EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the Privacy Policy.





<u>Course</u> > <u>Final exam (1 week)</u> > <u>Final Exam</u> > Problem 5

Problem 5

In this problem, we will do regression for data that are generated from a Gaussian Mixture Model. Let $X \in \mathbb{R}$ be the random variable for the features and $Y \in \mathbb{R}$ be the random variable for the output. We assume X is generated from a mixture of m Gaussian distributions, and Y is linearly correlated to X with some random noise. The generation process can be described as follows:

- 1. Sample a random variable T from a multinomial distribution on $\{1,2,\ldots,m\}$, where $P\left(T=t
 ight)=p_{t}$.
- 2. Sample X from the tth Gaussian distribution, with mean μ_t and variance σ_t^2 .
- 3. Given X from the tth Gaussian distribution, let $Y=w_tX+\epsilon$, where w_t is a fixed parameter for the tth Gaussian and ϵ is from an independent Gaussian with 0 mean and variance of 1.

5. (1)

1/1 point (graded)

Which of the following is the correct probability density of X?

$$\sum_{t=1}^m rac{1}{\sqrt{2\pi\sigma_t^2}} \mathrm{exp}\left(-rac{(X-\mu_t)^2}{2\sigma_t^2}
ight)$$

$$\sum_{t=1}^{m} \frac{p_t}{\sqrt{2\pi\sigma_t^2}} \exp\left(-\frac{(X-\mu_t)^2}{2\sigma_t^2}\right) \checkmark$$

$$\prod_{t=1}^{m}rac{1}{\sqrt{2\pi\sigma_{t}^{2}}}\mathrm{exp}\left(-rac{(X-\mu_{t})^{2}}{2\sigma_{t}^{2}}
ight).$$

$$\prod_{t=1}^m rac{p_t}{\sqrt{2\pi\sigma_t^2}} \mathrm{exp}\left(-rac{(X-\mu_t)^2}{2\sigma_t^2}
ight).$$

Solution:

$$p\left(X=x
ight) \ = \sum_{t=1}^{m} p\left(X=x|T=t
ight) P\left(T=t
ight)$$

$$=\sum_{t=1}^{m}rac{p_{t}}{\sqrt{2\pi\sigma_{t}^{2}}}\mathrm{exp}\left(-rac{(X-\mu_{t})^{2}}{2\sigma_{t}^{2}}
ight)$$

Submit

You have used 1 of 3 attempts

- **1** Answers are displayed within the problem
- 5. (2)

1/1 point (graded)

Now, given an observation of X=x, what is the likelihood that it is drawn from the tth Gaussian distribution?

$$rac{rac{1}{\sigma_t} \mathrm{exp}\left(-(x-\mu_t)^2/\left(2\sigma_t^2
ight)
ight)}{\sum_{i=1}^m rac{p_i}{\sigma_i} \mathrm{exp}\left(-(x-\mu_i)^2/\left(2\sigma_i^2
ight)
ight)}$$

$$rac{rac{p_t}{\sigma_t} \mathrm{exp}\left(-(x-\mu_t)^2/\left(2\sigma_t^2
ight)
ight)}{\sum_{i=1}^m rac{1}{\sigma_i} \mathrm{exp}\left(-(x-\mu_i)^2/\left(2\sigma_i^2
ight)
ight)}$$

$$\frac{\frac{p_{t}}{\sigma_{t}}\mathrm{exp}\left(-(x-\mu_{t})^{2}/\left(2\sigma_{t}^{2}\right)\right)}{\sum_{i=1}^{m}\frac{p_{i}}{\sigma_{i}}\mathrm{exp}\left(-(x-\mu_{i})^{2}/\left(2\sigma_{i}^{2}\right)\right)} \checkmark$$

$$rac{rac{1}{\sigma_t} \mathrm{exp}\left(-(x-\mu_t)^2/\left(2\sigma_t^2
ight)
ight)}{p_t \sum_{i=1}^m rac{p_i}{\sigma_i} \mathrm{exp}\left(-(x-\mu_i)^2/\left(2\sigma_i^2
ight)
ight)}$$

Solution:

Using Bayesian theorem, we have

$$egin{aligned} P\left(T=t|X=x
ight) &= rac{p_{t}}{p\left(X=x|T=t
ight)P\left(T=t
ight)} \ p\left(X=x
ight) \end{aligned} \ &= rac{rac{p_{t}}{\sqrt{2\pi\sigma_{t}^{2}}} \exp\left(-\left(X-\mu_{t}
ight)^{2}/2\sigma_{t}^{2}
ight)}{\sum_{i=1}^{m} rac{p_{i}}{\sqrt{2\pi\sigma_{i}^{2}}} \exp\left(-\left(X-\mu_{i}
ight)^{2}/2\sigma_{i}^{2}
ight)} \end{aligned} \ &= rac{rac{p_{t}}{\sigma_{t}} \exp\left(-\left(x-\mu_{t}
ight)^{2}/\left(2\sigma_{t}^{2}
ight)
ight)}{\sum_{i=1}^{m} rac{p_{i}}{\sigma_{i}} \exp\left(-\left(x-\mu_{t}
ight)^{2}/\left(2\sigma_{i}^{2}
ight)
ight)} \end{aligned}$$

Submit

You have used 1 of 3 attempts

1 Answers are displayed within the problem

5. (3)

1/1 point (graded)

The objective of regression is to find an optimal function $f^*: \mathbb{R} \to \mathbb{R}$ that minimizes to loss $\mathbb{E}\left[(Y-f(X))^2\right]$ over all choices of f. Suppose we know the generation process in prior (e.g. all the parameters for the multinomial and Gaussian distributions), which of the following is the explicit form of the solution f^* ?

$$f^{st}\left(X
ight) = rac{\sum_{t=1}^{m}rac{p_{t}}{\sigma_{t}}\mathrm{exp}\left(-\left(x-\mu_{t}
ight)^{2}/\left(2\sigma_{t}^{2}
ight)
ight)w_{t}X}{\sum_{i=1}^{m}rac{p_{i}}{\sigma_{i}}\mathrm{exp}\left(-\left(x-\mu_{i}
ight)^{2}/\left(2\sigma_{i}^{2}
ight)
ight)}$$

$$f^{st}\left(X
ight) = rac{\sum_{t=1}^{m}rac{p_{t}}{\sigma_{t}}\mathrm{exp}\left(-(x-\mu_{t})^{2}/\left(2\sigma_{t}^{2}
ight)
ight)}{\sum_{i=1}^{m}rac{p_{i}}{\sigma_{i}}\mathrm{exp}\left(-(x-\mu_{i})^{2}/\left(2\sigma_{i}^{2}
ight)
ight)w_{i}X}$$

$$f^{st}\left(X
ight) = rac{\sum_{t=1}^{m}rac{p_{t}}{\sigma_{t}}\mathrm{exp}\left(-\left(x-\mu_{t}
ight)^{2}/\left(2\sigma_{t}^{2}
ight)
ight)\left(w_{t}X+\epsilon
ight)}{\sum_{i=1}^{m}rac{p_{i}}{\sigma_{i}}\mathrm{exp}\left(-\left(x-\mu_{i}
ight)^{2}/\left(2\sigma_{i}^{2}
ight)
ight)}$$

$$f^{st}\left(X
ight) = rac{\sum_{t=1}^{m}rac{p_{t}}{\sigma_{t}}\mathrm{exp}\left(-(x-\mu_{t})^{2}/\left(2\sigma_{t}^{2}
ight)
ight)}{\sum_{i=1}^{m}rac{p_{i}}{\sigma_{i}}\mathrm{exp}\left(-(x-\mu_{i})^{2}/\left(2\sigma_{i}^{2}
ight)
ight)\left(w_{i}X+\epsilon
ight)}$$

Correction Note (Sept 9):

Solution:

To minimize the loss, we need $f^*\left(X\right)=\mathbb{E}\left[Y|X\right]$. As $p\left(Y|X\right)=\sum_t p\left(Y,T|X\right)=\sum_t p\left(Y|X,T\right)p\left(T|X\right)$, we have:

$$egin{aligned} f^*\left(X
ight) &= \mathbb{E}\left[Y|X
ight] \ &= \mathbb{E}_{T|X}\left[\mathbb{E}\left[Y|T,X
ight]
ight] \ &= \sum_{t}^{m} P\left(T=t|X
ight)\mathbb{E}\left[Y|T=t,X
ight] \ &= \sum_{t}^{m} rac{P\left(X|T=t
ight)P\left(T=t
ight)}{P\left(X
ight)} w_{t}X \ &= rac{\sum_{t=1}^{m} rac{p_{t}}{\sigma_{t}} \exp\left(-(x-\mu_{t})^{2}/\left(2\sigma_{t}^{2}
ight)
ight)w_{t}X}{\sum_{i=1}^{m} rac{p_{i}}{\sigma_{i}} \exp\left(-(x-\mu_{t})^{2}/\left(2\sigma_{i}^{2}
ight)
ight)} \end{aligned}$$

Submit

You have used 3 of 3 attempts

1 Answers are displayed within the problem

5. (4)

2/2 points (graded)

Now suppose we don't know the data generation process, but observe N datapoints (x_n, y_n) for $1 \le n \le N$. This time we would like to fit the function f^* , with the constraint that f^* is a linear function. In another word, we would like to find the optimal parameters a^* and b^* for $f^* = a^*X + b^*$ which minimize the empirical loss

$$\sum_{n=1}^N \left(y_n - (ax_n + b)
ight)^2$$

over $a \in \mathbb{R}$, $b \in \mathbb{R}$.

Recall in linear regression, we can derive a closed form solution for a^* and b^* by setting the derivative of the loss function to 0. Try to compute this closed form solution and think of the situation when $N\to\infty$, i.e. when we have infinite number of training examples, what is the value of a^* and b^* ?

Hint: When $N \to \infty$, $\bar{x} \to \mathbb{E}[X]$, $\bar{y} \to \mathbb{E}[Y]$, $\frac{1}{N} \sum_{i=1}^{N} (x_i - \bar{x}) (y_i - \bar{y}) \to \operatorname{Cov}(X, Y)$, $\frac{1}{N} \sum_{i=1}^{N} (x_i - \bar{x})^2 \to \operatorname{Var}(X)$, where \bar{x} represents the mean of the observed x, Cov refers to the variance.

$$a^* = rac{\mathrm{Cov}\left(X,Y
ight)}{\mathbb{E}\left[X
ight]}$$

$$a^{st} = rac{\mathrm{Var}\left(X
ight)}{\mathrm{Cov}\left(X,Y
ight)}$$

$$a^{st} = rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(Y
ight)}$$

$$a^* = rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(X
ight)}$$
 🗸

$$b^{st} = \mathbb{E}\left[Y
ight] - rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(X
ight)} \mathbb{E}\left(X
ight) oldsymbol{\checkmark}$$

$$b^{st} = \mathbb{E}\left[X
ight] - rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(X
ight)} \mathbb{E}\left(Y
ight)$$

$$b^{st} = \mathbb{E}\left[Y
ight] - rac{\mathrm{Var}\left(X
ight)}{\mathrm{Cov}\left(X,Y
ight)} \mathbb{E}\left(X
ight)$$

$$b^{st}=\mathbb{E}\left[X
ight]-rac{\mathrm{Var}\left(X
ight)}{\mathrm{Cov}\left(X,Y
ight)}\mathbb{E}\left(Y
ight)$$

Solution:

Taking the derivative of the loss function to a and b, we have

$$egin{align} rac{\partial loss}{\partial a} &= -2\sum_{n=1}^{N}\left(y_n - ax_n - b
ight)x_n \ rac{\partial loss}{\partial b} &= -2\sum_{n=1}^{N}\left(y_n - ax_n - b
ight) \end{aligned}$$

Set both of them to 0 and we can solve for a and b as

$$egin{aligned} a^* &= rac{rac{1}{N} \sum_{n=1}^N x_n y_n - (rac{1}{N} \sum_{n=1}^N x_n) \left(rac{1}{N} \sum_{n=1}^N y_n
ight)}{rac{1}{N} \sum_{n=1}^N x_n^2 - \left(rac{1}{N} \sum_{n=1}^N x_n
ight)^2} \ &= rac{rac{1}{N} \sum_{n=1}^N \left(x_n - ar{x}
ight) \left(y_n - ar{y}
ight)}{rac{1}{N} \sum_{n=1}^N \left(x_n - ar{x}
ight)^2} \ b^* &= rac{1}{N} \sum_{n=1}^N y_n - rac{a^*}{N} \sum_{n=1}^N x_n \end{aligned}$$

When $N o \infty$, we have:

$$egin{aligned} a^{*} &= rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(X
ight)} \ b^{*} &= \mathbb{E}\left[Y
ight] - rac{\mathrm{Cov}\left(X,Y
ight)}{\mathrm{Var}\left(X
ight)} \mathbb{E}\left(X
ight) \end{aligned}$$

Submit

You have used 1 of 3 attempts

• Answers are displayed within the problem

5. (5)

3/3 points (graded)

Now, let's consider a concrete example when m=2, $p_1=p_2=0.5$, $w_1=1$, $w_2=-1$, $\mu_1=2$, $\mu_2=-2$, and $\sigma_1=\sigma_2=1$, what is the value of $\mathbb{E}\left[X\right]$, $\mathbb{E}\left[Y\right]$, $\mathbb{E}\left[XY\right]$? Enter your solutions below.

$$\mathbb{E}\left[X
ight]=iggl[0]$$
 Answer: 0

$$\mathbb{E}\left[Y
ight]=igg|$$
 2 wo Answer: 2

$$\mathbb{E}\left[XY\right]= \boxed{egin{array}{c} 0 \end{array}}$$
 Answer: 0

Solution:

The probability density of X now is:

$$p\left(X
ight) = rac{1}{2\sqrt{2\pi}}[\exp{(-rac{(x-2)^2}{2})} + \exp{[-rac{(x+2)^2}{2}]}]$$

The expectation of X is therefore:

$$\begin{split} \mathbb{E}\left[X\right] &= \int X P\left(X\right) dX \\ &= \int X \frac{1}{2\sqrt{2\pi}} \left[\exp\left(-\frac{\left(X-2\right)^2}{2}\right) + \exp\left[-\frac{\left(X+2\right)^2}{2}\right]\right] dX \\ &= \frac{1}{2} \left[\int X \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\left(X-2\right)^2}{2}\right) dX + \int X \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\left(X+2\right)^2}{2}\right) dX\right] \\ &= \frac{1}{2} [2 + (-2)] \\ &= 0 \end{split}$$

The distribution of *Y* here is:

$$egin{aligned} p\left(Y
ight) &= p\left(Y | T=1
ight) P\left(T=1
ight) + p\left(Y | T=2
ight) P\left(T=2
ight) \ &= p_1 \mathcal{N} \left(w_1 \mu_1, w_1^2 \sigma_1^2 + 1
ight) + p_2 \mathcal{N} \left(w_2 \mu_2, w_2^2 \sigma_2^2 + 1
ight) \end{aligned}$$

$$egin{aligned} &=rac{1}{2}\mathcal{N}\left(2,2
ight)+rac{1}{2}\mathcal{N}\left(2,2
ight) \ &=\mathcal{N}\left(2,2
ight) \end{aligned}$$

Therefore, $\mathbb{E}\left[Y
ight]=2$

Similarly,

$$\mathbb{E}[XY] = \mathbb{E}_{T} [\mathbb{E}[XY|T]]$$

$$= \frac{1}{2} \mathbb{E}[XY|T = 1] + \frac{1}{2} \mathbb{E}[XY|T = 2]$$

$$= \frac{1}{2} \mathbb{E}[w_{1}X^{2} + \epsilon X] + \frac{1}{2} \mathbb{E}[w_{2}X^{2} + \epsilon X]$$

$$= \frac{1}{2} \mathbb{E}[X^{2}] + \frac{1}{2} \mathbb{E}[-X^{2}]$$

$$= 0$$

Submit

You have used 1 of 5 attempts

• Answers are displayed within the problem

5. (6)

1/2 points (graded)

Given the knowledge of $\operatorname{Cov}(X,Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$ and $\operatorname{Var}(X) = \operatorname{Cov}(X,X)$, what is the value of a^* and b^* in this concrete example, assuming we have infinite number of training data? Enter your solutions below

$$a^* = \boxed{f0}$$
 Answer: 0

$$b^* = \boxed{f 0}$$
 X Answer: 2

Solution:

$$egin{aligned} a^* &= rac{ ext{Cov}\left(X,Y
ight)}{ ext{Var}\left(X
ight)} \ &= rac{\mathbb{E}\left[XY
ight] - mathbbE\left[X
ight]\mathbb{E}\left[Y
ight]}{\mathbb{E}\left[X^2
ight] - \left(\mathbb{E}\left[X
ight]
ight)^2} \ &= 0 \ b^* &= \mathbb{E}\left[Y
ight] - a^*\mathbb{E}\left(X
ight) \ &= 2 \end{aligned}$$

Submit

You have used 2 of 5 attempts

1 Answers are displayed within the problem

5. (7)

1/1 point (graded)

Does this mean with infinite number of training data, the linear regression model is a good fit for this given scenario? In other words, is the linear regression model a good model for predicting Y from X for X for

Correction Note (Sept 3): An earlier version does not include the second sentence starting with "in other words".

O Yes				



Solution:

The linear regression gives us the solution $f(X) = a^*X + b^* = 2$ in this concrete example with infinite training data. Apparently this is not a good model to predict Y from X.

Submit

You have used 1 of 3 attempts

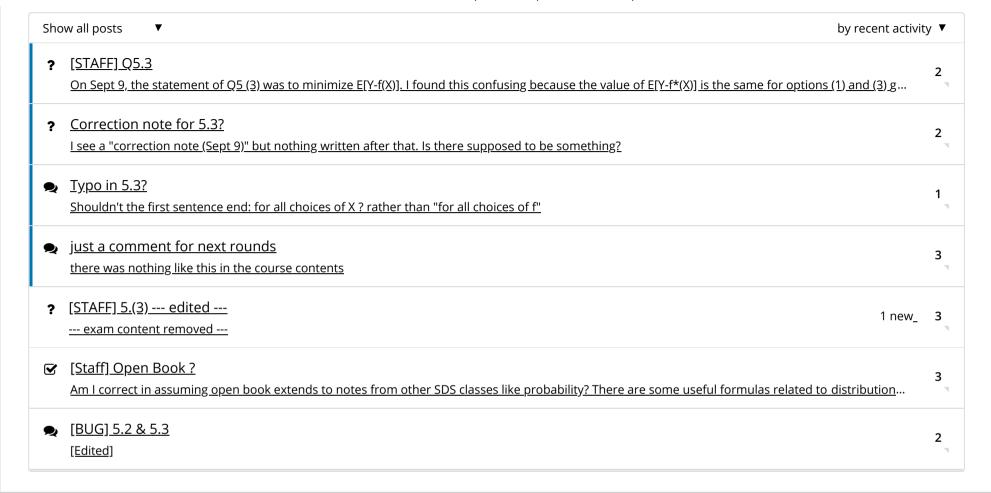
• Answers are displayed within the problem

Error and Bug Reports/Technical Issues

Topic: Final exam (1 week):Final Exam / Problem 5

Hide Discussion

Add a Post



© All Rights Reserved