Neural Networks

doc. RNDr. lveta Mrázová, CSc.

DEPARTMENT OF THEORETICAL COMPUTER SCIENCE AND MATHEMATICAL LOGIC FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY IN PRAGUE

Neural Networks:

Introduction to the Area

doc. RNDr. Iveta Mrázová, CSc.

DEPARTMENT OF THEORETICAL COMPUTER SCIENCE AND MATHEMATICAL LOGIC FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY IN PRAGUE

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm

Computer Versus Brain

- The speed of information processing
- The kind of information processing
 - serial × parallel
- The kind of information storage
- Redundancy
- Control

Computer Versus Brain

The Z1 computer designed by Konrad Zuse in 1936/7

Mark I Perceptron machine built by F. Rosenblatt in 1957 (equipped with a camera to Process 20x20 pixel images)



<u>Frontal lobes</u> plan for the future, movement control, speech formation.

<u>**Temporal lobes**</u> process and interpret music and speech.

<u>Cerebral cortex</u> covers all the lobes, that form together the left and right cerebral hemisphere. It is just a few millimeters thick.

Amygdala creates emotions from the perceptions and thoughts

<u>**Hippocampus</u>** consolidates recently acquired information and transforms the short-term memory to the long-term one</u>

Parietal lobes collect and process the data provided by the senses.

Occipital lobes specialize in visual information processing.

<u>Cerebellum</u> controls muscle coordination and learning of automated movements.

Brain stem controls automatically performer body functions such as breathing. It connects the brain with the spinal cord.

<u>**Thalamus</u>** collects the information from the senses and forwards it further to the brain cortex.</u>

Computer Versus Brain The Structure of the Human Brain

Neural Networks – a Brief History

- 1943 formal neuron (W. McCulloch, W. Pitts)
- 1949 mathematical notion of learning (D. Hebb)
- 1958 perceptron (F. Rosenblatt)
- 1962 Adaline and sigmoidal transfer function (B. Widrow, M. Hoff)
- 1969 The perceptrons (M. Minsky, S. Papert)
- 1980s a further development

Neural Networks – a Brief History

since the eighties – **further developments**:

- The back-propagation training algorithm (P. Werbos, D. Rumelhart, G. Hinton, Y. Le Cun)
- Kohonen self-organizing feature maps (T. Kohonen)
- RBF-networks (Radial Basis Function, J. Moody, C. Darken)
- GNG-model (Growing Neural Gas, B. Fritzke)
- Convolutional neural networks (Y. Le Cun)
- SVM-machines (Support Vector Machines, V. Vapnik)
- ELM-networks (Extreme Learning Machines, G.-B. Huang)

Neural Networks – 21st Century

2003 - Allen Brain Atlas (Allen Institute for Brain Science, USA)

HBP – Human Brain Project, EU (January 2013)

Goal: mimic the human brain and identify faults its function

Expected costs – 1.2 billions Euro /10 year

https://www.humanbrainproject.eu/

http://www.nature.com/news/brain-simulation-and-graphene-projects-win-billion-eurocompetition-1.12291

2013 – BigBrain (Montreal Neurological Institute and German Forschungszentrum Jülich)

https://bigbrain.loris.ca/main.php

Neural Networks – 21st Century

GRAND CHALLENGE BRAIN Initiative, President Obama, 2. 4. 2013, USA:

~ <u>Brain Research Through Advancing Innovative Neurotechnologies</u>

https://obamawhitehouse.archives.gov/BRAIN

Goal: understand, how we think, how we learn and how works our memory

Expected costs – 3 billions USD / 10 years

Participants: DARPA ~ Defense Advanced Research Projects Agency

NIH ~ National Institutes of Health

NSF ~ National Science Foundation

private sector

https://www.nimh.nih.gov/news/science-news/science-news-about-the-brain-initiative

http://www.nature.com/news/flashing-fish-brains-filmed-in-action-1.12621



New Technologies

Neurosynaptic chip



11. 8. 2024 – PinnacleBusiness Systems:Scientists from IBMunveiled the firstneurosynaptic chip toachieve a scale of:

- one million programmable neurons,
- 256 million programmable synapses, and
- 46 billion synaptic operations per second per Watt.

□ Light-sheet microscopy



http://www.nature.com/nmeth/journal/v10/n5/fig_tab/nmeth.243 4_SV4.html

Neural Networks – a General Introduction

Recent problems:

- Training strategies parallelization and efficiency, few-shot training, generative models
- Architecture generalization and robustness, non-standard forms of data
- Scalability GPU, TPU, Google Cloud Platform, Edge TPU, ...
- Convergence, (few-shot) training, and over-training
- Prediction and generative models
- Adversarial patterns

Applications:

- Data mining "black-box", "white-box"
- Clustering and classification
- Information processing NLP, speech, vision, olfactory, tactile, motoric
- Artworks
- Solutions of optimization tasks
- and many others

Neural Networks – a General Introduction

Recent problems:

- Training strategies parallelization and efficiency, few-shot training, generative models
- Architecture generalization and robustness, non-standard forms of data
- Scalability GPU, TPU, Google Cloud Platform, Edge TPU, ...
- Convergence, (few-shot) training, and over-training
- Prediction and generative models
- Adversarial patterns

Applications:

- Data mining "black-box", "white-box"
- Clustering and classification
- Information processing NLP, speech, vision, olfactory, tactile, motoric
- Artworks
- Solutions of optimization tasks
- and many others

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm

Biological Background (1)

Model of a neuron

- ~ basic "computational unit" of a more complex system
 - $_{\sim}$ neural network (contains cca 8.6× 10^{10} neurons)
- \sim biological neurons consist of:

body (soma), dendrites, axon and synapses





Biological Background (1) - Biological Neuron



Biological Background (2) - Biological Neuron

Body (soma):

- $^{\circ}$ summarizes signals transmitted by surrounding neurons \rightarrow **potential**
- inner potential leads to the excitation of the neuron
- $^{\circ}$ the size varies from several μ m to several tens of μ m (~ 10⁻⁶ m)

Dendrites:

- represent signal input to neuron body
- $^{\circ}$ their length varies around 2-3 mm (~ 10⁻³ m)

Biological Background (3) - Biological Neuron

Axon:

- the only output of a neuron, but branched out widely at its end
- transmits the signal given by the level of excitation to the synapses
- $^{\circ}\,$ its length can reach over 1 m

Synapse:

- represent the "output device" of the neurons, can the signal amplify or diminish and transmit it to other neurons
- $^{\circ}$ for each neuron, there are up to 10^{6} connections to other neurons

Neuron output:

• Depends on neuron inputs and their processing inside neuron body



Biological Background (4) - Biological Networks

Biological neural networks:

neurons are mutually interconnected into networks

 by means of axons connected to dendrites of other neurons via synapses



- density of the neurons:
 - reaches cca $70 80 \cdot 10^3$ / mm³ in the human brain
 - cca 10 · 10³ neurons die every day without replacement
 - synapses are formed on the dendrites during the whole life
 - new synapses are formed, resp. non-functioning synapses can be revived

=> <u>LEARNING</u>

Question:

What percentage of a human brain is lost while alive?

Assume the age at death to be 100 years.

Biological Background (5) – Memory Types

Memory types

Short-term memory mechanism

- based on cyclical circulation of signals in neural networks
- after cca 300 circulations, fixation of the information starts in mid-term memory – this takes cca 30 s

Mid-term memory mechanism

- based on the changes of "neural weights"
- the change of synaptic weight coefficients is caused by multiple actions of the same signal on the respective synapse

Biological Background (6) – Memory Types

Memory types

Mid-term memory mechanism

- some information stored in mid-term memory moves to long-term memory while sleeping
- information stays in mid-term memory for several hours or days

Long-term memory mechanism

- consists in copying the information from mid-term memory to proteins inside the neurons – in particular in their nuclei
- the stored information can remain in the organism for its entire life

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm

Adaptation and Learning

Adaptation:

ability to accommodate to the changes of the environment

Adaptive process: the process of the adjustment

- every adaptation represents for the system some costs (material, energy, ...)
- living organisms are capable of reducing these costs during multiply repeated adaptations to environment changes

LEARNING:

- its objective is to minimize the costs spent for adaptation
- is the result of a multiply repeated adaptation

Adaptation and Learning - the Formalism (1)

Manifestation of the environment: x

~ input pattern generated by the environment

- Feature description of the objects (~ input patterns) :
 - selection of n basic characteristics features x_1, \dots, x_n
 - $\boldsymbol{x} = (x_1, \dots, x_n)$
- Information about the desired system reaction to the manifested environment: Ω

~ e.g., the true class assignment function for the input patterns

• The system reacts to any manifestation of the environment (~ *input pattern generated* by the environment) x and information Ω (~ true class assignment) by yielding one of the symbols ω_r ; r = 1, ..., R at its output (~ actual class labels).

Adaptation and Learning - the Formalism (2)

- Every assignment $[\mathbf{x}, \Omega] \rightarrow \omega_r$ is accompanied by some costs given by the function $Q(\mathbf{x}, \Omega, \omega_r)$ for each time unit
- The goal of the system:
 - find for any x and Ω such an assignment $[x, \Omega] \rightarrow \omega_r$, for which the **cost is minimal:**

$$Q(x, \Omega, \omega_r) = \min_{\omega} Q(x, \Omega, \omega)$$

Adaptive Systems (1)

Adaptive system \sim a system with two inputs and one output determined by:

- 1) a set *X* of manifestations of the environment *x* (~ input patterns generated by the environment)
- 2) a set O_1 of information about the desired system reaction Ω (e.g., true class assignments for the input patterns)
- 3) a set O_2 of output symbols ω (~ actual class labels).
- 4) a set *D* of (parametrized class label) decision rules $\omega = d(\mathbf{x}, q)$
- 5) the cost $Q(\mathbf{x}, \Omega, q)$

For any pair $[x, \Omega]$ we seek such a parameter q^* , for which it holds:

$$Q(x, \Omega, q^*) = \min_{q} Q(x, \Omega, q)$$

Adaptive Systems (2)

- Initial assignment $[\mathbf{x}, \Omega] \rightarrow \omega_s$
- If the system stays for time T in its initial assignment, this will be associated with total costs corresponding to $T \cdot Q(x, \Omega, \omega_s)$
- If the system is able to change its behavior based on an ongoing cost assessment, it finds after the time τ necessary for evaluating ω_r , for which the cost is minimal

Adaptive Systems (3)

Total costs after time T :

 $\tau Q(x,\Omega,\omega_s) + (T-\tau)Q(x,\Omega,\omega_r)$

- are bigger than the least possible total costs $T Q(x, \Omega, \omega_r)$,
- but smaller than the total costs of a system, that cannot change its decision, $TQ(x, \Omega, \omega_s)$ Adaptive system

$$T Q(x, \Omega, \omega_r) < \tau Q(x, \Omega, \omega_s) + (T - \tau)Q(x, \Omega, \omega_r) < T Q(x, \Omega, \omega_s)$$
Learned system
Non-adaptive system

Learning Systems (1)

The result of adaptation is stored in the memory:

- Save the time τ necessary to find minimum costs for repeated manifestations of the environment (~ input pattern sets generated by the environment)
- Further, it is not necessary to evaluate the costs

 \rightarrow after training, the information Ω about the desired system reaction is not necessary anymore

• The total costs of a learning system after training $T Q(x, \Omega, \omega_r)$ are smaller than the total costs of an adaptive system

Learning Systems (2)

<u>Learning system</u> \sim a system with 2 inputs and 1 output determined by:

- 1) a set *X* of manifestations of the environment *x* (~ input patterns generated by the environment)
- 2) a set O_1 of information about the desired system reaction Ω (e.g., true class assignments for the input patterns; is not necessary after training)
- 3) a set O_2 of output symbols ω (~ actual class labels)
- 4) a set *D* of (parametrized class label) decision rules $\omega = d(x,q)$
- 5) The desired behavior $\Omega = T(x)$ (in addition to adaptive systems)
- 6) Mean costs J(q) evaluated over $X \times O_1$ (*different from adaptive systems*)

Learning Systems (3)

Learning system

• After presenting the pair elements from the sequence

 $\{ [x_k, \Omega_k] \}; 1 \le k \le \infty$, where $\Omega_k = T_k(x_k)$,

it finds such a parameter q^* , for which it holds:

 $J(q^*) = \min_q J(q)$

- Sequential ~ sequential presentation of the pairs $[x_k, \Omega_k]$
- Inductive ~ find after the evaluation of countably many pairs $[x_k, \Omega_k]$ the parameter q^* , that minimizes the mean costs over the entire set X

Efficiency of Adaptation and Learning

The efficiency of an adaptive system is the higher, the shorter is the

time $\, au\,$ necessary for its adaptation and the longer are the time

intervals T when the environment does not change:

 $\circ \tau/T \rightarrow 0$:

The efficiency of the AS is comparable with the efficiency of a learning system after training

° $\tau/T \rightarrow 1$ ($\tau/T < 1$):

AS has about the same efficiency as a non-adaptive system

• $\tau/T \geq 1$: no adaptation takes place

The efficiency of the (trained) learning system is the highest possible

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm

Data Standardization

- In the Euclidean space, standardization of attributes is recommended so that all attributes can have an equal impact on the computation of distances.
- Example: Consider the following pair of data points

•
$$\mathbf{x}_i$$
: (0,1; 20) and \mathbf{x}_j : (0,9; 720),
 $dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(0,9 - 0.1)^2 + (720 - 20)^2} = 700,000457$

- The distance is dominated by (720 20) = 700.
- Standardize the attributes to have a common value range

Interval-Scaled Attributes

- Their values are real numbers following a linear scale.
 - e.g., the difference in Age between 10 and 20 is the same as that between 40 and 50.
 - The key idea is that the intervals keep the same importance throughout the scale
- Decimal scaling to the interval [-1,1]: divides the attribute values by the smallest power of 10 that keeps all the transformed values within the interval [-1,1].
- Range standardization to the interval [0,1] normalizes attribute values in the following way (*f* is an attribute):

$$range(x_{if}) = \frac{x_{if} - \min(f)}{\max(f) - \min(f)}$$
Interval-Scaled Attributes (cont ...)

• Range standardization to the interval [-1,1] normalizes attribute values in the following way (using $range(x_{if})$ defined above):

 $range^{[-1,1]}(x_{if}) = 2 \cdot range(x_{if}) - 1$

- Standardization according to the mean absolute deviation (MAD) transforms the attribute values so that they have zero mean and mean absolute deviation equal to 1.
 - The mean absolute deviation of attribute f (denoted by ${}^{MAD}S_f$) is computed as follows:

$$m_{\rm f} = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf}),$$

^{MAD} $s_f = \frac{1}{n} (|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|),$

MAD-score: $z(x_{if}) = \frac{x_{if} - m_f}{MAD_{S_f}}$

Interval-Scaled Attributes (cont ...)

- Standardization according to standard deviation transforms the attribute values so that they have zero mean and (corrected) sample standard deviation equal to 1.
 - The (corrected) sample standard deviation of attribute f (denoted by $cors_f$) is computed as it follows:

$$m_f = \frac{1}{n} \left(x_{1f} + x_{2f} + \dots + x_{nf} \right) ,$$

$$cors_{f} = \sqrt{\frac{1}{n-1} \left((x_{1f} - m_{f})^{2} + (x_{2f} - m_{f})^{2} + \dots + (x_{nf} - m_{f})^{2} \right)},$$

std-score:
$$std(x_{if}) = \frac{x_{if} - m_f}{cor_{s_f}}$$

Question:

Let X be a dataset consisting of the following five elements:

$$X = \{ 0, 1; -2, 0; 1, 5; -0, 2; -1, 3 \}.$$

State its standardized forms.

Selection and Order of Features

Probability of a wrong decision

X

Information contained in the input patterns

Too many features:

- technical feasibility
- speed of processing
- danger of over-training
- the number of variables \times the number of training patterns
- correlated features

Selection of Informative Features

- Selection of the minimum number of features from the considered set of features
 - the chosen set is not guaranteed to contain really informative features
 - the choice depends on the actual task solved
- The order of features from the considered set of features
 - according to the amount of information contained
 - can be used, e.g., in the case of sequential classifiers

Karhunen-Loeve Transform (1)

Properties of the Karhunen-Loeve transform:

- For the given number of expansion members, it yields the least mean squared error between the original and the transformed patterns
- 2. After the application of the covariance matrix the approximated patterns are decorrelated
 - \rightarrow decorrelation of features

Karhunen-Loeve Transform (2)

- 3. Expansion members do not contribute equally to the approximation
 - The influence of each respective expansion member on the approximation accuracy falls with its index
 - → The impact of members with high indexes will be small and we can thus omit them
- 4. The magnitude of the approximation error does not influence the structure of the expansion
 - Changed demands on the approximation error do not require the recomputation of the entire expansion
 - \rightarrow It is sufficient to add or remove a few of the last members

Of advantage especially for sequential classification methods

Karhunen-Loeve Transform (3)

- The choice of a suitable mapping of patterns from X^m to X^p such that the patterns from X^p will represent the best approximation of the original patterns from X^m in the sense of the mean squared error
 - *K* patterns *m* features *p* orthonormal vectors e_i $(1 \le i \le p)$ in X^m $(p \le m)$

S

→ Approximate the vectors \mathbf{x}_k from X^m ($1 \le k \le K$) by a linear combination of \mathbf{e}_i : $\mathbf{y}_k = \sum_{i=1}^p c_{ki} e_i$

uch that the squared error
$$arepsilon_k^2 = \|\mathbf{x}_k - \mathbf{y}_k\|^2$$
, will be minimal

the data matrix $\mathbf{X} =$

 $\mathbf{v} = \mathbf{V}^{\mathsf{T}} \mathbf{x}$

1st pattern

Kth pattern

Karhunen-Loeve Transform (4)



the *p*-th EV

Karhunen-Loeve Transform (5)

Computation of the matrix V:

• Center the data: $\mu_j = \frac{1}{K} \sum_{k=1}^{K} x_{kj}$



$$w_{ij} = w_{ji} = \frac{1}{K} \sum_{k=1}^{K} (x_{ki} - \mu_i) \cdot (x_{kj} - \mu_j)$$

 The vectors defining the most important features correspond to the eigenvectors of the covariance matrix

Observations

are the rows

 $X = \begin{pmatrix} \cdots & x_{1j} & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & x_{Kj} & \cdots \end{pmatrix}$

 $\mu = (\cdots \mu_j \cdots)$

Karhunen-Loeve Transform (6)

The eigenvalues correspond to the variance of the most important features

- the 1st column of the matrix V is the eigenvector corresponding to the biggest eigenvalue, ...
- further columns of $\,V\,$ will be added until the following eigenvalues are too small and can be omitted
- covariance matrices are positive semi-definite (~ their eigenvalues are non-negative;
 λ_i ≥ 0 ∀ i, because they are computed from the sums of squares, which themselves are each non-negative).

Problem:

- The choice of an adequate number of eigenvalues (p)
- An optimal choice of p cannot be guaranteed as the expansion does not reflect the true importance of each respective feature

Question:

Let Σ be a covariance matrix of a two-dimensional dataset: $\Sigma = \begin{pmatrix} 2,0 & 0,8 \\ 0,0 & 0,8 \end{pmatrix}$

$$= \begin{pmatrix} 0,8 & 0,6 \end{pmatrix}$$

Find its principal components.

Karhunen-Loeve Transform (7)

Modifications:

- **1.** Centered most important features: $y = V^T (x \mu)$, where $\mu = (\mu_1, ...)$ is the vector of mean values
- 2. Normalized most important features: $y = L^{-\frac{1}{2}} V^T (x \mu)$,

where **L** is the matrix $p \times p$, diagonal elements are the eigenvalues corresponding to the columns of **V**, the other elements are zero

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm

Probability – Basic Notions (1)

Probability (of an event A from the space S):

- $\mathbf{P}(A) \ge \mathbf{0}$ ($\mathbf{P}(\emptyset) = \mathbf{0}$)
- P(S) = 1
- For a finite number of mutually exclusive events $A_1, A_2, ..., A_n$ the probability $P(A_1 \cup A_2 \cup \cdots \cup A_n) = \sum_{i=1}^n P(A_i)$
- For an infinite number of mutually exclusive events $A_1, A_2, ...$ the probability $P(A_1 \cup A_2 \cup \cdots) = \sum_{i=1}^{\infty} P(A_i)$

Probability – Basic Notions (2)

• Conditional probability of the event B given that the event A has occurred (P(A) > 0):

 $P(B|A) = \frac{P(A \cap B)}{P(A)}^{\text{probability of A and B}}$ probability of B given A

- Mutual independence of the events A and B: $P(A \cap B) = P(A) \cdot P(B)$
- Formula for the probability of *A*:

 $P(A) = \sum_{i} P(A|B_{i}) P(B_{i})$

Probability – Basic Notions (3)

Bayesian formula for the conditional probability:

probability of B given A $P(B|A) = \frac{P(A|B) P(B)}{P(A)}; P(A), P(B) > 0$ probability of A

- Random variable:
 - 'the name of an experiment with a probabilistic outcome'
 - its value is the outcome of the experiment
- Probability distribution (for the random variable Y):
 - probability $\mathbf{P}(Y = y_i)$, that Y will take on the value y_i
- Expected value (~mean) of a random variable Y:

 $\mu_Y = \mathbf{E}(Y) = \sum_i y_i \ \mathbf{P}(Y = y_i)$

Probability – Basic Notions (4)

Variance (of a random variable):

$$VAR(Y) = E[(Y - \mu_Y)^2]$$

- Characterizes the width (dispersion) of the distribution around its mean
- Standard deviation of $Y: \sigma_y = \sqrt{VAR(Y)}$
- Binomial distribution
 - The probability of observing *r* 'heads' in a series of *n* independent coin tosses
 - The probability of 'heads' in a single toss is p

Probability – Basic Notions (5)

Binomial distribution

- The probability of observing *r*
 `heads' in a series of *n* independent coin tosses
- The probability of 'heads' in a single toss is *p*
- Probability function (probability that X will take on the value r):

$$P(r) = \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r}$$



Probability – Basic Notions (6)

- Expected (mean) value of X: E[X] = np
- Variance: VAR(X) = n p (1-p)
- Standard deviation: $\sigma_X = \sqrt{np(1-p)}$
- For sufficiently large values of *n* the binomial distribution is closely approximated by a normal distribution with the same mean and variance
- <u>Recommendation</u>: use the normal approximation only when:
 $np(1-p) \ge 5$ (i.e., for $n \ge \frac{5}{p(1-p)}$)

Probability – Basic Notions (7)

Normal distribution

- also called Gaussian distribution
- Normal probability density function

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Probability that the value of the random variable X will fall into the interval (a, b):

$$\int_a^b p(x) dx$$



Probability – Basic Notions (8)

Normal distribution

- Suitable for a large number of natural phenomena
- Expected (mean) value of X: $E[X] = \mu$
- Variance: $VAR(X) = \sigma^2$
- Standard deviation: $\sigma_X = \sigma$

<u>Central limit theorem:</u>

'The distribution of the mean of a large number of independent random variables of the same distribution approximates the normal distribution.'

Probability – Basic Notions (9)

- Estimator ~ random variable Y
 - Used to estimate the parameter *p* from the tested population
- Estimation bias of Y for p: E[Y] p
 - 'unbiased' estimator for p: $\mathbf{E}[Y] = p$
- N% confidence interval for the parameter p
 - Interval that contains p with probability N%
- Test ~ procedure deciding on the correctness of a statistical hypothesis H
 - Significance level α corresponds to the probability of rejecting the true hypothesis \rightarrow usually set to $\alpha = 0.05$

Hypotheses Testing (1)

- Given the observed accuracy of a hypothesis over a limited sample of data
 → how well does this estimate its accuracy over additional examples?
- 2. Given that one hypothesis outperforms another over some sample of data \rightarrow how probable is it that this hypothesis is more accurate in general?
- 3. When the data is limited → how to best use this data to both learn a hypothesis and estimate its accuracy as well as to compare the performance of two learning algorithms?
 - → limit the difference between the accuracy observed on the given data and the actual accuracy of the whole data distribution

Hypotheses Testing (2)

Aim: 1) Understand whether to use the hypothesis or not

2) Evaluating hypotheses represents an integral component of many learning methods (e.g., when post-pruning decision trees to avoid overfitting)

Estimate future accuracy of a hypothesis given only a limited set of data:

- Bias in the estimate: over-training × unbiased estimate of future accuracy (mutually independent training and test sets)
- Variance in the estimate: the measured accuracy can vary from the true accuracy; bigger variance for fewer test examples

Hypotheses Testing (3)

Estimating hypothesis accuracy

- Space of possible instances *X*, e.g., the set of all people
- Various target functions may be defined over X, f: X → {0, 1},
 e.g., people who plan to purchase new skis this year
- Different instances $x \in X$ may be encountered with different frequencies, e.g., probability that x arrives at the ski resort
 - **D** ... probability of encountering the instances in **X**

Hypotheses Testing (4)

<u>Task</u>: learn the target function *f* from the space *H* of possible hypotheses

provided are training examples *x*, along with their correct target value *f*(*x*), drawn randomly from *X* according to the distribution *D*

Questions:

Given a hypothesis h and a data sample containing n examples drawn at random according to the distribution D:

- 1. What is the best estimate of the accuracy of *h* over future instances drawn from the same distribution?
- 2. What is the probable error in this accuracy estimate?

<u>The sample error on the training set $S \subset X$ </u>

 \sim the fraction of S, misclassified by h

$$ERROR_{S}(h) \equiv \frac{1}{n} \sum_{x \in S} \delta(f(x), h(x))$$

- **n** ... the number of examples in **S**
- $\delta(f(x), h(x)) = 1$ for $f(x) \neq h(x)$
- $\delta(f(x), h(x)) = 0$ for f(x) = h(x)
- Binomial distribution $ERROR_{s}(h)$: $ERROR_{s}(h) = \frac{r}{n}$
 - \circ *r* ... the number of examples from *S*, that were misclassified by *h*

Hypotheses Testing (6)

The true error of hypothesis h

- ∼ probability of misclassification for an instance $x \in X$ drawn at random according to D: $ERROR_D(h) \equiv \Pr_{x \in D} [f(x) \neq h(x)]$
- Binomial distribution: $ERROR_D(h) = p$ ($= \frac{r}{n}$... estimate for p)
 - *p* ... probability of misclassifying a single instance drawn from *D*
 - unbiased estimator $ERROR_D(h)$ (~ p = r/n)
 - The hypothesis *h* and the sample set *S* must be chosen independently.
 - The sample set **S** contains $n \ (\geq 30)$ examples drawn at random from **X** according to the probability distribution **D**

Hypotheses Testing (7)

Estimator variance

- An unbiased estimator with the least variance would yield the smallest expected squared error between the estimate and the true value of the parameter
- Given no other information, the most probable value of ERROR_D(h) is ERROR_S(h)
- With approximately 95% probability, the true error ERROR_D(h) lies in the interval

$$ERROR_{S}(h) \pm 1.96 \sqrt{\frac{ERROR_{S}(h) \cdot (1 - ERROR_{S}(h))}{n}}$$

→ for approximately 95% of experiments, the calculated interval will contain the true error value

The expression for general (N%) confidence intervals – constant z_N :

$$ERROR_{S}(h) \pm z_{N} \sqrt{\frac{ERROR_{S}(h) \cdot (1 - ERROR_{S}(h))}{n}}$$

	The values of z_N for two-sided N% confidence intervals						
N%	50%	68%	80%	90%	95%	98%	99%
$\boldsymbol{z_N}$	0.67	1.00	1.28	1.64	1.96	2.33	2.58

- Wider intervals for a higher probability
- Good approximation for $n \ge 30$, resp. $n \cdot ERRORS(h) (1 ERRORS(h)) \ge 5$

Hypotheses Testing (9)

General approach to derive the confidence intervals:

- Identify the underlying population parameter *p* to be estimated, e.g., *ERROR_p(h)*.
- **2.** Define the estimator Y (e.g., $ERROR_{S}(h)$).

- choose a minimum-variance, unbiased estimator

- 3. Determine the probability distribution D_Y that governs the estimator Y including its mean and variance.
- 4. Determine the *N%* confidence interval

- find the thresholds L and U such that N% of the mass in the probability distribution D_Y falls between L a U.

Hypotheses Testing (10)

Difference in error of two hypotheses **D** (with a discrete-valued target function):

- Hypothesis h_1 has been tested on a sample S_1 containing n_1 randomly drawn examples
- Hypothesis h_2 has been tested on an independent sample S_2 containing n_2 examples drawn from the same distribution
- We want to estimate the difference d between the true errors of these two hypotheses:

 $d = ERRORD(h_1) - ERRORD(h_2)$

Hypotheses Testing (11)

 \rightarrow Estimator \hat{d} : $\hat{d} \equiv ERROR_{S_1}(h_1) - ERROR_{S_2}(h_2)$

- \hat{d} (~ difference between sample errors) yields an unbiased estimate of d
- Normal distribution with the mean $E[\hat{d}] = d$ and variance $\sigma_{\hat{d}}^2$

$$\sigma_{\hat{d}}^{2} \approx \frac{ERROR_{S_{1}}(h_{1}) \cdot \left(1 - ERROR_{S_{1}}(h_{1})\right)}{n_{1}} + \frac{ERROR_{S_{2}}(h_{2}) \cdot \left(1 - ERROR_{S_{2}}(h_{2})\right)}{n_{2}}$$

• *N%* confidence interval:

$$\widehat{d} \pm z_N \sqrt{\frac{ERROR_{S_1}(h_1) \cdot (1 - ERROR_{S_1}(h_1))}{n_1}} + \frac{ERROR_{S_2}(h_2) \cdot (1 - ERROR_{S_2}(h_2))}{n_2}$$

Hypotheses Testing (12)

Comparing the learning algorithms:

- test for comparing the learning algorithms L_A and L_B
- statistical significance of the observed difference between the algorithms
 - → determine, which of the learning methods, L_A and L_B , is better for learning the target function f
- Consider the relative performance of the two algorithms averaged over all the training sets of size *n* that might be drawn from the distribution *D*

Hypotheses Testing (13)

Comparing learning algorithms:

ightarrow estimate the expected value of the difference in the errors

 $\mathop{\mathbb{E}}_{S \subseteq D} \left[ERROR_D \left(L_A(S) \right) - ERROR_D \left(L_B(S) \right) \right]$

- L(S) ... hypothesis obtained by the learning algorithm L on the training set S
- S ⊂ D ... the expected value is taken over the samples S drawn according to the underlying instance distribution D
- \rightarrow in practice, just a limited number of training data D_0 is available to compare the considered learning algorithms
Hypotheses Testing (14)

- Divide the set D_0 into a training set S_0 and a disjoint test set T_0
 - Training data are used to train both L_A and L_B
 - Test data are used to compare the accuracy of the two learned hypotheses:

 $ERROR_{T_0}(L_A(S_0)) - ERROR_{T_0}(L_B(S_0))$

- $ERROR_{T_0}(h)$ approximates the true error $ERROR_D(h)$
- The difference in errors is measured only for the training set S_0 (rather than taking the expected value of this difference over all samples S that might be drawn from the distribution D)

k-Fold Cross Validation (1)

- 1. Partition the available data D_0 into k disjoint subsets $T_1, T_2, ..., T_k$ of equal size (≥ 30).
- **2.** FOR i := 1 TO k DO

use T_i for the test set, and the remaining data to build the training set S_i

- $S_i \leftarrow D_0 \setminus T_i$
- $h_A \leftarrow L_A(S_i)$
- $h_B \leftarrow L_B(S_i)$
- $\delta_i \in ERROR_{T_i}(h_A) ERROR_{T_i}(h_B)$
- 3. Return the value $\overline{\delta}$, where $\overline{\delta} = \frac{1}{k} \sum_{i=1}^{k} \delta_i$

k-Fold Cross Validation (2)

• N % - confidence interval: $\overline{\delta} \pm t_{N,k-1} \frac{\delta_{\overline{\delta}}}{\sqrt{k}}$

• $\sigma_{\overline{\delta}}$... estimate of the standard deviation: $\sigma_{\overline{\delta}} \equiv \sqrt{\frac{1}{k-1} \cdot \sum_{i=1}^{k} (\delta_{i} - \overline{\delta})^{2}}$

- $t_{N,k-1}$... constant (values of $t_{N,\nu}$ for two-sided confidence intervals approach the values of z_N with $\nu \to \infty$)
- **N** the desired confidence level
- $\boldsymbol{\nu}$ Nr. of degrees of freedom (nr. of independent random events that influence the value of $\bar{\delta}$; $\nu = k 1$ in the current setting)

k-Fold Cross Validation (3)

Values of $t_{N,\nu}$ for two-sided confidence intervals

	Confidence level N			
	90%	95%	98%	99%
$\nu = 2$	2.92	4.30	6.96	9.92
$\nu = 4$	2.13	2.78	3.75	4.60
$\nu = 5$	2.02	2.57	3.36	4.03
$\nu = 9$	1.83	2.26	2.82	3.25
$\nu = 10$	1.81	2.23	2.76	3.17
$\nu = 20$	1.72	2.09	2.53	2.84
$\nu = 30$	1.70	2.04	2.46	2.75
$\nu = 120$	1.66	1.98	2.36	2.62
$\boldsymbol{\nu}=\infty$	1.64	1.96	2.33	2.58

N ... desired confidence level

 $oldsymbol{
u}$ Nr. of degrees of freedom

k-Fold Cross Validation (4)

Testing has to be done on identical test sets!

• in contrast to comparing hypotheses that requires independent test sets

→ Paired tests

 typically produce tighter confidence intervals because any differences in observed errors are due to differences between the hypotheses and not due to differences in the makeup of the sampled data

A Two-Sided Test

(looks for the change in the estimated parameter and has two critical regions)



A One-Sided × a Two-Sided Test



A one-sided test specifies the direction of change in the estimated parameter (e.g., "precision is higher than") and has only one critical region.

Neural Networks:

Contents:

- Introduction to the Field
- Perceptron and Linear Separability

Contents:

- Introduction to the Field
 - Motivation and a Brief History
 - Biological Background
 - Adaptation and Learning
 - Feature Selection and Ordering
 - Probability and Hypotheses Testing (Review)
- Perceptron and Linear Separability
 - A Formal Neuron
 - Perceptron and Linear Separability
 - Perceptron Learning Algorithm
 - Convergence of Perceptron Learning
 - The Pocket Algorithm