

- You have 3 hours for the exam.
- The exam is closed book, closed notes except your one-page (two sides) or two-page (one side) crib sheet.
- Please use non-programmable calculators only.
- Mark your answers ON THE EXAM ITSELF. If you are not sure of your answer you may wish to provide a *brief* explanation. All short answer sections can be successfully answered in a few sentences AT MOST.
- For true/false questions, fill in the *True/False* bubble.
- For multiple-choice questions, fill in the bubbles for **ALL CORRECT CHOICES** (in some cases, there may be more than one). For a question with p points and k choices, every false positive will incur a penalty of $p/(k - 1)$ points.
- For short answer questions, **unnecessarily long explanations and extraneous data will be penalized**. Please try to be terse and precise and do the side calculations on the scratch papers provided.
- Please **draw a bounding box around your answer** in the Short Answers section. A missed answer without a bounding box will not be regraded.

First name	
Last name	
SID	

For staff use only:

Q1.	True/False	/23
Q2.	Multiple Choice Questions	/36
Q3.	Short Answers	/26
	Total	/85

Q1. [23 pts] True/False

- (a) [1 pt] Solving a non linear separation problem with a hard margin Kernelized SVM (Gaussian RBF Kernel) might lead to overfitting.
 True False
- (b) [1 pt] In SVMs, the sum of the Lagrange multipliers corresponding to the positive examples is equal to the sum of the Lagrange multipliers corresponding to the negative examples.
 True False
- (c) [1 pt] SVMs directly give us the posterior probabilities $P(y = 1|x)$ and $P(y = -1|x)$.
 True False
- (d) [1 pt] $V(X) = E[X]^2 - E[X^2]$
 True False
- (e) [1 pt] In the discriminative approach to solving classification problems, we model the conditional probability of the labels given the observations.
 True False
- (f) [1 pt] In a two class classification problem, a point on the Bayes optimal decision boundary x^* always satisfies $P(y = 1|x^*) = P(y = 0|x^*)$.
 True False
- (g) [1 pt] Any linear combination of the components of a multivariate Gaussian is a univariate Gaussian.
 True False
- (h) [1 pt] For any two random variables $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $Y \sim \mathcal{N}(\mu_2, \sigma_2^2)$, $X + Y \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.
 True False
- (i) [1 pt] Stanford and Berkeley students are trying to solve the same logistic regression problem for a dataset. The Stanford group claims that their initialization point will lead to a much better optimum than Berkeley's initialization point. Stanford is correct.
 True False
- (j) [1 pt] In logistic regression, we model the odds ratio ($\frac{p}{1-p}$) as a linear function.
 True False
- (k) [1 pt] Random forests can be used to classify infinite dimensional data.
 True False
- (l) [1 pt] In boosting we start with a Gaussian weight distribution over the training samples.
 True False
- (m) [1 pt] In Adaboost, the error of each hypothesis is calculated by the ratio of misclassified examples to the total number of examples.
 True False
- (n) [1 pt] When $k = 1$ and $N \rightarrow \infty$, the kNN classification rate is bounded above by twice the Bayes error rate.
 True False
- (o) [1 pt] A single layer neural network with a sigmoid activation for binary classification with the cross entropy loss is exactly equivalent to logistic regression.
 True False

- (p) [1 pt] The loss function for LeNet5 (the convolutional neural network by LeCun et al.) is convex.
 True False
- (q) [1 pt] Convolution is a linear operation i.e. $(\alpha f_1 + \beta f_2) * g = \alpha f_1 * g + \beta f_2 * g$.
 True False
- (r) [1 pt] The k-means algorithm does coordinate descent on a non-convex objective function.
 True False
- (s) [1 pt] A 1-NN classifier has higher variance than a 3-NN classifier.
 True False
- (t) [1 pt] The single link agglomerative clustering algorithm groups two clusters on the basis of the maximum distance between points in the two clusters.
 True False
- (u) [1 pt] The largest eigenvector of the covariance matrix is the direction of minimum variance in the data.
 True False
- (v) [1 pt] The eigenvectors of AA^T and $A^T A$ are the same.
 True False
- (w) [1 pt] The non-zero eigenvalues of AA^T and $A^T A$ are the same.
 True False

Q2. [36 pts] Multiple Choice Questions

(a) [4 pts] In linear regression, we model $P(y|x) \sim \mathcal{N}(w^T x + w_0, \sigma^2)$. The irreducible error in this model is _____.

- σ^2
 $E[(y - E[y|x])|x]$
 $E[(y - E[y|x])^2|x]$
 $E[y|x]$

(b) [4 pts] Let S_1 and S_2 be the set of support vectors and w_1 and w_2 be the learnt weight vectors for a linearly separable problem using hard and soft margin linear SVMs respectively. Which of the following are correct?

- $S_1 \subset S_2$
 S_1 may not be a subset of S_2
 $w_1 = w_2$
 w_1 may not be equal to w_2 .

(c) [4 pts] Ordinary least-squares regression is equivalent to assuming that each data point is generated according to a linear function of the input plus zero-mean, constant-variance Gaussian noise. In many systems, however, the noise variance is itself a positive linear function of the input (which is assumed to be non-negative, i.e., $x \geq 0$). Which of the following families of probability models correctly describes this situation in the univariate case?

- $P(y|x) = \frac{1}{\sigma\sqrt{2\pi x}} \exp\left(-\frac{(y-(w_0+w_1x))^2}{2x\sigma^2}\right)$
 $P(y|x) = \frac{1}{\sigma\sqrt{2\pi x}} \exp\left(-\frac{(y-(w_0+(w_1+\sigma^2)x))^2}{2\sigma^2}\right)$
 $P(y|x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y-(w_0+w_1x))^2}{2\sigma^2}\right)$
 $P(y|x) = \frac{1}{\sigma x\sqrt{2\pi}} \exp\left(-\frac{(y-(w_0+w_1x))^2}{2x^2\sigma^2}\right)$

(d) [3 pts] The left singular vectors of a matrix A can be found in _____.

- Eigenvectors of AA^T
 Eigenvectors of A^2
 Eigenvectors of $A^T A$
 Eigenvalues of AA^T

(e) [3 pts] Averaging the output of multiple decision trees helps _____.

- Increase bias
 Increase variance
 Decrease bias
 Decrease variance

(f) [4 pts] Let A be a symmetric matrix and S be the matrix containing its eigenvectors as column vectors, and D a diagonal matrix containing the corresponding eigenvalues on the diagonal. Which of the following are true:

- $AS = SD$
 $SA = DS$
 $AS = DS$
 $AS = DS^T$

(g) [4 pts] Consider the following dataset: $A = (0, 2)$, $B = (0, 1)$ and $C = (1, 0)$. The k-means algorithm is initialized with centers at A and B . Upon convergence, the two centers will be at

- A and C
 C and the midpoint of AB
 A and the midpoint of BC
 A and B

(h) [3 pts] Which of the following loss functions are convex?

- Misclassification loss
- Logistic loss
- Hinge loss
- Exponential Loss ($e^{-yf(x)}$)

(i) [3 pts] Consider T_1 , a decision stump (tree of depth 2) and T_2 , a decision tree that is grown till a maximum depth of 4. Which of the following is/are correct?

- $Bias(T_1) < Bias(T_2)$
- $Bias(T_1) > Bias(T_2)$
- $Variance(T_1) < Variance(T_2)$
- $Variance(T_1) > Variance(T_2)$

(j) [4 pts] Consider the problem of building decision trees with k -ary splits (split one node into k nodes) and you are deciding k for each node by calculating the entropy impurity for different values of k and optimizing simultaneously over the splitting threshold(s) and k . Which of the following is/are true?

- The algorithm will always choose $k = 2$
- The algorithm will prefer high values of k
- There will be $k - 1$ thresholds for a k -ary split
- This model is strictly more powerful than a binary decision tree.

Q3. [26 pts] Short Answers

- (a) [5 pts] Given that (x_1, x_2) are jointly normally distributed with $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and $\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$ ($\sigma_{21} = \sigma_{12}$), give an expression for the mean of the conditional distribution $p(x_1|x_2 = a)$.

This can be solved by writing $p(x_1|x_2 = a) = \frac{p(x_1, x_2 = a)}{p(x_2 = a)}$. x_2 being a component of a multivariate Gaussian is a univariate Gaussian with $x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$. Write out the Gaussian densities and simplify (complete squares) to see the following:

$$x_1|x_2 = a \sim \mathcal{N}(\bar{\mu}, \bar{\sigma}^2), \quad \bar{\mu} = \mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(a - \mu_2)$$

- (b) [4 pts] The logistic function is given by $\sigma(x) = \frac{1}{1+e^{-x}}$. Show that $\sigma'(x) = \sigma(x)(1 - \sigma(x))$.

$$\sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{(1+e^{-x})} \cdot \frac{e^{-x}}{(1+e^{-x})} = \left(\frac{1}{1+e^{-x}} \right) \left(1 - \frac{1}{1+e^{-x}} \right) = \sigma(x)(1 - \sigma(x))$$

- (c) Let X have a uniform distribution

$$p(x; \theta) = \begin{cases} \frac{1}{\theta} & 0 \leq x \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

Suppose that n samples x_1, \dots, x_n are drawn independently according to $p(x; \theta)$.

- (i) [5 pts] The maximum likelihood estimate of θ is $x_{(n)} = \max(x_1, x_2, \dots, x_n)$. Show that this estimate of θ is biased.

Biased estimator: $\hat{\theta}$ (the sample estimate) is a biased estimator of θ (the population distribution parameter) if $E[\hat{\theta}] \neq \theta$.

Here $\hat{\theta} = x_{(n)}$. And $E[x_{(n)}] = \frac{n}{n+1}\theta \neq \theta$. The steps for finding $E[x_{(n)}]$ are given in the solutions of Homework 2, problem 5(c).

- (ii) [2 pts] Give an expression for an unbiased estimator of θ .

$$\hat{\theta}_{unbiased} = \frac{n+1}{n}x_{(n)}$$

$$E[\hat{\theta}_{unbiased}] = E\left[\frac{n+1}{n}x_{(n)}\right] = \frac{n+1}{n}E[x_{(n)}] = \frac{n+1}{n} \times \frac{n}{n+1}\theta = \theta$$

- (d) [5 pts] Consider the problem of fitting the following function to a dataset of 100 points $\{(x_i, y_i)\}, i = 1 \dots 100$:

$$y = \alpha \cos(x) + \beta \sin(x) + \gamma$$

This problem can be solved using the least squares method with a solution of the form:

$$\begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} = (X^T X)^{-1} X^T Y$$

What are X and Y ?

$$X = \begin{bmatrix} \cos(x_1) & \sin(x_1) & 1 \\ \cos(x_2) & \sin(x_2) & 1 \\ \vdots & \vdots & \vdots \\ \cos(x_{100}) & \sin(x_{100}) & 1 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{100} \end{bmatrix}$$

- (e) [5 pts] Consider the problem of binary classification using the Naive Bayes classifier. You are given two dimensional features (X_1, X_2) and the categorical class conditional distributions in the tables below. The entries in the tables correspond to $P(X_1 = x_1 | C_i)$ and $P(X_2 = x_2 | C_i)$ respectively. The two classes are *equally likely*.

$X_1 =$ \backslash Class	C_1	C_2
-1	0.2	0.3
0	0.4	0.6
1	0.4	0.1

$X_2 =$ \backslash Class	C_1	C_2
-1	0.4	0.1
0	0.5	0.3
1	0.1	0.6

Given a data point $(-1, 1)$, calculate the following posterior probabilities:

$$P(C_1 | X_1 = -1, X_2 = 1) = \text{Using Bayes' Rule and conditional independence assumption of Naive Bayes}$$

$$\frac{P(X_1=-1, X_2=1 | C_1) P(C_1)}{P(X_1=-1, X_2=1)} = \frac{P(X_1=-1 | C_1) P(X_2=1 | C_1) P(C_1)}{P(X_1=-1 | C_1) P(X_2=1 | C_1) P(C_1) + P(X_1=-1 | C_2) P(X_2=1 | C_2) P(C_2)} = 0.1$$

$$P(C_2 | X_1 = -1, X_2 = 1) = 1 - P(C_1 | X_2 = -1, X_1 = 1) = 0.9$$

SCRATCH PAPER

SCRATCH PAPER