

CUB Spring 2025. Machine Learning.
Mid-Term Exam 28.03.2025

(write with printed letters):

Number:

Duration time: 120 minutes.

Written A4 page cheat sheet (two-sided) is allowed.

Number of points is 27. The exam grade is computed as a sum of points. For getting 100% it is enough to get in total 22 points out of 27. It is needed to get in total 10 points.

With given set of answers the number of correct answers may vary (depending on the set of answers).

Give your own examples (not from the lectures) of nominal features.

Examples can be colors of picture; red, blue, green.
Examples can be ordered list of buying property: rooms, rooms, rooms.

Suppose that you have an overfitted ML model. Choose methods to reduce overfitting.

Trees increase the maximal depth of a tree; ← incorrect

Trees reduce the maximal depth of a tree;

Increase σ parameter in RBF kernel $K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2))$

Task 4 (2 pts) Let's consider the following function:

①

$$f(x) = \text{tr}[(A + xy^T)^{-1}].$$

Here $x, y \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$. Using the techniques of differentials, find gradient of the function $f(x)$ w.r.t. x .

Task 5 (2 pts) Let's consider the following constrained optimization problem:

①

$$x^2 - y^2 \rightarrow \min_{x,y}$$

$$x^2 + y^2 \leq 1.$$

$$\begin{aligned} &\Rightarrow x^2 + y^2 - 1 = 0 \\ &\Rightarrow x^2 + y^2 - 1 < 0 \end{aligned}$$

$$F(x) = x^2 - y^2 - \delta(x^2 + y^2 - 1)$$

Find all stationary points (the points satisfying the KKT theorem).

$$\frac{dF(x)}{dx} = 2x - \delta 2x = 0$$

$$\frac{dF(y)}{dy} = -2y - \delta 2y = 0$$

Task 6 (2 pts) Suppose we have independent samples from discrete distribution, where random variable can take values 1, 2, 3 with the following probabilities:

①

$$p(x|\alpha): \begin{matrix} 1 & 2 & 3 \\ \alpha & \alpha^2 & 1 - \alpha - \alpha^2 \end{matrix}$$

Find maximal likelihood estimate α_{ML} if we have in the samples in total 20 ones, 10 twos and 10 threes.

$$\frac{20 \cdot 1 \cdot \alpha + 10 \cdot \alpha^2 + 10 \cdot 3(1 - \alpha - \alpha^2)}{20 + 10 + 10} = \frac{20\alpha + 10\alpha^2 + 30 - 30\alpha - 30\alpha^2}{70} = \frac{3 - \alpha - \alpha^2}{7}$$

I know it's wrong :)

Task 17 (1 pts) Choose correct statements about Target Encoding procedure;

- ☐ it is some procedure for working with missing/unknown values;
- ☐ in this procedure we need to train some regression model;
- ☐ in contrast to one-hot-encoding this procedure doesn't create additional features;
- ☐ this procedure may lead to leakage of target variable into training set features;
- ☐ this procedure is usually applied for both numerical and discrete features.

Task 18 (2 pts) Let's consider the following optimization problem:

$$\begin{aligned} \frac{1}{2} \mathbf{x}^T A \mathbf{x} &\rightarrow \min_{\mathbf{x}}, \\ x_i &\leq b \quad \forall i \end{aligned}$$

Here $\mathbf{x} \in \mathbb{R}^n$, b is some scalar and $A \in \mathbb{R}^{n \times n}$ is some symmetric positively defined matrix. Construct the dual optimization problem.

Task 19 (2 pts) Negative Binomial distribution is a discrete probability distribution where a random variable takes values $0, 1, 2, 3, \dots$ with the following probabilities:

$$p(x = k | q, r) = \binom{k+r-1}{k} (1-q)^k q^r.$$

Here $q \in (0, 1)$ and $r > 0$ are parameters of the distribution and $\binom{M}{k}$ – binomial coefficient. Suppose we have independent samples from this distribution:

$$x_1, x_2, \dots, x_N \sim p(x|q, r).$$

Find maximal likelihood estimate q_{ML} for fixed r .

Task 7 (1 pts) Suppose we have two ML models for bank loans approval (scoring models). The first model approves 100 loans, among which 80 were returned and 20 were not returned. The second model approves 50 loans, among which 48 were returned and 2 were not returned. Also, you know that in the dataset there are in total 200 loans that were returned. Which of these two models is better in terms of F_1 measure?

I - precision - $\frac{80}{100} = 0.8$
recall - $\frac{80}{200} = 0.4$

$F_1 = \frac{2 \cdot 0.8 \cdot 0.4}{0.8 + 0.4} = \frac{1.6}{1.2} = \frac{4}{3}$

← first one is better.

$\frac{0.8 \cdot 0.4}{0.8 + 0.4} = \frac{0.32}{1.2} = \frac{8}{30} = \frac{4}{15}$

II - precision - $\frac{48}{50} = 0.96$
recall - $\frac{48}{200} = 0.24$

$F_1 = \frac{2 \cdot 0.96 \cdot 0.24}{0.96 + 0.24} = \frac{0.96 \cdot 0.48}{1.2} = \frac{0.4608}{1.2} = \frac{1.152}{3}$

Task 8 (1 pts) Suppose that for some classifier, the precision is 0.9 and the recall is 0.8. The training set is balanced and has in total 5000 objects. Find false positive rate. ← FPR

precision - 0.9
recall - 0.8
total - 5000

$TP = 0.9TP + 0.9FP$

$0.1TP = 0.9FP$

$FP = \frac{TP}{9}$

$TP = 0.8TP + 0.8FN$

$0.2TP = 0.8FN$

$TP = 4FN$

$FPR = \frac{FP}{FP + TN} = \frac{FP}{5000 - \frac{4}{9}FP} = \frac{400}{5000} = 0.08$

$\frac{TP}{FP} = \frac{FN}{TN}$

$TN = FP \cdot \frac{FN}{TP} = FP \cdot \frac{0.8}{0.9} = \frac{8}{9}FP$

$TP + FP + FN + TN = 5000$

$9FP + FP + \frac{9}{4}FP + TN = 5000$

$TN = 5000 - \frac{13}{4}FP$

$\frac{9}{4}FP = 5000 - \frac{13}{4}FP \rightarrow FP = 400 \rightarrow TN = 1000$

Task 9 (1 pts) Stochastic gradient descent optimizer with momentum comparing to standard stochastic gradient descent:

0.62

☐ uses information about second derivatives;

☒ makes some kind of averaging of gradients from previous steps; ← correct

☒ computes gradient not for the current parameters but for some averaged version of parameters; ← in correct

☐ uses different scaling for every optimizing parameter.

Task 10 (1 pts) Choose correct statements about MAE loss function (mean absolute error) for regression problem:

0.62

☐ it is a differentiable function;

☐ weights of linear regression with this loss function can be found by analytical solution;

☒ it allows to find solutions less sensitive to the presence of outliers in dataset; ← correct

☒ it is usually faster to compute comparing to standard MSE loss function. ← in correct

Task 11 (1 pts) Let's consider support vector machine approach without using kernels. Suppose that D - number of features, N - number of training objects, M - number of support vectors. What is computational complexity for making inference on test object?

0

Task 12 (2 pts) Suppose that objects in some dataset after sorting them by score value of some two-class classifier have the following class labels (the first object has the highest score):

②

$[1, 1, -1, -1, 1, -1]$

$a_1, a_2, a_3, a_4, a_5, a_6$

Find F_1 measure for this classifier for the best threshold.

1) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=3 \ FP=3 \ FN=0$
 $Precision=0.5$
 $Recall=1$
 $F_1 = \frac{2 \cdot 0.5 \cdot 1}{0.5 + 1} = \frac{2}{3}$

3) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=2 \ FP=2 \ FN=1$
 $Precision=0.5$
 $Recall=\frac{2}{3}$
 $F_1 = \frac{2 \cdot 0.5 \cdot \frac{2}{3}}{0.5 + \frac{2}{3}} = \frac{4}{7}$

5) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=2 \ FP=0 \ FN=1$
 $Precision=1$
 $Recall=\frac{2}{3}$
 $F_1 = \frac{2 \cdot 1 \cdot \frac{2}{3}}{1 + \frac{2}{3}} = \frac{4}{5}$

2) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=3 \ FP=2 \ FN=0$
 $Precision=0.6$
 $Recall=1$
 $F_1 = \frac{2 \cdot 0.6 \cdot 1}{0.6 + 1} = \frac{3}{4}$

4) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=2 \ FP=1 \ FN=1$
 $Precision=\frac{2}{3}$
 $Recall=\frac{2}{3}$
 $F_1 = \frac{2 \cdot \frac{2}{3} \cdot \frac{2}{3}}{\frac{2}{3} + \frac{2}{3}} = \frac{2}{3}$

6) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=1 \ FP=0 \ FN=2$
 $Precision=1$
 $Recall=\frac{1}{3}$
 $F_1 = \frac{2 \cdot 1 \cdot \frac{1}{3}}{1 + \frac{1}{3}} = \frac{1}{2}$

Task 13 (1 pts) Let's consider training a two-class linear classification model with exponential loss function $\exp(-M)$, where $M = yw^T x$. Write down one step of stochastic gradient descent algorithm for training weights w with constant stepsize α and one object in a mini-batch. *I'm not gonna solve 13, so makepace for 12.*

3) $1 \ 1 \ -1 \ -1 \ 1 \ -1$
 $TP=0 \ FP=0 \ FN=3$
 $Precision=1$
 $Recall=0$
 $F_1=0$

⑦

Task 14 (1 pts) Let's consider a two-class classification problem and exponential loss function

⑦

$$L(y, z) = \exp(-yz).$$

Here $y \in \{-1, +1\}$. Does this function allow to predict correct class probabilities? Justify your answer.

⑦ Task 15 (1 pts) For convex loss functions stochastic gradient descent comparing to full batch gradient descent allows:

- ☐ find model parameters with lower value of loss function;
- ☐ find model parameters faster due to acceleration of one optimization iteration;
- ☒ find model parameters faster due to reduction of number of iterations needed for convergence.

the answer
the correct

Task 16 (1 pts) Choose the correct statements about precision-recall curve:

0.72

- ☐ it is monotonic curve w.r.t. increasing of recall value;
- ☐ The area under the curve is more suitable value for the case of unbalanced classes comparing to the area under ROC curve; *the answer*
- ☐ The starting and ending point of the curve does not depend on a training set; *this depends*
- ☒ The area under the curve takes values between 0 and 1. *correct*