

CS221 Problem Workout Problems Solutions

Week 2

1) Problem 1: Non-linear features

Consider the following two training datasets of (x, y) pairs:

- $\mathcal{D}_1 = \{(-1, +1), (0, -1), (1, +1)\}$.
- $\mathcal{D}_2 = \{(-1, -1), (0, +1), (1, -1)\}$.

Observe that neither dataset is linearly separable if we use $\phi(x) = [x]$, so let's fix that.

Define a two-dimensional feature function $\phi(x)$ such that:

- There exists a weight vector \mathbf{w}_1 that classifies \mathcal{D}_1 perfectly (meaning that $\mathbf{w}_1 \cdot \phi(x) > 0$ if x is labeled $+1$ and $\mathbf{w}_1 \cdot \phi(x) < 0$ if x is labeled -1); and
- There exists a weight vector \mathbf{w}_2 that classifies \mathcal{D}_2 perfectly.

Note that the weight vectors can be different for the two datasets, but the features $\phi(x)$ must be the same.

Solution One option is $\phi(x) = [1, x^2]$, and using $\mathbf{w}_1 = [-1, 2]$ and $\mathbf{w}_2 = [1, -2]$.

Then in \mathcal{D}_1 :

- For $x = -1$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 1] = 1 > 0$
- For $x = 0$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 0] = -1 < 0$
- For $x = 1$, $\mathbf{w}_1 \cdot \phi(x) = [-1, 2] \cdot [1, 1] = 1 > 0$

In \mathcal{D}_2 :

- For $x = -1$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 1] = -1 < 0$
- For $x = 0$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 0] = 1 > 0$
- For $x = 1$, $\mathbf{w}_2 \cdot \phi(x) = [1, -2] \cdot [1, 1] = -1 < 0$

Note that there are many options that work, so long as -1 and 1 are separated from 0 .

Some additional food for thought: Is every dataset linearly separable in some feature space? In other words, given pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, can we find a feature extractor ϕ such that we can perfectly classify $(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_n), y_n)$ for some linear model \mathbf{w} ? If so, is this a good feature extractor to use?

Solution In theory, yes we can. If we assume that our inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$ are distinct, then we can construct a feature map $\phi : x_i \mapsto y_i$ for $i = 1, \dots, n$. By setting $\mathbf{w}^* = [1]$, it's clear that

$$y_i \mathbf{w}^* \cdot \phi(x_i) = y_i * y_i = 1 > 0, \quad i = 1, \dots, n, \quad (1)$$

so \mathbf{w}^* correctly classifies all the points in the dataset.

Hopefully, it's clear that this is a poor choice of feature map. For one, this feature extractor is undefined for any points outside of the training set! But even more broadly, this process is not at all *generalizeable*. We are essentially just memorizing our dataset instead of learning patterns and structures within the data that will allow us to accurately predict new points in the future. While minimizing training loss is an important part of the machine learning process (the aforementioned procedure gives you zero training loss!), it does not guarantee you good performance in the future.

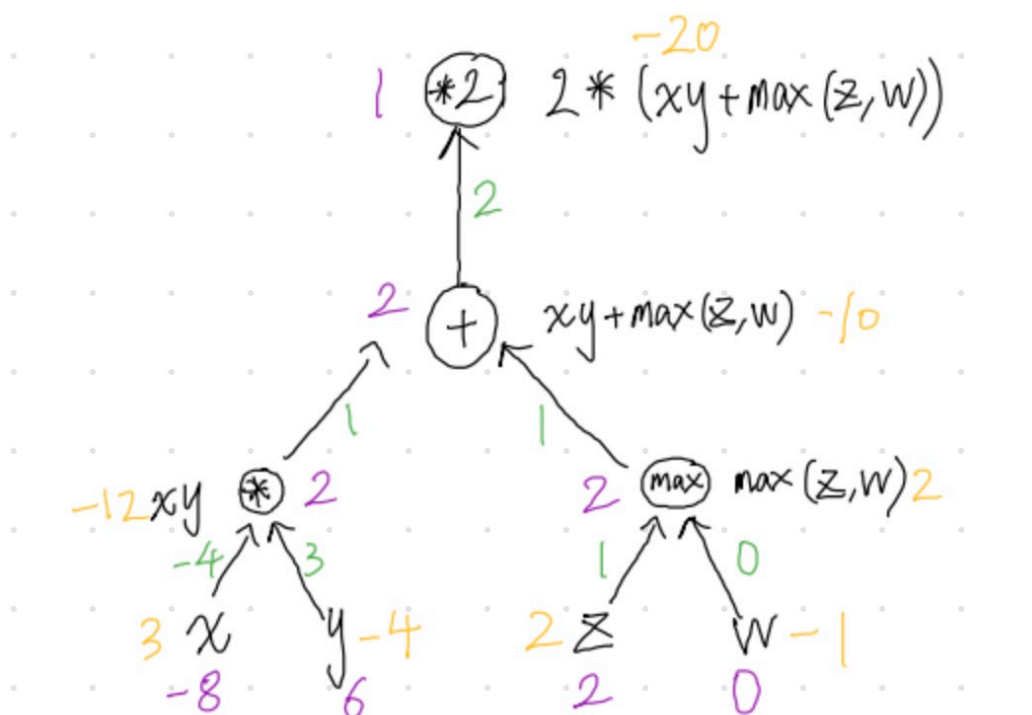
2) Problem 2: Backpropagation

Consider the following function

$$\text{Loss}(x, y, z, w) = 2(xy + \max\{w, z\})$$

Run the backpropagation algorithm to compute the four gradients (each with respect to one of the individual variables) at $x = 3$, $y = -4$, $z = 2$ and $w = -1$. Use the following nodes: addition, multiplication, max, multiplication by a constant.

Solution When calculating the gradients, we run backpropagation from the root node to the leaves nodes. As shown on the computation graph below, the purple values are the gradients of Loss with respect to each node.



3) Problem 3: K-means

Consider doing ordinary K -means clustering with $K = 2$ clusters on the following set of 3 one-dimensional points:

$$\{-2, 0, 10\}. \tag{2}$$

Recall that K -means can get stuck in local optima. Describe the precise conditions on the initialization $\mu_1 \in \mathbb{R}$ and $\mu_2 \in \mathbb{R}$ such that running K -means will yield the global optimum of the objective function. Notes:

- Assume that $\mu_1 < \mu_2$.
- Assume that if in step 1 of K -means, no points are assigned to some cluster j , then in step 2, that centroid μ_j is set to ∞ .
- Hint: try running K -means from various initializations μ_1, μ_2 to get some intuition; for example, if we initialize $\mu_1 = 1$ and $\mu_2 = 9$, then we converge to $\mu_1 = -1$ and $\mu_2 = 10$.

Solution The objective is minimized for $\mu_1 = -1$ and $\mu_2 = 10$. First, note that if all three points end up in one cluster, K -means definitely fails to recover the global optimum. Therefore, -2 must be assigned to the first cluster, and 10 must be assigned to the second cluster. 0 can be assigned to either: If 0 is assigned to cluster 1, then we're done. If it is assigned to cluster 2, then we have $\mu_1 = -2, \mu_2 = 5$; in the next iteration, 0 will be assigned to cluster 1 since it's closer. Therefore, the condition on the initialization written formally is $|-2 - \mu_1| < |-2 - \mu_2|$ and $|10 - \mu_1| > |10 - \mu_2|$.

4) Problem 4: Non-linear decision boundaries

Suppose we are performing classification where the input points are of the form $(x_1, x_2) \in \mathbb{R}^2$. We can choose any subset of the following set of features:

$$\mathcal{F} = \left\{ x_1^2, x_2^2, x_1x_2, x_1, x_2, \frac{1}{x_1}, \frac{1}{x_2}, 1, \mathbf{1}[x_1 \geq 0], \mathbf{1}[x_2 \geq 0] \right\} \quad (3)$$

For each subset of features $F \subseteq \mathcal{F}$, let $D(F)$ be the set of all decision boundaries corresponding to linear classifiers that use features F .

For each of the following sets of decision boundaries E , provide the minimal F such that $D(F) \supseteq E$. If no such F exists, write ‘none’.

- E is all lines:

 (4)

- E is all circles centered at the origin:

 (5)

- E is all circles:

 (6)

- E is all axis-aligned rectangles:

 (7)

- E is all axis-aligned rectangles whose lower-right corner is at $(0, 0)$:

 (8)

Solution

- Lines: $x_1, x_2, 1$ ($ax_1 + bx_2 + c = 0$)
- Circles centered at the origin: $x_1^2, x_2^2, 1$ ($x_1^2 + x_2^2 = r^2$)
- Circles centered anywhere in the plane: $x_1^2, x_2^2, x_1, x_2, 1$ ($(x_1 - a)^2 + (x_2 - b)^2 = r^2$)
- Axis aligned rectangles: not possible (need features of the form $\mathbf{1}[x_1 \leq a]$)
- Axis aligned rectangles with lower right corner at $(0, 0)$: not possible

5) Problem 5: k-means Initialization

Consider performing k -means clustering with $k = 2$ clusters on the following set of 3 one-dimensional data points:

$$\{-8, 0, 14\} \quad (9)$$

- (a) What is the globally optimal clustering and the associated reconstruction loss for the above dataset? Your answer should specify the following:
- Centroids μ_1 and μ_2 for clusters 1 and 2 respectively.
 - Assignments of points $x_1 = -8, x_2 = 0$ and $x_3 = 14$ to clusters.
 - Reconstruction loss (sum of squared distances) of this globally optimal clustering.

Solution $\mu_1, \mu_2 = -4, 14$

Assignments of points: $x_1 = -8$ and $x_2 = 0$ to cluster centered at -4 . $x_3 = 14$ to cluster centered at 14 .

Reconstruction loss: $(4)^2 + (4)^2 + 0^2 = \mathbf{32}$

- (b) Consider the following algorithm (henceforth referred to as **k-means++**) to initialize the k-means centroids intelligently rather than at random:¹
- Choose one centroid μ_1 , uniformly at random from among the data points.
 - For each data point x , compute $D^*(x)$, the Euclidean distance between x and the nearest centroid that has already been chosen. In other words $D^*(x) = \min_{\mu_i} D(x, \mu_i)$
 - Choose one new data point at random as a new centroid, using a weighted probability distribution where a point x is chosen with probability proportional to $D^*(x)$. In other words,

$$P[x_i \text{ is chosen as } \mu] = \frac{D^*(x_i)}{\sum_{j \in J} D^*(x_j)}$$

where J is the set of all points that haven't been chosen as a centroid yet.

- Repeat Steps II and III until k centroids have been chosen.
- Now that the initial centroids have been chosen, proceed using standard k -means clustering.

Now consider the dataset specified in equation (9) with $k = 2$.

Fill in the table given below. For each row, you should specify (1) the probability of the centroids being chosen in this order with the k-means++ algorithm, (2) the **final** clusters that would result from these initial centroids (at the conclusion of k-means clustering), and (3) whether or not this clustering is optimal. You may leave your probabilities as unsimplified fractions and/or sums of fractions,

¹Adapted from the following work: Arthur, David, and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Stanford, 2006.

e.g. $\frac{1}{1+2+3}$.

Note that, in the table on the following page, μ_1 is the first centroid chosen and μ_2 is the second centroid chosen.

μ_1	μ_2	Probability	Final Clusters	Optimal (Y/N)
-8	0			
-8	14			
0	-8			
0	14			
14	-8			
14	0			

What is the overall probability that using the k-means++ algorithm returns the globally optimal k -means cluster for this dataset with $k = 2$? (You may leave your answer as an unsimplified fraction or a sum of fractions.)

Solution

μ_1	μ_2	Probability	Final Clusters	Optimal (Y/N)
-8	0	$\frac{4}{45}$	$\{-8\}$ at centroid -8 and $\{0, 14\}$ at centroid 7	No
-8	14	$\frac{11}{45}$	$\{-8, 0\}$ at centroid -4 and $\{14\}$ at centroid 14	Yes
0	-8	$\frac{4}{33}$	$\{-8\}$ at centroid -8 and $\{0, 14\}$ at centroid 7	No
0	14	$\frac{7}{33}$	$\{-8, 0\}$ at centroid -4 and $\{14\}$ at centroid 14	Yes
14	-8	$\frac{11}{54}$	$\{-8, 0\}$ at centroid -4 and $\{14\}$ at centroid 14	Yes
14	0	$\frac{7}{54}$	$\{-8, 0\}$ at centroid -4 and $\{14\}$ at centroid 14	Yes

Probability of obtaining the optimal clustering = $\frac{11}{45} + \frac{7}{33} + \frac{11}{54} + \frac{7}{54} = \frac{391}{495} = 0.7899$

Detailed explanation of solution:

- Since the first centroid is assigned randomly,
 $Pr[\mu_1 = -8] = Pr[\mu_1 = 0] = Pr[\mu_1 = 14] = \frac{1}{3}$
- If $\mu_1 = -8$: $Pr[\mu_2 = 0|\mu_1 = -8] = \frac{8}{8+22} = \frac{4}{15}$
 $Pr[\mu_2 = 14|\mu_1 = -8] = \frac{22}{8+22} = \frac{11}{15}$
- if $\mu_1 = 0$: $Pr[\mu_2 = -8|\mu_1 = 0] = \frac{8}{8+14} = \frac{4}{11}$
 $Pr[\mu_2 = 14|\mu_1 = 0] = \frac{14}{8+14} = \frac{7}{11}$
- if $\mu_1 = 14$: $Pr[\mu_2 = -8|\mu_1 = 14] = \frac{22}{14+22} = \frac{11}{18}$
 $Pr[\mu_2 = 0|\mu_1 = 14] = \frac{14}{14+22} = \frac{7}{18}$

For initial centroids $\{-8, 0\}$, $x_3 = 14$ is assigned to cluster centered at $\mu = 0$ at iteration 1. Centroids are then updated to $\mu_1 = -8, \mu_2 = 7$. There is no change in assignments at the second iteration which leads to a final clustering of $\{-8\}$ at centroid $\mu_1 = -8$ and $\{0, 14\}$ at centroid $\mu_2 = 7$. This is suboptimal (optimal clustering is described in (a)).

For initial centroids either $\{-8, 14\}$ or $\{0, 14\}$, the optimal clustering described in (a) is obtained.

- (c) For part (c), consider an unspecified dataset in \mathbb{R}^d with n datapoints, and with $k \in \{2, 3, \dots, n\}$.

An alternative initialization method discussed in class is "multiple random initialization": (1) run multiple k-means clustering runs with random initializations of centroids on each run, and (2) choose the clustering with the lowest reconstruction loss.

Under what circumstances might you choose: (i) k-means++ initialization over multiple random initialization, (ii) multiple random initialization over k-means++ initialization? **Describe no more than one circumstance in each case.**

The circumstances you describe can include: the problem or application that clustering is being applied to, specific prior knowledge of the dataset, compute constraints, and/or the importance of getting the globally optimal clustering.

Solution We would choose k-means++ over multiple random initialization when compute is constrained, because multiple runs of k-means clustering is compute intensive.

We would choose multiple random initialization over k-means++ when compute costs are not all that important and accuracy is very highly valued (e.g., clustering to determine the likelihood of a malignant tumour or a similar medical application where getting an accurate clustering is very important). This is because a faulty k-means++ initialization - like any other one-time initialization method - could yield a poor (locally optimum but not globally optimum) clustering.

6) Problem 6: Two views

Alice and Bob are trying to predict movie ratings. They cast the problem as a standard regression problem, where $x \in \mathcal{X}$ contains information about the movie and $y \in \mathbb{R}$ is the real-valued movie rating.

To split up the effort, Alice will create a feature extractor $A : \mathcal{X} \rightarrow \mathbb{R}^d$ using information from the text of the movie reviews and Bob will create a feature extractor $B : \mathcal{X} \rightarrow \mathbb{R}^d$ from the movie metadata (genre, cast, number of reviews, etc.) These features will be used to do linear regression as follows: If $\mathbf{u} \in \mathbb{R}^d$ is a weight vector associated with A and $\mathbf{v} \in \mathbb{R}^d$ is a weight vector associated with B , then we can define the resulting predictor as:

$$f_{\mathbf{u}, \mathbf{v}}(x) \stackrel{\text{def}}{=} \mathbf{u} \cdot A(x) + \mathbf{v} \cdot B(x). \quad (10)$$

We are interested in the squared loss:

$$\text{Loss}(x, y, \mathbf{u}, \mathbf{v}) \stackrel{\text{def}}{=} (f_{\mathbf{u}, \mathbf{v}}(x) - y)^2. \quad (11)$$

Let $\mathcal{D}_{\text{train}}$ be the training set of (x, y) pairs, on which we can define the training loss:

$$\text{TrainLoss}(\mathbf{u}, \mathbf{v}) \stackrel{\text{def}}{=} \sum_{(x, y) \in \mathcal{D}_{\text{train}}} \text{Loss}(x, y, \mathbf{u}, \mathbf{v}). \quad (12)$$

(a) Compute the gradient of the training loss with respect to Alice's weight vector \mathbf{u} .

$$\nabla_{\mathbf{u}} \text{TrainLoss}(\mathbf{u}, \mathbf{v}) = \underline{\hspace{15cm}}$$

Solution

$$\nabla_{\mathbf{u}} \text{TrainLoss}(\mathbf{u}, \mathbf{v}) = 2 \sum_{(x, y) \in \mathcal{D}_{\text{train}}} (f_{\mathbf{u}, \mathbf{v}} - y) A(x) \quad (13)$$

$$= 2 \sum_{(x, y) \in \mathcal{D}_{\text{train}}} (\mathbf{u} \cdot A(x) + \mathbf{v} \cdot B(x) - y) A(x). \quad (14)$$

Suppose Alice and Bob have found a current pair of weight vectors (\mathbf{u}, \mathbf{v}) where the gradient with respect to \mathbf{u} is zero: $\nabla_{\mathbf{u}} \text{TrainLoss}(\mathbf{u}, \mathbf{v}) = 0$. Mark each of the following statements as true or false (no justification is required).

i. Even if Bob updates his weight vector \mathbf{v} to \mathbf{v}' , \mathbf{u} is still optimal for the new \mathbf{v}' ; that is, $\nabla_{\mathbf{u}} \text{TrainLoss}(\mathbf{u}, \mathbf{v}') = 0$ for any \mathbf{v}' .

—

True / False

ii. The gradient of the loss on each example is also zero: $\nabla_{\mathbf{u}} \text{Loss}(x, y, \mathbf{u}, \mathbf{v}) = 0$ for each $x, y \in \mathcal{D}_{\text{train}}$.

—

True / False

Solution Neither statement is true:

- i. Alice's weight vector \mathbf{u} is only optimal with respect to the particular \mathbf{v} . As another example, think of alternating minimization (k -means), where the centroids are only optimal *given* the current assignments, not for all assignments; otherwise, there would be no point in iterating.
 - ii. The loss gradient is only zero *on average*, definitely not for each example.
- (b) Alice and Bob got into an argument and now aren't on speaking terms. Each of them therefore trains his/her weight vector separately: Alice computes

$$\mathbf{u}_A = \arg \min_{\mathbf{u} \in \mathbb{R}^d} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} (\mathbf{u} \cdot A(x) - y)^2$$

and Bob computes

$$\mathbf{v}_B = \arg \min_{\mathbf{v} \in \mathbb{R}^d} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} (\mathbf{v} \cdot B(x) - y)^2.$$

Just for reference, if they had worked together, then they would have gotten:

$$(\mathbf{u}_C, \mathbf{v}_C) = \arg \min_{\mathbf{u}, \mathbf{v}} \text{TrainLoss}(\mathbf{u}, \mathbf{v}).$$

Let L_A, L_B, L_C be the training losses of the following predictors:

- (L_A) Alice works alone: $x \mapsto \mathbf{u}_A \cdot A(x)$
- (L_B) Bob works alone: $x \mapsto \mathbf{v}_B \cdot B(x)$
- (L_C) They work together from the beginning: $x \mapsto \mathbf{u}_C \cdot A(x) + \mathbf{v}_C \cdot B(x)$

For each statement, mark it as necessarily true, necessarily false, or neither:

- i. $L_A \leq L_B$ True / False / Neither
- ii. $L_C \leq L_A$ True / False / Neither
- iii. $L_C \leq L_B$ True / False / Neither

Solution We can view L_A as the resulting training loss from optimizing over all weight vectors $(\mathbf{u}, 0)$, whereas L_C optimizes over all (\mathbf{u}, \mathbf{v}) ; therefore $L_C \leq L_A$. The same logic applies to conclude $L_C \leq L_B$. However, we cannot make a comparison between L_A and L_B ; it's possible that either Alice or Bob's features get lower training error.

Now, let L'_A, L'_B, L'_C be the corresponding losses on the *test set*. For each statement, mark it as necessarily true, necessarily false, or neither:

- i. $L'_A \leq L'_B$ True / False / Neither
- ii. $L'_C \leq L'_A$ True / False / Neither
- iii. $L'_C \leq L'_B$ True / False / Neither

Solution Neither for all three! We might hope that $L'_C \leq L'_A$ and $L'_C \leq L'_B$, but we might have overfit on the training set, and the same trends might not hold.

7) Problem 7: Movie Genre(s) (20 points)

You are interested in classifying the genre(s) of a movie given its plot summary. For simplicity, assume there are only three genres in the universe: action, romance, and comedy. You decide that you will model this as three binary classification problems, one for each genre.

For each plot summary x , you represent the corresponding movie M 's genre(s) with $y \in \{0, 1\}^3$, where

$$\begin{aligned}y_1 &= \mathbf{1}[M \text{ is an action movie}], \\y_2 &= \mathbf{1}[M \text{ is a romance movie}], \\y_3 &= \mathbf{1}[M \text{ is a comedy movie}].\end{aligned}$$

a. (12 points) Sharing is Caring

You decide to use a feature extractor ϕ that maps each word in the English vocabulary to the number of occurrences of that word in the plot summary. Assume our English vocabulary D contains only the top 10,000 words in English (by frequency). Any out-of-vocabulary word is ignored.

Recall that the logistic function $\sigma(z) = (1 + e^{-z})^{-1}$ takes in a real number and outputs a probability in $(0, 1)$. After some research, you learned that the logistic function is often used with the *logistic loss* function for binary classification. Given a label $y \in \{0, 1\}$ and predicted probability $p \in [0, 1]$,

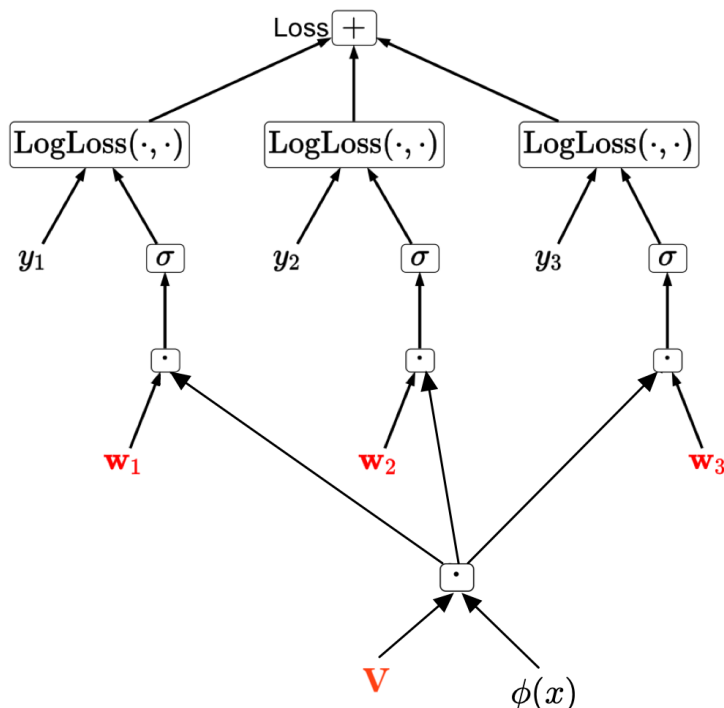
$$\text{LogLoss}(y, p) = -(y \log p + (1 - y) \log(1 - p)).$$

You decide to use these shiny new tools to solve your problem.

- (i) [1 point] One way to solve your problem is to build three binary classifiers. For a plot summary x and genre $k \in \{1, 2, 3\}$, the predicted probability that $y_k = 1$ is $p_k = \sigma(\mathbf{w}_k \cdot \phi(x))$. How many parameters will there be in the three classifiers combined? No justification required.

Solution 30,000. Each classifier has 10,000 parameters, so there are $3 \times 10,000 = 30,000$ parameters.

- (ii) [6 points] You've heard that having too many parameters can lead to bad generalization and want to find a way to reduce it. Since the three binary classification tasks are all about predicting movie genre and are thus related in nature, why not let the three classifiers share parameters? You come up with the following computation graph:



where \mathbf{V} is a trainable weight matrix shared among all three classifiers such that $\mathbf{V}\phi(x) \in \mathbb{R}^2$ for all x .

Now, for any $k \in \{1, 2, 3\}$, the predicted probability that $y_k = 1$ is $p_k = \sigma(\mathbf{w}_k \cdot \mathbf{V}\phi(x))$.

- In this setup, how many parameters are there in the three classifiers combined? No justification required.

Solution 20,006. Since $\mathbf{V}\phi(x) \in \mathbb{R}^2$, \mathbf{V} must be a $2 \times |D|$ matrix, which has $2 \times 10,000$ parameters. For the dimensionality to match, w_1, w_2, w_3 are all 2-dimensional vectors. So there are 20,006 parameters in total.

- The pointwise loss function $\text{Loss}(x, y, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3)$ is defined as follows:

$$\begin{aligned}\text{Loss}(x, y, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) &= \sum_{k=1}^3 \text{LogLoss}(y_k, p_k) \\ &= \sum_{k=1}^3 -\left(y_k \log \sigma(\mathbf{w}_k \cdot \mathbf{V}\phi(x)) + (1 - y_k) \log (1 - \sigma(\mathbf{w}_k \cdot \mathbf{V}\phi(x)))\right).\end{aligned}$$

Derive the gradient of the pointwise loss with respect to \mathbf{w}_1 , i.e., $\nabla_{\mathbf{w}_1} \text{Loss}(x, y, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3)$. You can use p_1 in your final expression. Show your work.

Hint: Feel free to use $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ directly without showing the intermediate steps.

Solution

$$\nabla_{\mathbf{w}_1} \text{Loss}(x, y, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) = (p_1 - y_1) \mathbf{V}\phi(x)$$

Derivation:

$$\nabla_{\mathbf{w}_1} \text{Loss}(x, y, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) \tag{1}$$

$$= \nabla_{\mathbf{w}_1} \text{LogLoss}(y_1, p_1)$$

$$= \nabla_{\mathbf{w}_1} - \left(y_1 \log \sigma(\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) + (1 - y_1) \log (1 - \sigma(\mathbf{w}_1 \cdot \mathbf{V}\phi(x))) \right) \tag{2}$$

$$= - \left(y_1 \frac{1}{p_1} \sigma'(\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) \nabla_{\mathbf{w}_1} (\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) \right.$$

$$\left. + (1 - y_1) \frac{1}{1 - p_1} \left(- \sigma'(\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) \right) \nabla_{\mathbf{w}_1} (\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) \right) \tag{3}$$

$$= - \left(y_1 \frac{p_1(1 - p_1)}{p_1} \nabla_{\mathbf{w}_1} (\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) - (1 - y_1) \frac{p_1(1 - p_1)}{1 - p_1} \nabla_{\mathbf{w}_1} (\mathbf{w}_1 \cdot \mathbf{V}\phi(x)) \right) \tag{4}$$

$$= - (y_1(1 - p_1) \mathbf{V}\phi(x) - (1 - y_1)p_1 \mathbf{V}\phi(x)) \tag{5}$$

$$= - \mathbf{V}\phi(x)(y_1 - p_1) = (p_1 - y_1) \mathbf{V}\phi(x) \tag{6}$$

- (iii) [5 points] You plan to use regular gradient descent to train your classifiers. Suppose you have a dataset of N points $(x^{(i)}, y^{(i)})$. You define the overall training loss as follows:

$$\text{TrainLoss}(\mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) = \frac{1}{N} \sum_{i=1}^N \text{Loss}(x^{(i)}, y^{(i)}, \mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3).$$

Below is your gradient descent algorithm:

Algorithm 1: Gradient Descent

```
1: Randomly shuffle the training data
2: Initialize  $\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}_3 = \mathbf{0}$ , initialize  $\mathbf{V}$  with random non-zero matrix
3: for  $t = 1, 2, \dots, T$  do
4:   // calculate gradients
5:   for  $k = 1, 2, 3$  do
6:      $\mathbf{v}_k \leftarrow \nabla_{\mathbf{w}_k} \text{TrainLoss}(\mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3)$  // get gradient w.r.t  $\mathbf{w}_k$ 
7:   end for
8:    $\mathbf{U} \leftarrow \nabla_{\mathbf{V}} \text{TrainLoss}(\mathbf{V}, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3)$  // get gradient w.r.t  $\mathbf{V}$ 
9:
10:  // update weights
11:  for  $k = 1, 2, 3$  do
12:     $\mathbf{w}_k \leftarrow \mathbf{w}_k - \eta \mathbf{v}_k$  // update  $\mathbf{w}_k$ 
13:  end for
14:   $\mathbf{V} \leftarrow \mathbf{V} - \eta \mathbf{U}$  // update  $\mathbf{V}$ 
15: end for
```

Consider a training dataset \mathcal{D} with two datapoints $(x^{(1)}, y^{(1)})$ and $(x^{(2)}, y^{(2)})$. You run gradient descent on \mathcal{D} and initialize $\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}_3 = \mathbf{0}$ per Line 2 in the algorithm described above. Upon initializing \mathbf{V} randomly, we see that

$$\begin{aligned}\mathbf{V}\phi(x^{(1)}) &= [5, -2]^T, & y^{(1)} &= [0, 1, 1]^T \\ \text{and } \mathbf{V}\phi(x^{(2)}) &= [-1, 2]^T, & y^{(2)} &= [1, 0, 1]^T.\end{aligned}$$

You run gradient descent on this dataset with $\eta = 0.1$. What is the value of \mathbf{w}_1 at the end of one iteration? Show your work.

Note: We are expecting the value of \mathbf{w}_1 only, **not** \mathbf{w}_2 or \mathbf{w}_3 .

Solution For all k ,

$$\begin{aligned}\mathbf{v}_k &= \nabla_{\mathbf{w}_k} \text{TrainLoss}(\mathbf{V}_0, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) = \frac{1}{2} \sum_{i=1}^2 \nabla_{\mathbf{w}_k} \text{Loss}(x^{(i)}, y^{(i)}, \mathbf{V}_0, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3) \\ &= \frac{1}{2} ((p_k^{(1)} - y_k^{(1)}) \mathbf{V}_0 \phi(x^{(1)}) + (p_k^{(2)} - y_k^{(2)}) \mathbf{V}_0 \phi(x^{(2)})).\end{aligned}$$

Note that for all k , since $\mathbf{w}_k = \mathbf{0}$, $p_k^{(1)} = p_k^{(2)} = \sigma(\mathbf{0} \cdot \mathbf{V}_0 \phi(x)) = \sigma(0) = (1 + e^0)^{-1} = 0.5$. And since $\eta = 0.1$, after the first iteration, $\mathbf{w}_k = \mathbf{0} - \eta \mathbf{v}_k = -0.1 \mathbf{v}_k$.

We see that

$$\mathbf{v}_1 = \frac{1}{2} ((0.5 - 0)[5, -2]^T + (0.5 - 1)[-1, 2]^T) = [1.5, -1]^T.$$

Therefore,

$$\mathbf{w}_1 = [-0.15, 0.1]^T.$$

Misses that the teaching staff saw include:

- Not averaging the gradient update across the dataset (i.e., dropping the $\frac{1}{N}$ term)
- Calculating gradient but forgetting to update
- Mixing up the subscript used in \mathbf{w}_1 with the superscripts used in the datapoints

b. (8 points) Less is More

You tried running gradient descent, but since \mathbf{V} is a large matrix, your computer does not have enough memory to store the gradient with respect to \mathbf{V} , so you couldn't train your model properly.

Luckily, your friend Alice has worked on a similar problem (e.g., classifying fiction genres given synopsis) and trained a model with the same architecture as yours. You ask Alice to share with you the weights of her *trained* classifier, which includes a weight matrix \mathbf{V}_0 that has exactly the same shape as \mathbf{V} .

Since Alice's problem is similar to yours, you believe the solution should be similar as well. Hence, you decide to initialize \mathbf{V} with \mathbf{V}_0 and never update it in order to get around the memory constraint. In other words, **you run gradient descent only on \mathbf{w}_1 , \mathbf{w}_2 , and \mathbf{w}_3 .**

You tell Alice about your new gradient descent algorithm, and the idea to make a new memory-efficient classifier that shares \mathbf{V}_0 .

Alice responds: "Actually, what you have here isn't a new kind of classifier. It's actually three classifiers with a shared feature extractor." Though shocked at first, after some thought, you agree with her as well.

(i) [2 points] Show that Alice is right by filling in the blank below:

Given a plot summary x , for each genre $k \in \{1, 2, 3\}$,
the predicted probability that $y_k = 1$ is $p_k = \sigma(\mathbf{w}_k \cdot \tau(x))$,
where $\tau(x) = \underline{\hspace{2cm}}$ is the shared feature extractor.

No justification required.

Solution $\tau(x) = \mathbf{V}_0\phi(x)$.

(ii) [1 point] In this setup, how many parameters are there in the three classifiers combined?
No justification required.

Solution 6. The only parameters are \mathbf{w}_1 , \mathbf{w}_2 , and \mathbf{w}_3 , which are all 2-dimensional.

(iii) [5 points] Between the following:

- (A) three classifiers that share $\phi(x)$ as feature extractor (i.e., your idea in a(i))
- (B) three classifiers that share $\tau(x)$ as feature extractor

Which setup has higher approximation error? Which one has higher estimation error?
Justify your answer with 1-2 sentences.

Solution (B) has higher approximation error. (A) has higher estimation error. This is because (A) has a much larger hypothesis class. Each of the three weight vectors in (A) has much higher dimensionality than each of the three weight vectors in (B). More specifically, as we saw in 1a(i), each weight vector in (A) is 10,000-dimensional; as we saw in 1b(ii), each weight vector in (B) is 2-dimensional.

A commonly seen incorrect justification relates estimation/approximation error to the accuracy/inaccuracy of having a pre-trained weight matrix rather than the number of parameters.