

## Q1. Kernel Perceptron

Remember that the perceptron update rule looks like:

$$w \leftarrow w + yx$$

for input feature vectors  $x$  and class labels  $y \in \{-1, 1\}$ . If  $w \cdot x = 0$ , we predict positive label.

- (a) Suppose  $w = [1, 1]$  initially, and we observe the following training examples:

$x_0$	$x_1$	$y$
1	2	-1
3	1	1
1	1	-1
1	0	1

What is the final value of  $w$ ?

- (b) Notice that because of the update rule, the final weight vector  $w^*$  is just a linear combination of training examples and the initializer. Suppose we iterate over the training set following the order in the table until all the training samples are classified correctly, fill in the coefficients below:

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \text{---} \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \text{---} \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} + \text{---} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \text{---} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

This means that instead of *explicitly* representing  $w$  as a vector of feature weights, we can *implicitly* represent the decision rule with a vector  $v$  with one weight per example.

- (c) Now suppose  $x$  and  $w$  are  $D$ -dimensional, and we have  $N$  training examples. How many numbers do I need to represent  $w$  *explicitly*?
- (d) How many numbers do I need to represent  $w$  *implicitly*? (Assume that the initial value for  $w$  is public information so you do not need to consume any memory to store it.)
- (e) Write the update rule for the implicit representation if you pick the  $i^{\text{th}}$  training sample (use  $e_i$  to represent a vector whose  $i^{\text{th}}$  position is 1 and all the other positions are 0s):

$$v \leftarrow$$

- (f) Write the prediction rule for the implicit representation if the initial  $w$  is a zero vector. You can use  $v_i$  to represent the  $i^{\text{th}}$  value from  $v$  and  $x_i$  to represent the  $i^{\text{th}}$  training sample:

$$\text{pred}(x) =$$

- (g) When is it more space-efficient to use the implicit representation? (Your answer should be at most three words/symbols, and be expressed in terms of  $D$  and  $N$ .)

(h) Now suppose  $x$  is two-dimensional, and we introduce a feature transformation

$$f(x) = [x_0^2, x_1^2, \sqrt{2}x_0, \sqrt{2}x_1, \sqrt{2}x_0x_1, 1]$$

It is not too hard to show that for any vectors  $a$  and  $b$ ,

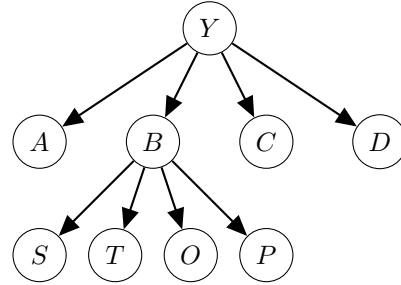
$$f(a) \cdot f(b) = (a \cdot b + 1)^2$$

How can we take advantage of this fact when working with the implicit representation? (Your answer should be 1 sentence.)

## Q2. A Nonconvolutional Nontrivial Network

You have a robotic friend MesutBot who has trouble passing Recaptchas (and Turing tests in general). MesutBot got a 99.99% on the last midterm because he could not determine which squares in the image contained stop signs. To help him ace the final, you decide to design a few classifiers using the below features.

- $A = 1$  if the image contains an octagon, else 0.
- $B = 1$  if the image contains the word STOP, else 0.
  - $S = 1$  if the image contains the letter S, else 0.
  - $T = 1$  if the image contains the letter T, else 0.
  - $O = 1$  if the image contains the letter O, else 0.
  - $P = 1$  if the image contains the letter P, else 0.
- $C = 1$  if the image is more than 50% red in color, else 0.
- $D = 1$  if the image contains a post, else 0.



(a) First, we use a Naive Bayes-inspired approach to determine which images have stop signs based on the features and Bayes Net above. We use the following features to predict  $Y = 1$  if the image has a stop sign anywhere, or  $Y = 0$  if it doesn't.

(i) Which expressions would a Naive Bayes model use to predict the label for  $B$  if given the values for features  $S = s, T = t, O = o, P = p$ ? Choose all valid expressions.

$b = \arg \max_b P(b)P(s|b)P(t|b)P(o|b)P(p|b)$

$b = \arg \max_b P(s|b)P(t|b)P(o|b)P(p|b)$

$b = \arg \max_b P(b|s, t, o, p)$

$b = \arg \max_b P(b, s, t, o, p)$

$b = \arg \max_b P(s, t, o, p|b)$

None

(ii) [Optional] Which expressions would we use to predict the label for  $Y$  with our Bayes Net above? Assume we are given all features except  $B$ . So  $A = a, S = s, T = t$ , etc. For the below choices, the underscore means we are dropping the value of that variable. So  $y, \_ = (0, 1)$  would mean  $y = 0$ .

$y, \_ = \arg \max_{y, b} P(y)P(a|y)P(b|y)P(c|y)P(d|y)P(s|b)P(t|b)P(o|b)P(p|b)$

$y, \_ = \arg \max_{y, b} P(s)P(t)P(o)P(p)P(a)P(b|s, t, o, p)P(c)P(d)P(y|a, b, c, d)$

First compute  $b' = \arg \max_b$  of the formula chosen in part (ii).

Then compute  $y = \arg \max_y P(y)P(a|y)P(b'|y)P(c|y)P(d|y)$

First compute  $b' = \arg \max_b$  of the formula chosen in part (ii).

Then compute  $y = \arg \max_y P(y|a, b', c, d)$

$y = \arg \max_y \sum_{b'} P(y)P(a|y)P(b'|y)P(c|y)P(d|y)P(s|b')P(t|b')P(o|b')P(p|b')$

None

(iii) [Optional] One day MesutBot got allergic from eating too many cashews. The incident broke his

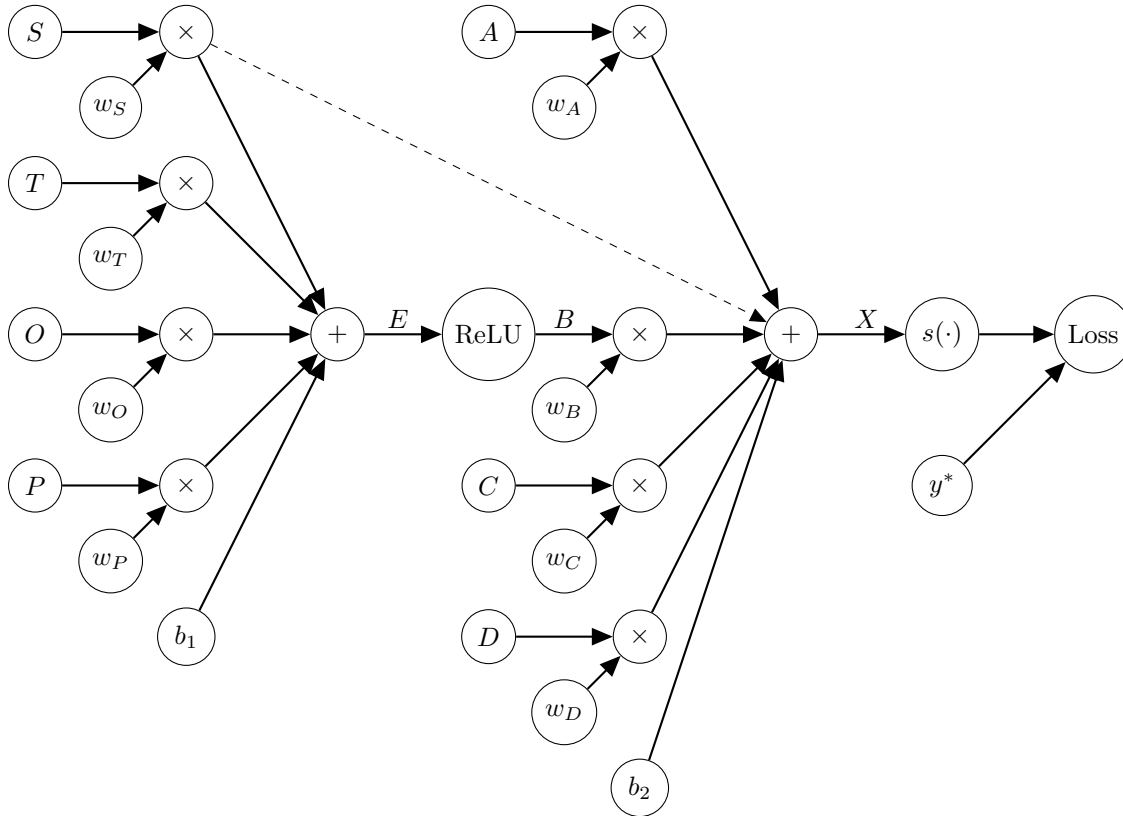
letter  $S$  detector, so that he no longer gets reliable  $S$  features. Now what expressions would we use to predict the label for  $Y$ ? Assume all features except  $B, S$  are given. So  $A = a, T = t, O = o$ , etc.

- $y = \arg \max_y P(y)P(a|y)P(c|y)P(d|y)$
- $y, \_-, \_ - = \arg \max_{y,b,s} P(y)P(a|y)P(b|y)P(c|y)P(d|y)P(s|b)P(t|b)P(o|b)P(p|b)$
- $y, \_ - = \arg \max_{y,s} P(y)P(a|y)P(b|y)P(c|y)P(d|y)P(s|b)P(t|b)P(o|b)P(p|b)$
- $y, \_ - = \arg \max_{y,b} P(y)P(a|y)P(b|y)P(c|y)P(d|y)P(t|b)P(o|b)P(p|b)$
- $y, \_ - = \arg \max_{y,b} P(y)P(a|y)P(b|y)P(c|y)P(d|y)P(s|b)P(t|b)P(o|b)P(p|b)$
- $y, \_ - = \arg \max_{y,b} P(y|a, b, c, d)$
- $y = \arg \max_y P(y)P(a|y)P(c|y)P(d|y) \sum_{b',s'} P(b'|y)P(s'|b')P(t|b')P(o|b')P(p|b')$
- None

(b) You decide to try to output a probability  $P(Y|features)$  of a stop sign being in the picture instead of a discrete  $\pm 1$  prediction. We denote this probability as  $P(Y|\vec{f}(x))$ . Which of the following functions return a **valid** probability distribution for  $P(Y = y|\vec{f}(x))$ ? Recall that  $y \in \{-1, 1\}$ .

- $\frac{e^{y \cdot \vec{w}^T \vec{f}(x)}}{e^{-y \cdot \vec{w}^T \vec{f}(x)} + e^{y \cdot \vec{w}^T \vec{f}(x)}}$
- $\frac{1}{2}$
- $\frac{0.5}{1 + e^{-\vec{w}^T \vec{f}(x)}}$
- $\frac{-1}{1 + e^{\vec{w}^T \vec{f}(x)}} + 1$
- None

You note that features are inputs into a neural network and the output is a label, so you modify the Bayes Net from above into a Neural Network computation graph. Recall the logistic function  $s(x) = \frac{1}{1+e^{-x}}$  has derivative  $\frac{\partial s(x)}{\partial x} = s(x)[1 - s(x)]$



(c) For this part, ignore the dashed edge when calculating the below.

(i) What is  $\frac{\partial Loss}{\partial w_A}$ ?

- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot A$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot A$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot 2A + 1$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot 2A$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot A + 1$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot A + 1$
- None

(ii) What is  $\frac{\partial Loss}{\partial w_s}$ ? Keep in mind we are still ignoring the dotted edge in this subpart.

- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot 2S + S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot 2S$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S + S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S + S$
- None

(d) [Optional] MesutBot is having trouble paying attention to the S feature because sometimes it gets zeroed out by the ReLU, so we connect it directly to the input of  $s(\cdot)$  via the dotted edge. For the below, treat the dotted edge as a regular edge in the neural net.

(i) Which of the following is equivalent to  $\frac{\partial Loss}{\partial w_A}$ ?

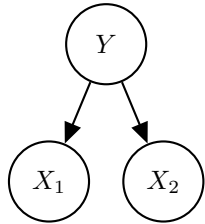
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot A$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot A$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot 2A + A$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot 2A$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot A + A$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot A + A$
- None

(ii) Which of the following is equivalent to  $\frac{\partial Loss}{\partial w_s}$ ? Keep in mind we are still treating the dotted edge as a regular edge.

- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot 2S + S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot 2S$
- $2(s(X) - y^*) \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S + S$
- $\frac{\partial Loss}{\partial s(X)} \cdot [s(X) \cdot (1 - s(X))] \cdot w_B \cdot \begin{pmatrix} 1 & E \geq 0 \\ 0 & E < 0 \end{pmatrix} \cdot S + S$
- None

### Q3. [Optional] Naive Bayes

You are given a naive bayes model, shown below, with label  $Y$  and features  $X_1$  and  $X_2$ . The conditional probabilities for the model are parameterized by  $p_1$ ,  $p_2$  and  $q$ .



$X_1$	$Y$	$P(X_1 Y)$
0	0	$p_1$
1	0	$1 - p_1$
0	1	$1 - p_1$
1	1	$p_1$

$X_2$	$Y$	$P(X_2 Y)$
0	0	$p_2$
1	0	$1 - p_2$
0	1	$1 - p_2$
1	1	$p_2$

$Y$	$P(Y)$
0	$1 - q$
1	$q$

Note that some of the parameters are shared (e.g.  $P(X_1 = 0|Y = 0) = P(X_1 = 1|Y = 1) = p_1$ ).

- (a) Given a new data point with  $X_1 = 1$  and  $X_2 = 1$ , what is the probability that this point has label  $Y = 1$ ? Express your answer in terms of the parameters  $p_1, p_2$  and  $q$  (you might not need all of them).

$$P(Y = 1|X_1 = 1, X_2 = 1) = \underline{\hspace{4cm}}$$

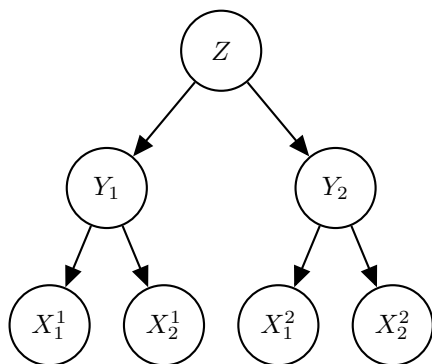
The model is trained with the following data:

sample number	1	2	3	4	5	6	7	8	9	10
$X_1$	0	0	1	0	1	0	1	0	1	1
$X_2$	0	0	0	0	0	0	0	1	0	0
$Y$	0	0	0	0	0	0	0	1	1	1

- (b) What are the maximum likelihood estimates for  $p_1, p_2$  and  $q$ ?

$$p_1 = \underline{\hspace{4cm}} \quad p_2 = \underline{\hspace{4cm}} \quad q = \underline{\hspace{4cm}}$$

- (c) For the next part, the model you are given is no longer simple naive bayes. Now there are two distinct label variables  $Y_1, Y_2$ , and there is a super label  $Z$  which conditions all of these labels, thus giving us this hierarchical naive bayes model. The conditional probabilities for the model are parametrized by  $p_1, p_2, q_0, q_1$  and  $r$ . Note that some of the parameters are shared as in the previous part.



$X_1^i$	$Y_i$	$P(X_1^i Y_i)$
0	0	$p_1$
1	0	$1 - p_1$
0	1	$1 - p_1$
1	1	$p_1$

$Y_i$	$Z$	$P(Y_i Z)$
0	0	$1 - q_0$
1	0	$q_0$
0	1	$1 - q_1$
1	1	$q_1$

$X_2^i$	$Y_i$	$P(X_2^i Y_i)$
0	0	$p_2$
1	0	$1 - p_2$
0	1	$1 - p_2$
1	1	$p_2$

$Z$	$P(Z)$
0	$1 - r$
1	$r$

- (i) What is the probability that  $Z = 1$  given the partial data point  $X_1^2 = 1, X_2^2 = 1, Y_1 = 1$ ? Simplify your answer as much as possible and express it in terms of the parameters  $p_1, p_2, q_0, q_1$  and  $r$  (you might not need all of them).

$$P(Z = 1 | X_1^2 = 1, X_2^2 = 1, Y_1 = 1) = \underline{\hspace{4cm}}$$

- (ii) Now we are given a partial data point with  $X_1^2 = 1, X_2^2 = 1, Y_1 = 1$ . What is the probability that  $Y_2 = 1$ . Simplify your answer as much as possible and express it in terms of the parameters  $p_1, p_2, q_0, q_1$  and  $r$  (you might not need all of them).

$$P(Y_2 = 1 | X_1^2 = 1, X_2^2 = 1, Y_1 = 1) = \underline{\hspace{4cm}}$$

- (d) Let  $L_{nb}$  and  $L_{hnb}$  be the likelihood of the training data under the naive bayes model and the hierarchical naive bayes model, respectively. Assume each of the models use their respective maximum likelihood parameters. Which of the following properties are guaranteed to be true?

- $L_{nb} \leq L_{hnb}$
- $L_{nb} \geq L_{hnb}$
- $L_{nb} = L_{hnb}$
- Insufficient information, the above relationships rely on the particular training data.
- None of the above.