

Bachelor Thesis, Mathematics

# **“Machine Learning using Bayesian statistics and Neural networks”**

*A thesis submitted in fulfillment of the requirements  
for the bachelor's degree in Mathematics*

Author: **Lukas Prokop**  
Supervisor: 高山・信毅 ※  
Supervisor: Bredies Kristian †

※ University of Kobe, Japan

† University of Graz, Austria

Version: August 2, 2017



# Abstract

Machine Learning is a vivid research area. Machine Learning fundamentally changes the idea that programmers mechanically write programs in order to perform tasks such as classification, regression, clustering, density estimation, and model selection. In supervised learning, machines learn by observing test vectors and their desired output. They successively adapt their estimation of the input to return desired outputs for actual data. Validation data is used to verify whether this estimate performs good on input, the machine has not learned about. Recent efforts such as Google DeepMind's AlphaGo or Neural Algorithms of Artistic Style (by Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge) show the success of Machine Learning in real-world applications. This thesis will cover Bayesian and Neural Networks as two branches of Machine Learning.

Bayesian networks consider probabilities as quantification of belief. Bayes' Theorem is used as a very generic tool to establish a relation between the prior, a likelihood, and the posterior belief. For example, consider a set of random variables modelling symptoms of a disease. Probabilistic dependencies between these symptoms exist inherently. For a set of sample patients, the belief in experiencing certain symptoms and a certain disease is specified. The machine can now learn these input-output relationships. In the following for a new patient, the belief in symptoms is provided and the machine returns a belief whether the patient suffers from a certain disease.

Neural networks are another branch of Machine Learning. Neural networks consist of multiple layers of neurons. Input signals traverse these layers and using sophisticated algorithms, neuron weights are adjusted so that neurons learn which input signals belong to which output class. This research area is prominent since the 1970s, when Paul Werbos described the Backpropagation algorithm in his Ph.D. thesis.

This bachelor thesis sums up fundamental theorems and theoretical background of the aforementioned fields. Furthermore it covers technical details of my implementation to recognize mathematical expressions.

**Keywords:** Machine learning, Bayesian network, Bayesian statistics, Bayes' Theorem, probability theory, Neural networks

# Abstract

抽象は日本語で仕上がります…...

**Keywords:** Machine learning, Bayesian statistics, Bayes' Theorem, Probability theory, Neural networks

## Acknowledgements

I would like to thank my advisor 高山先生 of Kobe University for his continued effort during my studies. Thank you for the valuable input and one year of academic support.

I would also like to thank my Austrian advisor Bredies Kristian for a final revision and grading at University of Graz.

During my year abroad in Japan, I met lot of new people and gained valuable experiences. But I wouldn't have taken the challenge without Martina. Thank you for staying with me and sharing your experiences continuously.

どうもありがとうございました。

All source codes are available at [lukas-prokop.at/proj/bakk\\_kobe](https://lukas-prokop.at/proj/bakk_kobe) and published under terms and conditions of Free/Libre Open Source Software. This document was printed with Xe<sub>La</sub>TeX in the Andada typeface.

# Contents

<b>1</b>	<b>Bayesian theory</b>	<b>1</b>
1.1	Probability theory and statistics . . . . .	1
1.1.1	$\sigma$ -algebra . . . . .	1
1.1.2	Basic definitions . . . . .	1
1.1.3	Average Value and Expected Value . . . . .	2
1.1.4	Continuous probability model . . . . .	3
1.1.5	Variance and standard deviation . . . . .	4
1.1.6	Covariance . . . . .	4
1.1.7	Law of Large Numbers . . . . .	5
1.1.8	Union and intersection of events . . . . .	8
1.1.9	Marginalization . . . . .	8
1.1.10	Joint distribution . . . . .	8
1.1.11	Independence . . . . .	8
1.1.12	Conditional independence . . . . .	9
1.1.13	Bayes' Theorem . . . . .	9
1.2	Probability distributions . . . . .	10

<b>CONTENTS</b>	<b>v</b>
1.2.1 Normal distribution . . . . .	10
1.2.2 Gaussian distribution . . . . .	10
1.2.3 Independent and identically distributed . . . . .	10
1.3 Graphical models . . . . .	10
1.4 Example: Polynomial curve fitting problem . . . . .	11
1.4.1 The problem . . . . .	11
1.4.2 Overfitting . . . . .	11
1.4.3 Regularization as countermeasure . . . . .	12
1.4.4 Maximum Likelihood Estimator . . . . .	12
<b>2 Neural networks</b>	<b>16</b>
2.1 Structure of neural networks . . . . .	16
2.2 Neural networks in practice . . . . .	16
2.2.1 Curve Fitting Problem in Neural Networks . . . . .	16
2.2.2 Backpropagation . . . . .	16
2.2.3 Gradient Descent . . . . .	17
2.3 Google TensorFlow . . . . .	17
2.4 Terminology . . . . .	17
2.5 A problem statement . . . . .	17
2.6 Current state . . . . .	18
2.7 Implementation . . . . .	18
2.7.1 Segmentation . . . . .	18
2.7.2 Recognition . . . . .	19

<i>CONTENTS</i>	vi
2.7.3 Annotation . . . . .	21
2.7.4 Correction . . . . .	21
2.7.5 Output format . . . . .	21

## Chapter 1

# Bayesian theory

## Probability theory and statistics

### $\sigma$ -algebra

The following mathematical object is necessary in order to define probability properly:

**Definition 1** Let  $X$  be a set. A  $\sigma$ -algebra is a set  $Y$  of subsets of  $X$  satisfying:

1.  $\emptyset \in Y$
2.  $Z \in Y \implies Z^C \in Y$  where  $Z^C$  denotes the complement of  $Z$ ,  $X \setminus Z$ .
3.  $(\bigcup_{i=1}^n Z_i) \in Y$  where  $n \in \mathbb{N}$  and  $Z_i \in Y$  where  $1 \leq i \leq n$ .

**Example 1** Let  $X := \{a, b, c, d\}$  and  $Z := \{\{a\}\}$ . We extend  $Z$  to a  $\sigma$ -algebra  $Y$ :

$$Y = \{\{\}, \{a, b, c, d\}, \{a\}, \{b, c, d\}\}$$

### Basic definitions

Probability theory is concerned with random experiments and random phenomena. Probability in its basic form is the fraction of events with a certain outcome to the total number of events.

**Definition 2** A *probability space*  $(\Omega, \mathcal{A}, P)$  denotes the set of possible outcomes, a set of events, and a map from an element of  $\mathcal{A}$  to a real value in  $[0, 1]$  respectively. An *event* is an arbitrary set of elements in  $\Omega$  satisfying the conditions of a  $\sigma$ -algebra. A *random variable* can realize one of the values in  $\mathcal{A}$ .



**Example 2** We toss a coin two times with possible outcomes heads ( $h$ ) or tails ( $t$ ). Then  $\Omega = \{h, t\}$  and  $\mathcal{A} = \{(h, h), (h, t), (t, h), (t, t)\}$ . Let  $R$  be our random variable. Our first experiment yields  $(h, t)$  as result. Our second experiment yields  $(h, h)$ . Then we denote  $\mathbb{P}[R = (h, t)] := 0.5$ ,  $\mathbb{P}[R = (h, h)] := 0.5$ ,  $\mathbb{P}[R = (t, h)] := 0$ , and  $\mathbb{P}[R = (t, t)] := 0$ .

## Average Value and Expected Value

The *average value*  $\mathbb{E}[R]$  of a random variable  $R$  is defined as,

$$\mathbb{E}[R] := \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} a \quad (1.1)$$

If all outcomes in the probability space are considered, we call  $\mathbb{E}[R]$  the *population mean*  $\mu$ , also denoting the expected value of random variable  $R$ .

$$\mathbb{E}[R] := \mu = \sum_{a \in \mathcal{A}} \mathbb{P}[R = a] \cdot a \quad (1.2)$$

So, average and expected values are only defined for numeric events. Our example event  $(h, t)$  does not satisfy this property. Thus, no example is given.

Let  $R$  and  $S$  be two random variables and  $c \in \mathbb{R}$ . The following properties are satisfied:

$$\mathbb{E}[c] := c \quad (1.3)$$

$$\begin{aligned} \mathbb{E}[R + c] &:= \sum_{a \in \mathcal{A}} (\mathbb{P}[R = a] \cdot (a + c)) \\ &= \sum_{a \in \mathcal{A}} (\mathbb{P}[R = a] \cdot a) + \sum_{a \in \mathcal{A}} (\mathbb{P}[R = a] \cdot c) \\ &= \mathbb{E}[R] + c \cdot \sum_{a \in \mathcal{A}} \mathbb{P}[R = a] \\ &= \mathbb{E}[R] + c \cdot 1 = \mathbb{E}[R] + c \end{aligned} \quad (1.4)$$

$$\begin{aligned} \mathbb{E}[R + S] &:= \sum_{a \in (\mathcal{A}_R \cup \mathcal{A}_S)} \begin{cases} \mathbb{P}[R = a] \cdot a & \text{if } a \in \mathcal{A}_R \\ \mathbb{P}[S = a] \cdot a & \text{if } a \in \mathcal{A}_S \end{cases} \\ &= \sum_{a \in \mathcal{A}_R} \mathbb{P}[R = a] \cdot a + \sum_{a \in \mathcal{A}_S} \mathbb{P}[S = a] \cdot a \\ &= \mathbb{E}[R] + \mathbb{E}[S] \end{aligned} \quad (1.5)$$

$$\begin{aligned} \mathbb{E}[c \cdot R] &:= \sum_{a \in \mathcal{A}} \mathbb{P}[R = a] \cdot (c \cdot a) \\ &= c \cdot \sum_{a \in \mathcal{A}} \mathbb{P}[R = a] \cdot a = c \cdot \mathbb{E}[R] \end{aligned} \quad (1.6)$$

### Continuous probability model

Let  $f$  be a continuous function defined in  $(-\infty, \infty) \subseteq \mathbb{R}$ . If  $f$  satisfies the properties 1.7 and 1.8,  $f$  is called a Probability Density Function (PDF):

$$f(x) \geq 0 \quad \forall x \in (-\infty, \infty) \quad (1.7)$$

$$1 = \int_{-\infty}^{\infty} f(x) dx \quad (1.8)$$

#### Example 3

$$f(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The normal distribution function depending on parameters  $\mu$  and  $\sigma^2$  is one example for a probability density function. The function is introduced in Section 1.2.1 in detail.

*Probability* is defined as follows:

$$\mathbb{P}[a \leq R \leq b] := \int_a^b f(x) dx \quad (1.9)$$

The *expected value* in the continuous model is defined as,

$$\mathbb{E}[R] := \int_{-\infty}^{\infty} (\mathbb{P}[R = x] \cdot x) dx = \int_{\mathbb{R}} (\mathbb{P}[R = x] \cdot x) dx \quad (1.10)$$

Let  $R$  and  $S$  be two random variables and  $c \in \mathbb{R}$ . The expected value satisfies:

$$\mathbb{E}[c] := c \quad (1.11)$$

$$\begin{aligned} \mathbb{E}[R + c] &:= \int_{\mathbb{R}} \mathbb{P}[R = x] \cdot (x + c) dx \\ &= \int_{\mathbb{R}} (\mathbb{P}[R = x] \cdot x + \mathbb{P}[R = x] \cdot c) dx \\ &= \int_{\mathbb{R}} \mathbb{P}[R = x] \cdot x dx + c \cdot \int_{\mathbb{R}} \mathbb{P}[R = x] dx \\ &= \mathbb{E}[R] + c \cdot 1 = \mathbb{E}[R] + c \end{aligned} \quad (1.12)$$

$$\begin{aligned} \mathbb{E}[R + S] &:= \int_{\mathbb{R}} (\mathbb{P}[R = x] \cdot x + \mathbb{P}[S = x] \cdot x) dx \\ &= \int_{\mathbb{R}} \mathbb{P}[R = x] \cdot x dx + \int_{\mathbb{R}} \mathbb{P}[S = x] \cdot x dx \\ &= \mathbb{E}[R] + \mathbb{E}[S] \end{aligned} \quad (1.13)$$

$$\begin{aligned} \mathbb{E}[c \cdot X] &:= \int_{\mathbb{R}} \mathbb{P}[R = x] \cdot (x \cdot c) dx \\ &= c \cdot \int_{\mathbb{R}} \mathbb{P}[R = x] \cdot x dx = c \cdot \mathbb{E}[X] \end{aligned} \quad (1.14)$$

Because the expected value operator  $\mathbb{E}$  satisfies the same properties in the discrete and continuous case, we combine those cases and denote  $\mathbb{E}[R]$  for both cases with a random variable  $R$ .

## Variance and standard deviation

*Variance* quantifies how strong values are spread out from their population mean:

$$\sigma^2 = \mathbb{V}[R] := \mathbb{E}[(R - \mu)^2]$$

In the discrete case, this is equivalent to,

$$\mathbb{V}[R] = \mathbb{E} \left[ \sum_{a \in \mathcal{A}} (\mathbb{P}[R = a] \cdot (a - \mu)^2) \right]$$

and in the continuous case, we have:

$$\mathbb{V}[R] = \mathbb{E} \left[ \int_{\mathbb{R}} (\mathbb{P}[R = x] \cdot (x - \mu)^2) dx \right]$$

The *standard deviation* is defined as its second root:

$$\sigma = \sqrt{\mathbb{V}[R]}$$

## Covariance

*Covariance* measures the joint variability of two given random variables. In general, it is defined as:

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

We will exploit the following properties:

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[(Y - \mathbb{E}[Y])(X - \mathbb{E}[X])] = \text{Cov}(Y, X) \quad (1.15)$$

$$\text{Cov}(X, X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{V}[X] \quad (1.16)$$

Let  $X$  be a set of  $n$  independent variables  $X_{1 \leq i \leq n}$ .

$$\begin{aligned} \mathbb{V} \left[ \sum_{i=1}^n X_i \right] &= \mathbb{E} \left[ \left( \sum_{i=1}^n X_i - \mathbb{E} \left[ \sum_{i=1}^n X_i \right] \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \cdot \sum_{j=1}^n (X_j - \mathbb{E}[X_j]) \right] \\ &= \mathbb{E} \left[ \sum_{j=1}^n \left( \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right) (X_j - \mathbb{E}[X_j]) \right] \\ &= \mathbb{E} \left[ \sum_{i,j \in [1,n]} (X_i - \mathbb{E}[X_i]) (X_j - \mathbb{E}[X_j]) \right] \end{aligned}$$

$$\begin{aligned}
&= \sum_{i,j \in [1,n]} \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] \\
&= \sum_{i,j \in [1,n]} \text{Cov}(X_i, X_j) \\
&= \sum_{\substack{i \neq j \\ i,j \in [1,n]}} \text{Cov}(X_i, X_i) + \sum_{\substack{i \neq j \\ i,j \in [1,n]}} \text{Cov}(X_i, X_j) \\
&= \sum_{i=1}^n \mathbb{V}[X_i] + \sum_{\substack{i \neq j \\ i,j \in [1,n]}} \text{Cov}(X_i, X_j) \\
&= \sum_{i=1}^n \mathbb{V}[X_i] + 2 \cdot \sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j) \tag{1.17}
\end{aligned}$$

## Law of Large Numbers

The Law of Large Numbers stresses the practical importance of the expected value.

**Theorem 1 (Law of Large Numbers)** First, we define the notion of the average value over a sample  $A_i$  of size  $n$ :

$$\bar{R}_n := \frac{1}{n} \sum_{i=0}^n A_i$$

Then, the Law of Large Numbers states that,

$$\lim_{n \rightarrow \infty} \bar{R}_n = \mathbb{E}[R]$$

In order to prove this theorem, we use Chebyshev's Inequality, the Weak Law of Large numbers and the Borel-Cantelli Lemma.

Consider the continuous case. Let  $f$  be a PDF,  $a \in \mathbb{R}_{\geq 0}$ ,  $n \in \mathbb{N}$  and let  $\mathbb{E}[R^n]$  be defined as follows:

$$\begin{aligned}
\mathbb{E}[R^n] &= \int_{\mathbb{R}} x^n \cdot f(x) dx \\
&\geq \int_{x \geq a} x^n \cdot f(x) dx \\
&\geq a^n \int_{x \geq a} f(x) dx \\
&= a^n \cdot \mathbb{P}[R \geq a]
\end{aligned}$$

The discrete case follows immediately. This concludes the correctness of the following theorem:

**Theorem 2 (Chebyshev's Inequality Theorem)** Let  $R$  be a random variable,  $a \in \mathbb{R}_{\geq 0}$  and  $n \in \mathbb{N}$ . Then it holds that,

$$\mathbb{P}[R \geq a] \leq \frac{1}{a^n} \mathbb{E}[R^n]$$

The next theorem is called Weak Law of Large Numbers.

**Theorem 3 (Weak Law of Large Numbers, Bernoulli's Theorem)** Let  $R_i$  be a sequence of independent and identically distributed random variables (see section 1.2 for a definition of i.i.d.) with common mean  $\mu$  and variance  $\sigma^2$ . Let

$$S_n := \sum_{i=1}^n R_i \quad T_n := \frac{S_n}{n} - \mu$$

Then for any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}[T_n \geq \varepsilon] = 0$$

**Proof 1**

$$\mathbb{E}[S_n] = \mathbb{E}\left[\sum_{i=1}^n R_i\right] = \sum_{i=1}^n \mathbb{E}[R_i] = \sum_{i=1}^n \mu = n \cdot \mu \quad (1.13) \quad (1.18)$$

$$\begin{aligned} \mathbb{E}[T_n] &= \mathbb{E}\left[\frac{1}{n} (R_1 + R_2 + \dots + R_n) - \mu\right] \\ &= \frac{1}{n} (\mathbb{E}[R_1] + \mathbb{E}[R_2] + \dots + \mathbb{E}[R_n]) - \mathbb{E}[\mu] \\ &= \frac{n \cdot \mu}{n} - \mu = 0 \end{aligned} \quad (1.14) \quad (1.19)$$

In the continuous and discrete case, it holds that  $a^2 \cdot \mathbb{V}[X] = \mathbb{V}[a \cdot X]$ :

$$\begin{aligned} a^2 \cdot \mathbb{V}[X] &= a^2 \cdot \int (x - \mu)^2 \cdot f(x) dx = \int (x \cdot a - \mu \cdot a)^2 \cdot f(x) dx \\ a^2 \cdot \mathbb{V}[X] &= a^2 \cdot \sum_{i=1}^n p_i \cdot (x_i - \mu)^2 = \sum_{i=1}^n p_i \cdot a^2 \cdot (x_i - \mu)^2 = \sum_{i=1}^n p_i \cdot (x_i \cdot a - \mu \cdot a)^2 \end{aligned}$$

The relation  $\mathbb{E}[X^2] = \mathbb{V}[X] + \mathbb{E}[X]^2$  holds as well,

$$\begin{aligned} \mathbb{V}[X] &= \mathbb{E}[(X - \mu)^2] = \mathbb{E}[X^2] - \mathbb{E}[2X\mu] + \mathbb{E}[\mu^2] \\ &= \mathbb{E}[X^2] - 2 \cdot \mathbb{E}[X \cdot \mathbb{E}[X]] + \mathbb{E}[\mathbb{E}[X]^2] \end{aligned} \quad (1.13) \quad (1.20)$$

$$= \mathbb{E}[X^2] - 2 \cdot \mathbb{E}[X]^2 + \mathbb{E}[X]^2 \quad (1.14)$$

$$= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (1.21)$$

$$\begin{aligned} \mathbb{V}[T_n] &= \mathbb{E}[(T_n - \mathbb{E}[T_n])^2] \\ &= \mathbb{E}[T_n^2] \\ &= \mathbb{E}\left[\left(\frac{S_n}{n} - \mu\right)^2\right] \\ &= \mathbb{E}\left[\frac{S_n^2}{n^2} - 2\frac{S_n}{n}\mu + \mu^2\right] \end{aligned} \quad (1.19)$$

$$= \mathbb{E} \left[ \left( \frac{S_n}{n} \right)^2 \right] - \mathbb{E} \left[ \frac{2\mu}{n} S_n \right] + \mathbb{E} [\mu^2] \quad (1.13)$$

$$\begin{aligned} &= \mathbb{V} \left[ \frac{S_n}{n} \right] + \mathbb{E} \left[ \left( \frac{S_n}{n} \right)^2 \right] - 2\mu \cdot \mathbb{E} \left[ \frac{S_n}{n} \right] + \mu^2 \\ &= \frac{1}{n^2} \cdot \mathbb{V} [S_n] + \left( \frac{1}{n} \cdot \mathbb{E} [S_n] \right)^2 - \frac{2\mu}{n} \cdot (n \cdot \mu) + \mu^2 \\ &= \frac{1}{n^2} \cdot (\mathbb{V} [S_n] + \mathbb{E} [S_n]) - 2\mu^2 + \mu^2 \\ &= \frac{1}{n^2} \left( \mathbb{V} \left[ \sum_{i=1}^n R_i \right] + \mathbb{E} [S_n]^2 \right) - 2\mu^2 + \mu^2 \\ &= \frac{1}{n^2} \left( \sum_{i=1}^n \mathbb{V} [R_i] + 2 \sum_{\substack{i < j \\ i, j=1}}^n \text{Cov} [R_i, R_j] + \mathbb{E} [S_n]^2 \right) - \mu^2 \quad (1.17) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{n^2} \left( \sum_{i=1}^n \sigma^2 + 0 + (n \cdot \mu)^2 \right) - \mu^2 \\ &= \frac{1}{n^2} (n \cdot \sigma^2 + n^2 \cdot \mu^2) - \mu^2 = \frac{\sigma^2}{n} + \mu^2 - \mu^2 = \frac{\sigma^2}{n} \quad (1.22) \end{aligned}$$

This was proved the so-called Bienaymé formula. A shorter approach uses the definition of the mean:

$$\mathbb{V} [\overline{T_n}] = \mathbb{V} \left( \frac{1}{n} \sum_{i=1}^n X_i \right) = \frac{1}{n^2} \sum_{i=1}^n \mathbb{V} [X_i] = \frac{\sigma^2}{n}$$

Furthermore, the following equation holds:

$$\begin{aligned} \mathbb{E} [|T_n|^2] &= \mathbb{E} [T_n^2] - 2 \cdot \mathbb{E} [T_n] \cdot \mathbb{E} [|T_n|] + \mathbb{E} [\mathbb{E} [T_n]^2] \\ &= \mathbb{E} [(|T_n| - \mu)^2] = \mathbb{V} [|T_n|] \end{aligned}$$

Now we can apply Chebyshev's Inequality Theorem to  $T_n$  ( $R = T_n$ ,  $a = \varepsilon$ ,  $n = 2$ ):

$$\mathbb{P} [|T_n| \geq \varepsilon] \leq \frac{1}{\varepsilon^2} \mathbb{E} [T_n^2] = \frac{1}{\varepsilon^2} \mathbb{V} [|T_n|]$$

For any  $\varepsilon > 0$  with  $n \rightarrow \infty$ , it holds that

$$\begin{aligned} &\mathbb{P} [|T_n| \geq \varepsilon] \rightarrow 0 \\ &\Leftrightarrow \forall \varepsilon > 0 : \lim_{n \rightarrow \infty} \mathbb{P} [|T_n| \geq \varepsilon] = 0 \end{aligned}$$

This concludes the proof of the Weak Law of Large Numbers. In order to finish our proof of the Strong Law of Large numbers, we will use the Borel-Cantelli Lemma:

**ToDo:** Borel-Cantelli Lemma definition

**ToDo:** Borel-Cantelli Lemma: show that it applies to our case

## Union and intersection of events

We denoted  $\mathcal{A}$  as the set of events. A element of  $\mathcal{A}$  is called *event* and satisfies the properties of a  $\sigma$ -algebra.

We define the union  $U$  of events  $a, b \in \mathcal{A}$  as union of the sets;  $c = a \cup b$ . The semantics following implicitly. Random variable  $R$  realizes  $c$  if event  $a$  or  $b$  or both occur.

**ToDo:** Better illustration, better consistency with marginalization, independence, conditional independence

## Marginalization

$$\mathbb{P}[R = r] = \sum_{s \in S} \mathbb{P}[R = r, S = s]$$

**ToDo:** Explanation

**ToDo:** Illustrative example

## Joint distribution

$$\mathbb{P}[R = r, S = s] = \mathbb{P}[R = r] \cdot \mathbb{P}[S = s]$$

**ToDo:** Explanation

**ToDo:** Illustrative example

## Independence

$$\mathbb{P}[R = r] = \sum_{s \in S} \mathbb{P}[R = r, S = s]$$

**ToDo:** Explanation

**ToDo:** Illustrative example

## Conditional independence

Conditional independence concerns two or more random variables  $R_i$ . Conditional dependence is given if the outcome of a random variable  $R_i$  changes the probability of the outcome of a random variable  $R_j$  with  $i \neq j$ . Let  $A$  and  $B$  be two random variables and  $a$  and  $b$  be two distinctive outcomes. Trivially, we can first observe that:

$$\mathbb{P}[A = a|A = a] = 1$$

$$\mathbb{P}[A = b|A = a] = 0$$

The notation  $\mathbb{P}[A = b|A = a]$  signifies the probability that the event  $A = b$  occurs if event  $A = a$  as condition occurred. In this particular case, the probability is obviously 0, because a random variable cannot take two values simultaneously. The first case evaluates to 1, because if we are assured that random variable  $A$  will be realized with  $a$ , then random variable  $A$  will take value  $a$  with probability 1.

Formally, conditional independence is defined the following way: Random variables  $A$  and  $B$  are conditionally independent if and only if

$$\mathbb{P}[A, B] = \mathbb{P}[A] \cdot \mathbb{P}[B]$$

An interesting question arises for three or more variables.

**ToDo:** definition 3+ variables case

**ToDo:** mutual independence  $\neq$  pairwise independence

## Bayes' Theorem

**Theorem 4 (Bayes' Theorem)** Let  $A$  and  $B$  be two events and  $\mathbb{P}[B] \neq 0$ . Then:

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[B|A] \cdot \mathbb{P}[A]}{\mathbb{P}[B]}$$

**Proof 2** In section 1.1.11, we showed the following relation between marginal and conditional probability:

$$\mathbb{P}[A, B] = \mathbb{P}[B|A] \cdot \mathbb{P}[A] = \mathbb{P}[A|B] \cdot \mathbb{P}[B]$$

Bayes' theorem follows immediately:

$$\frac{\mathbb{P}[B|A] \cdot \mathbb{P}[A]}{\mathbb{P}[B]} = \mathbb{P}[A|B]$$



Bayes' Theorem is fundamental to theory we will cover in the following.  $\mathbb{P}[A]$  is called *prior probability* and  $\mathbb{P}[A|B]$  is called *posterior probability* in the Bayesian interpretation. The names derive from the fact, that  $\mathbb{P}[A]$  is known beforehand in most applications and  $\mathbb{P}[A|B]$  is the degree of belief in  $A$  after  $B$  happened.

## Probability distributions

Probability distributions are templates for probability density functions satisfying the criteria mentioned in Section 1.1.4. They are parameterized by one or more variables and can be continuous or discrete.

### Normal distribution

**ToDo:** definition

**ToDo:** visualization

### Gaussian distribution

**ToDo:** definition

**ToDo:** visualization

### Independent and identically distributed

**ToDo:** definition

**ToDo:** illustrative example

## Graphical models

**ToDo:** Show symmetry, decomposition, weak union and contraction, via Dawid et al.

## Example: Polynomial curve fitting problem

### The problem

**ToDo:** visualization

In the following, we introduce the curve fitting problem similar to [2, p. 4 ff.]. The problem is defined as follows:

**Problem 1 (Polynomial curve fitting problem)** Consider a polynomial of arbitrary degree.

Given  $x = (x_1, \dots, x_n)^N$  as a vector of  $N$  x-values and  $t = (t_1, \dots, t_n)^N$  as the corresponding y-values drawn from the polynomial. Furthermore let  $E(w)$  be an error function for given polynomial coefficients  $w$ .

Find a polynomial with coefficients  $w$  which approximates values  $t$  minimizing  $E(w)$ .

The degree of the polynomial is unknown on purpose. *Model selection* is a branch of Machine Learning dedicated to finding appropriate models for given problems. So for polynomial degree choice for our curve fitting problem, we refer to research literature in Model Selection. **ToDo:** provide useful references for Curve Fitting

Popular error functions include

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2 \quad (\text{Mean squared error, MSE})$$

$$E(w) = \sqrt{\frac{1}{N} \sum_{n=1}^N (y(x_n, w) - t_n)^2} \quad (\text{Root mean square, RMS})$$

$$E(w) = \frac{1}{N} \sum_{n=1}^N (y(x_n, w) - t_n) \quad (\text{Mean signed deviation, MSD})$$

### Overfitting

Machine Learning distinguishes between a *training* and *validation* dataset as input. It uses the training set to learn which output is desired for some given input. Therefore all elements of the training set are labelled such that the error in the output can be quantified. *Overfitting* describes the situation, when the learning algorithm approximates the output with little error, but input from the validation set (which contains different inputs) is computed with high error.

**ToDo:** visualization

## Regularization as countermeasure

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{\lambda}{2} |w|^2$$

We now model the problem from a probabilistic view:

## Maximum Likelihood Estimator

The Maximum Likelihood Estimator (MLE) is a technique to estimate the parameters of a probability distribution. It maximizes the likelihood that the given data actually occurs.

**ToDo:** Illustrate that the Curve Fitting problem is considered Bayesian here

**Theorem 5** Consider input data  $x$ , mean  $\mu$  and variance  $\sigma^2$ :

$$\ln \mathbb{P}[x|\mu, \sigma^2] = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

Then  $\mu_{\text{ML}} = \frac{1}{N} \cdot \sum_{n=1}^N x_n$  for maximized  $\mu$  and  
 $\sigma_{\text{ML}} = \frac{1}{N} \cdot \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$  for maximized  $\sigma^2$

So we want to determine the 2 parameters of a Gaussian distribution, namely  $\mu$  and  $\sigma^2$ , in the maximum likelihood case. We begin with  $\mu$ :

**Proof 3** 1. Derive  $\ln \mathbb{P}[x|\mu, \sigma^2]$  for  $\mu$

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln \mathbb{P}[x|\mu, \sigma^2] &= \frac{\partial}{\partial \mu} \left( -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \right) \\ &= \frac{\partial}{\partial \mu} \left( -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (x_n^2 - 2x_n\mu + \mu^2) - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \right) \\ &= -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (-2x_n + 2\mu) \\ &= -\frac{1}{\sigma^2} \cdot \sum_{n=1}^N (\mu - x_n) \end{aligned}$$

2. Set result zero

$$0 = -\frac{1}{\sigma^2} \cdot \sum_{n=1}^N (\mu - x_n) = \sum_{n=1}^N (\mu - x_n) = N \cdot \mu - \sum_{n=1}^N x_n$$

$$\implies \mu_{\text{ML}} = \frac{1}{N} \cdot \sum_{n=1}^N x_n \quad \text{commonly called "sample mean"}$$

We continue with  $\sigma^2$  and use the same approach:

**Proof 4** 1. Derive  $\ln \mathbb{P}[x|\mu, \sigma^2]$  for  $\sigma^2$

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \ln \mathbb{P}[x|\mu, \sigma^2] &= \frac{\partial}{\partial \sigma^2} \left( -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \right) \\ &= \frac{1}{2\sigma^4} \cdot \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \cdot \frac{1}{\sigma^2} \\ &= \frac{1}{2\sigma^2} \left( \frac{1}{\sigma^2} \cdot \sum_{n=1}^N (x_n - \mu)^2 - N \right) \end{aligned}$$

2. Set result zero

$$\begin{aligned} 0 &= \frac{1}{2\sigma^2} \left( \frac{1}{\sigma^2} \cdot \sum_{n=1}^N (x_n - \mu)^2 - N \right) \\ N \cdot \sigma^2 &= \sum_{n=1}^N (x_n - \mu)^2 \\ \sigma_{\text{ML}}^2 &= \frac{1}{N} \cdot \sum_{n=1}^N (x_n - \mu)^2 \quad \text{commonly called "sample variance"} \end{aligned}$$

And now we derive the precision parameter  $\beta$  in the maximum likelihood case:

**Theorem 6** Given

$$\ln \mathbb{P}[t|x, w, \beta] = -\frac{\beta}{2} \cdot \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

then find

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \cdot \sum_{n=1}^N (y(x_n, w_{\text{ML}}) - t_n)^2$$

by maximizing  $\beta$

**Proof 5** 1. Derive  $\ln \mathbb{P}[t|x, w, \beta]$  with  $\beta$

$$\begin{aligned} \frac{\partial}{\partial \beta} \ln \mathbb{P}[t|x, w, \beta] &= \frac{\partial}{\partial \beta} \left( -\frac{\beta}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \right) \\ &= -\frac{1}{2} \cdot \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{N}{2} \cdot \frac{1}{\beta} \end{aligned}$$

2. Set result zero

$$0 = -\frac{1}{2} \cdot \sum_{n=1}^N (y(x_n, w) - t_n)^2 + \frac{N}{2\beta}$$

$$\frac{N}{\beta} = \sum_{n=1}^N (y(x_n, w) - t_n)^2$$

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \cdot \sum_{n=1}^N (y(x_n, w) - t_n)^2$$

The maximum of the logarithm of an expression corresponds to the minimum of the negative logarithm of the same expression. **ToDo:** so why do we minimize and not maximize?

$$-\log \mathbb{P}[\omega|x, t, \alpha, \beta] \propto -\log [\mathbb{P}[t|x, \omega, \beta] \cdot \mathbb{P}[\omega|\alpha]] \quad (1.23)$$

due to proportionality,  $\exists c \in \mathbb{R}$  such that

$$(1.24)$$

$$= -\log \mathbb{P}[t|x, \omega, \beta] - \log \mathbb{P}[\omega|\alpha] - \log c \quad (1.25)$$

insert formula Bishop 1.62

$$(1.26)$$

$$= \frac{\beta}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi - \log \left( \left( \frac{\alpha}{2\pi} \right)^{\frac{M+1}{2}} \cdot \exp \left( -\frac{\alpha}{2} \omega^T \omega \right) \right) - \log c \quad (1.27)$$

$$= \frac{\beta}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi - \frac{M+1}{2} \log \alpha + \frac{M+1}{2} \log 2\pi + \frac{\alpha}{2} \omega^T \omega - \log c \quad (1.28)$$

Let  $f$  be any function with a minimum. Then  $\arg \min_{\omega} f(\omega) = \arg \min_{\omega} c \cdot f(\omega) + a$  for any  $c, a \in \mathbb{R}$ . This applies also to our case:

$$\arg \min_{\omega} -\log \mathbb{P}[\omega|x, t, \alpha, \beta] = \arg \min_{\omega} \left( \frac{\beta}{2} \cdot \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\alpha}{2} \omega^T \omega \right) \quad (1.29)$$

$$= \arg \min_{\omega} \beta \left( \frac{1}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\frac{\alpha}{\beta}}{2} \omega^T \omega \right) \quad (1.30)$$

$$= \arg \min_{\omega} \frac{1}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\frac{\alpha}{\beta}}{2} \omega^T \omega \quad (1.31)$$

$$(1.32)$$

Hence, the coefficients maximizing the probability that the coefficients correspond to our model parameters  $(x, t, \alpha, \beta)$  are given in the last line. Considering we determine the best coefficients by minimizing the error function, it is justified to consider these coefficients as optimum. Let  $\lambda = \frac{\alpha}{\beta}$ , then ...

$$\tilde{E}(\omega) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\lambda}{2} \|\omega\|^2$$

## Chapter 2

# Neural networks

Neural networks are our tool of choice for our implementation project. We want to recognize mathematical expressions and also consider Kanji recognition. However, due to the ambiguity of handwriting and implicit conventions applied, this implementation cannot compete with industrial applications. Furthermore this implementation only covers online recognition and can therefore not recognize handwriting from images.

## Structure of neural networks

**ToDo:** Explain layers

**ToDo:** Discuss the properties of a neural network of two inputs, one hidden layer, and one output

## Neural networks in practice

### Curve Fitting Problem in Neural Networks

**ToDo:** Illustrate the exercises given by Takayama-sensei

## Backpropagation

**ToDo:** Explain the algorithm

## Gradient Descent

**ToDo:** Explain the algorithm and exercises

**ToDo:** Give convergence proof

## Google TensorFlow

Google TensorFlow is an open-source software library for Machine Intelligence. It provides users a choice between a high-level and a low-level API. This enables them to use classification or regression models to implement Machine Learning applications.

For our application, it is the tool of choice. We used the Python API to implement our handwriting recognition system.

## Terminology

A *glyph* is a visual unit drawn on any surface visible to a potential reader. A *character* is the perceived unit of writing. *Unicode* [5] is an encoding covering most of the world's writing systems. A *Unicode code point* is a unique number assigned to a character.

## A problem statement

The problem of handwriting recognition dates back to the early days of machine learning even before 1960 [1]. Since the beginning, it was considered a more convenient and natural way to input text into a machine. This is because most people learn to write with pen and paper before they learn to type text on a keyboard. The requirement of a keyboard for typing is impractical for mobile applications. Instead, the industry developed devices, where a pen or your finger is used as an input interface on a surface tracking your motion. Handwriting recognition has gained attention on handheld devices such as personal digital assistants (PDA). With the recent advent of smartphones after the year 2000, new input devices were used for handwriting. The screens got large enough such that the user can draw the desired characters and applications recognize the glyphs. Even though the quality has improved tremendously since the 1960s, most users still prefer a virtual keyboard over character drawing interfaces.

In this thesis, we want to focus on two tasks. The first task is a generic version of the second task.



1. Given the stroke traces of drawn characters, recognize the syntax and semantics of the characters and forward the data to a domain-specific application in an appropriate format.
2. Given the stroke traces of mathematical characters, recognize the syntax and semantics of the characters and return its representation in mathematical notation formats such as  $\text{\LaTeX}$  or MathML [6].

has gained a lot of momentum since the increase of the Internet [3].

**ToDo:** Add the wonderful pathological examples mentioned in the papers, I read

## Current state

**ToDo:** Sum up the papers, you read

**ToDo:** Show some existing implementations

## Implementation

Our application implements 5 stages. Namely,

**segmentation** splitting into individual characters,

**recognition** given a drawn character, return the corresponding Unicode point,

**annotation of spatial information** spatial relationship between characters is added,

**correction** using the context information, we apply corrections to the recognized characters,

**output format** we generate the desired MathML/ $\text{\LaTeX}$  output.

### Segmentation

We decided in favor of a very simple implementation for this step. Every stroke is considered as individual character unless it intersects with another stroke. Two or more intersecting strokes are considered as one character. This unit is passed on to the next stage.

## Recognition

First, we define the glyphs, we want to recognize. In general, mathematical notation can be intermingled with text, but we restrict our recognition system to the latin script combined with mathematical symbols. This restriction might be insufficient for other recognition applications. But the task to recognize text in languages such as Japanese, Arabic and Hebrew goes beyond the scope of the problem tackled in this thesis.

We define four properties:

1. We want to recognize mathematical symbols, latin script and numbers. We use Unicode points to reference these characters.
2. A character  $X$  is *visually similar* to another character  $Y$ , if at least one of the following conditions is met <sup>1</sup>:
  - if  $X$  and  $Y$  are commonly written in the same stroke layout by handwriters,
  - $X$  is a larger version of  $Y$  or vice versa,
  - $X$  has the same layout like  $Y$  and only distinguishes itself by the position within the character boundaries.
3. We want to recognize most characters individually. Some characters are visually difficult to distinguish. Two or more visually similar glyphs are represented in a *character group* (a set of Unicode points). We distinguish the characters later context-sensitively.
4. A character group is represented by its *representative* (one Unicode point), i.e. the member with the smallest Unicode point value.

In order to retrieve the list of character groups and representatives, we use the following algorithm:

1. We consider
  - the Latin alphabet of the Basic Latin block (U+0041-U+005A, U+0061-U+007A),
  - the Mathematical Operators block (U+2200-U+22FF),
  - the Supplemental Mathematical Operators block (U+2A00-U+2AFF),
  - extended by Hindu-Arabic digits, and
  - some more characters considered important for mathematical character recognition.

---

<sup>1</sup> Visual similarity as defined above is a reflexive, symmetric and transitive binary relation (hence, an equivalence relation).

Our sources specify 706 Unicode points.

2. These characters are taken and Unicode-normalized [4] using the form NFKD (decomposition by compatibility followed by recomposition respecting canonical equivalence). One exception is (FORKING). This character is not normalized (see below). The resulting characters from the normalization constitute our new list of Unicode points. 8 characters get lost.
3. A custom list defines visually similar characters. Thus some characters are merged into the same character group. 648 Unicode point representatives are left.

This process is required, because character properties such as size are relative. Therefore we defined visual similarity, character classes and their representatives above. Consider, for example, a glyph drawn to represent digit 2. We cannot determine whether it corresponds to 2 (DIGIT TWO) or <sub>2</sub> (SUBSCRIPT TWO), because size is only difference between the characters. Size vastly depends on the size of the drawing area. As such we want to recognize both inputs as 2 (DIGIT TWO) and use context-sensitive information later on to potentially replace it with unicode point <sub>2</sub> (SUBSCRIPT TWO). Unicode normalization provides us a standardized way to normalize the characters.

We want to illustrate this process with examples.

- The character A (LATIN CAPITAL LETTER A) is passed through the whole algorithm and is added in its original form.
- The character II (ROMAN NUMERAL TWO) is normalized to II (LATIN CAPITAL LETTER I, LATIN CAPITAL LETTER I). Because these characters are equivalent to the existing individual symbol I (LATIN CAPITAL LETTER I), nothing is actually changed in the database by adding II (ROMAN NUMERAL TWO).
- However, ||| (LARGE TRIPLE VERTICAL BAR OPERATOR) is not decomposed into three individual code points.
- Sometimes the opposite process takes place. Character (FORKING) is identified as one unicode point point. But normalization returns two characters. Namely, (NONFORKING) and / (COMBINING LONG SOLIDUS OVERLAY). This exception has been removed explicitly in step 2 of the algorithm above.
- The character = (EQUAL TO BY DEFINITION) is a combination of the character = (EQUALS SIGN) and def (LATIN SMALL LETTER D, LATIN SMALL LETTER E, LATIN SMALL LETTER F). Symbol = (EQUALS SIGN) itself is a composition of two characters - (HYPHEN-MINUS). However, Unicode normalization retains the same character.

Because of this process, a hierarchical structure of glyph components is implied. A glyph can be decomposed if two or more non-intersecting strokes are part of a glyph and those strokes occur in at least one other glyph.

**ToDo:** describe stroke simplification algorithm

**ToDo:** Describe the database we used to compare characters with

### **Annotation**

**ToDo:** describe how the spatial information is encoded

### **Correction**

**ToDo:** describe how correction is done

**ToDo:** HMM desired, but was too complex for this thesis

### **Output format**

**ToDo:** Show MathML/LaTeX example

## Bibliography

- [1] C. E. G. Bailey. “Introductory lecture on character recognition”. In: *Proceedings of the IEE - Part B: Electronic and Communication Engineering* 106.29 (1959), pp. 444–445. issn: 0369-8890. doi: [10.1049/pi-b-2.1959.0286](https://doi.org/10.1049/pi-b-2.1959.0286).
- [2] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. isbn: 0387310738.
- [3] Kam-Fai Chan and Dit-Yan Yeung. “Mathematical expression recognition: a survey”. In: *International Journal on Document Analysis and Recognition* 3.1 (2000), pp. 3–15.
- [4] Unicode Consortium. *UAX #15: Unicode Normalization Forms*. <http://www.unicode.org/reports/tr15/>. [Online; accessed 02-August-2016]. 2017.
- [5] Unicode Consortium. *Unicode Consortium*. <http://www.unicode.org/>. [Online; accessed 02-August-2016]. 2017.
- [6] The W3C Math Working Group. *W3C Math Home*. since 1996. url: <https://www.w3.org/TR/MathML3/>.