

1. **Short Questions.** Answer each questions with at most two sentences, and/or a picture. For the (true/false) questions, if the answer is true, provide a short justification. If false, explain why or give a small counter example. No point if no proper explanation is provided.

a. ( 2pts) (true/false) We can use linear regression to learn the parameters for the model  $y = w_1x_1x_2 + w_2x_2 \log x_1$ .

b. (2 pts) (true/false) We can use linear regression to learn the parameters for the model  $y = x_1^{w_1} + x_2^{w_2}$ .

d. (3 pts) Why does the kernel trick allow us to solve SVMs with high dimensional feature spaces, without significantly increasing the running time?

d. (4 pts) Consider the following data set:

$\bigcirc$	$+$
$+$	$\bigcirc$

Circle all the classifiers that will achieve zero training error on this data set. (maybe more than one).

- i. Logistic Regression
- ii. Depth-2 Decision tree
- iii. Support vector machine with linear kernel
- iv. Support vector machine with quadratic kernel

e. (4pts) Assume we computed the parameters for a Naive Bayes classifier. How can we use these parameters to compute  $P(x)$  for a given input vector  $x = [x_1, x_2, \dots, x_d]^T$ ?

2. **Linear Regression and regularization.** We have a regression problem with the target  $y$  and a single input feature  $x$ . We know that  $y$  is a polynomial function of  $x$  but do not know the exact order of the polynomial except that it is  $\leq 4$ . Answer the following questions.

a. (4pts) Given a training set, consider the following strategy for learning the polynomial function. For each  $k \in \{1, 2, 3, 4\}$ , we learn a linear regression model  $\hat{y} = w_0 + w_1x + \dots + w_kx^k$  by minimizing the Sum of Squared Error (SSE)  $\sum_{i=1}^N (\hat{y}_i - y_i)^2$  on the training data. We then pick the model that achieves the lowest training SSE. Will this strategy find the correct order  $k$ ? Briefly justify your answer.

b. (4pts) Now we will focus on the 4-th order model  $\hat{y} = w_0 + w_1x + \dots + w_4x^4$  and learn the parameters by minimizing the following objective with  $L_1$  regularization:

$$\sum_{i=1}^N (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^4 |w_j|$$

As the value for  $\lambda$  increases, what will happen to the training SSE of the learned model and why?

c. (4pts) Continue from [b], the true model is  $y = x + x^2$ . As we gradually increase  $\lambda$  to a very large value, which of the following do you expect to happen to the learned parameters  $w_1, \dots, w_4$ ? (may have multiple correct ones.) Briefly justify your answer.

- i. The changes to the weights are unpredictable.
- ii. All weights will become smaller at the same time.
- iii.  $w_3$  and  $w_4$  will decrease to zero first.
- iv. Eventually all weights will decrease to zero as  $\lambda$  gets to be large enough.

3. Support vector machines. For the data set shown below, we apply soft margin SVM.

a (2.5pts) Where would the decision boundary be for a very small  $c$  value, e.g.,  $c = 0.00001$ ?

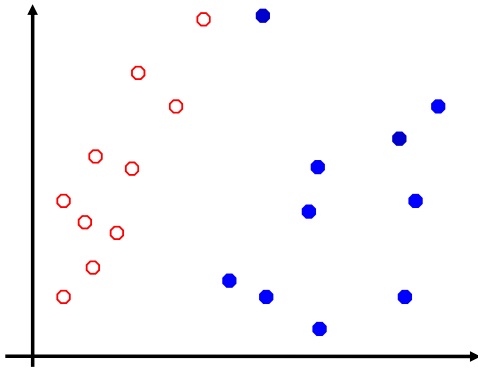
b (2.5pts) Where would the decision boundary be for a very large  $c$  value, e.g.,  $c = 100000$ ?

Now we will add another **positive** example (solid dot) to the training set.

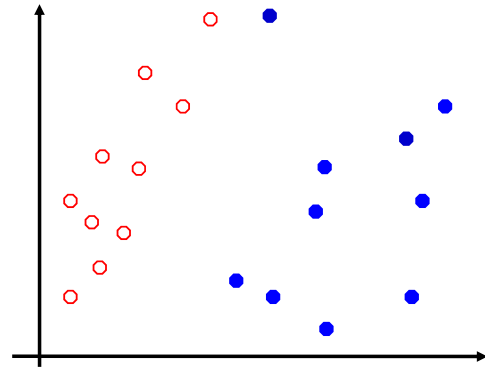
c (2.5pts) Where could you place it so that it will not influence the learned decision boundary for a very large  $c$  value.

d (2.5pts) Where could you place it so that it will significantly influence the learning decision boundary for a very large  $c$  value.

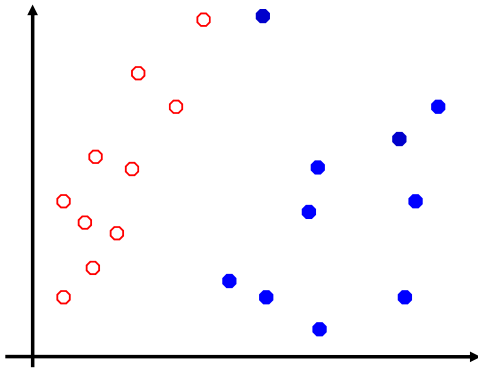
Use the following figures to answer the above questions.



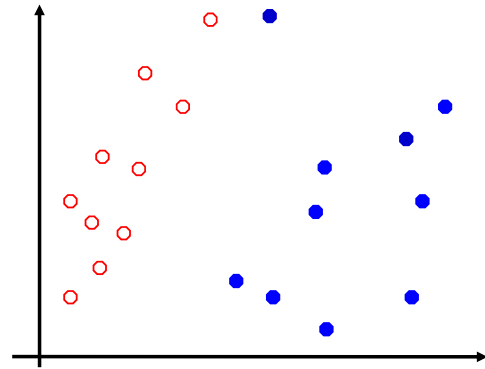
(a)



(b)



(c)



(d)

4. Use the following data to build a decision tree to predict whether a student will be lazy (L) or diligent (D) based on weight (Normal or Underweight), eye color (Amber or Violet) and the number of eyes(2, 3 or 4).

Weight	Eye Color	Num. Eyes	Output
N	A	2	L
N	V	2	L
N	V	2	L
U	V	3	L
U	V	3	L
U	A	4	D
N	A	4	D
N	V	4	D
U	A	3	D
U	A	3	D

- a. (3pts) What is the conditional entropy  $H(\text{Eyecolor}|\text{Weight} = N)$ ?

- b. (4pts) What attribute will information gain select to place at the root of the tree? (Note: we do not worry about the bias for multi-way split here. Just evaluate multiway split by its information gain.)

c. (4pts) Draw the full decision tree learned from this data (no pruning or early stopping).

d. (2pts) What is the training set error of this tree?