# Questions

1. Consider the logistic regression setting taught in the lecture. Here, the learning task is a <u>binary classification</u> problem [[*Fill this blank with either "regression" or "multi-class classification" or "binary classification". 1/4 mark*]]. Let $\phi : \mathcal{X} \mapsto \mathbb{R}^n$ be a given feature map. The model employed in logistic regression is the <u>linear model</u> [[*Fill this blank with the appropriate proper noun. 1/4 mark*]]. The mathematical definition of this model is given by the expression:

$$\mathcal{L}_{n,\phi} = \left\{ f \mid \exists w \in \mathbb{R}^n \ni f(x) = w^T \phi(x) \ast x \in \mathcal{X} \right\}.$$
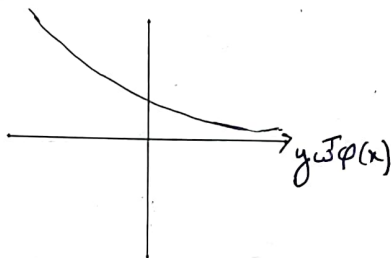
[[*In the above blank, use 'w' to denote the parameter of this model. 1/2 mark*]]

The mathematical expression for the loss function used here is:

$$l(w, x, y) \equiv \underline{Log\left(1 + e^{-y w^T \phi(x)}\right)}.$$

[1 mark]

This loss function can be visualized using the plot below:

[[*Fill the plot appropriately to roughly depict the graph of logistic loss. Clearly label the axes, without which no marks will be given. 1/2 mark*]].

If $p^*$ is the underlying (unknown) likelihood relating the inputs and labels, then, the Bayes optimal, restricted to the functions in the model, is given by the mathematical expression:

$$f^* = \underset{w \in \mathbb{R}^n}{\operatorname{argmin}} \ \underset{x, y \sim p^*}{\mathbb{E}}\left[ Log\left(1 + e^{-y w^T \phi(x)}\right) \right]$$

2

---

*Handwritten margin notes:*

You may also write "linear classifier" model

$\{f \mid \exists w \in \mathbb{R}^n \ni$
$f(x) = \operatorname{sign}(w^T \phi(x))$
$\ast x \in \mathcal{X}\}$ is also

OK if you wrote linear classifiers model

[[*No marks will be given if you write general expressions for Bayes optimal. You need to write the specific expression for the logistic regression set-up. 1 mark*]]. Let the training data be $D \equiv \{(x_1, y_1), \ldots, (x_m, y_m)\}$. The name of the important assumption that relates $p^*$ to $D$ is <u>Supervised Learning</u>. [[*Fill this blank with the appropriate proper noun. 1/2 mark*]]. This assumption formally means the following:

*Any other equivalent sentence is fine.* ← $\boxed{D \text{ is net } \textit{i.i.d} \text{ samples from } p^*}$

[1/2 mark]

The ERM problem in this case is the following mathematical optimization problem:

*argmin is also fine.* ← $$\min_{w \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 + e^{-y_i \, w^T \varphi(x_i)}\right)$$

[1/2 mark]

If the ERM solution is denoted by $\hat{w}_m$, then the label for any $x \in \mathcal{X}$ shall be computed using the formula: <u>$\text{sign}\left(\hat{w}_m^T \varphi(x)\right)$</u>.

[1/2 mark]

Now, say, the training data actually is

$$D = \left\{ \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, 1 \right), \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}, -1 \right), \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}, -1 \right), \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, 1 \right) \right\}$$

and the feature map $\phi$ is defined by

$$\phi\left( \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right) \equiv ((z_1 + z_2)\%2) - \frac{1}{2} \ \forall \ z_1, z_2 \in \mathbb{R}.$$

Here, $a\%b$ is the remainder when $a$ is divided by $b$. For this case, in the box below, write down the optimal solution[1] of the ERM problem along with justification:

---

[1]You are welcome to solve this optimization problem in any way you prefer. For e.g., analytically, manually iterating through gradient descent etc.

No Solution. Because $\mathcal{D}$ is
linearly separable in feature space.

[1 mark]

For this specific training data and feature map, suppose we wish to perform linear classification using the 0-1 loss. Then, run the perceptron algorithm in rough using manual calculations and write all parameter iterates until convergence including initialization in the box below. For each iterate write down the update equation too. No other details are required

$$w^{(0)} = 0$$
$$w^{(1)} = 0 + (-1)\left(\tfrac{1}{2}\right) = -\tfrac{1}{2}.$$

$\longrightarrow w^{(1)} = 0 + (1)\left(-\tfrac{1}{2}\right) = -\tfrac{1}{2}$ *is also correct.*

[1 mark]

2. Consider the linear regression setting taught in lectures with training data as: $\mathcal{D} = \{(2,1),(4,5)\}$ (usual convention of set of input,label pairs). Consider the feature map $\phi(x) = x$. Analytically solve the ERM problem in rough work and write down the final ERM solution in this blank: $w_\phi^{ERM} = \underline{1.1}$ [[1/2 mark]]. With this solution, the explained variance computed on the training set is $\underline{0.775}$ [[Fill the blank with appropriate number. 1/2 mark]]. Now, consider another feature map, $\psi(x) \equiv \begin{bmatrix} x \\ 1 \end{bmatrix}$. With this feature map, analytically solve the ERM problem in rough work and write down the final ERM solution in this blank: $w_\psi^{ERM} = \begin{bmatrix} 2 \\ -3 \end{bmatrix}$ [[1 mark]]. With this solution, the explained variance computed on the training set is $\underline{1}$ (one) [[Fill the blank with appropriate number. 1/4 mark]]. Now, consider another feature map, $\nu(x) \equiv \begin{bmatrix} x \\ 1 \\ -1 \end{bmatrix}$. With this feature map, analytically solve the ERM problem in rough work and write down the final ERM solution you obtained in this blank: $w_\nu^{ERM} = \begin{bmatrix} 2 \\ -3 \\ 0 \end{bmatrix}$

*Any other where $w_2 - w_3 = -3$ is also fine.*
*For e.g.*
$$\begin{bmatrix} 2 \\ 45 \\ 48 \end{bmatrix}$$

4

[[*1 mark*]]. With this solution, the explained variance computed on the training set is **0** [[*Fill the blank with appropriate number. 1/4 mark*]]. (one)

3. In the lectures you were taught how to model the Bayes optimal in a binary classification task using linear functions (over input feature space). Now suppose you have a multi-class classification problem with 3 classes: 'Ξ', '‡', and 'ς'. However, still you are only allowed to use the linear model taught in lectures. Think about how you can model the Bayes optimal in a 3-class classification task using these linear functions. With this way of modelling the Bayes optimal in mind, according to you, the loss function, $l$, appropriate for this task would be defined by $l(w, x, \Xi) \equiv \underline{\mathbb{1}_{\{w^T\varphi(x) > -1\}}}$ , $l(w, x, \ddagger) \equiv$ $\underline{\mathbb{1}_{\{w^T\varphi(x) \notin (-1, 1)\}}}$ , $l(w, x, \varsigma) \equiv \underline{\mathbb{1}_{\{w^T\varphi(x) < 1\}}}$ . Here, $w$ denotes the parameter of the linear model.

*[handwritten left: Any cyclic permutation of these is also ok]*

*[handwritten right: Any other tie break is also fine → Any thresholds instead of 1, -1 are also oks.]*

[1.5 Marks]

Observe that your way of modelling the Bayes optimal with linear functions has an inherent ('wrong'?) bias. More specifically, if the parameter changes a little then the label for a fixed $x$ changes preferentially to one of the other two classes. In this sense, there is am implicit (unequal) nearness between different class pairs.

*[handwritten left: "logistic" & "hinge versions" are also fine.]*

Now, suppose you are allowed to model functions of the form $f(x) = W^T\phi(x)$, where $W$ is $n \times 3$, where $\phi$ is a feature map. You may use the notation $W = [w_1 \ w_2 \ w_3]$, where $w_i \in \mathbb{R}^n$. Think about how you can model the Bayes optimal in a 3-class classification task using these '3-dimensional linear functions'. With this way of modelling the Bayes optimal in mind, according to you, the loss function, $l$, appropriate for this task would be defined by $l(W, x, \Xi) \equiv$ $\underline{\mathbb{1}_{\{w_1^T\varphi(x) < w_2^T\varphi(x)\}} + \mathbb{1}_{\{w_1^T\varphi(x) < w_3^T\varphi(x)\}}}$ $l(W, x, \ddagger) \equiv \underline{\mathbb{1}_{\{w_2^T\varphi(x) < w_1^T\varphi(x)\}} + \mathbb{1}_{\{w_2^T\varphi(x) < w_3^T\varphi(x)\}}}$ $l(W, x, \varsigma) \equiv \underline{\mathbb{1}_{\{w_3^T\varphi(x) < w_1^T\varphi(x)\}} + \mathbb{1}_{\{w_3^T\varphi(x) < w_2^T\varphi(x)\}}}$.

[2.5 Marks]

*[handwritten: Many alternatives (see above part).]*

5

① Perceptron question.

At iteration $k$

$~~~~~~~~$ ~~$\omega^{(k)}$~~ ~~$\omega^{(k-1)}$~~ if $\left( y_i \omega^T \varphi(u_i) \leq 0 \text{ for some } i \right)$

$~~~~~~~~$ then, $\omega^{(k)} = \omega^{(k-1)} + y_i \varphi(u_i)$

$\omega^{(0)} = 0$

$\varphi(u_1) = \varphi(x_4) = -\frac{1}{2}$

$\varphi(u_3) = \varphi(x_2) = \frac{1}{2}$.

$\omega^{(1)}$

$y_1 \omega^{(0)T} \varphi(u_1) = y_4 \omega^{(0)T} \varphi(u_4) = 0 ~~\cancel{\geq} \leq 0$

$y_3 \omega^{(0)T} \varphi(u_3) = y_2 \omega^{(0)T} \varphi(x_2) = 0 ~~~~~\leq 0$

We can choose any $i = 1$ to $4$.

Let's choose $i = 3$. $~~\omega^{(1)} = 0 + (-1)\frac{1}{2} = -\frac{1}{2}$

---

② $\omega_\varphi^{ERM} = \underset{\omega \in \mathbb{R}}{\text{argmin}} ~\frac{1}{2}\left[ \left( \omega(2) - 1 \right)^2 + \left( \omega(4) - 5 \right)^2 \right]$

$~~~~~~~~ = \underset{\omega \in \mathbb{R}}{\text{argmin}}$

$~~~~~~~~~~~~~~~~~~~ 10\omega^2 - 22\omega ~=~ \frac{22}{20} = 1.1$

$~~~~~~~~ = ~~$ ~~argmin~~ ~~$\omega \in \mathbb{R}$~~

$MSE\left(\omega_\varphi^{ERM}\right) = \frac{1}{2}\left[ \left( 1.1 \times 2 - 1 \right)^2 + \left( 1.1 \times 4 - 5 \right)^2 \right]$

$~~~~~~~~~~~~~~ = \frac{1.44}{2} + \frac{0.36}{2} = 0.9$

$MSE\left(\bar{\omega}\right) = \frac{1}{2}\left[ (3-1)^2 + (3-5)^2 \right] = 4$

$\bar{\omega}^T \varphi(u) = \frac{5+1}{2} = 3$

$~~~~ \forall x$

$\therefore$ exp. val $= 1 - \frac{0.9}{4} = 1 - 0.225 = 0.775$

(2)

$$W_\psi^{ERM} = \underset{w_1 \in \mathbb{R}, w_2 \in \mathbb{R}}{\operatorname{argmin}} \; \frac{1}{2} \underbrace{\left[ (2w_1 + w_2 - 1)^2 + (4w_1 + w_2 - 5)^2 \right]}_{g(w)}$$

$$\frac{\partial g(w)}{\partial w_1} = 2(2w_1 + w_2 - 1) + 4(4w_1 + w_2 - 5) = 0 \iff \begin{array}{l} 20w_1 + 6w_2 - 22 = 0 \\ \iff w_2 = \frac{22}{6} - \frac{20w_1}{6} \end{array}$$

$$\frac{\partial g(w)}{\partial w_2} = (2w_1 + w_2 - 1) + (4w_1 + w_2 - 5) = 0 \iff 6w_1 + 2w_2 - 6 = 0$$

$$18w_1 + (22 - 20w_1) - 18 = 0$$

$$\iff w_1 = 2$$

$$w_2 = \frac{-18}{6} = -3$$

$$MSE\left(W_\psi^{ERM}\right) = \frac{1}{2}\left[ (2 \times 2 - 3 - 1)^2 + (4 \times 2 + 3 - 5)^2 \right]$$

$$= 0 \iff \text{equivariance} = 1$$

(2)

$$W_\gamma^{ERM} = \underset{w_1 \in \mathbb{R}, w_2 \in \mathbb{R}, w_3 \in \mathbb{R}}{\operatorname{argmin}} \; \frac{1}{2}\left[ (2w_1 + w_2 - w_3 - 1)^2 + (4w_1 + w_2 - w_3 - 5)^2 \right]$$

∴ equations will be same as above except $w_2 \to w_2 - w_3$.

∴ $w_1 = 2$, $w_2 - w_3 = -3$ are optimality conditions.

Any choice of $w$ that ⌢satisfies⌢ is fine.

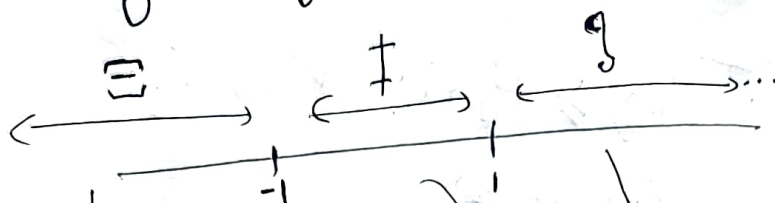My choice is $\begin{bmatrix} 2 \\ -3 \\ 0 \end{bmatrix}$.

expvariance will remain 1.

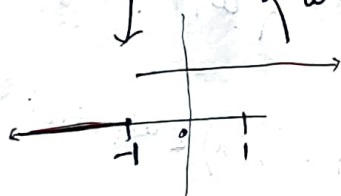I part

I'll bin map reals onto three bins
instead of two (as)

Two bins → require one threshold
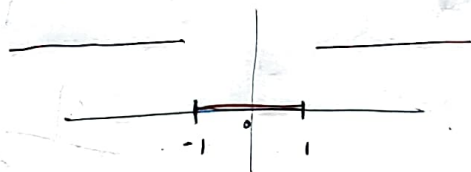Three bins → require two thresholds.

My choice for thresholds is $-1, 1$

$$\xleftarrow{\quad \Xi \quad} \qquad \xleftrightarrow{\quad \mp \quad} \qquad \xrightarrow{\quad 9 \quad} \cdots$$

$$\ell(\omega, x, \Xi) = \mathbb{1}_{\{\omega^T \varphi(x) > -1\}}$$

$$\ell(\omega, x, 9) = \mathbb{1}_{\{\omega^T \varphi(x) < 1\}}$$

$$\ell(\omega, x, \mp) = \mathbb{1}_{\{\omega^T \varphi(x) \geq 1\}} + \mathbb{1}_{\{\omega^T \varphi(x) \leq -1\}}$$

$$= \mathbb{1}_{\{\omega^T \varphi(x) \geq 1 \,\&\, \omega^T \varphi(x) \leq -1\}}$$

$$= \mathbb{1}_{\{\omega^T \varphi(x) \notin (-1, 1)\}}$$

3) II part

We may say high $w_i^T \phi(x) \iff i^{th}$ class.

Let threshold be "$h$".

Then $w_i^T \phi(x) \geq h$, $w_j^T \phi(x) < h \implies i^{th}$ class

$(j \neq i)$

This is same as

$$w_i^T \phi(x) \geq w_j^T \phi(x) \quad \forall j \neq i.$$

loss for $1^{st}$ class

$$\mathbb{1}\left\{w_1^T \phi(x) < w_2^T \phi(x)\right\} + \mathbb{1}\left\{w_1^T \phi(x) < w_3^T \phi(x)\right\}$$

$$= \mathbb{1}\left\{w_1^T \phi(x) < w_2^T \phi(x) \text{ or } w_1^T \phi(x) < w_3^T \phi(x)\right\}$$

$$= \mathbb{1}\left\{w_1^T \phi(x) < \min_{j=2,3} w_j^T \phi(x)\right\}.$$