





TOPIC PLAN				
Partner organization	Belgrade Metropolitan University			
Торіс	Artificial intelligence			
Lesson title	Classification in machine learning			
Learning objectives	Students can interpret basic concepts of classification. Students can project Bayesian classifier. Students are able to solve practical classification problems on artificial data and real-world datasets.	Methodology Modeling Collaborative learning Project based learning Problem based learning Strategies/Activitie s Graphic Organizer Think/Pair/Share Discussion questions		
Aim of the lecture / Description of the practical problem	The aim of this lecture is to learn basic concepts of classification, to gain knowledge for projecting Bayesian classifier and for solving practical classification problems on artificial data as well as on real-world datasets. In this lecture classification with Bayesian classifier is explained on two problems. In first problem there are artificial data where covariance matrices of classes are same and mathematical model for solving of this type of classification problem is presented. In second problem there are real data (famous Iris dataset) where covariance matrices of classes are different and mathematical model for solving of this type of classification problem is presented. Problem is presented. Problem is presented. Problems are practically solved in Python.			
Previous knowledge assumed:	Basics of calculus. Basics of linear algebra. Basics of statistics. Basic of Python programming.	Assessment for learning Observations Conversations Work sample Conference Check list Diagnostics		



Co-funded by the Erasmus+ Programme of the European Union



Introduction /	RANDOM VECTORS AND THEIR PROPERTIES	Assessment as
Theoretical		learning
basics	In theory of statistical pattern recognition, measurements from one pattern or object are treated as random vectors. We will use term random variable in context of pattern recognition because it is clear that every time when we repeat the measurement, we will get some other value, which leads us to the idea of characterizing the measurement of physical quantity as some random variable that gives a page different value in each realization	Self-assessment Peer-assessment Presentation Graphic Organizer Homework
	However, these values are not completely random, but subject to some laws. In order to formalize these laws, we will introduce two functions that will be joined to this random variable. First function is called <i>distribution function</i> , noted as $F_X(x)$ and second is called <i>probability density function</i> , noted as $f_X(x)$.	Iearning □Test □Quiz □Presentation □Project □Published work
	Distribution function is defined as:	
	$F_X(x)=P_r\ \{X\leq x\}$	
	and it represents probability that random variable X will take value that is less or equal to argument x . It has three main characteristics:	
	$F_X(\infty) = 1 \Leftrightarrow P_r \{X \le \infty\} = 1$	
	$F_X(-\infty)=0\Leftrightarrow P_r\ \{X\leq -\infty\}=0$	
	$x_1 \leq x_2 \Rightarrow F_X(x_1) \leq F_X(x_2)$	
	Probability density function is defined as:	
	$f_{\mathcal{X}}(x) = rac{\partial F(x)}{\partial x}$	
	Distribution function can be written as:	
	$F_{oldsymbol{\mathcal{X}}}(x) = \int_{-\infty}^x f_{oldsymbol{\mathcal{X}}}(au) d au$	
	Probability density function has two constraints:	



Co-funded by the Erasmus+ Programme of the European Union



 $F_{\mathcal{X}}(\infty) = 1 \Rightarrow \int_{-\infty}^{\infty} f_{\mathcal{X}}(x) = 1$ $(x_1 \le x_2 \Rightarrow F_X(x_1) \le F_X(x_2)) \Rightarrow (\forall x) f_X(x) \ge 0$ Random vector X is completely described with distribution function or with probability density function. However, in many practical situations, these functions can't be determined or they are too complex. In those cases, we must choose some other parameters that are less informative, but much more convenient in numerical sense. First and most important parameter is mathematical expectation or mean value of random vector X: $M_X = E\{X\} = \int_{\mathcal{T}_X} x f_X(x) dx$ where integration goes through whole space of random vector X. Conditional mathematical expectation of random vector X for class w_i is integral: $M_i = E\{X/w_i\} = \int xf(x/w_i)dx$ Second very important parameter that characterize random vector X is covariance matrix: $\Sigma = E\left\{ (X - M_X)(X - M_X)^T \right\}$ Component c_{ii} of this matrix is: $c_{ij} = E\{(X_i - m_i)(X_j - m_j)\}; (i, j = 1, ..., n)$ Thus, the diagonal elements of the covariance matrix form the variances of individual random variables in random vector, while non-diagonal elements represent covariances between random variables X_i and X_i . Although mathematical expectation and the covariance matrix are very important parameters which describe the distribution of

[&]quot;The European Commission's support for the production of this publication does not constitute an endorsement of the contents, which reflect the views only of the authors, and the Commission cannot be held responsible for any use which may be made of the information contained therein."





a random vector, they are mostly unknown in practice at need to be estimated based on measured samples. The procedure is called sample estimation technique. The technique is most commonly used to estimate the mean and variance of a random variable. Namely, if it necessary to estimate the mean value of the random variable Y whose realizations Y_i ; $i=1, N$ are known, the simplest way is to determine arithmetic mean: $\hat{m_Y} = \frac{1}{N} \sum_{i}^{N} Y_i$ Similarly, estimation of variance can be written as: $\hat{\sigma_Y}^2 = \frac{1}{N} \sum_{i}^{N} (Y_i - \hat{m_Y})^2$	nd nis nis an is m he
HYPOTHESIS TESTING	
The main goal of pattern recognition is to decide to white category observed sample belongs. Based on observation or measurement a measurement vector is formed. The vector serves as an input to the decision rule through white this vector joins one of the analyzed classes. <i>Hypothes</i> <i>testing</i> is a whole family of methods solving of this type a problem. These methods are very powerful, but the assume knowing of joined probability density functions from all classes and this information is often unknown practice.	ch on is ch <i>is</i> of ey om in
In pattern recognition theory, we deal with random vector gained from different classes and each of them characterized with its distribution function and probabil density function. These functions are called condition functions, and according to that, conditional probabil density function for <i>i</i> -th class is noted as:	is ity al ity
$f(x/w_i) or f_i(x); i=1,2,\ldots,L$	
where w_i denotes class <i>i</i> , and <i>L</i> is overall number classes. Unconditional probability density function random vector X, which is often called mixed dens function is given as:	of of ity



Co-funded by the Erasmus+ Programme of the European Union





[&]quot;The European Commission's support for the production of this publication does not constitute an endorsement of the contents, which reflect the views only of the authors, and the Commission cannot be held responsible for any use which may be made of the information contained therein."





Because mixed (a priori) probability density function is positive and joint for both a posteriori probability, decision rule can be written as follows: $P_1 f_1(X) > P_2 f_2(X) \Rightarrow X \in w_1$ $P_1 f_1(X) < P_2 f_2(X) \Rightarrow X \in w_2$ or we can write above equations as follows: $l(X) = \frac{f_1(X)}{f_2(X)} > \frac{P_2}{P_1} \Rightarrow X \in w_1$ $l(X) = \frac{f_1(X)}{f_2(X)} < \frac{P_2}{P_1} \Rightarrow X \in w_2$ Expression I(X) is called likelihood ratio, and that is very important quantity in pattern recognition theory. Ratio $P_2/$ P1 is called threshold value in decision making. It is common practice to apply negative logarithm on likelihood ratio, and then decision rule has form: $h(X) = -ln(l(X)) = -ln(f_1(X)) + ln(f_2(X)) < ln(\frac{P_1}{P_2}) \Rightarrow X \in w_1$ $h(X) = -ln(l(X)) = -ln(f_1(X)) + ln(f_2(X)) > ln(\frac{P_1}{P_2}) \Rightarrow X \in w_2$ A sign of inequality changed direction because of negative logarithm use. Expression h(X)is called discrimination function. Further on we will consider that $P_1 = P_2 = 0.5 \Rightarrow \ln\left(\frac{P_1}{P_2}\right) = 0$. Stated rules above are called Bayesian rule or minimal error decision test. For analysis of stated rule, it is very important to determine probability of decision error. This classification rule can't ensure perfect classification (as well as other rules). When we say probability error, we consider probability of event that rule will bring wrong decision about measurement vector belonging to a class. Conditional probability of error for measurement vector, noted as r(X), is equal to smaller of probabilities $q_1(X)$ and $q_2(X)$, i.e. $r(X) = min[q_1(X), q_2(X)]$

"The European Commission's support for the production of this publication does not constitute an endorsement of the contents, which reflect the views only of the authors, and the Commission cannot be held responsible for any use which may be made of the information contained therein."





Total error which is called Bayesian error, noted as ϵ can be calculated as follows:

$$egin{aligned} \epsilon &= E\left\{r(X)
ight\} = \int r(X)f(X)dX = \int min[q_1(X),q_2(X)]f(X)dX \ &= \int min[P_1f_1(X),P_2f_2(X)]dX = P_1\int_{L_2}f_1(X)dX + P_2\int_{L_1}f_2(X)dX \ &= P_1\epsilon_1 + P_2\epsilon_2 \end{aligned}$$

where

$$\epsilon_1=\int_{L_2}f_1(X)dX; \epsilon_2=\int_{L_1}f_2(X)dX$$

Probability ε_1 is called **Type I probability error** and it represents probability that sample which comes from first class is wrongly classified. Similarly, probability ε_2 is called Type II probability error and it represents probability that sample which comes from second class is wrongly classified. In a total error relation first equality represents definition of this error, while other equality is applied Bayesian theorem. Integration area L_1 is the area from which the decision rule joins the measurement vector X to the class w_1 and analogously, integration area L_2 corresponds to those vectors X which the decision rule classifies into a class w_2 . Consequently, these areas are often called w_1 -area and w_2 -area, respectively. For measurement vectors from L_1 area holds the relation $P_1f_1(X) > P_2f_2(X)$ and according to that conditional probability error is $r(X) = P_2 f_2(X)/f(X)$. Analogously, for vectors from L_2 area holds the relation $r(X) = P_1 f_1(X) / P_2 f_2(X)$ f(X). Based on that, we can say that Bayesian probability error consists of two terms. One of them refers to wrongly classified vectors from class w_1 , while other refers to wrongly classified vectors from the class w_2 .













	Where $f_h(h/w_i)$ is a posteriori probability density function of discrimination function h for samples that come from class w_i .		
	For our problem for artificially generated data random vector X is normally distributed. If random vector X is normally distributed, its probability density function can be written as:		
	$N_{m{X}}(M,\Sigma) = rac{1}{(2\pi)^{n/2} \Sigma ^{1/2}} e^{-rac{1}{2} d^2(m{X})} onumber \ 1 = -rac{1}{2} (X-M)^T \Sigma^{-1} (X-M)$		
	$=rac{1}{(2pi)^{n/2} \Sigma ^{1/2}}e^{-\frac{1}{2}(n-m)}\Sigma^{-(n-m)}$		
	Where $N_X(M, \Sigma)$ is shorted notation for normal distribution with mathematical expectation M and covariance matrix Σ . Function $d^2(X)$ is called statistical distance (or d^2 curve) of vector X from mathematical expectation vector M.		
Action	After introducing theoretical concepts, discussion with students will be performed for solving two posted problems. First problem is binary classification of artificially generated data. Data are normally distributed and data from first and second class have same covariance matrix. Second problem is classification of Iris dataset where data from classes have different covariance matrix.		
	SOLVING OF FIRST CLASSIFICATION PROBLEM		
	For our problem of artificially generated data, we have a posteriori probability density functions $f_i(X)$, where <i>i</i> =1,2. These functions are normal, with mathematical expectation vector M_i and covariance matrix Σ_i . Bayesian minimal error decision rule can be written in form:		
	$h(x) = -ln(l(X)) = -ln(f_1(X)) + ln(f_2(X)) < ln(rac{P_1}{P_2}) \Rightarrow X \in w_1$		
	Now we can replace $f_1(X)$ and $f_2(X)$ in upper equation and rewrite expression for $h(x)$:		







$$\begin{split} h(x) &= -\ln\left(\frac{1}{(2pi)^{n/2}|\Sigma_1|^{1/2}}e^{-\frac{1}{2}(X-M_1)^T\Sigma_1^{-1}(X-M_1)}\right) + \\ &+ \ln\left(\frac{e^{-\frac{1}{2}(X-M_2)^T\Sigma_1^{-1}(X-M_2)}}{(2\pi)^{n/2}|\Sigma_2|^{1/2}}\right) \end{split}$$
We can write expression for discrimination function in final form:

$$h(x) &= \frac{1}{2}(X-M_1)^T\Sigma_1^{-1}(X-M_1) - \\ &- \frac{1}{2}(X-M_2)^T\Sigma_2^{-1}(X-M_2) + \frac{1}{2}\ln\left(\frac{|\Sigma_1|}{|\Sigma_2|}\right) \end{aligned}$$
Last equation shows that decision boundary is quadratic function of X. If two classes have same covariance matrix, i.e. if $\Sigma_1 = \Sigma_2 = \Sigma$, decision boundary becomes linear function of X and can be written in next form:

$$h(x) = (M_2 - M_1)^T\Sigma^{-1}X + \frac{1}{2}(M_1^T\Sigma^{-1}M_1 - M_2^T\Sigma^{-1}M_2)$$
In our classification problem of artificially generated data classes have same covariance matrix, and a priori probabilities of classes appearance are same. So, we can write expression for classification in the first class:

$$h(x) = (M_2 - M_1)^T\Sigma^{-1}X + \frac{1}{2}(M_1^T\Sigma^{-1}M_1 - M_2^T\Sigma^{-1}M_2) < 0 \Rightarrow X \in w_1$$
Similarly, we can write expression for classification in second class:

$$h(x) = (M_2 - M_1)^T\Sigma^{-1}X + \frac{1}{2}(M_1^T\Sigma^{-1}M_1 - M_2^T\Sigma^{-1}M_2) < 0 \Rightarrow X \in w_2$$
As we can see our classifier is line (we have linear function of X) it we and to draw our line, we can do that as follows. Decision of houndary is the game of the classes is abseed on value of decision boundary is the game of the classes is appearance are same. So, we can a classe is the draw our line, we can see our classification in the first class:

$$h(x) = (M_2 - M_1)^T\Sigma^{-1}X + \frac{1}{2}(M_1^T\Sigma^{-1}M_1 - M_2^T\Sigma^{-1}M_2) > 0 \Rightarrow X \in w_2$$
As we can see our classifier is line (we have linear function of X) it we want to draw our line, we can do that as follows. Decision of belonging to one of two classes is based on value of decision boundary (is th bigger or smaller than zero), so if we want to draw our line, we can seel or smaller than zero), so if we want to draw boundary is the draw function of h(x) with zero:

$$h(x) = (M_2 - M_1)^T\Sigma^{-1}X + \frac{1}{2}(M_1^T\Sigma^{-1}M_1 - M_2^T\Sigma^{-1}M_2) = 0$$

[&]quot;The European Commission's support for the production of this publication does not constitute an endorsement of the contents, which reflect the views only of the authors, and the Commission cannot be held responsible for any use which may be made of the information contained therein."

















Consolidatio n	 The teacher's discussion with the students through different questions; Individually solving of simple tasks by the students under the supervision of the teacher; Given of examples by the teacher for introducing a new concepts and creative discussion with the students; Given homework by the teacher with a time limit until the next class. 			
Reflections and next steps				
Activities that	worked	Parts to be revisited		
Introducing to cl Solving problem	assificiation basics. Is in Python programming language.	Classification theoretical background.		
References				
 [1] Emilija Kisić, CS374 – Artificial Intelligence, Authorized Lectures on Metropolitan University Belgrade eLearning platform – LAMS, 2022. [2] Artificial Intelligence: A Modern Approach. 4th Edition, S. Russell and P. Norvig. Prentice Hall, 2021. [3] Artificial Intelligence: A New Synthesis, Nils Nilsson, Morgan Kaufmann, 1998 				