



# Inhaltsverzeichnis

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Overview . . . . .	1
1.2	Chapter 1.1 - What is machine learning . . . . .	1
1.3	Chapter 1.2 - What is deep learning . . . . .	2
1.3.1	What is a neural network (NN) . . . . .	3
1.4	Chapter 1.3 - Examples for Deep Learning . . . . .	3
<b>2</b>	<b>Tools for Deep Learning</b>	<b>3</b>
2.1	Datasets . . . . .	3
<b>3</b>	<b>Machine learning basics</b>	<b>3</b>
3.1	Linear Algebra . . . . .	4
3.2	Chapter 3.2 Random variable and probability distribution . . . . .	4
3.2.1	3.2.1 One random vector . . . . .	4
3.3	Chapter 3.3 - Multiple random vectors . . . . .	6
3.3.1	Chapter 3.2.3 - Kernel basaed density estimation . . . . .	6
3.4	Chapter 3.4 Kullback-Leibler divergence and cross entropy . . . . .	7
3.4.1	E3.5 KLD between normal and Laplace distribution . . . . .	8
3.5	Chapter 3.5 Probabilistic framework . . . . .	9
3.5.1	Role of a NN . . . . .	11
<b>4</b>	<b>Dense Neural Networks</b>	<b>11</b>
4.0.1	4.1 Fully connected neural networks - Neuron . . . . .	12
4.0.2	Chapter 4.2 Layer of Nurons . . . . .	12
4.1	4.3 Feedforward neural network . . . . .	13

# Abbildungsverzeichnis

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Mathematical notations . . . . .	1
1.2	Three steps of machine learning . . . . .	2
1.3	Neural Network . . . . .	3
<b>2</b>	<b>Tools for Deep Learning</b>	<b>3</b>
<b>3</b>	<b>Machine learning basics</b>	<b>3</b>
3.1	Vector Normes . . . . .	4
3.2	One random vector . . . . .	4
3.3	Moments of a vector . . . . .	5

3.4	Multivariate normal (Gaussian) distribution . . . . .	5
3.5	Kernel function . . . . .	6
3.6	Estimated PDF function . . . . .	7
3.7	Forward vs. Backward KL divergence . . . . .	8
3.8	Probabilistic framework of supervised learning . . . . .	9
3.9	The data generating distribution . . . . .	9
3.10	Bayes Rule in DL . . . . .	10
3.11	Bayes decision theorem . . . . .	10
3.12	Supervised learning . . . . .	10
3.13	Calc part1 . . . . .	11
3.14	Calc part2 . . . . .	11
<b>4</b>	<b>Dense Neural Networks</b>	<b>11</b>
4.1	Neuron . . . . .	12
4.2	Layer of Neurons . . . . .	12
4.3	Meanings of phi . . . . .	13
4.4	Feedforward multilayer neural network . . . . .	13
4.5	Network Parameters . . . . .	14
4.6	. . . . .	14

# List of Equations

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Tools for Deep Learning</b>	<b>3</b>
<b>3</b>	<b>Machine learning basics</b>	<b>3</b>
<b>4</b>	<b>Dense Neural Networks</b>	<b>11</b>

# Lecture 1: Introduction

Date: 03/05/2020

Lecturer: David Silver

By: Nithish Moudhgalya

## 1.1. Overview

• p.4 expert talks are not relevant for the exam

• **scalar**:  $x, A, \alpha$

• column **vector**:  $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = [x_i] \in \mathbb{R}^N$  with element  $x_i = [\underline{x}]_i$

• **matrix**:  $\mathbf{A} = [a_{ij}]_{1 \leq i \leq M, 1 \leq j \leq N} = [a_{ij}] \in \mathbb{R}^{M \times N}$  with element  $a_{ij} = [\mathbf{A}]_{ij}$

• 3D, 4D **tensor**:  $\mathbf{A} = [a_{ijk}], \mathbf{A} = [a_{ijkl}]$

• **transpose** of vector and matrix:  $\underline{x}^T = [x_1, \dots, x_N], \mathbf{A}^T = [a_{ji}]_{ij}$

• **determinant** of a square matrix  $\mathbf{A}$ :  $|\mathbf{A}|$

• **inverse** of a square matrix  $\mathbf{A}$ :  $\mathbf{A}^{-1}$

• **trace** of a square matrix  $\mathbf{A}$ :  $\text{tr}(\mathbf{A}) = \sum_i a_{ii}$

Abbildung 1.1: Mathematical notations

• Element notation (element of vector):  $[x_i = \underline{a} + \underline{b}]_i$

## 1.2. Chapter 1.1 - What is machine learning

• Signal processing is not a subset of ML or the other way around, they are identical in the task

• Basic difference is in how to design the processing rule

• Regression means to calculate (SP, ML) a continuous-valued output in the real numbers from a signal, can be one-dimensional

• Classification means to calculate (SP, ML) a discrete-valued output from the natural numbers

• p. 1-5 filter in the middle is the view of a pixel and its 8 neighbours

## Three steps of machine learning (1)

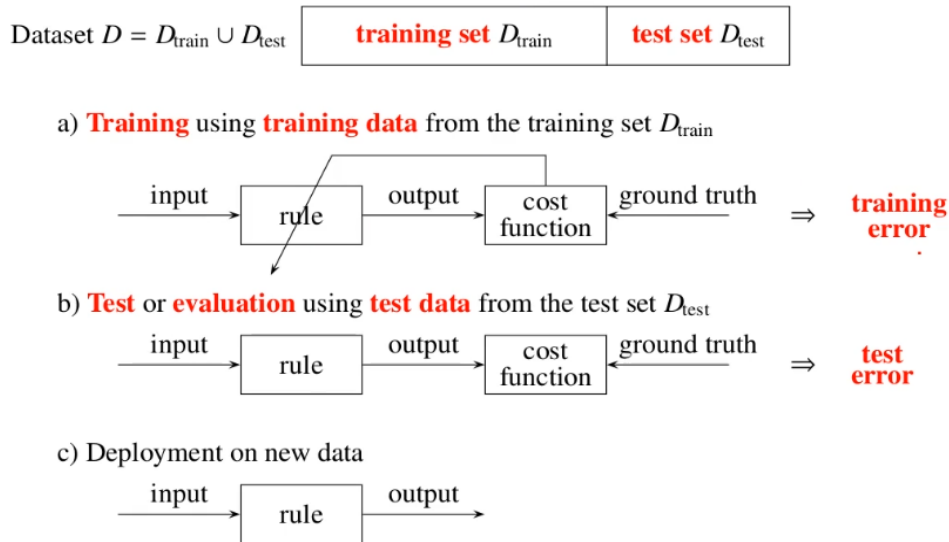


Abbildung 1.2: Three steps of machine learning

- over-fitting just memorizes the training set so test set is needed
- test error rate is similar to training error rate  $\rightarrow$  no over-fitting
- Test set and training set need to be disjoint concatenated
- TODO: Summary supervised learning and unsupervised learning

## 1.3. Chapter 1.2 - What is deep learning

- feature extraction: a feature is a clustered subset of information for recognition
- fish example: fish length, color etc.
- Needs human experience, can't be calculated
- DNN solves the problem without feature extraction

### 1.3.1. What is a neural network (NN)

- Cascaded pipeline of layers

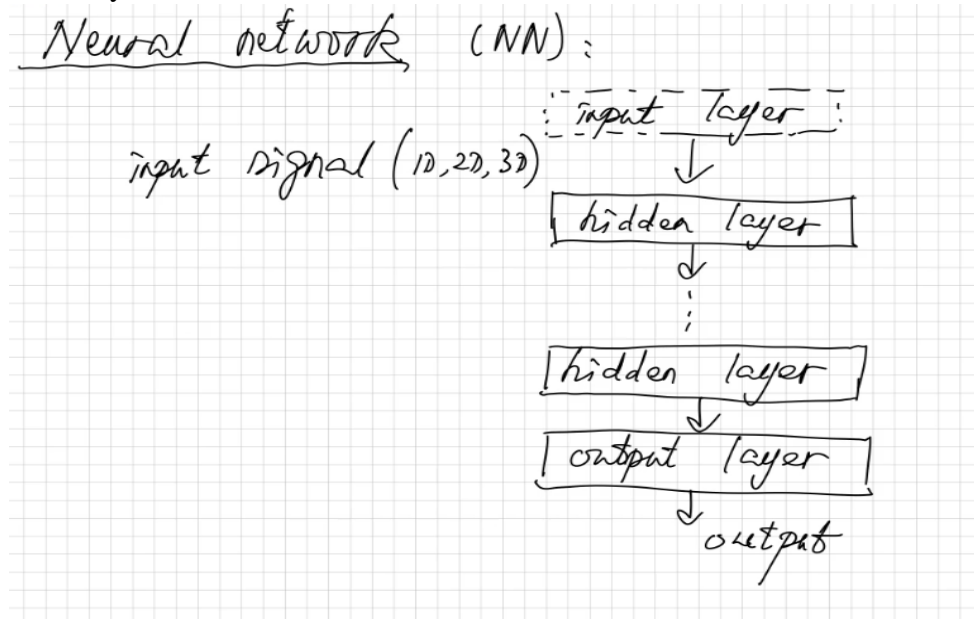


Abbildung 1.3: Neural Network

## 1.4. Chapter 1.3 - Examples for Deep Learning

- Semantic image segmentation is pixelwise classification

### Lecture 2: Tools for Deep Learning

Date: 02/05/2020

Lecturer: David Silver

By: Nithish Moudhgalya

## 2.1. Datasets

- Overview for Training Datasets up to ImageNet they are for teaching

### Lecture 3: Machine learning basics

Date: 02/05/2020

Lecturer: David Silver

By: Nithish Moudhgalya

## 3.1. Linear Algebra

•TODO get Math nicely into the Context book

3.1 Linear algebra

3-4

### Vector norms

Given a vector  $\underline{x} = [x_i] \in \mathbb{R}^M$ . There are different definitions for the vector norm:

- **2-norm** or  **$l_2$ -norm** or **Euclidean norm**:  $\|\underline{x}\|_2 = \sqrt{\sum_{i=1}^M x_i^2} = \sqrt{\underline{x}^T \underline{x}}$
- **1-norm** or  **$l_1$ -norm**:  $\|\underline{x}\|_1 = \sum_{i=1}^M |x_i|$
- **0-norm** or  **$l_0$ -norm**:  $\|\underline{x}\|_0 = \text{number of non-zero elements in } \underline{x}$

Comments:

- $\|\underline{x}\|_2^2$  represents the energy of  $\underline{x}$ .
- $\|\underline{x}\|_0$  measures the sparsity of  $\underline{x}$ .
- Different vector norms have different unit-norm contour lines  $\{\underline{x} \mid \|\underline{x}\|_p = 1\}$ .

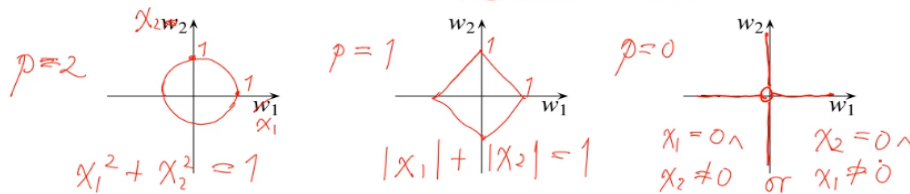


Abbildung 3.1: Vector Normes

## 3.2. Chapter 3.2 Random variable and probability distribution

### 3.2.1. 3.2.1 One random vector

- **scalar**:  $x, A, \alpha$
- **column vector**:  $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = [x_i] \in \mathbb{R}^N$  with element  $x_i = [\underline{x}]_i$
- **matrix**:  $\mathbf{A} = [a_{ij}]_{1 \leq i \leq M, 1 \leq j \leq N} = [a_{ij}] \in \mathbb{R}^{M \times N}$  with element  $a_{ij} = [\mathbf{A}]_{ij}$
- **3D, 4D tensor**:  $\mathbf{A} = [a_{ijk}], \mathbf{A} = [a_{ijkl}]$
- **transpose** of vector and matrix:  $\underline{x}^T = [x_1, \dots, x_N], \mathbf{A}^T = [a_{ji}]_{ij}$
- **determinant** of a square matrix  $\mathbf{A}$ :  $|\mathbf{A}|$
- **inverse** of a square matrix  $\mathbf{A}$ :  $\mathbf{A}^{-1}$
- **trace** of a square matrix  $\mathbf{A}$ :  $\text{tr}(\mathbf{A}) = \sum_i a_{ii}$

Abbildung 3.2: One random vector

PDF for a discrete-valued RV:

$p(\underline{x}) = \sum_i p_i \delta(\underline{x} - \underline{x}_i)$  Dirac function

**cumulative distribution function CCDF**

$$F(\underline{x}) = \int_{-\infty}^{\infty} p(\underline{z}) d\underline{z}, \quad p(\underline{x}) = \frac{\partial^d F(\underline{x})}{\partial x_1 \dots \partial x_d}$$

- **scalar**:  $x, A, \alpha$
- column **vector**:  $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = [x_i] \in \mathbb{R}^N$  with element  $x_i = [\underline{x}]_i$
- **matrix**:  $\mathbf{A} = [a_{ij}]_{1 \leq i \leq M, 1 \leq j \leq N} = [a_{ij}] \in \mathbb{R}^{M \times N}$  with element  $a_{ij} = [\mathbf{A}]_{ij}$
- 3D, 4D **tensor**:  $\mathbf{A} = [a_{ijk}], \mathbf{A} = [a_{ijkl}]$
- **transpose** of vector and matrix:  $\underline{x}^T = [x_1, \dots, x_N], \mathbf{A}^T = [a_{ji}]_{ij}$
- **determinant** of a square matrix  $\mathbf{A}$ :  $|\mathbf{A}|$
- **inverse** of a square matrix  $\mathbf{A}$ :  $\mathbf{A}^{-1}$
- **trace** of a square matrix  $\mathbf{A}$ :  $\text{tr}(\mathbf{A}) = \sum_i a_{ii}$

Abbildung 3.3: Moments of a vector

special case  $d = 1$ :  $\underline{X} \rightarrow X \in \mathbb{R}$

$\mu \rightarrow \mu \in \mathbb{R}$

$\underline{C} \rightarrow$  variance of  $X = \text{Var}(X) = \delta^2 = E[(X - \mu)^2] = \dots = E(X^2) - \mu^2$

$\delta = \sqrt{\text{Var}(X)}$ : standard deviation

For any function  $g(\underline{X})$  of  $\underline{X}$ :  $E[g(\underline{x})] = \int g(x) \cdot p(\underline{x}) d\underline{x} = (d.v.) \sum_i g(\underline{x}_i) \cdot P(\underline{x}_i)$

### Multivariate normal (Gaussian) distribution

#### PDF

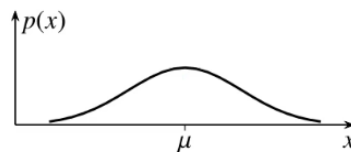
$$\underline{X} \in \mathbb{R}^d \sim N(\underline{\mu}, \mathbf{C}) \quad \boxed{\underline{x}} \cdot \boxed{\underline{\mu}} \cdot \boxed{\mathbf{C}}$$

$$p(\underline{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{C}|^{1/2}} e^{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \mathbf{C}^{-1}(\underline{x} - \underline{\mu})}$$

$$\ln(p(\underline{x})) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\mathbf{C}|) - \frac{1}{2}(\underline{x} - \underline{\mu})^T \mathbf{C}^{-1}(\underline{x} - \underline{\mu})$$

$N(0, \mathbf{I})$  is called the **standard normal distribution**.

#### 1D-Visualization



$$\underline{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

identity matrix

#### Moments

$$E(\underline{X}) = \underline{\mu}, \quad \text{Cov}(\underline{X}) = \mathbf{C}$$

Abbildung 3.4: Multivariate normal (Gaussian) distribution

**one-hot coding**: Only one bit is 1 e.g. 0100 one cold coding is the inverse

Coding for class label, random vector  $\mathbf{y}$  of length  $c$  so the identity matrix with dimension  $c$  is used for class labels

Reformulation of the PMF (categorical distribution) by one-hot coding of the classes

$\underline{x} = [x_i] \in \{\underline{e}_1, \underline{e}_2, \dots, \underline{e}_c\}$ , i.e. all  $x_i = 0$  except for one single element equal to 1

$$\text{PMF: } P(\underline{X} = \underline{x}) = P(\underline{x}) = \begin{cases} P_i \text{ if } \underline{x} = \underline{e}_1 \text{ or } x_1 = 1 \\ \dots \\ P_c \text{ if } \underline{x} = \underline{e}_c \text{ or } x_c = 1 \end{cases} = p_1^{x_1} \cdot p_2^{x_2} \cdot \dots \cdot p_c^{x_c} = \prod_{i=1}^c p_i^{x_i}$$



$$\ln(P(\underline{x})) = \sum_{i=1}^c x_i \cdot \ln(p_i) = [x_1, \dots, x_c] \cdot \begin{bmatrix} \ln(P_1) \\ \dots \\ \ln(P_c) \end{bmatrix} = \underline{x}^T \cdot \ln(\underline{P})$$

In function applied element wise

### 3.3. Chapter 3.3 - Multiple random vectors

•3-16 Table for distributions

•product rule for probability  $p(\underline{x}, \underline{y}) = p(\underline{x}) \cdot p(\underline{y})$

•Bayes rule  $p(\underline{y}|\underline{x}) = p(\underline{y}) \cdot \frac{p(\underline{x}|\underline{y})}{p(\underline{x})}$

•Independent and identically distributed

$\underline{x}$  and  $\underline{y}$  are independent if:

$$p(\underline{x}, \underline{y}) = p(\underline{x}) \cdot p(\underline{y}) \leftrightarrow p(\underline{x}|\underline{y}) = p(\underline{x}) \text{ or } p(\underline{x}|\underline{y}) = p(\underline{y})$$

$\underline{x}_1, \dots, \underline{x}_N$  are independent and identically distributed (i.i.d)

$$P(\underline{x}_1, \dots, \underline{x}_N) = \prod_{i=1}^N p_i(\underline{x}_i), \underline{X}_i \sim p_i(\underline{x}_i)$$

$$p_i(\underline{x}_i) = p(\underline{x}_i) \rightarrow p(\underline{x}_1, \dots, \underline{x}_N) = \prod_{i=1}^N p(\underline{x}_i)$$

#### 3.3.1. Chapter 3.2.3 - Kernel based density estimation

PDF:  $p(\underline{x})$  of  $\underline{X} \in \mathbf{R}^d$  unknown, only i.i.d samples  $\underline{x}(n), 1 \leq n \leq N$

kernel-based estimate  $\hat{p}(\underline{x})$  of  $p(\underline{x})$  from  $\underline{x}(n)$

kernel function  $k(\underline{x})$ , like a PDF

$$1. k(\underline{x}) \geq 0 \forall \underline{x}$$

$$2. \int k(\underline{x}) d\underline{x} = 1$$

$$\hat{p}(\underline{x}) = \frac{1}{N} \sum_{n=1}^N k(\underline{x} - \underline{x}(n))$$

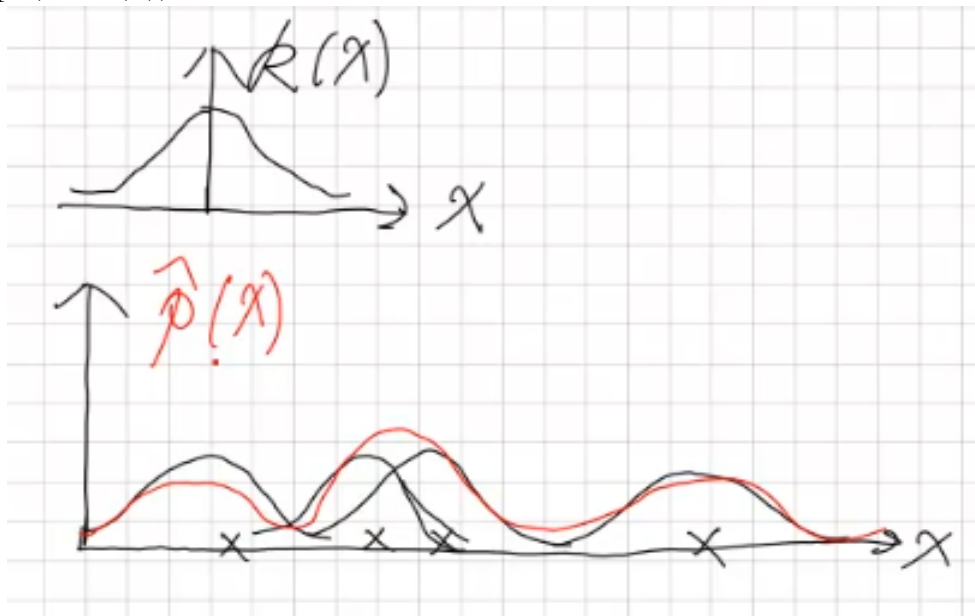


Abbildung 3.5: Kernel function

**Smooth Gaussian Kernel:**

$$N(0, I) : k(\underline{x}) = \frac{1}{2\pi^{\frac{d}{2}}} \cdot e^{-\frac{1}{2} \|\underline{x}\|^2}$$

**Dirac Kernel**

$k(\underline{x}) = \delta(\underline{x})$  : Dirac function

$$\delta(\underline{x}) = \begin{cases} \infty, \underline{x} = \underline{0} \\ 0, \underline{x} \neq \underline{0} \end{cases}$$

$$\int \delta(\underline{x}) d\underline{x} = 1$$

sampling property :  $\int \delta(\underline{x} - \underline{x}_0) f(\underline{x}) d\underline{x} = f(\underline{x}_0)$

**empirical distribution**

$$\hat{p}(\underline{x}) \cdot \frac{1}{N} \sum_{n=1}^N \delta(\underline{x} - \underline{x}(n))$$

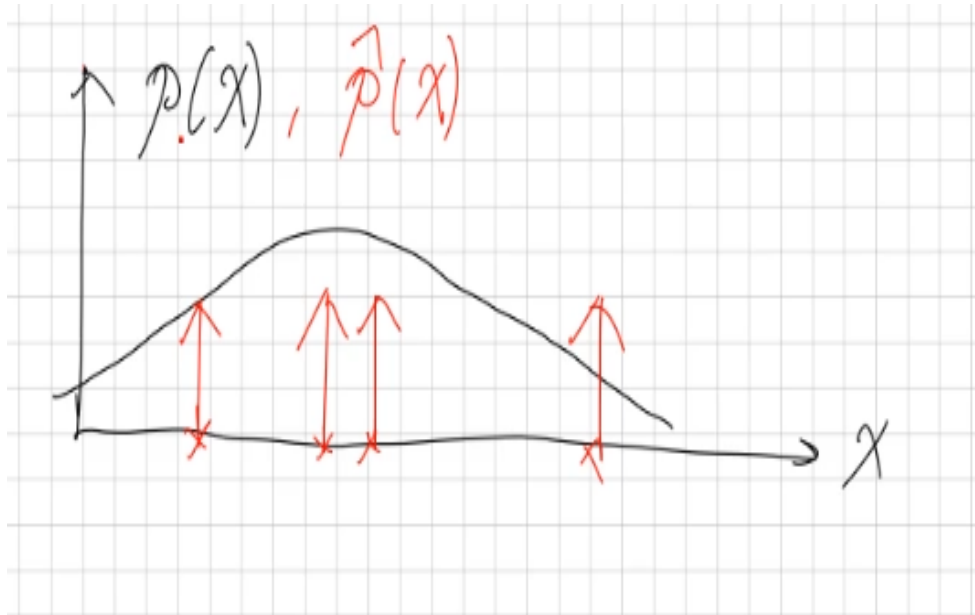


Abbildung 3.6: Estimated PDF function

### 3.4. Chapter 3.4 Kullback-Leibler divergence and cross entropy

Dissimilarity measure between 2 distributions:

**Case A: continuous-valued random variables: PDF**

$\underline{X} \sim p(\underline{x})$  : true statistical distribution of  $\underline{X}$

$q(\underline{x})$ : approximation for  $p(\underline{x})$ , e.g. by DNN

**KL divergence (KLD) between p and q:**

$$D_{KL}(p||q) = \int p(\underline{x}) \cdot \ln\left(\frac{p(\underline{x})}{q(\underline{x})}\right) d\underline{x}$$

$$= E_{\underline{X} \sim p}[\ln\left(\frac{p(\underline{X})}{q(\underline{X})}\right)]$$

expectation over  $p(\underline{x})$

DKL is real valued scalar positive or negative or 0

**Case B: discrete-valued random vector: PMF**

$\underline{X} \in \{\underline{x}_1, \dots, \underline{x}_c\} \sim$ : true PMF of  $\underline{X} \sim Q(\underline{x})$ : approximation for  $P(\underline{x})$

$$D_{KL}(P||Q) = \sum_{i=1}^c P(\underline{x}_i) \cdot \ln\left(\frac{P(\underline{x}_i)}{Q(\underline{x}_i)}\right) = E_{\underline{X} \sim P}[\ln\left(\frac{P(\underline{x})}{Q(\underline{x})}\right)]$$

Properties of the KL divergence: P1) Nonnegative  $D_{KL}(P||Q) \geq 0 \forall p, q$

P2) Equality  $D_{KL}(P||Q) = 0$  iff (if and only if)  $p(\underline{x}) = q(\underline{x})$

proof for sufficient":  $\ln\left(\frac{p(\underline{x})}{q(\underline{x})}\right) = 0 \forall \underline{x}$

P1 and P2:  $D_{KL}(p||1)$  is a suitable metric for approximation p by q

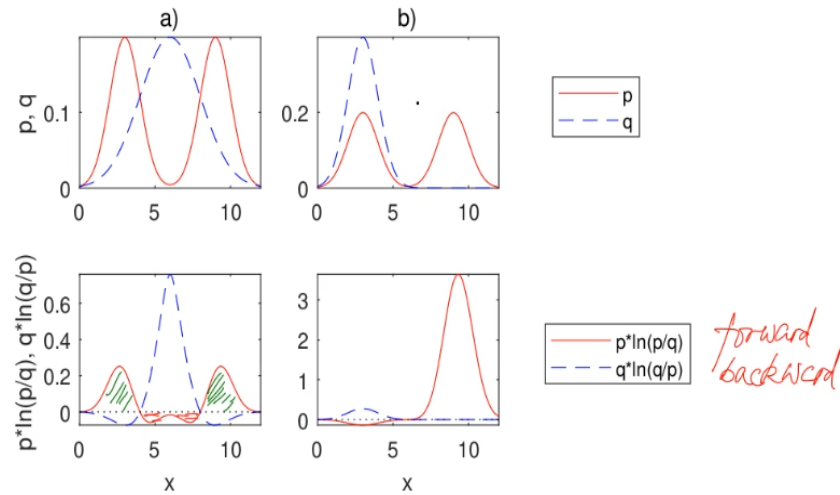
P3 Asymmetry

$$D_{KL}(p||q) = E_{\underline{X} \sim p}[\ln\left(\frac{p(\underline{x})}{q(\underline{x})}\right)] \neq D_{KL}(q||p) = E_{\underline{X} \sim q}[\ln\left(\frac{q(\underline{x})}{p(\underline{x})}\right)]$$

forward KLD

backward KLD

$D_{KL}$  is not a true distance measure with  $D(\underline{x}, \underline{y}) = D(\underline{y}, \underline{x})$



- When minimizing  $D_{KL}(p||q)$ , a) is better than b) because  $q(x)$  is broad.
- When minimizing  $D_{KL}(q||p)$ , b) is better than a) because  $q(x)$  is narrow.

Abbildung 3.7: Forward vs. Backward KL divergence

### 3.4.1. E3.5 KLD between normal and Laplace distribution

$$p(x) \sim N(0, \sigma^2), p(x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{x^2}{2\sigma^2}}$$

$$q(x) \sim \text{Laplace}: (0, b) q(x) = \frac{1}{2b} e^{-\frac{|x|}{b}}$$

Task: choose b to best approximate p by q.

$$\frac{p(x)}{q(x)} = \sqrt{\frac{2}{\pi}} \cdot \frac{b}{\sigma} \cdot \exp\left(-\frac{x^2}{2\sigma^2} + \frac{|x|}{b}\right)$$

$$D_{KL}(p||q) = E_{\underline{X}} \text{simp} = \ln\left(\frac{p(\underline{x})}{q(\underline{x})}\right) = \ln\left(\sqrt{\frac{2}{\pi}} \cdot \frac{b}{\sigma}\right) + E_{X \sim p}\left(\left(-\frac{x^2}{2\sigma^2} + \frac{|x|}{b}\right)\right)$$

$$E_{\underline{X}} \text{simp} = \sigma^2$$

$$E_{\underline{X}} \text{simp}(|x|) = \int_{-\infty}^{\infty} |x| \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = 2 \int_0^{\infty} | \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = \sqrt{\frac{2}{\pi}} \cdot \sigma$$

$$\text{Let } \alpha = \frac{\sigma}{b}, D_{KL}(p||q) = \dots = \sqrt{\frac{2}{\pi}} \cdot \alpha - \ln(\alpha) + \ln\left(\sqrt{\frac{2}{\pi}}\right) - \frac{1}{2} \frac{dD_{KL}(p||q)}{d\alpha} = \sqrt{\frac{2}{\pi}} - \frac{1}{\alpha} = 0 \rightarrow \alpha = \sqrt{\frac{2}{\pi}}, \text{ i.e. } b \approx 0,86$$

$$D_{KL, \min}(p||q) = D_{KL}(p||q)|_{\alpha} = \sqrt{\frac{2}{\pi}} = \dots = \frac{1}{2} - \ln\left(\frac{\pi}{2}\right) \approx 0,048$$

• Probable exam question calculate this for 2 distributions

P4) Additive :

$\underline{X} = (\underline{x}_1, \underline{x}_2)$ ,  $\underline{x}_1$  and  $\underline{x}_2$  are independent, i.e.

$$p(\underline{x}) = p_1(\underline{x}_1) \cdot p_2(\underline{x}_2), q(\underline{x}) = q_1(\underline{x}_1) \cdot q_2(\underline{x}_2)$$

$$\text{Then: } D_{KL}(p||q) = D_{KL}(p_1||q_1) + D_{KL}(p_2||q_2)$$

P5) Relation to cross entropy :

Definition of entropy 3-24 probability always greater than 0 but smaller than 1

$$D_{KL}(p||q) = \int p \ln\left(\frac{p}{q}\right) d\underline{x} = \int p \ln(p) d\underline{x} - \int p \ln(q) d\underline{x} = -H(p) + H(p, q) \text{ or}$$

$$\text{cross entropy: } H(p, q) = D_{KL}(p||q) + H(p) \geq H(p) \geq 0$$

For a given (fixed)  $p(\underline{x})$  :  $H(p)$  fixed

**Hence:**  $\min D_{KL}(p||q) \leftrightarrow \min H(p, q)$

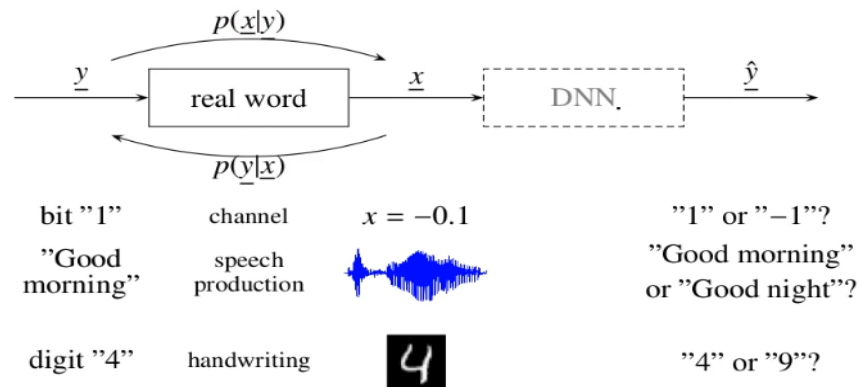
Not the case for backward KLD  $D_{KL}(q||p)$ !

Minimizing is not the same anymore because then it is  $H(q)$  and that's what we are trying to optimize

## 3.5. Chapter 3.5 Probabilistic framework

valid for both SP and ML

valid for both regression problem and classication problem



- $y$ : **latent variable**, hidden, not directly measurable, quantity of interest
- $x$ : **observed variable**, measurement, input for DNN
- $\hat{y}$ : output of DNN as estimate for  $y$
- real world: describes how is  $x$  generated from  $y$
- DNN: describes how to estimate  $y$  from  $x$

Abbildung 3.8: Probabilistic framework of supervised learning

Both  $y$  and  $x$  are modeled as random variables. They are described by the joint **data generating distribution**

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x).$$

$p(y)$  **prior** PDF of  $y$  or **prior**, available before any measurement of  $x$

$p(x|y)$  **likelihood**. It describes the real world, the generation of  $x$  from  $y$ . It is a kind of channel-sensor model, e.g.

- bit: communication channel + receiver
- speech: speech production system + microphone
- digit: handwriting + camera

$p(x)$  prior PDF of  $x$ , also called **evidence**

$p(y|x)$  **posterior** PDF of  $y$  after a measurement  $x$  or **posterior**

Abbildung 3.9: The data generating distribution

Bayes Rule:

Bayes rule :

$$p(\underline{y}|\underline{x}) = p(\underline{x}|\underline{y}) \cdot \frac{p(\underline{y})}{p(\underline{x})}$$

posterior      likelihood      prior      evidence

i.e.  $p(\underline{y}|\underline{x})$  contains both  $p(\underline{x}|\underline{y})$  and  $p(\underline{y})$

Training in ML : calculate/model  $p(\underline{y}|\underline{x}) \forall \underline{x}, \underline{y}$  from  $D_{\text{train}}$

Inference "... : Draw conclusion about  $\underline{y}$  for a given  $\underline{x}$  based on  $p(\underline{y}|\underline{x})$

Abbildung 3.10: Bayes Rule in DL

a) **Maximum a posterior (MAP)** inference:

$$\max_{\underline{y}} p(\underline{y}|\underline{x}) = p(\underline{x}|\underline{y}) \frac{p(\underline{y})}{p(\underline{x})}$$

b) **Minimum Bayesian risk (MBR)** inference:

$$\min_{\hat{\underline{y}}} E_{\underline{x}, \underline{y}}[l(\underline{y}, \hat{\underline{y}}(\underline{x}))] = \int l(\underline{y}, \hat{\underline{y}}(\underline{x})) p(\underline{x}, \underline{y}) d\underline{x} d\underline{y}$$

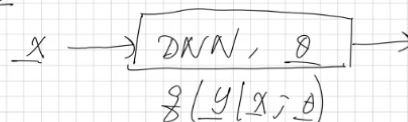
see course DPR and SASP for more details.

Abbildung 3.11: Bayes decision theorem

Supervised learning :

\*  $p(\underline{x}|\underline{y})$ ,  $p(\underline{y})$  unknown  $\rightarrow p(\underline{y}|\underline{x})$  unknown

\* approximate  $p(\underline{y}|\underline{x})$  by a parametric posterior model  $g(\underline{y}|\underline{x}; \underline{\theta})$ , given by a DNN with parameter vector  $\underline{\theta}$



▷ function  $g(\cdot)$  known  $\leftarrow$  DNN architecture

▷  $\underline{\theta}$  unknown  $\leftarrow$  " coefficient

\* learning  $\underline{\theta}$  from  $D_{\text{train}}$

Abbildung 3.12: Supervised learning

## Learning Criterion:

Learning criterion:

$$\min_{\Theta} D_{KL}(p(x, y) \| q(x, y; \Theta)) \xleftrightarrow{p(x, y) \text{ fixed}}$$

$$\min_{\Theta} H(\quad, \quad, \quad)$$

Since  $q(x, y; \Theta) = q(y|x; \Theta) \cdot q(x)$ ,

$$\begin{aligned} \text{CE } H(p(x, y), q(x, y; \Theta)) &= -\int p(x, y) \ln q(x, y; \Theta) dx dy \\ &= -\int p(x, y) \cdot \ln q(y|x; \Theta) dx dy \\ &\quad - \int p(x, y) \cdot \ln q(x) dx dy \\ &\quad \underbrace{\hspace{10em}}_{\text{independent of } \Theta, \text{ const}} \end{aligned}$$

Abbildung 3.13: Calc part1

$$= \int p(x, y) \cdot \underbrace{[-\ln q(y|x; \Theta)]}_{\text{loss}} dx dy + \text{const.}$$

average loss, Bayesian risk, see DPR

In practice:  $p(x, y)$  replaced by the empirical distrib.

$$\hat{p}(x, y) = \frac{1}{N} \sum_{n=1}^N \delta(x - x(n), y - y(n))$$

3.2.3: sampling property of  $\delta(\cdot)$ :

$$\min_{\Theta} H(\hat{p}, q) = \text{const} + \frac{1}{N} \sum_{n=1}^N \underbrace{\left[ -\ln q(y(n)|x(n); \Theta) \right]}_{\substack{\text{loss } L(x(n), y(n); \Theta) \\ \text{cost function } L(\Theta)}}$$

Abbildung 3.14: Calc part2

### 3.5.1. Role of a NN

1. approximate true posterior  $p(y|x)$  by  $q(y|x; \Theta)$
2. learn  $\Theta$  from  $D_{\text{train}}$

## Lecture 4: Dense Neural Networks

Date: 05/05/2020

Lecturer: David Silver

By: Nithish Moudhgalya

A general model for  $q(y|x; \Theta)$

\*) can learn any linear or nonlinear mapping

\*) suitable for both regression and classification problems

artificial NN: mimic biological NN (brain)

#### 4.0.1. 4.1 Fully connected neural networks - Neuron

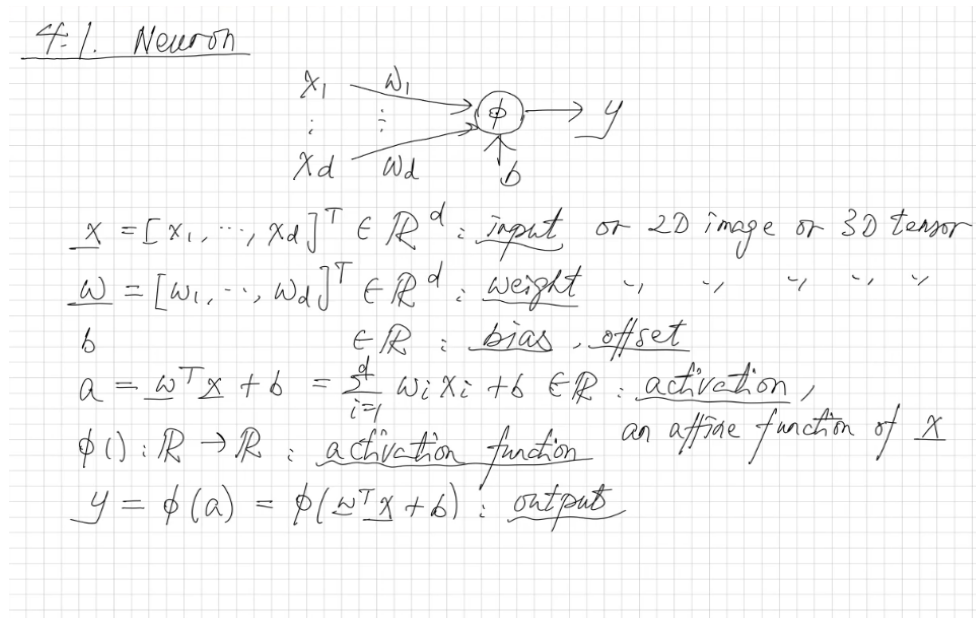


Abbildung 4.1: Neuron

#### 4.0.2. Chapter 4.2 Layer of Neurons

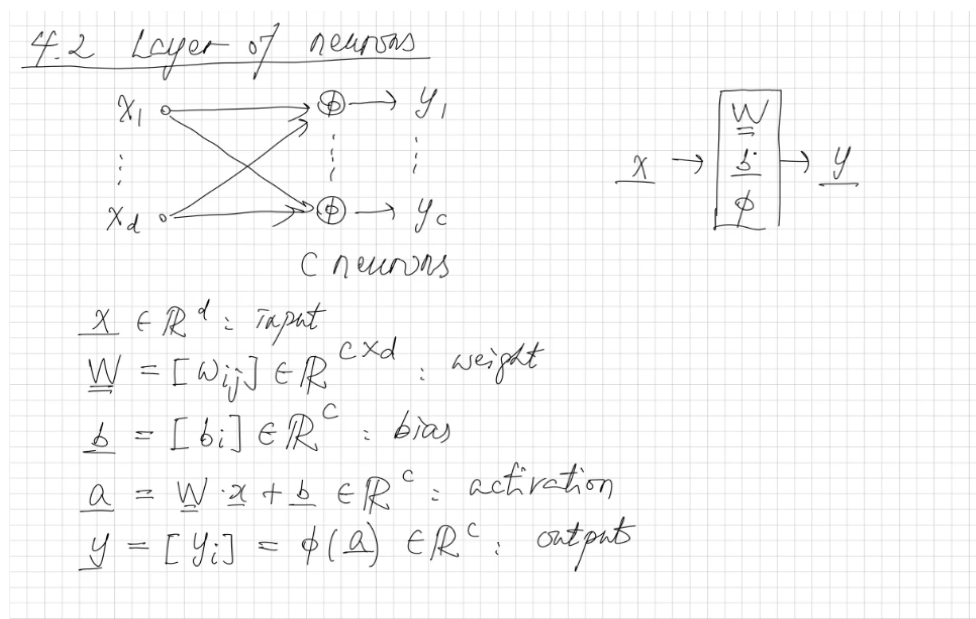


Abbildung 4.2: Layer of Neurons

2 meanings of  $\phi()$ :



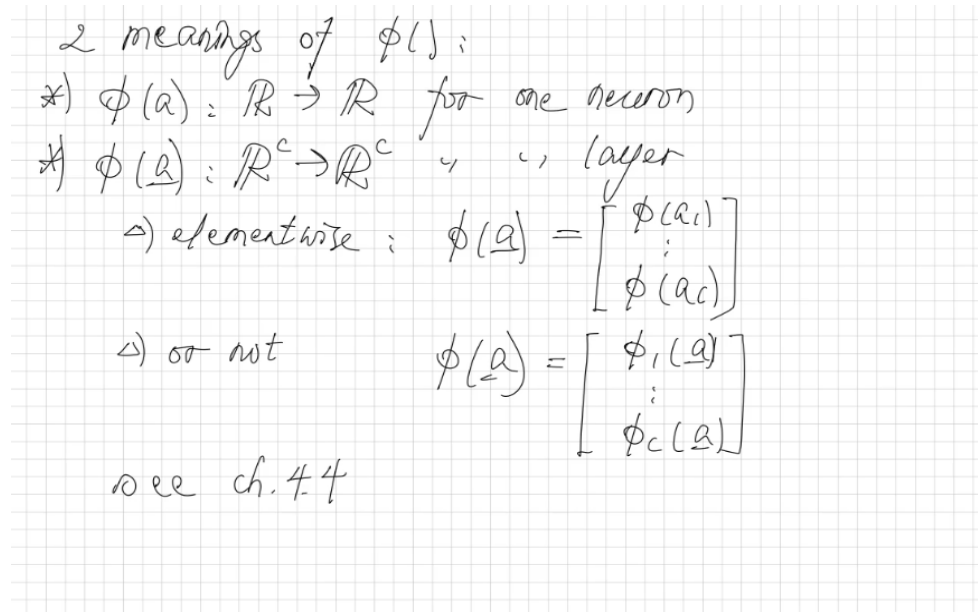


Abbildung 4.3: Meanings of phi

Comments :

- no interconnections between neurons in the same layer
  - dense layer, fully connected layer:
    - each input  $x_j$  connected to each neuron  $i$
- $\rightarrow c \cdot d$  weights and  $w_{ij}$  and  $c$  biases  $b_i$ ,  $1 \leq i \leq c$ ,  $1 \leq j \leq d$ ,  
 $\rightarrow c \cdot (d + 1)$  parameters

## 4.1. 4.3 Feedforward neural network

A cascade of dense layers

**Layer**  $1 \leq l \leq L$  :

$M_l - 1$  number of the input neurons

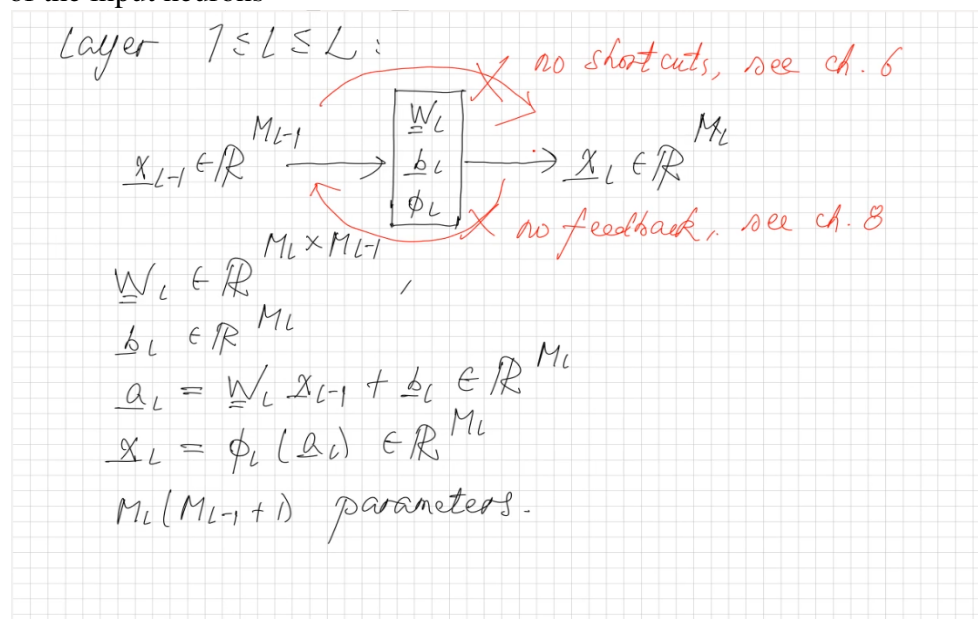


Abbildung 4.4: Feedforward multilayer neural network



Network:

$$x_L = f(x_0; \theta) : \mathbb{R}^{M_0} \rightarrow \mathbb{R}^{M_L}$$

\* parameter vector

learned from  $\mathcal{D}_{\text{train}}$

$$\theta = \begin{bmatrix} \text{vec}(\underline{W}_1) \\ \underline{b}_1 \\ \vdots \\ \text{vec}(\underline{W}_L) \\ \underline{b}_L \end{bmatrix} \in \mathbb{R}^{N_p}$$

$$\text{Number of parameters: } N_p = \sum_{l=1}^L M_l (M_{l-1} + 1)$$

$$\text{multiplications: } N_x = \sum_{l=1}^L M_l M_{l-1}$$

\*  $L, \{M_1, \dots, M_L\}$  } hyperparameters chosen by you

\*  $\{\phi_1, \dots, \phi_L\}$

Abbildung 4.5: Network Parameters