






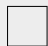








Profs. Nicolas Flammarion and Martin Jaggi
Machine Learning – CS-433 - MA
16.01.2025 from 15h15 to 18h15 in STCC
Duration : 180 minutes

Student 1

SCIPER : 999000

Wait for the start of the exam before turning to the next page. This document is printed double sided, 20 pages. Do not unstaple.

- This is a closed book exam. No electronic devices of any kind.
- Place on your desk: your student ID, writing utensils, one double-sided A4 page cheat sheet if you have one; place all other personal items below your desk or on the side.
- You each have a different exam.
- This exam has many questions. We do *not* expect you to solve all of them even for the best grade
- Only answers in this booklet count. No extra loose answer sheets. You can use the last two pages as scrap paper.
- For the **multiple choice** questions, we give +2 points if your answer is correct, and 0 points for incorrect or no answer.
- For the **true/false** questions, we give +1.5 points if your answer is correct, and 0 points for incorrect or no answer.
- Use a **black or dark blue ballpen** and clearly erase with **correction fluid** if necessary.
- If a question turns out to be wrong or ambiguous, we may decide to nullify it.

Respectez les consignes suivantes Observe this guidelines Beachten Sie bitte die unten stehenden Richtlinien		
choisir une réponse select an answer Antwort auswählen	ne PAS choisir une réponse NOT select an answer NICHT Antwort auswählen	Corriger une réponse Correct an answer Antwort korrigieren
  		 
ce qu'il ne faut PAS faire what should NOT be done was man NICHT tun sollte		
     		



First part: multiple choice questions

For each question, mark the box corresponding to the correct answer. Each question has **exactly one correct answer**.

Question 1 Recall that given a model $f : \mathcal{X} \rightarrow \{-1, 1\}$, an adversarial example can be found by solving the following *maximization* problem:

$$\max_{\hat{x}: \|\hat{x} - x\|_p \leq \varepsilon} 1_{f(\hat{x}) \neq y},$$

where (x, y) is the original example, $1_{f(\hat{x}) \neq y}$ is the zero-one loss and $\|\cdot\|_p$ is the ℓ_p norm. Often, the zero-one loss is replaced by a surrogate loss function.

Which of the following statements is **FALSE**?

- ☐ The zero-one loss is usually replaced by a surrogate loss function $\ell(yg(\hat{x}))$ where g is the output of the model before classification, i.e., $f(x) = \text{sign}(g(x))$.
- ☐ The maximization problem can be approximated by running projected gradient descent algorithm on a surrogate loss.
- ☐ The Fast gradient sign method (FGSM) approximates the maximization problem by linearizing a surrogate loss for ℓ_2 perturbations.
- ☐ There exist closed-form solutions to relaxations of the maximization problem with surrogate losses when the model and loss are simple.
- ☐ The zero-one loss is not differentiable; hence, the original maximization problem is generally hard to optimize.

Question 2 Recall that the output of self-attention is given by the following formula:

$$Z = \text{softmax} \left(\frac{XW_QW_K^\top X^\top}{\sqrt{d_k}} \right) XW_V,$$

where $X \in \mathbb{R}^{n \times d}$, $W_Q \in \mathbb{R}^{d \times d_k}$, $W_K \in \mathbb{R}^{d \times d_k}$, $W_V \in \mathbb{R}^{d \times d}$ are the data, query, key, and value matrices, respectively. Here, d is the dimensionality of the tokens, d_k is the hidden dimensionality, n is the sequence length. The softmax is applied row-wise. Additionally, the dimensions satisfy $d_k < d < n$. Consider the following linear self-attention modification:

$$Z = (XW_QW_K^\top X^\top) XW_V.$$

Which of the following statements is **FALSE**?

- ☐ The linear self-attention modification can be used to reduce the complexity of the self-attention mechanism to be *linear* in the sequence length.
- ☐ Replacing the product $W_QW_K^\top$ with a single matrix $W_{QK} \in \mathbb{R}^{d \times d}$ does not change the class of functions that can be expressed by both formulations of self-attention.
- ☐ The original self-attention mechanism is used in the Transformer model has a *quadratic* complexity in the sequence length.
- ☐ Both the original self-attention and the linear self-attention modification can be used with *causal masking*.
- ☐ The linear self-attention modification does not include normalization and thus does not benefit from a probabilistic interpretation.



Question 3 Which of the following statements is **always** true for a real-valued data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$ with $d > n$ where n is the number of data points and d represent the feature dimension?

- ☐ The sample points are linearly separable.
- ☐ The data points lie in an at most n -dimensional subspace of \mathbb{R}^d , so there is at least one non-zero direction in \mathbb{R}^d that is orthogonal to all these points.
- ☐ The weights \mathbf{w} can be computed with least-squares linear regression as follows: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.
- ☐ $\mathbf{X}^\top \mathbf{X}$ has exactly $d - n$ eigenvectors with eigenvalue zero.

Question 4 If $X \sim N(\mu, \sigma^2)$ and $Y = aX + b$, then which expression below gives the variance of Y ?

- ☐ $a^2\sigma^2 + b$
- ☐ $a\sigma^2 + b$
- ☐ $a\sigma^2$
- ☐ $a^2\sigma^2$

Question 5 Given a dataset $\mathcal{X} = \{1, 2, 4, 9, 16\}$, which choice of b minimizes the mean average error (MAE): $\mathcal{L}(b) = \frac{1}{5} \sum_{x \in \mathcal{X}} |x - b|$?

- ☐ π
- ☐ 4.0
- ☐ $\sqrt{358/5}$
- ☐ 6.4
- ☐ 2.0
- ☐ 9.0
- ☐ $\sqrt[3]{16}$
- ☐ 2π
- ☐ 8.5

Question 6 Imagine you are designing a convolutional neural network (CNN) for image classification with 10 classes in total. Each image is of size $(32, 32, 3)$, and the layers are of the following configurations. We denote a convolutional layer by (kernel height, kernel width, number of filters).

- (3, 3, 16) Convolutional layer with stride 1 and no padding
- Max-pooling layer with pooling size (2, 2) and stride 2
- (4, 4, 32) Convolutional layer with stride 1 and no padding
- Max-pooling layer with pooling size (2, 2) and stride 2
- Flatten and connect to a fully connected layer with output dimension 128
- Fully connected layer with output dimension 10

Which of the following statements is true?

- ☐ The max pooling layers have learnable parameters.
- ☐ After the first convolutional layer, the output size is $(30, 30, 3)$.
- ☐ After the second convolutional layer, the output size is $(12, 12, 32)$.
- ☐ The second convolutional layer has the most learnable parameters among all the layers.



Question 7 Which of the following statements is true?

- ☐ For a batch-normalized NN, the batch statistics are calculated from the current batch during inference time.
- ☐ Data augmentation can make CNNs less sensitive to image rotations.
- ☐ The computational complexity of training a fully connected neural network is quadratic in the number of layers for a single iteration.
- ☐ When normalization is applied to ensure stable training, the neural networks can be arbitrarily initialized and still achieve near-optimal performance.

Question 8 Which of the following statements is TRUE regarding the application of logistic regression and linear regression on binary classification tasks?

- ☐ Logistic regression models the probability of an instance belonging to a particular class, while linear regression predicts class labels directly.
- ☐ Logistic regression is able to generate non-linear decision boundaries with respect to original input features when the data are not linearly separable.
- ☐ Both logistic regression and linear regression have a closed-form solution for the optimal parameters.
- ☐ For linearly separable data, logistic regression via gradient descent converges in a finite number of steps.

Question 9 Suppose we are using a logistic regression model to predict the classes of y with a single feature x . Let θ_0 be the bias and θ_1 be the coefficient for x . The predicted probability that $y = 1$ given $x = x_0$ is:

$$\hat{p}(y = 1 \mid x = x_0) = \sigma(\theta_0 + \theta_1 x_0) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_0)}}.$$

Suppose we increase x_0 by one unit, resulting in:

$$\hat{p}(y = 1 \mid x = x_0 + 1) = \sigma(\theta_0 + \theta_1(x_0 + 1)).$$

Which of the following statements best describes how the predicted probability for $y = 1$ changes?

- ☐ The predicted probability is multiplied by a fixed amount of $\sigma(\theta_1)$.
- ☐ The predicted probability is multiplied by a fixed amount equal to θ_1 .
- ☐ The amount by which the predicted probability changes depends on its current value of x_0 .
- ☐ The predicted probability is multiplied by a fixed amount of e^{θ_1} .

Question 10 Which of the following statements is true about the K-means algorithm with the following objective function?

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|x_n - \mu_k\|^2 \text{ subject to } z_{nk} \in \{0, 1\} \text{ and } \sum_{k=1}^K z_{nk} = 1 \text{ for all } n.$$

- ☐ The objective value might increase during training ($J_{t+1} > J_t$) due to non-convexity of the objective function or poor initialization.
- ☐ K-means converges to a local minimum because the objective function is convex.
- ☐ Increasing the number of clusters (K) always decreases (or keeps constant) the optimal objective value.
- ☐ K-means is guaranteed to converge to the global minimum of the objective function.



Question 11 Suppose you are using the K-means++ algorithm to initialize the centroids for a clustering problem. The dataset consists of the following 1D points: $\{1, 2, 3, 10, 11, 12\}$. The first centroid is chosen at random as 2. Which point is more likely to be selected as the next centroid?

- ☐ All points are equally likely to be selected as the next centroid.
- ☐ 1 or 3 (equally likely)
- ☐ 11
- ☐ 10
- ☐ 12

Expectation Maximization

We consider a mixture of two exponential distributions, where a random variable X is constructed as follows: First, the class $Z \in \{0, 1\}$ is chosen according to:

$$P(Z = 0) = \pi_0, \quad P(Z = 1) = \pi_1 = 1 - \pi_0.$$

Then, the random variable X is sampled from an exponential distribution with parameter λ_0 or λ_1 , depending on the class Z . The conditional densities are:

$$p(x | Z = 0) = \begin{cases} \lambda_0 e^{-\lambda_0 x}, & x \geq 0 \\ 0, & \text{otherwise,} \end{cases} \quad p(x | Z = 1) = \begin{cases} \lambda_1 e^{-\lambda_1 x}, & x \geq 0 \\ 0, & \text{otherwise.} \end{cases}$$

We are given n i.i.d. observations $x_{1:n} = (x_1, \dots, x_n)$ from the generative model. Let $z_{1:n} = (z_1, \dots, z_n)$ denote the hidden variables, where $z_i = 0$ if the i -th observation belongs to class 0, and $z_i = 1$ otherwise. Denote the parameter vector as $\theta = (\lambda_0, \lambda_1, \pi_0, \pi_1)$.

Question 12 Which of the following represents the complete-data likelihood $p(x_{1:n}, z_{1:n} | \theta)$?

- ☐ $\prod_{i=1}^n [(\pi_0 \lambda_0 e^{-\lambda_0 x_i})^{1-z_i} (\pi_1 \lambda_1 e^{-\lambda_1 x_i})^{z_i}]$
- ☐ $\prod_{i=1}^n [(\pi_0 \lambda_0 e^{-\lambda_0 x_i}) + (\pi_1 \lambda_1 e^{-\lambda_1 x_i})]$
- ☐ $\prod_{i=1}^n [(\lambda_0 e^{-\lambda_0 x_i})^{1-z_i} (\lambda_1 e^{-\lambda_1 x_i})^{z_i}]$
- ☐ $\prod_{i=1}^n [(\lambda_0 e^{-\lambda_0 x_i}) + (\lambda_1 e^{-\lambda_1 x_i})]$

Question 13 We run the EM algorithm on this data. The vector $\theta^{(t)}$ denotes the parameters at the t -th iteration of the algorithm. At iteration $t + 1$, we perform the E-step by computing:

$$Q(\theta; \theta^{(t)}) := \mathbb{E}_{z_{1:n}} [\log p(x_{1:n}, z_{1:n} | \theta) | x_{1:n}, \theta^{(t)}].$$

Denote $q_1^{(t)}(x_i) = P(z_i = 1 | x_i, \theta^{(t)})$, and define $q_0^{(t)}(x_i) = 1 - q_1^{(t)}(x_i)$.

What is the value of $Q(\theta; \theta^{(t)})$?

- ☐ $\sum_{i=1}^n [(\lambda_0 e^{-\lambda_0 x_i}) q_0^{(t)}(x_i) + (\lambda_1 e^{-\lambda_1 x_i}) q_1^{(t)}(x_i)]$
- ☐ $\sum_{i=1}^n [q_0^{(t)}(x_i) (\log(\lambda_0) - \lambda_0 x_i) + q_1^{(t)}(x_i) (\log(\lambda_1) - \lambda_1 x_i)]$
- ☐ $\sum_{i=1}^n [q_0^{(t)}(x_i) (\log(\pi_0 \lambda_0) - \lambda_0 x_i) + q_1^{(t)}(x_i) (\log(\pi_1 \lambda_1) - \lambda_1 x_i)]$
- ☐ $\sum_{i=1}^n [(\pi_0 \lambda_0 e^{-\lambda_0 x_i}) q_0^{(t)}(x_i) + (\pi_1 \lambda_1 e^{-\lambda_1 x_i}) q_1^{(t)}(x_i)]$



Question 14 What is the value of $q_1^{(t)}(x_i) = P(z_i = 1 \mid x_i, \theta^{(t)})$?

- ☐ $\frac{\lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}}{\lambda_0^{(t)} e^{-\lambda_0^{(t)} x_i} + \lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}}$
- ☐ $\pi_1^{(t)} \lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}$
- ☐ $\frac{\pi_1^{(t)} \lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}}{\pi_0^{(t)} \lambda_0^{(t)} e^{-\lambda_0^{(t)} x_i} + \pi_1^{(t)} \lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}}$
- ☐ $\pi_0^{(t)} \lambda_0^{(t)} e^{-\lambda_0^{(t)} x_i} + \pi_1^{(t)} \lambda_1^{(t)} e^{-\lambda_1^{(t)} x_i}$

Question 15 In the M-step, what is the updated value of $\lambda_0^{(t+1)}$?

- ☐ $\frac{\sum_i x_i}{\sum_i q_0^{(t)}(x_i)}$
- ☐ $\frac{\sum_i q_0^{(t)}(x_i)}{\sum_i x_i q_0^{(t)}(x_i)}$
- ☐ $\frac{\sum_i q_0^{(t)}(x_i)}{\sum_i x_i}$
- ☐ $\frac{\sum_i x_i q_0^{(t)}(x_i)}{\sum_i q_0^{(t)}(x_i)}$

Question 16 In matrix factorization for recommender systems, the goal is to approximate a partially observed user-item rating matrix X by decomposing it as the product of two lower-dimensional matrices: W (item features) and Z (user features). The observed ratings x_{dn} represent the rating of the n^{th} user for the d^{th} item, with Ω as the set of indices for observed ratings. Which of the following is true regarding the regularized matrix factorization objective, given by:

$$\min_{W, Z} \frac{1}{2} \sum_{(d, n) \in \Omega} (x_{dn} - (WZ^T)_{dn})^2 + \frac{\lambda_W}{2} \|W\|_F^2 + \frac{\lambda_Z}{2} \|Z\|_F^2$$

- ☐ The optimization requires filling in all missing entries in X before training.
- ☐ The Frobenius norm regularization ensures the sparsity (the number of non-zero entries) of W and Z .
- ☐ The cost function is jointly convex in W and Z .
- ☐ Regularization terms $\frac{\lambda_W}{2} \|W\|_F^2$ and $\frac{\lambda_Z}{2} \|Z\|_F^2$ are used to prevent overfitting to the observed ratings.

Question 17 Let $N \in \mathbb{N}$. Which of the following sequence-to-sequence functions $f : \{0, 1\}^N \rightarrow \{0, 1\}^N$ **cannot** be arbitrarily well approximated by a decoder-only transformer without causal masking and without positional embeddings? In the choices below, $f_i(x_1, \dots, x_N)$ refers to the i^{th} coordinate of $f(x_1, \dots, x_N)$.

- ☐ $f(x) = \mathbf{0}$, $\forall x \in \{0, 1\}^N$, where $\mathbf{0} \in \{0, 1\}^N$ is the zero vector.
- ☐ $f : \{0, 1\}^N \rightarrow \{0, 1\}^N$ such that $f_i(x) = x_1 \oplus \dots \oplus x_{i-1} \oplus x_{i+1} \oplus \dots \oplus x_N$, where \oplus is the addition modulo 2 operation.
- ☐ $f(x) = v$, $\forall x \in \{0, 1\}^N$, where $v \in \{0, 1\}^N$ is the vector of alternating 0's and 1's: $v = (0, 1, 0, 1, \dots)$.
- ☐ $f : \{0, 1\}^N \rightarrow \{0, 1\}^N$ such that $f_i(x) = 1 - x_i$.

**Question 18**

Why would you use LASSO over Ridge regression?

- A) It can help us identify which features are important.
- B) It is faster to learn the weights for LASSO than for Ridge.
- C) LASSO usually achieves a lower generalization error than Ridge.
- D) If there are many features, the model learned using LASSO can make predictions more efficient.

- ☐ C and D
- ☐ A
- ☐ C
- ☐ A and B
- ☐ D
- ☐ B and C
- ☐ A and D.
- ☐ B
- ☐ A and C
- ☐ B and D

Question 19 Suppose we fit “LASSO Regression” to a data set, which has 100 features (X_1, X_2, \dots, X_{100}). Now, we rescale one of these features by multiplying with 10 (say that feature is X_1), and then refit LASSO regression with the same regularization parameter. Which of the following options will be **TRUE**?

- ☐ It is more likely for X_1 to be included in the model.
- ☐ Can't say.
- ☐ It is more likely for X_1 to be excluded from the model.
- ☐ None of these.

Question 20 We are given the following dataset as illustrated, and assume that "+" and "-" denotes the label of each data point:

+	+	-	-
	-		-
+	+	-	-

The Leave-One-Out cross validation errors using 1-Nearest-Neighbor and 3-Nearest-Neighbor classifiers are:

- ☐ 0, 5/10.
- ☐ 5/10, 1/10.
- ☐ 0, 1/10.
- ☐ 1/10, 5/10.



Question 21 Consider a contrastive learning setup with an encoder f_θ that outputs normalized embeddings of images. During training, each image x is augmented twice to create views x_i, x'_i . The contrastive loss is then defined as

$$\mathcal{L}(\theta) = - \sum_{i=1}^N \log \frac{\exp(f_\theta(x_i)^\top f_\theta(x'_i)/\tau)}{\sum_{k=1}^N \mathbb{1}_{k \neq i} \exp(f_\theta(x_i)^\top f_\theta(x_k)/\tau)},$$

where τ is the temperature parameter and N denotes the batch of size. Which of the following statements is FALSE?

- ☐ To achieve optimal performance, all possible image augmentations should be applied with equal probability during training.
- ☐ Two different random crops of the same image should produce similar embeddings, while crops from different images should produce dissimilar embeddings.
- ☐ The quality of learned representations depends heavily on the choice of image augmentations, which should preserve semantically meaningful features while varying nuisance factors.
- ☐ The temperature parameter τ controls the sensitivity of the model to hard vs. easy negative examples.

Question 22 Consider a diffusion model with forward process $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1 - \alpha_t)I)$ where α_t decreases with t . Which of the following statements is FALSE?

- ☐ The forward process gradually adds Gaussian noise to the data until reaching a standard normal distribution in the limit of $t \rightarrow \infty$.
- ☐ The backward process can be implemented by training a neural network to predict the noise component that was added to the clean data.
- ☐ Similar to GANs, diffusion models require training two separate models: one for the forward process and one for the backward process of diffusion.
- ☐ For training, solving the regression problem of predicting the noise at each timestep is equivalent to solving a denoising score matching problem.

Question 23 Consider four points $A = (-1, 0), B = (-1, 1), C = (1, 0), D = (1, 1)$ in \mathbb{R}^2 . Each point can be assigned a label of -1 or 1 . Which of the following statements are TRUE?

- P) For any assignment of labels to the points A, B, C, D , there always exists a linear classifier that perfectly separates them.
- Q) Consider the cases where A, C are labelled $-1, 1$, respectively, and B and D have arbitrary labels. In such cases, no linear classifier can achieve a margin of 1.25 over these points.

- ☐ P
- ☐ Both P and Q
- ☐ Q
- ☐ Neither P nor Q

Question 24 Recall that the soft-margin SVM corresponds to the following optimization problem:

$$\min_{w, \xi} \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \sum_{i=1}^N \xi_i, \quad \text{subject to: } y_i(w^T x_i) \geq 1 - \xi_i \text{ for all } 1 \leq i \leq N, \text{ and } \xi_i \geq 0$$

Which of the following statements is TRUE?

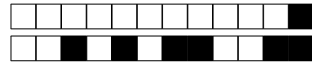
- ☐ For the optimal solution, it is always true that $\xi_i \leq 1$ for all i .
- ☐ The soft-margin SVM problem is equivalent to the minimization of a regularized logistic loss function.
- ☐ There cannot be a solution to the soft-margin SVM where $\xi_i \geq 1$ for all i .
- ☐ The given optimization problem admits a solution if and only if the data is linearly separable.



Question 25 Which of the following statements about the k -nearest neighbors algorithm is TRUE?
Here, N denotes the number of data points.

- ☐ The variance of the k -nearest neighbors method is 0 when $k = N$.
- ☐ The bias of the k -nearest neighbors method is 0 when $k = N$.
- ☐ The k -nearest neighbors method is sensitive to the choice of distance metric.
- ☐ The k -nearest neighbors approach is only applicable to classification tasks.

DRAFT



Second part: true/false questions

For each question, mark the box (without erasing) TRUE if the statement is **always true** and the box FALSE if it is **not always true** (i.e., it is sometimes false).

Question 26 (Large Language Models) Suppose we have a decoder-only transformer T with a context window of size N . We want to use T to compute an infinite binary function $f : \{0, 1\}^\infty \rightarrow \{0, 1\}^\infty$, where we will use f_i to refer to the i^{th} coordinate of f . Now, for every $x = (x_1, x_2, \dots) \in \{0, 1\}^\infty$, we have that $f_i(x) = 0$ for $i < N$ and $f_i(x) = g_i(x_{i-N+1}, \dots, x_i)$ for $i \geq N$, where $g_i : \{0, 1\}^N \rightarrow \{0, 1\}$. For every sequence $x \in \{0, 1\}^\infty$, we will apply T with a sliding context window. That is, $T(x) = y$ will mean that $y_i = T(x_{\max(1, i-N+1)}, \dots, x_i)$. Suppose it is known that each g_i requires at least linear time, $\Omega(i)$, to be computed. Then, T cannot compute f .

☐ TRUE ☐ FALSE

Question 27 (Adversarial robustness) Adversarial training with the Fast Gradient Sign Method roughly has twice the computational cost of training the model without adversarial training.

☐ TRUE ☐ FALSE

Question 28 (Masked Language Models) BERT-style masked language models cannot be used for autoregressive generation of text.

☐ TRUE ☐ FALSE

Question 29 (Text Representation Learning) Suppose we have trained a **FastText** model to learn sentence representations for the task of determining whether the subject in a given sentence is human or not. It is possible that our model could **incorrectly** classify both the sentences "Francesco devoured the unsuspecting octopus." and "The unsuspecting octopus devoured Francesco.".

☐ TRUE ☐ FALSE

Question 30 (Overfitting) When using a train/test split, if the split is not random, this will primarily cause underfitting.

☐ TRUE ☐ FALSE

Question 31 (Gaussian Mixture Models) We fit a GMM to a dataset utilizing the (soft) EM algorithm. Let L_t denote the log-likelihood of the data at iteration t . During this process, the log-likelihood L_t may decrease in some iterations, i.e., there may exist t such that $L_{t+1} < L_t$.

☐ TRUE ☐ FALSE

Question 32 (Nearest Neighbours) In a binary classification problem with a training set of size 2, where the points have different labels, the decision boundary of the 1-nearest neighbor (1-NN) classifier with ℓ_1 -distance coincides with that of the maximum margin linear classifier.

☐ TRUE ☐ FALSE

Question 33 (Generative Adversarial Networks) In GAN training, if the discriminator achieves perfect classification accuracy on both real and fake samples, this means we have reached the training optimum and training should stop.

☐ TRUE ☐ FALSE

Question 34 (Overfitting) Training your model until it achieves a low loss value on your test data is a good way to prevent overfitting.

☐ TRUE ☐ FALSE



Question 35 (Masked Language Models) In BERT-style masked language modeling, if we mask multiple tokens in a sentence during pre-training, the prediction of a masked token is independent of the predictions of other masked tokens, conditioned on the context.

☐ TRUE ☐ FALSE

Question 36 (Fairness) The independence criterion of fairness requires that the false positive rate (FPR) and false negative rate (FNR) must be equal across all groups.

☐ TRUE ☐ FALSE

Question 37 (Batch Normalization) Batch normalization is a technique that reduces the processing time of a single batch, and thus leads to much faster convergence.

☐ TRUE ☐ FALSE

Question 38 (Optimization) Given a continuous, strictly convex, loss function $\mathcal{L}(w)$, a gradient descent step $w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t)$ with learning rate $\eta > 0$, results in a decrease in the loss, i.e. $\mathcal{L}(w_{t+1}) \leq \mathcal{L}(w_t)$.

☐ TRUE ☐ FALSE

Question 39 (Optimization) The 0–1 loss is **not** a suitable loss function for training a neural network classifier with stochastic gradient descent.

☐ TRUE ☐ FALSE



Third part, open questions

Answer in the empty space below. Your answer should be carefully justified, and all the steps of your argument should be discussed in details. Leave the check-boxes empty, they are used for the grading.

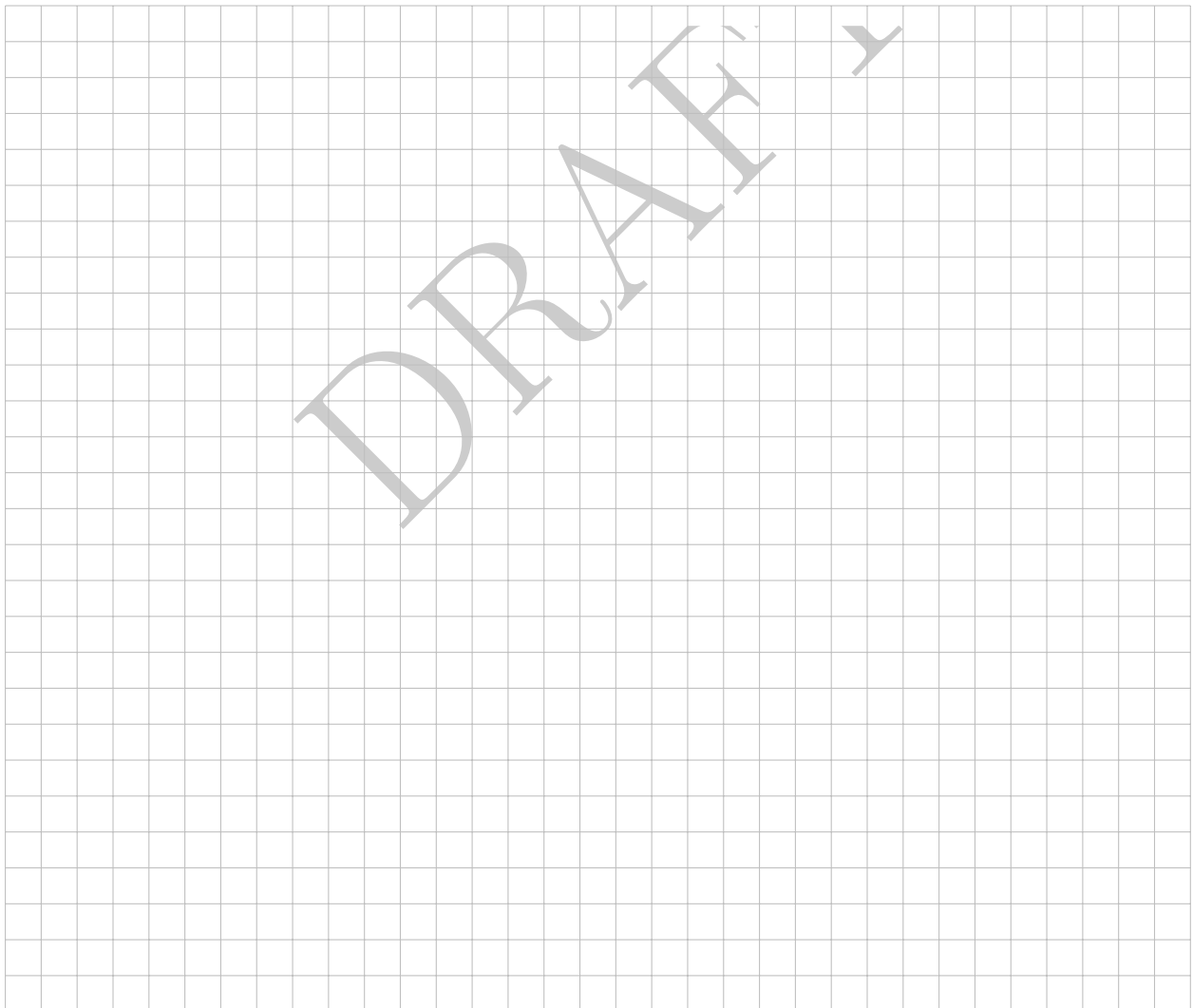
Answer in the space provided! Your answer must be justified with all steps. Leave the check-boxes empty, they are used for the grading.

1 Convexity

Question 40: (3 points) Let $f_1, f_2 : \mathbb{R} \rightarrow \mathbb{R}$ be convex functions. Consider the following statements and determine whether each statement is true or false. If the statement is true, provide a proof. If the statement is false, either provide a counterexample or disprove it.

- (a) The function $g(x) = \min(f_1(x), f_2(x))$ is convex.
- (b) The function $h(x) = \max(f_1(x), f_2(x))$ is convex.

☐₀ ☐₁ ☐₂ ☐₃





2 Logistic Regression Loss

Question 41: (3 points) Consider the logistic regression loss $L : \mathbb{R}^d \rightarrow \mathbb{R}$ for a binary classification task with data $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \{0, 1\}$ for $i \in \{1, \dots, N\}$:

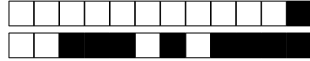
$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (\log(1 + e^{\mathbf{x}_i^\top \mathbf{w}}) - y_i \mathbf{x}_i^\top \mathbf{w}).$$

Questions:

- (a) Compute the gradient $\nabla_{\mathbf{w}} L(\mathbf{w})$ of the loss with respect to \mathbf{w} .
- (b) Prove that the loss function $L(\mathbf{w})$ is convex with respect to \mathbf{w} .

☐ 0 ☐ 1 ☐ 2 ☐ 3





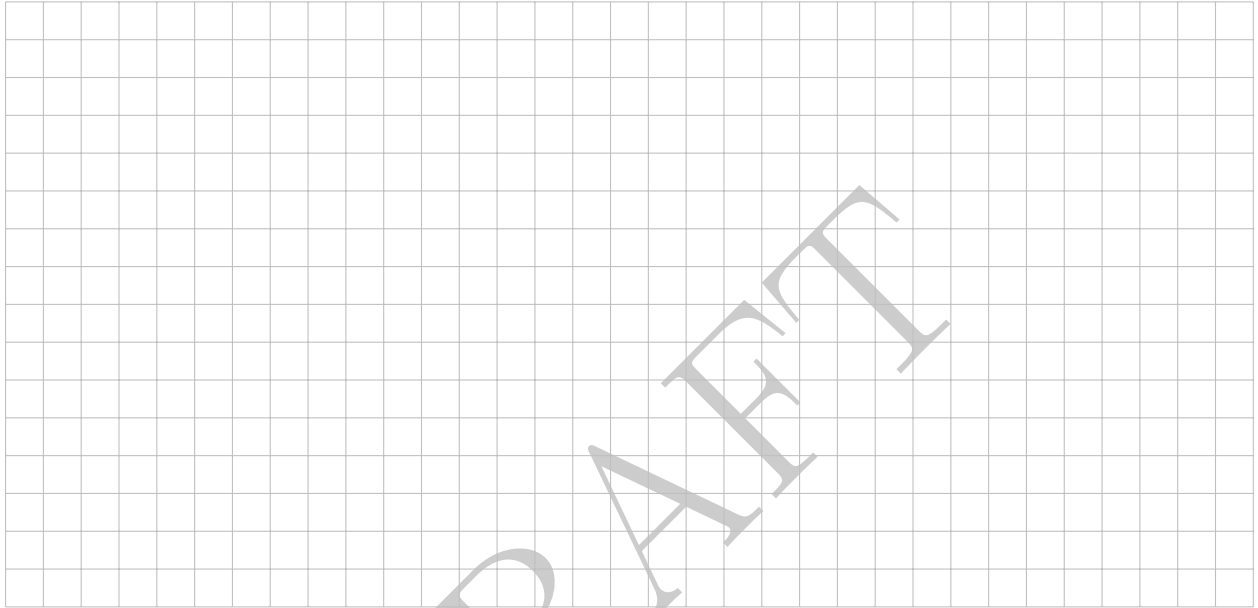
3 Kernels

Let $m \in \mathbb{N}^+$ be a positive integer, and let $c \in \mathbb{R}$ satisfy $c < 0$. For any positive integer $d \in \mathbb{N}^+$, consider the function $k : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ defined by:

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^m.$$

Question 42: (1 points) **Definition of a valid kernel:** State what it means for a general function $k : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ to be a *valid kernel*.

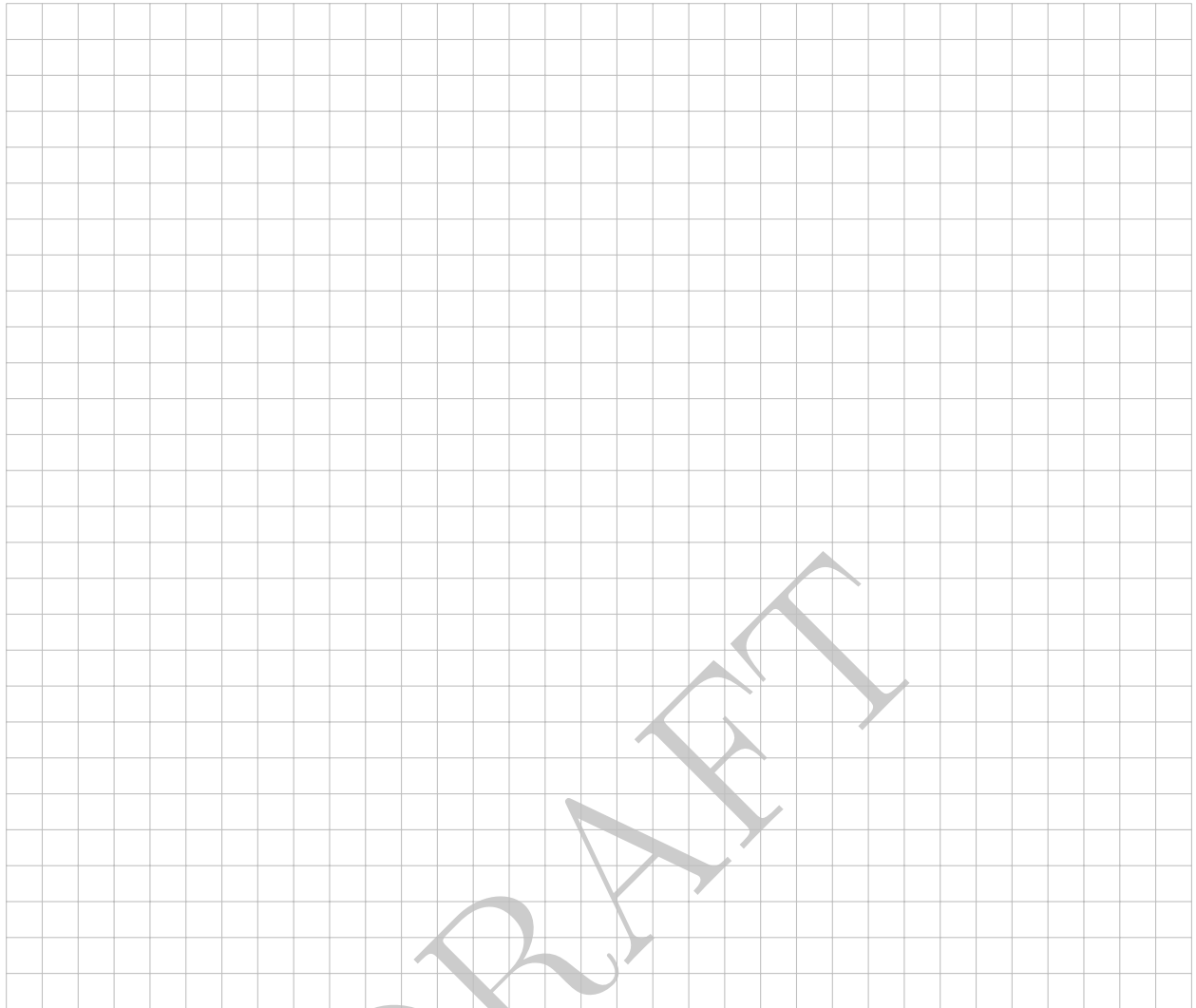
☐_0 ☐_1



Question 43: (4 points) **Proof for $m = 1$:** Show that for $c < 0$ and any $d \in \mathbb{N}^+$, the function $(xy + c)^d$ (with $x, y \in \mathbb{R}$) **fails** to be a valid kernel. (Note: If you give a single proof covering all m , that is sufficient for full points, see next question.)

☐_0 ☐_1 ☐_2 ☐_3 ☐_4

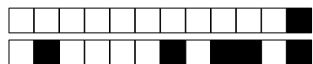




Question 44: (2 points) **Proof for any m :** Show that for any arbitrary dimension m , the function $(\mathbf{x}^\top \mathbf{x}' + c)^d$ (with $c < 0$) **fails** to be a valid kernel.

☐ ₀ ☐ ₁ ☐ ₂





DRAFT



4 Neural networks

We consider single-hidden-layer networks with the ReLU activation function. Recall that

$$\text{ReLU}(z) = \max\{0, z\}.$$

A *single-hidden-layer ReLU network* (with no skip connections) taking an input $\mathbf{x} \in \mathbb{R}^d$ and producing a real output can be written as

$$F(\mathbf{x}) = v^\top \text{ReLU}(W\mathbf{x} + b) + c,$$

where $W \in \mathbb{R}^{m \times d}$, $b \in \mathbb{R}^m$, $v \in \mathbb{R}^m$, and $c \in \mathbb{R}$. The ReLU function is applied elementwise to its argument. The “no skip connections” constraint means \mathbf{x} does not directly feed into the output neuron; all entries of \mathbf{x} must pass through the hidden layer. We wish to study how certain elementary functions can be represented by such networks.

Question 45: (*4 points*) **Implementing the identity function:** consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$ defined by

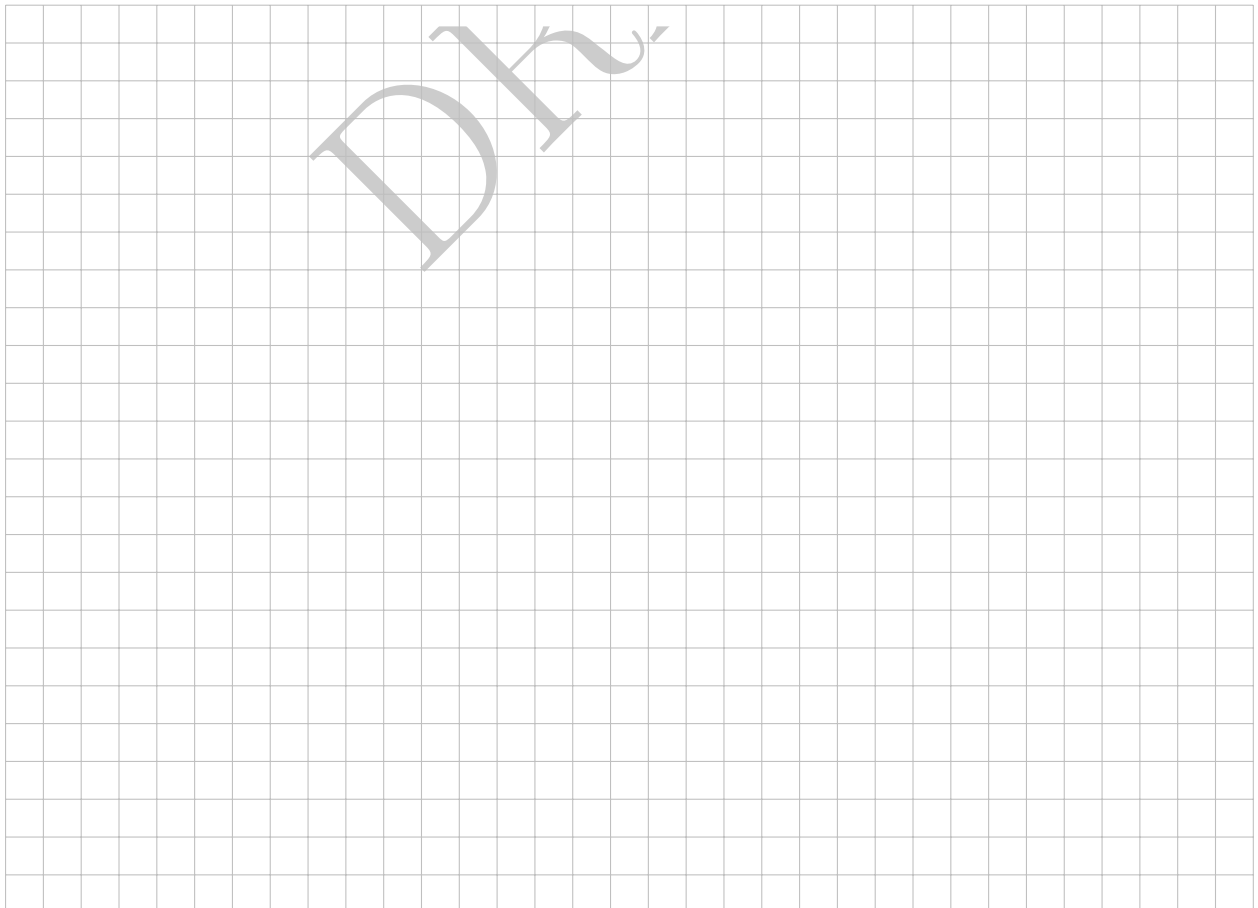
$$f(x) = x.$$

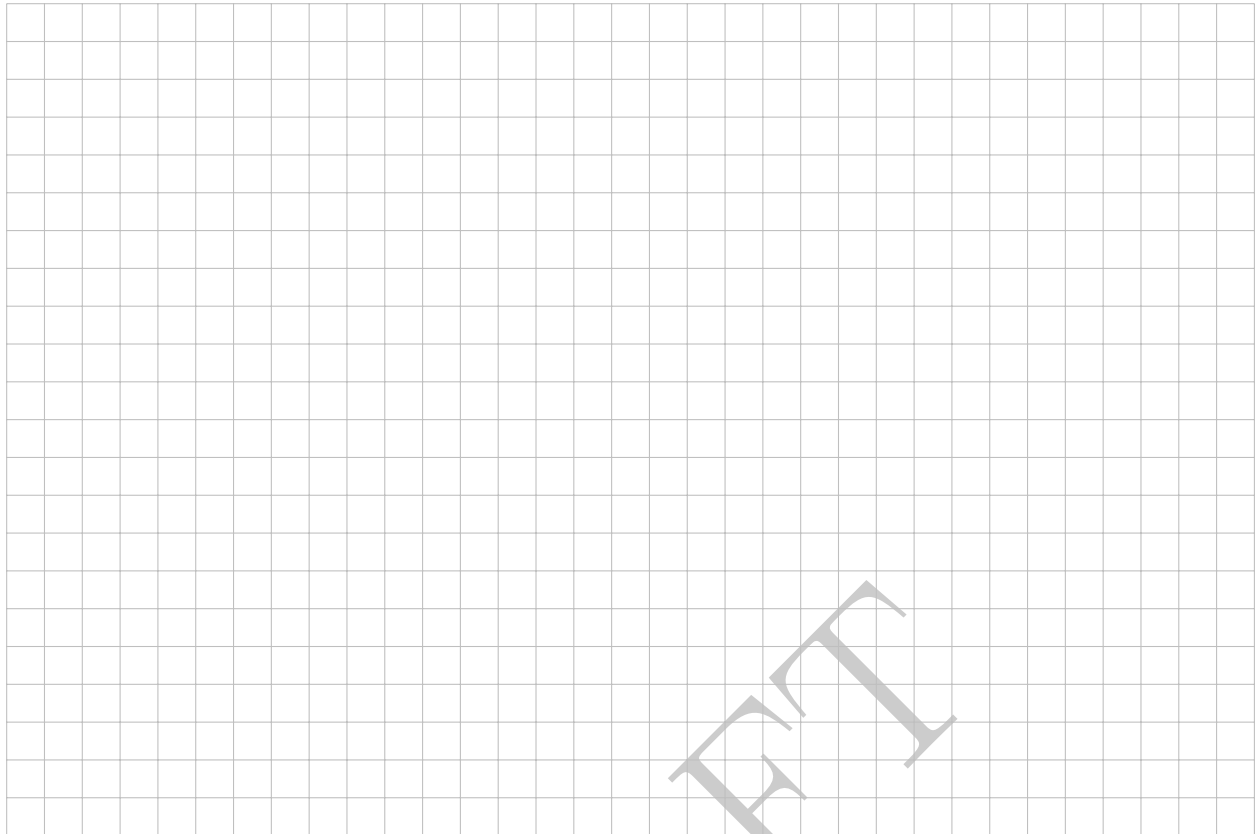
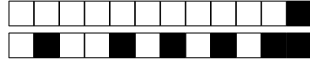
Show how to construct a single-hidden-layer ReLU network (with no skip connections) that exactly represents $f(x)$. In particular, specify:

- The number of hidden units m .
- The weight matrix $W \in \mathbb{R}^{m \times 1}$ and bias vector $b \in \mathbb{R}^m$ for the hidden layer.
- The vector $v \in \mathbb{R}^m$ and scalar $c \in \mathbb{R}$ for the output unit.

show that for every $x \in \mathbb{R}$, this network’s output equals x .

₀ ₁ ₂ ₃ ₄





Question 46: (6 points) **Implementing the maximum function:** consider the function $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by:

$$g(x_1, x_2) = \max\{x_1, x_2\}.$$

design a single-hidden-layer ReLU network (again, with no skip connections) whose output is $\max(x_1, x_2)$:

- Specify the number of hidden units m .
- Provide the matrix $W \in \mathbb{R}^{m \times 2}$ and the bias vector $b \in \mathbb{R}^m$.
- Provide the output-layer weight vector $v \in \mathbb{R}^m$ and scalar bias $c \in \mathbb{R}$.
- Show that for every $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$, the network's output equals $\max\{x_1, x_2\}$.

Hint: use ideas from the construction in the previous question.

₀ ₁ ₂ ₃ ₄ ₅ ₆





DRAFT



DRAFT