

Problem 1.**(a)**

Let's assume that the observation set uses the following notation:

$$X_{N \times p} = (x_1, x_2, \dots, x_N)^T, x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$$

This notation indicates that there are N samples, and each sample is a p -dimensional vector. Each observation is generated by $p(x|\theta)$.

For discriminative model, it assumes that The θ in $p(x|\theta)$ is a constant. For N observations, the probability of the observation set is $p(X|\theta) = \prod_{iid, i=1}^N p(x_i|\theta)$. In order to find the size of θ , we use the maximum log-likelihood MLE method:

$$\theta_{MLE} = \underset{\theta}{argmax} \log p(X|\theta) \underset{iid}{=} \underset{\theta}{argmax} \sum_{i=1}^N \log p(x_i|\theta)$$

For generative model, it assumes The θ in $p(x|\theta)$ is not a constant. This θ satisfies a preset prior distribution $\theta \sim p(\theta)$. So according to Bayes' theorem, the posterior dependent on the parameters of the observation set can be written as:

$$p(\theta|X) = \frac{p(X|\theta) \cdot p(\theta)}{p(X)} = \frac{p(X|\theta) \cdot p(\theta)}{\int_{\theta} p(X|\theta) \cdot p(\theta) d\theta}$$

In order to find the value of θ , we have to maximize this parameter posterior MAP:

$$\theta_{MAP} = \underset{\theta}{argmax} p(\theta|X) = \underset{\theta}{argmax} p(X|\theta) \cdot p(\theta)$$

In general, generative model learns the joint probability distribution $p(X, \theta)$ and a discriminative model learns the conditional probability distribution $p(X|\theta)$ - which is the probability of X given θ .

(b)

I would choose Logistic regression.

The requirements of the generative model are relatively high, and the joint distribution of the input data x and the category needs to be learned. For the person who only has limited knowledge of the process generating the data, the Naive Bayes would be hard to compute. The generative model needs to

calculate the class conditional probability and the prior probability, which increases the amount of calculation. Furthermore, the generative model has higher bias, because it focuses on a smaller set of models

On the contrary, the discriminant model is to solve directly, learn the optimal classification surface between different categories, and reflect the differences between different types of data. The discriminant model focuses on a bigger set of models and has lower bias.

Problem 2.

(1)

False

(2)

True

(3)

False

Problem 3.

(a)

1. For gradient descent, we compute the full gradient. But for stochastic gradient descent, we compute the random gradient. This means that GD is $O(n)$ slower than SGD.

2. Stochastic gradient descent is memory efficient, extended to mini-batches.

(b)

Total runtime of gradient descent: $O(\frac{n}{\epsilon}) + O(\frac{n}{\epsilon}) = O(\frac{n}{\epsilon})$

Total runtime of stochastic gradient descent: $O(\frac{1}{\epsilon^2}) + O(\frac{1}{\epsilon^2}) = O(\frac{1}{\epsilon^2})$

(c)

Total runtime of sub-gradient descent: $O(\frac{n}{\epsilon}) + O(\frac{n}{\epsilon^2}) = O(\frac{n}{\epsilon^2})$

Total runtime of stochastic sub-gradient descent: $O(\frac{1}{\epsilon^2}) + O(\frac{1}{\epsilon^2}) = O(\frac{1}{\epsilon^2})$

Problem 4.**(a)**

For perceptron, the prediction on x_i is incorrect if $y_i w^T x_i < 0$. Let the $I(w)$ be the set of points on which prediction is incorrect. The objective function would be:

$$E(w) = -\sum_{i \in I(w)} y_i w^T x_i = \sum_{i=1}^n \max(0, -y_i w^T x_i)$$

The objective function of perceptron is hinge loss, which has positive value when the data point is misclassified and 0 when the data point is correctly classified. So there exists a hinge point that makes the hinge loss non-smooth.

(b)

For this question, $\eta = 1$, $w_0 = 0$. Assuming we have n data points and let the $I(w)$ be the set of points on which prediction is incorrect.

For perceptron, if prediction is correct, w does not change.

If prediction is wrong:

$$w_{new} = w_{old} - \eta \nabla E_i(w) = w_{old} + \eta y_i x_i$$

For the first round, n iterations of running perceptron:

$$\begin{aligned} w_n &= w_0 + \sum_{i \in I(w)} y_i x_i \\ &= w_0 + \sum_{i=1}^n n_i y_i x_i \\ &= \sum_{i=1}^n n y_i x_i, \text{ where } n = 1\{i \in I(w)\} \end{aligned}$$

If we keep running the perceptron until it converges, w_n would be updated every time the point (x_i, y_i) is misclassified and the corresponding number of mistakes made on (x_i, y_i) would increase each time. Since the data points are independent and identically distributed:

$$w^* = \sum_{i=1}^n n_i y_i x_i, \text{ where } n_i \text{ is the number of mistakes made on } (x_i, y_i).$$

Because we let the final vector after convergence $w^* = \sum_{i=1}^n \alpha_i y_i x_i$, $\alpha_i = n_i$ is the number of mistakes made on (x_i, y_i) during the perceptron iterations before convergence.

(c)

Linear separability means that a linear function can be used to separate two classes without error, such as a straight line in two-dimensional space, a plane in three-dimensional space, and a linear function in high-dimensional space. Linear inseparable refers to the situation where some samples are divided by linear classification planes that will produce classification errors.

For the classes are not linearly separable, there must exist a data point (x_i, y_i) make $y_i w^T x_i < 0$ for all w . Thus, the perceptron with a fixed learning rate η cannot converge.

I cannot modify the learning rate η so that the iterations converge in expectation. Once the iterations can converge, it means we can find a w that can clearly separate the two classes for all data points, which contradict the assumption that the classes are not linearly separable.

Extra Credit 1:

$$\begin{aligned}
\theta^* &= \operatorname{argmax}_{\theta} P(x_1, x_2, \dots, x_n | \theta) \\
&= \operatorname{argmax}_{\theta} \prod_{i=1}^n P(x_i | \theta) \text{ (i.i.d. samples)} \\
&= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln(P(x_i | \theta)) \text{ (take the logarithm)} \\
&= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln(2\theta x_i e^{-\theta x_i^2}) \\
&= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln(2\theta x_i) + \ln(e^{-\theta x_i^2}) \\
&= \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln(2\theta x_i) - \theta x_i^2 \\
&= \operatorname{argmin}_{\theta} \sum_{i=1}^n \theta x_i^2 - \ln(2\theta x_i)
\end{aligned}$$

To get the optimal θ , we take the derivative of θ for the above equation and let the derivative equals to 0.

$$\begin{aligned}
\frac{\partial(\sum_{i=1}^n \theta x_i^2 - \ln(2\theta x_i))}{\partial \theta} &= 0 \\
\frac{\partial(\sum_{i=1}^n \theta x_i^2)}{\partial \theta} - \frac{\partial(\sum_{i=1}^n \ln(2\theta x_i))}{\partial \theta} &= 0 \\
(\sum_{i=1}^n x_i^2) - (\sum_{i=1}^n \frac{2x_i}{2\theta x_i}) &= 0 \\
(\sum_{i=1}^n x_i^2) - (\sum_{i=1}^n \frac{1}{\theta}) &= 0 \\
(\sum_{i=1}^n x_i^2) &= \frac{n}{\theta} \\
\theta &= \frac{n}{(\sum_{i=1}^n x_i^2)}
\end{aligned}$$

Extra Credit 2:

$$\frac{\partial L(y_i|x_i, w)}{\partial w_j} = \frac{\partial y_i \log(\frac{1}{1+e^{-w^T x_i}}) + (1-y_i) \log(\frac{1}{1+e^{w^T x_i}})}{\partial w_j}$$

For the first part:

$$\begin{aligned} & \frac{\partial \log(\frac{1}{1+e^{-w^T x_i}})}{\partial w_j} \\ &= -\frac{\partial \log(1+e^{-w^T x_i})}{\partial w_j} \text{ using } \log(a^b) = b \log(a) \end{aligned}$$

Now, using the chain rule, set $1 + e^{-w^T x_i} = y$

$$\begin{aligned} &= -\frac{\partial \log y}{\partial y} \frac{\partial y}{\partial w_j} \\ &= -\frac{1}{y} \frac{\partial y}{\partial w_j} \end{aligned}$$

Now, for $\frac{\partial 1+e^{-w^T x_i}}{\partial w_j}$, we do not need to consider the constant term. So, using the chain rule again, set $-w^T x_i = y$

$$\begin{aligned} &= -\frac{1}{y} \frac{\partial e^y}{\partial y} \frac{\partial y}{\partial w_j} \\ &= -\frac{1}{y} * (e^{-w^T x_i} * (-x_{ij})) \\ &= \frac{e^{-w^T x_i} * x_{ij}}{1+e^{-w^T x_i}} \\ &= (1 - \frac{1}{1+e^{-w^T x_i}}) x_{ij} \end{aligned}$$

For the second part:

$$\begin{aligned} & \frac{\partial \log(\frac{1}{1+e^{w^T x_i}})}{\partial w_j} \\ &= -\frac{\partial \log(1+e^{w^T x_i})}{\partial w_j} \text{ using } \log(a^b) = b \log(a) \end{aligned}$$

Now, using the chain rule, set $1 + e^{w^T x_i} = y$

$$\begin{aligned} &= -\frac{\partial \log y}{\partial y} \frac{\partial y}{\partial w_j} \\ &= -\frac{1}{y} \frac{\partial y}{\partial w_j} \end{aligned}$$

Now, for $\frac{\partial 1+e^{w^T x_i}}{\partial w_j}$, we do not need to consider the constant term. So, using the chain rule again, set $w^T x_i = y$

$$\begin{aligned}
&= -\frac{1}{y} \frac{\partial e^y}{\partial y} \frac{\partial y}{\partial w_j} \\
&= -\frac{1}{y} * (e^{w^T x_i} * (x_{ij})) \\
&= -\frac{e^{w^T x_i} * x_{ij}}{1 + e^{w^T x_i}} \\
&= -\frac{x_{ij}}{1 + e^{-w^T x_i}}
\end{aligned}$$

The original formula:

$$\begin{aligned}
&= y_i \left(\left(1 - \frac{1}{1 + e^{-w^T x_i}} \right) x_{ij} \right) + (1 - y_i) \left(-\frac{x_{ij}}{1 + e^{-w^T x_i}} \right) \\
&= \left(y_i - \frac{y_i}{1 + e^{-w^T x_i}} - \frac{1}{1 + e^{-w^T x_i}} + \frac{y_i}{1 + e^{-w^T x_i}} \right) x_{ij} \\
&= \left(y_i - \frac{1}{1 + e^{-w^T x_i}} \right) x_{ij} \\
&= (y_i - \sigma(w^T x_i)) x_{ij}
\end{aligned}$$

Extra Credit 3:

Suppose K is a gram matrix which $K_{ij} = K(x_i, x_j) \in R^{n \times n}$.

For any vector $\alpha \in R^n$:

$$\begin{aligned}
&\alpha^T K \alpha \\
&= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j (K(x_i, x_j)) \\
&= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \\
&\because K(x, y) = \min\{x, y\}, x \text{ and } y \text{ are positive numbers.} \\
&\therefore \forall K(x_i, x_j), \text{ it is a positive number. Let's assume it is } m_{ij}. \\
&= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \Phi(\sqrt{m_{ij}})^T \Phi(\sqrt{m_{ij}}) \\
&\because \alpha_i \text{ and } \alpha_j \text{ are real numbers} \\
&= \sum_{i=1}^n \alpha_i \Phi(\sqrt{m_{ij}})^T \sum_{j=1}^n \alpha_j \Phi(\sqrt{m_{ij}}) \\
&= \sum_{k=1}^m \left\{ \sum_{i=1}^n \alpha_i \Phi(\sqrt{m_{ij}}) \right\}^T \sum_{j=1}^n \alpha_j \Phi(\sqrt{m_{ij}})
\end{aligned}$$

\therefore inner product always ≥ 0

$\therefore \alpha^T K \alpha \geq 0$ and K is positive semi-definite function.

For symmetry:

$$\begin{aligned} K(x, y)^T &= (\min\{x, y\})^T \\ &= (\min\{x, y\}) \\ &= K(x, y) \end{aligned}$$

$K(x, y) = \min\{x, y\}$ is a valid kernel function because it is symmetric and positive semi-definite function.