- Please do not open the exam unless you are instructed to do so.

- This is a closed book and closed notes exam.

- Everything you need in order to solve the problems is supplied in the body of this exam OR in a cheatsheet at the end of the exam.

- Mark your answers ON THE EXAM ITSELF. If you make a mess, clearly indicate your final answer (box it).

- For true/false questions, CIRCLE True OR False <u>and</u> provide a brief justification for full credit.

- Unless otherwise instructed, for multiple-choice questions, CIRCLE ALL CORRECT CHOICES (in some cases, there may be more than one) <u>and</u> provide a brief justification if the question asks for one.

- If you think something about a question is open to interpretation, feel free to ask the instructor or make a note on the exam.

- If you run out of room for your answer in the space provided, please use the blank pages at the end of the exam and indicate clearly that you've done so.

- DO NOT put answers on the back of any page of the exam.

- DO NOT detach ANY pages.

- You may use scratch paper if needed (provided at the end of the exam).

- You have 2 hours 45 minutes.

- **Besides having the correct answer, being concise and clear is very important. For full credit, you must show your work and explain your answers.**

**Good Luck!**

Legibly write your name and UID in the space provided below to earn 2 points.

# Name:

# UID:

| | | |
|---|---|---|
| Name and UID | | /2 |
| True/False | | /20 |
| Multiple choice | | /32 |
| Short questions | | /16 |
| Poisson regression | | /10 |
| Kernelized K-means | | /10 |
| SVM | | /10 |
| **Total** | | /100 |

# 1 True or False (20 pts)

Choose either True or False for each of the following statements. (If your answer is incorrect, partial credit may be given if your explanation is reasonable.)

1. (2 pts) After mapping the instances into a high dimensional space, a Perceptron may be able to achieve better classification performance on instances it was not able to classify before.

            True                                  False

   **Solution:** True

2. (2 pts) A neural network with all linear activation functions can learn non-linear decision boundaries.

            True                                  False

   **Solution:**  False

3. (2 pts) In the AdaBoost algorithm, at each iteration, we increase the weight for misclassified examples.

            True                                  False

   **Solution:** True

4. (2 pts) The training error of 1-Nearest Neighbor classifier is zero even if the dataset is not linearly separable.

            True                                  False

   **Solution:** True

5. (2 pts) Given a function $k(\boldsymbol{x}_i, \boldsymbol{x}_j) \geq 0$ for all $\boldsymbol{x}_i, \boldsymbol{x}_j$, $k$ is a valid kernel function.

            True                                  False

   **Solution:**  False.

6. (2 pts) Ridge regression is <u>more</u> likely to overfit when we increase $\lambda$ where $\lambda$ is the non-negative weight for the regularization term.

            True                                  False

   **Solution:**  False

7. (2 pts) The hard-margin SVM can <u>sometimes</u> fail to find a solution.

            True                                  False

8. (2 pts) The minimum value of the K-means objective function can be zero for large enough $K$.

<div align="center">True          False</div>

9. (2 pts) A sigmoid function, $\sigma(x) = 1/(1 + \exp(-x))$, can map the hidden layer output of a neural network to a Boolean/binary output.

<div align="center">True          False</div>

10. (2 pts) Given a sentence consisting of $T$ words $(y_1, y_2, \ldots, y_T)$, we would like to use a hidden Markov model (HMM) to predict parts of speech (POS) tags $x_t$ for each word $y_t, t \in \{1, \ldots, T\}$. We have $x_t \in \{1, \ldots, A\}, y_t \in \{1, \ldots, B\}$ for $t \in \{1, \ldots, T\}$ where the total number of possible words is $B$ and the total number of POS tags is $A$. You use the Viterbi algorithm to compute the most probable sequence $(x_1, \ldots, x_T)$. The computational complexity of the algorithm scales as $\mathcal{O}(B^2 T)$.

<div align="center">True          False</div>

# Multiple choice (32 pts)

CIRCLE ALL CORRECT CHOICES (in some cases, there may be more than one)

11. (4 pts) Which of these is an example of an unsupervised learning problem ?

    (a) From emails labeled as spam/not-spam, learn to predict if an email is spam.

    (b) From documents labeled by their topic, learn to classify a document into topics.

    (c) From images of handwritten digits labeled with the digit, learn a handwritten digit classifier.

    (d) From a set of documents, group documents according to their topics.

    **Solution:** d

12. (4 pts) Your classifier has substantially higher test error than training error. Which of the following statements could explain this observation?

    (a) The hypothesis space from which the classifier was chosen is too complex.

    (b) The hypothesis space from which the classifier was chosen is too simple.

    (c) The distribution of test data differs from the distribution of training data.

    (d) The size of the training dataset is too small.

    **Solution:** a,c,d

13. (4 pts) You are given a training dataset for a binary classification problem: $\{(x_1, y_1), \ldots, (x_N, y_N)\}$ where $(x_n, y_n), n \in \{1, \ldots, N\}$ is an instance-label pair and $y_n \in \{0, 1\}$ where 1 refers to the positive class. Your classifier predicts 1 with probability $0 < \theta < \frac{1}{2}$ and 0 otherwise, independent of the features of an instance. Assume the sample size $N$ is large. Which of following statements is true ?

    (a) The sensitivity of the classifier is $\theta$.

    (b) The specificity of the classifier is $\theta$.

    (c) The false positive rate of the classifier is $\theta$.

    (d) The ROC curve for this classifier is the diagonal line.

    **Solution:** a,c,d

14. (4 pts) Given the same training data consisting of $N$ instances and $D$ features, we fit linear regression and obtain the optimal parameters $\boldsymbol{\theta}_{OLS}$. We also fit ridge regression with regularization parameter $\lambda > 0$ and obtain the optimal parameters $\boldsymbol{\theta}_{Ridge}$. Let $RSS(\boldsymbol{\theta})$ denote the residual sum of squares cost function evaluated on the training set for the model associated with the parameter $\boldsymbol{\theta}$. Which of the following will always hold ?

   (a) $RSS(\boldsymbol{\theta}_{Ridge}) \geq RSS(\boldsymbol{\theta}_{OLS})$

   (b) $RSS(\boldsymbol{\theta}_{Ridge}) \leq RSS(\boldsymbol{\theta}_{OLS})$

   (c) $RSS(\boldsymbol{\theta}_{Ridge}) = RSS(\boldsymbol{\theta}_{OLS})$

   (d) $RSS(\boldsymbol{\theta}_{OLS}) \geq 0$.

   **Solution:** a,d

15. (4 pts) Let $\boldsymbol{X} \in \mathbb{R}^{N \times D}$ be the design matrix with each row corresponding to the features of an example and $\boldsymbol{y} \in \mathbb{R}^{N}$ be a vector of all the labels. The OLS solution is $\boldsymbol{\theta}_{OLD} = (\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{y}$. Given a new test data point $\boldsymbol{x}$, our prediction for this data point is given by $\hat{y}_{OLD} = \boldsymbol{\theta}_{OLD}{}^{\mathrm{T}} \boldsymbol{x}$. We then scale each feature in the training and the test data by 2 and compute the OLS solution $\boldsymbol{\theta}_{NEW}$ to make our prediction $\hat{y}_{NEW} = \boldsymbol{\theta}_{NEW}{}^{\mathrm{T}} \boldsymbol{x}$. What is the relation between $\hat{y}_{NEW}$ and $\hat{y}_{OLD}$?

   (a) $\hat{y}_{NEW} = 2\hat{y}_{OLD}$

   (b) $\hat{y}_{NEW} = 4\hat{y}_{OLD}$

   (c) $\hat{y}_{NEW} = \frac{1}{2}\hat{y}_{OLD}$

   (d) $\hat{y}_{NEW} = \hat{y}_{OLD}$

   **Solution:** d

16. (4 pts) Random variables $(X_1, X_2, X_3, X_4)$ are distributed according to a Markov model. Which of the following statements is true of the distributions of these random variables?

   (a) $P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2)P(x_3)P(x_4)$

   (b) $P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)P(x_4|x_1, x_2, x_3)$

   (c) $P(x_1, x_2, x_3, x_4) = P(x_1)P(x_2|x_1)P(x_3|x_2)P(x_4|x_3)$

   (d) $P(x_4|x_1, x_2, x_3) = P(x_4|x_3)$

   **Solution:** b,c,d

17. (4 pts) We would like to run PCA on a dataset with 100 samples and five features. The eigenvalues of the covariance matrix are $(80, 15, 5, 0, 0)$. What is the minimum number $K$ of principal components needed so that the transformed $K$ features explain at least 90% of the variance?

   (a) 1
   (b) 2
   (c) 3
   (d) 4

   **Solution:** b

18. (4 pts) We would like to apply the EM algorithm to estimate the parameters of Gaussian Mixture Models (GMMs). Which of the following are true of the EM algorithm applied to GMMs?

   (a) In the E-step, we compute the probability that a data point is drawn from each mixture component.
   (b) In the E-step, each data point is assigned to one of the mixture components.
   (c) In the M-step, we update the mixture weights of the GMM.
   (d) In the M-step, we update the mean and covariance matrix of each mixture component.

   **Solution:** a,c,d.

# Short Answer Questions (16 pts)

Most of the following questions can be answered in one or two sentences. Please make your answer concise and to the point.

19. (4 pts) Describe the difference between *maximum likelihood estimation* (MLE) and *maximum a posteriori estimation* (MAP), and state under what condition MAP is equivalent to MLE?

    **Solution:** MLE maximizes the likelihood function, while MAP computes the value of the parameters that maximizes the posterior distribution (i.e,. P(parameter | data)) or P(parameter) P(data | parameter). When the prior distribution on the parameters is a uniform distribution.

20. (6 pts) Given vectors $x$ and $z$ in $R^3$, define the kernel $K_\beta(x; z) = (\beta + x \cdot z)^2$ for any value $\beta > 0$. Find the corresponding feature map $\phi_\beta(\cdot)$.

    **Solution:**  $x = [x_1, x_2, x_3]^T$, $z = [z_1, z_2, z_3]^T$.

    $$\begin{aligned}
    \text{Now, } K_\beta(x; z) &= (\beta + x \cdot z)^2 \\
    &= (\beta + (x_1 z_1 + x_2 z_2 + x_3 z_3))^2 \\
    &= \beta^2 + 2\beta(x_1 z_1 + x_2 z_2 + x_3 z_3) + (x_1 z_1 + x_2 z_2 + x_3 z_3)^2 \\
    &= \beta^2 + 2\beta x_1 z_1 + 2\beta x_2 z_2 + 2\beta x_3 z_3 + x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 \\
    &\quad + x_2^2 z_2^2 + 2x_1 z_1 x_3 z_3 + 2x_2 z_2 x_3 z_3 + x_3^2 z_3^2
    \end{aligned}$$

    Comparing this with $K_\beta(x; z) = \phi_\beta(x) \cdot \phi_\beta(z)$,
    $\phi_\beta(x) = [\beta, \sqrt{2\beta}x_1, \sqrt{2\beta}x_2, \sqrt{2\beta}x_3, \sqrt{2}x_1 x_2, \sqrt{2}x_1 x_3, \sqrt{2}x_2 x_3, x_1^2, x_2^2, x_3^2]^T$

21. (6 pts) We are given the confusion matrix for a binary classifier.

| Actual class | Negative | Positive |
|---|---|---|
| Predicted class | | |
| Negative | 80 | 50 |
| Positive | 20 | 50 |

Compute the following quantities related to its accuracy:

(a) (1 pts) True positives **Solution:** TP=50

(b) (1 pts) False positives **Solution:** FP=20

(c) (1 pts) Recall **Solution:** Recall=$\frac{TP}{Positives} = \frac{50}{100}$

(d) (1 pts) False positive rate **Solution:** False positive rate=$\frac{FP}{Negatives} = \frac{20}{100}$

(e) (1 pts) Precision **Solution:** Precision=$\frac{TP}{TP+FP} = \frac{50}{70}$

(f) (1 pts) Accuracy **Solution:** Accuracy=$\frac{TP+TN}{Positives+Negatives} = \frac{130}{200}$

# 22 Kernelized K-means (10 pts)

K-means with Euclidean distance metric assumes that each pair of clusters is linearly separable. This may not be the case. We have seen that we can use kernels to obtain a non-linear version of an algorithm that is linear by nature and K-means is no exception. Recall that there are two main aspects of kernelized algorithms: (i) the solution is expressed as a linear combination of training examples, (ii) the algorithm relies only on inner products between data points rather than their explicit representation. We will show that these two aspects can be satisfied in K-means.

1. (3 pts) Let $z_{nk}$ be an indicator that is equal to 1 if the $x_n$ is currently assigned to the $k^{th}$ cluster and 0 otherwise ($1 \leq n \leq N$ and $1 \leq k \leq K$). Show that the $k^{th}$ cluster center $\boldsymbol{\mu}_k$ can be updated as $\sum_{n=1}^{N} \alpha_{nk} \boldsymbol{x}_n$. Specifically, show how $\alpha_{nk}$ can be computed given all $z$'s.

   **Solution:**

   $$\alpha_{nk} = \frac{z_{nk}}{\sum_{m=1}^{N} z_{mk}}$$

2. (3 pts) Given two data points $\boldsymbol{x}_1$ and $\boldsymbol{x}_2$, show that the square distance $\|\boldsymbol{x}_1 - \boldsymbol{x}_2\|^2$ can be computed using only (linear combinations of) inner products.

   **Solution:**

   $$\|\boldsymbol{x}_1 - \boldsymbol{x}_2\|^2 = \boldsymbol{x}_1^T \boldsymbol{x}_1 + \boldsymbol{x}_2^T \boldsymbol{x}_2 - 2\boldsymbol{x}_1^T \boldsymbol{x}_2$$

3. (4 pts) Given the results of the above two parts, show how to compute the squared distance $\|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|^2$ using only (linear combinations of) inner products between the data points. $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ (You can leave your answer in terms of $\alpha_{nk}$ and inner product of $\boldsymbol{x}_n$ and $\boldsymbol{x}_k$.

**Solution:**

$$
\begin{aligned}
\|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|^2 &= \|\boldsymbol{x}_n - \sum_{n=1}^{N} \alpha_{nk}\boldsymbol{x}_n\|^2 \\
&= \boldsymbol{x}_n^T\boldsymbol{x}_n + \sum_i \sum_j \alpha_{ik}\alpha_{jk}\boldsymbol{x}_i^T\boldsymbol{x}_j - 2\sum_i \alpha_{ik}\boldsymbol{x}_n^T\boldsymbol{x}_i
\end{aligned}
$$

Note: This means that given a kernel $K$, we can run Lloyd's algorithm. We begin with some initial data points as centers and use the answer to part c) to find the closest center for each data point, giving us the initial $z_{nk}$'s. We then repeatedly use the answer to part a) to reassign the cluster centers and use the answer to part c) to reassign points to centers and update the $z_{nk}$'s.

# 23 Poisson Regression (10 pts)

We want to predict the number of user clicks on a website. The number of clicks takes on values in $\{0, 1, 2, \ldots\}$. In class, we showed how linear regression (ordinary least squares) can be interpreted as a probabilistic model where the output is real-valued. Logistic regression is a probabilistic model where the output takes values in $\{0, 1\}$. In this problem, we want to model outputs in $\{0, 1, 2, \ldots\}$.

In this example, we have a training set $\{(\boldsymbol{x}_n, y_n)\}_{n=1}^{N}$ where $\boldsymbol{x}_n \in \mathbb{R}^D$ and $y_n \in \{0, 1, 2, \ldots\}$. Now we model our target $y_n$ as distributed according to a Poisson distribution. Specifically

$$p(y_n|\boldsymbol{x}_n; \boldsymbol{\theta}) = \frac{1}{y_n!} \exp\left(y_n \boldsymbol{\theta}^T \boldsymbol{x}_n\right) \exp\left(-\exp(\boldsymbol{\theta}^T \boldsymbol{x}_n)\right)$$

1. (4 pts) Write the log likelihood of the parameters $l(\boldsymbol{\theta})$. Express your answer in terms of $y_n$, $\boldsymbol{x}_n$, $\boldsymbol{\theta}$.

   **Solution:**

$$\begin{aligned} l(\boldsymbol{\theta}) &= \sum_{n=1}^{N} \log p(y_n|\boldsymbol{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^{N} \left(y_n \boldsymbol{\theta}^T \boldsymbol{x}_n - \exp(\boldsymbol{\theta}^T \boldsymbol{x}_n) - \log(y_n!)\right) \end{aligned}$$

   [2pt for the log likelihood as a sum of log probabilities; 2pt for expressing in terms of the model parameters] [Full credit if the final answer is correct (even if line 1 is missing)]

2. (6 pts) Show that the gradient can be written in the form $\nabla l(\boldsymbol{\theta}) = \sum_{n=1}^{N} \epsilon_n \boldsymbol{x}_n$ for some $\epsilon_n$. Write $\epsilon_n$ in terms of $\boldsymbol{x}_n, y_n$ and $\boldsymbol{\theta}$.

**Solution:**

$$
\begin{aligned}
\nabla l(\boldsymbol{\theta}) &= \sum_{n=1}^{N} \left( y_n \boldsymbol{x}_n - \exp(\boldsymbol{\theta}^T \boldsymbol{x}_n) \boldsymbol{x}_n \right) \\
&= \sum_{n=1}^{N} \left( y_n - \exp(\boldsymbol{\theta}^T \boldsymbol{x}_n) \right) \boldsymbol{x}_n \\
&= \sum_{n=1}^{N} \epsilon_n \boldsymbol{x}_n
\end{aligned}
$$

Thus,

$$
\epsilon_n := y_n - \exp\left( \theta^\top x_n \right)
$$

# 24 SVM (10 pts)

We are attempting to use hard-margin SVM to solve a binary classification problem given a dataset $\mathcal{D}$ that has two samples $\{(x_1, y_1), (x_2, y_2)\}$ with $x_i \in R$ and $y_i \in \{-1, +1\}$, $(x_1 = 0, y_1 = -1)$ and $(x_2 = \sqrt{2}, y_2 = 1)$. To obtain a non-linear classifier, consider mapping the data points by $\boldsymbol{\phi}(x) = [1, \sqrt{2}x, x^2]^T$ to a 3-dimensional space. The hard-margin SVM has the form

$$
\begin{aligned}
min_{\boldsymbol{w},b} \quad & \frac{1}{2}||\boldsymbol{w}||_2^2 \\
\text{s.t.} \quad & y_1(\boldsymbol{w}^T\boldsymbol{\phi}(x_1) + b) \geq 1 \\
& y_2(\boldsymbol{w}^T\boldsymbol{\phi}(x_2) + b) \geq 1
\end{aligned} \tag{1}
$$

1. (2 pts) Write a vector that is parallel to the optimal vector $\boldsymbol{w}^*$ and justify your answer.
   **Solution:** $\boldsymbol{w}^*$ is parallel to the vector from $\boldsymbol{\phi}(x_1)$ to $\boldsymbol{\phi}(x_2)$, that is $\boldsymbol{\phi}(x_2) - \boldsymbol{\phi}(x_1)$. Remember that the vector $\boldsymbol{w}$ is orthogonal to the hyperplane. The optimal hyperplane should pass through the midpoint of the two samples and be perpendicular to the segment connecting them. Thus, we know that the optimal $\boldsymbol{w}^*$ should be parallel to $\boldsymbol{\phi}(x_2) - \boldsymbol{\phi}(x_1) = [0, 2, 2]^T$.

2. (2 pts) Write down the value of the margin achieved by the optimal $\boldsymbol{w}^*$. **Solution:** As stated above, the optimal hyperplane will pass through the midpoint and the midpoint should be the closest point on the hyperplane to either sample. Thus, you can answer this question by finding the midpoint and finding the distance from the midpoint to either sample.

   Alternatively, you could note that the hyperplane should be equidistant from both samples. Thus you could calculate the distance between the two samples and divide by 2.
   $\gamma = \frac{||[1,0,0]^T - [1,2,2]^T||_2}{2} = \frac{\sqrt{8}}{2} = \sqrt{2}$

3. (2 pts) Solve for $\boldsymbol{w}^*$ using the fact that the margin is equal to $\frac{1}{||\boldsymbol{w}^*||_2}$. **Solution:** From part 1 we know that $\boldsymbol{w}^*$ takes the form of $[0, a, a]^T$ for some positive scalar $a$. From part 2 we know that $\gamma = \sqrt{2}$. Thus,

$$
\begin{aligned}
\frac{1}{||\boldsymbol{w}^*||_2} &= \gamma = \sqrt{2} \\
||\boldsymbol{w}^*||_2^2 &= \left(\frac{1}{\sqrt{2}}\right)^2 \\
\boldsymbol{w}^{*\top}\boldsymbol{w}^* &= \frac{1}{2} \\
2a^2 &= \frac{1}{2} \\
a &= \sqrt{\frac{1}{4}} = \frac{1}{2}
\end{aligned}
$$

Thus, $\boldsymbol{w}^* = [0, \frac{1}{2}, \frac{1}{2}]$

4. (2 pts) Solve for $b^*$ and write down the decision boundary $f(x) = \boldsymbol{w}^{*T}\boldsymbol{\phi}(x) + b^*$.
   **Solution:** We know that $y_1\left(\boldsymbol{w}^{*\top}\boldsymbol{\phi}(x_1) + b^*\right) = 1$, and from part 3 we know the value of $\boldsymbol{w}^*$. Thus, we can solve for $b^*$.

$$
\begin{aligned}
1 \quad &= y_1\left(\boldsymbol{w}^{*\top}\boldsymbol{\phi}(x_1) + b^*\right) \\
-1 \quad &= [0, \tfrac{1}{2}, \tfrac{1}{2}]\left([1, 0, 0]^\top\right)) + b^* \\
-1 \quad &= b^*
\end{aligned}
$$

$$
\begin{aligned}
f(x) \quad &= [0, \tfrac{1}{2}, \tfrac{1}{2}]\left([1, \sqrt{2}x, x^2]^\top\right)) - 1 \\
&= \tfrac{x^2}{2} + \tfrac{\sqrt{2}x}{2} - 1
\end{aligned}
$$

5. (2 pts) Given a new data point $x_{\text{new}}$, what is the prediction of the label for $x_{\text{new}}$ using the above parameters? **Solution:**

$$\hat{y}_{\text{new}} = sign\left(f(x_{\text{new}})\right)$$
$$= sign\left(\frac{x_{\text{new}}^2}{2} + \frac{\sqrt{2}x_{\text{new}}}{2} - 1\right)$$

# Identities

## Probability density/mass functions for some distributions

$$\text{Normal} \quad : \quad P(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{Multinomial} \quad : \quad P(\boldsymbol{x}; \boldsymbol{\pi}) = \prod_{k=1}^{K} \pi_k{}^{x_k}$$

$\boldsymbol{x}$ is a length $K$ vector with exactly one entry equal to 1
and all other entries equal to 0

$$\text{Poisson} \quad : \quad P(x; \lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$$

## Matrix calculus

Here $\boldsymbol{x} \in \mathbb{R}^n, \boldsymbol{b} \in \mathbb{R}^n, \boldsymbol{A} \in \mathbb{R}^{n\times n}$. $\boldsymbol{A}$ is symmetric.

$$\begin{aligned}
\nabla \boldsymbol{x}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{x} &= 2\boldsymbol{A}\boldsymbol{x} \\
\nabla \boldsymbol{b}^{\mathrm{T}} \boldsymbol{x} &= \boldsymbol{b}
\end{aligned}$$

You may use this page for scratch space.

You may use this page for scratch space.