
112                success.append(l<=miss and miss<=u)
113                if conv_time != float('inf'):
114                    tid.append(conv_time)
115                end_epoch.append(conv_epoch)
116                means.append(m.tolist())
117                std.append(s.tolist())
118
119                # Save results in a json file
120                results = [means, std, sum(tid)/len(tid), sum(end_epoch)/(len(end_epoch)),
121                    float(100*sum(success)/len(success))]
122                with open(f"Poly_results/poly_results_{length}_y_sqr.json", "w") as outfile:
123                    outfile.write(json.dumps(results))

```

## C.8 Implementation av linjärkombination av stokastiska variabler

```
1 import torch
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import time
5
6 class Optimizer():
7     def __init__(self, x, y, order=3, n=3, ax = False):
8         self.order = order
9         self.n = n
10        self.x = x
11        self.y = y
12        if ax == False:
13            self.h = [self.create_lambda(i) for i in range(self.order)]
14        else:
15            self.h = [self.create_lambda(1)]
16        h_all_data = [h(x) for h in self.h]
17        self.h_all = torch.transpose(torch.stack(h_all_data, 0), 0, 1)
18        self.mu = torch.tensor([0. for _ in range(order)], dtype=float, requires_grad=True)
19        self.sigma = torch.tensor([1. for _ in range(order)], dtype=float, requires_grad=True)
20        self.beta = [self.mu, self.sigma]
21
22    def create_lambda(self, i):
23        if i == 0:
24            return lambda u: u*0+1
25        return lambda u:u**i
26
27    def optimize(self, epochs=300, lr=0.1, print_frequency = 10, test = False):
28        optimizer = torch.optim.Adam(self.beta,lr=lr, maximize=True)
29        old_loss = float('inf')
30        cur_epoch = epochs
31        t1 = time.time()
32        t2 = float('inf') # In case optimization doesnt terminate before max epochs are reached
33        for epoch in range(epochs):
34            optimizer.zero_grad()
35            loss = self.log_lik(self.y, self.beta, self.h_all)
36            loss.backward()
37            optimizer.step()
38            if epoch%print_frequency==0:
39                print(epoch, "mu:", self.mu.detach().numpy(), "sigma:", self.sigma.detach().numpy(), "loss:", loss)
40            if torch.abs(loss-old_loss) < 1e-15:
41                t2 = time.time()
42                cur_epoch = epoch
43                break
44            old_loss = loss
45        self.mu_optim = self.beta[0].detach().numpy()
46        self.sigma_optim = self.beta[1].detach().numpy()
47        print(epoch, "mu:", self.mu_optim, "sigma:", self.sigma_optim)
48        if test == False:
49            return [self.m(self.mu, self.h_all[i,:]).detach().numpy().tolist() for i in range(self.x.size(dim=0))],
50                [(self.sigma_2(self.sigma,
51                    self.h_all[i,:])**0.5).detach().numpy().tolist() for i in range(self.x.size(dim=0))]
52        else:
53            if t2 != float('inf'):
54                return [self.m(self.mu, self.h_all[i,:]).detach().numpy().tolist() for i in range(self.x.size(dim=0))],
55                    [(self.sigma_2(self.sigma,
56                        self.h_all[i,:])**0.5).detach().numpy().tolist() for i in range(self.x.size(dim=0))],
57                    cur_epoch, t2-t1
58            else:
59                return [self.m(self.mu, self.h_all[i,:]).detach().numpy().tolist() for i in range(self.x.size(dim=0))],
60                    [(self.sigma_2(self.sigma,
61                        self.h_all[i,:])**0.5).detach().numpy().tolist() for i in range(self.x.size(dim=0))],
62                    cur_epoch, None
63
64
```

```

65     def m(self, mu, h_x):
66         return (mu * h_x).sum()
67
68     def sigma_2(self, sigma, h_x):
69         return (sigma * h_x).pow(2).sum()
70
71     def sigma_2_regular_exp(self, sigma, h_x):
72         return ((sigma * h_x)**2).sum()
73
74     def log_normal_pdf(self, y, mu, sigma_2):
75         return -1/2*(torch.log(2*np.pi*(sigma_2)) + ((y-mu)**2/sigma_2))
76
77     def log_lik(self, y, beta, h_all):
78         return sum([self.log_normal_pdf(y_i, self.m(beta[0], h_all[i]),
79                                     self.sigma_2(beta[1], h_all[i])) for i, y_i in enumerate(y)])
80
81     def get_plot_CI(self, mu, sigma, n=3):
82         y_new = []
83         sigma_upper = []
84         sigma_lower = []
85         sample_x = []
86         sample_y = []
87
88         for i in range(len(self.x)):
89             h_x = [h_i.item() for h_i in self.h_all[i]]
90             y_new.append(self.m(mu, h_x))
91             s = self.sigma_2_regular_exp(sigma, h_x)**0.5
92             sigma_upper.append(y_new[-1] + 2*s)
93             sigma_lower.append(y_new[-1] - 2*s)
94             sample = torch.normal(float(y_new[-1]), float(s), size = (1, n))
95             for j in range(n):
96                 sample_x.append(self.x[i].item())
97                 sample_y.append(sample[0][j].item())
98         return y_new, sigma_upper, sigma_lower
99
100    def visualize(self):
101        y_new, sigma_upper, sigma_lower = self.get_plot_CI(self.mu_optim, self.sigma_optim)
102        markersize = 20
103        plt.scatter(self.x, self.y, sizes=[markersize]*len(self.x))
104        plt.plot(self.x, y_new, 'r-')
105        plt.plot(self.x, sigma_upper, 'r--')
106        plt.plot(self.x, sigma_lower, 'r--')
107        plt.show()

```