

Assignment Two

ECE 4200

- Provide credit to **any sources** other than the course staff that helped you solve the problems. This includes **all students** you talked to regarding the problems.
- You can look up definitions/basics online (e.g., wikipedia, stack-exchange, etc)
- Submission rules are the same as previous assignment.

Problem 1 (15 points). In class we said that for a generative model (eg Naive Bayes), the optimal estimator will be the maximum a posteriori probability (MAP) estimator that when given a feature \vec{X} , outputs the label that satisfies:

$$\arg \max_{y \in \mathcal{Y}} p(y|\vec{X}).$$

The maximum likelihood (ML) estimator outputs the label satisfying:

$$\arg \max_{y \in \mathcal{Y}} p(\vec{X}|y).$$

In this problem we will see that this is the predictor with the least error probability if the underlying data is generated from the model.

We will simplify the setting by considering a binary classification task, where the labels have two possible values, say $\mathcal{Y} = \{-1, +1\}$. Suppose the model that generates the data is $p(\vec{X}, y)$.

1. What is the distribution of y when we observe a feature \vec{X} ?
2. Suppose we predict the label -1 upon seeing \vec{X} . Show that the probability of error is $p(y = +1|\vec{X})$.
3. Use this to argue that for any \vec{X} the prediction to minimize the error probability is

$$\max_{y \in \{-1, +1\}} p(y|\vec{X}).$$

This shows that the MAP estimator is the optimal estimator for the binary task. This also extends to larger \mathcal{Y} .

4. Show that if the distribution over the labels is uniform, namely $p(y = -1) = p(y = +1) = 0.5$, then the MAP estimator and ML estimator are the same.
5. Construct *any* generative model where the MAP and ML estimator are not the same.

Solution:

1. It is the conditional distribution $p(y|\vec{X})$.

2. Since $\hat{y} = -1$

$$\Pr(\hat{y} \neq y|\vec{X}) = p(y = +1|\vec{X})$$

3. If $\hat{y} = -1$, the error probability is $p(y = +1|\vec{X}) = 1 - p(y = -1|\vec{X})$ and if $\hat{y} = 1$, the error probability is $1 - p(y = 1|\vec{X})$. Hence the error is minimized if we take

$$\hat{y} = \arg \min_{y \in \{-1, 1\}} (1 - p(y|\vec{X})) = \arg \max_{y \in \{-1, 1\}} p(y|\vec{X}).$$

4.

$$\arg \max_{y \in \{-1, 1\}} p(y|\vec{X}) = \arg \max_{y \in \{-1, 1\}} \frac{p(\vec{X}|y)p(y)}{\sum_{y \in \{-1, 1\}} p(\vec{X}|y)p(y)} = \arg \max_{y \in \{-1, 1\}} p(\vec{X}|y)$$

5. Suppose you have two coins A with bias 0.9 and B with bias 0.1. Your prior is $p(A) = 0.01, p(B) = 0.99$. You toss a coin and observe it is a head.

Then by ML principle, we have:

$$\arg \max\{p(H|A), p(H|B)\} = \arg \max\{0.9, 0.1\} = A$$

By MAP principle, we have:

$$\begin{aligned} \arg \max\{p(H|A), p(H|B)\} &= \arg \max\left\{\frac{p(H|A)p(A)}{p(H|A)p(A) + p(H|B)p(B)}, \frac{p(H|B)p(B)}{p(H|A)p(A) + p(H|B)p(B)}\right\} \\ &= \arg \max\{p(H|A)p(A), p(H|B)p(B)\} = \arg \max\{0.9 \times 0.01, 0.1 \times 0.99\} = B \end{aligned}$$

Problem 2. (10 points). ML vs MAP and add constant smoothing. Suppose you generate n independent coin tosses using a coin with bias μ . What you get is n_H heads and $n_T = n - n_H$ tails. Show the following:

1. According to maximum likelihood principle, show that your estimate for μ should be:

$$\hat{\mu} = \frac{n_H}{n_H + n_T}.$$

2. Suppose that $p(\mu)$ is the *prior* distribution on the value of μ on the bias of the coin. Show that

$$\arg \max_{\mu} p(\mu|n_H, n_T) = \arg \max_{\mu} p(\mu)p(n_H, n_T|\mu)$$

3. Let the prior of the bias be a *Beta* distribution, which is a distribution over $[0, 1]$

$$p(\mu) = \frac{\mu^{\alpha}(1-\mu)^{\beta}}{\int_0^1 x^{\alpha}(1-x)^{\beta} dx}$$

show that:

$$\arg \max_{\mu} p(\mu|n_H, n_T) = \frac{n_H + \alpha}{n_H + \alpha + n_T + \beta}$$

Remark: This shows that add constant smoothing is equivalent to inducing a certain prior on the parameter of the generating model.

Solution:

1.

$$\Pr(n_H, n_T | \mu) = \binom{n_H + n_T}{n_H} \mu^{n_H} (1 - \mu)^{n_T}.$$

Take the derivative with respect to μ , we get:

$$\frac{\partial \Pr(n_H, n_T | \mu)}{\partial \mu} = \binom{n_H + n_T}{n_H} \mu^{n_H-1} (1 - \mu)^{n_T-1} (-n_H(1 - \mu) + n_T \mu)$$

Set it to be 0, we get:

$$\mu = \frac{n_H}{n_H + n_T}.$$

2.

$$\arg \max_{\mu} p(\mu | n_H, n_T) = \arg \max_{\mu} \frac{p(\mu) p(n_H, n_T | \mu)}{\sum_{\mu} p(\mu) p(n_H, n_T | \mu)} = \arg \max_{\mu} p(\mu) p(n_H, n_T | \mu)$$

3.

$$\begin{aligned} \arg \max_{\mu} p(\mu | n_H, n_T) &= \arg \max_{\mu} \frac{p(\mu) p(n_H, n_T | \mu)}{\sum_{\mu} p(\mu) p(n_H, n_T | \mu)} \\ &= \arg \max_{\mu} \binom{n_H + n_T}{n_H} \mu^{n_H} (1 - \mu)^{n_T} \frac{\mu^{\alpha} (1 - \mu)^{\beta}}{\int_0^1 x^{\alpha} (1 - x)^{\beta} d\mu} \\ &= \arg \max_{\mu} \mu^{n_H} (1 - \mu)^{n_T} \mu^{\alpha} (1 - \mu)^{\beta} \\ &= \arg \max_{\mu} \mu^{n_H + \alpha} (1 - \mu)^{n_T + \beta} \end{aligned}$$

Take the derivative and set it to be zero, we get:

$$\mu^{n_H + \alpha - 1} (1 - \mu)^{n_T + \beta - 1} (-(n_H + \alpha)(1 - \mu) + (n_T + \beta)\mu) = 0.$$

Hence

$$\mu = \frac{n_H + \alpha}{n_H + \alpha + n_T + \beta}$$

Problem 3 (10 points). Consider the Tennis data set.

1. For $\beta = 1$, write down the probabilities of all the features conditioned on the labels. The total number of probabilities you need to compute should not be more than twenty.

Solution:

When $\beta = 1$,

For the feature outlook,

$$\begin{aligned} P(\text{sunny} | \text{no}) &= \frac{3 + \beta}{5 + 3\beta} = \frac{1}{2}, & P(\text{sunny} | \text{yes}) &= \frac{2 + \beta}{9 + 3\beta} = \frac{1}{4}. \\ P(\text{over} | \text{no}) &= \frac{0 + \beta}{5 + 3\beta} = \frac{1}{8}, & P(\text{over} | \text{yes}) &= \frac{4 + \beta}{9 + 3\beta} = \frac{5}{12}. \\ P(\text{rain} | \text{no}) &= \frac{2 + \beta}{5 + 3\beta} = \frac{3}{8}, & P(\text{sunny} | \text{yes}) &= \frac{3 + \beta}{9 + 3\beta} = \frac{1}{3}. \end{aligned}$$

For the feature temperature,

$$\begin{aligned} P(\text{hot}|\text{no}) &= \frac{2+\beta}{5+3\beta} = \frac{3}{8}, & P(\text{hot}|\text{yes}) &= \frac{2+\beta}{9+3\beta} = \frac{1}{4}. \\ P(\text{mild}|\text{no}) &= \frac{2+\beta}{5+3\beta} = \frac{3}{8}, & P(\text{mild}|\text{yes}) &= \frac{4+\beta}{9+3\beta} = \frac{5}{12}. \\ P(\text{cool}|\text{no}) &= \frac{1+\beta}{5+3\beta} = \frac{1}{4}, & P(\text{cool}|\text{yes}) &= \frac{3+\beta}{9+3\beta} = \frac{1}{3}. \end{aligned}$$

For the feature humidity,

$$\begin{aligned} P(\text{high}|\text{no}) &= \frac{4+\beta}{5+2\beta} = \frac{5}{7}, & P(\text{high}|\text{yes}) &= \frac{3+\beta}{9+2\beta} = \frac{4}{11}. \\ P(\text{normal}|\text{no}) &= \frac{1+\beta}{5+2\beta} = \frac{2}{7}, & P(\text{normal}|\text{yes}) &= \frac{6+\beta}{9+2\beta} = \frac{7}{11}. \end{aligned}$$

For the feature wind

$$\begin{aligned} P(\text{weak}|\text{no}) &= \frac{2+\beta}{5+2\beta} = \frac{3}{7}, & P(\text{weak}|\text{yes}) &= \frac{6+\beta}{9+2\beta} = \frac{7}{11}. \\ P(\text{strong}|\text{no}) &= \frac{3+\beta}{5+2\beta} = \frac{4}{7}, & P(\text{strong}|\text{yes}) &= \frac{3+\beta}{9+2\beta} = \frac{4}{11}. \end{aligned}$$

2. What are the probabilities of the labels?

Solution:

$$P(\text{no}) = \frac{5}{14}, P(\text{yes}) = \frac{9}{14}.$$

3. For a new label (*Overcast, Hot, High, Strong*), what does the Naive Bayes classifier predict for $\beta = 0$, $\beta = 1$, and for $\beta \rightarrow \infty$?

Solution:

$$\begin{aligned} P(\text{over}|\text{no}) &= \frac{0+\beta}{5+3\beta}, & P(\text{over}|\text{yes}) &= \frac{4+\beta}{9+3\beta} \\ P(\text{hot}|\text{no}) &= \frac{2+\beta}{5+3\beta}, & P(\text{hot}|\text{yes}) &= \frac{2+\beta}{9+3\beta} \\ P(\text{high}|\text{no}) &= \frac{4+\beta}{5+2\beta}, & P(\text{high}|\text{yes}) &= \frac{3+\beta}{9+2\beta} \\ P(\text{strong}|\text{no}) &= \frac{3+\beta}{5+2\beta}, & P(\text{strong}|\text{yes}) &= \frac{3+\beta}{9+2\beta} \end{aligned}$$

$$\frac{P(\text{no}|\text{over,hot,high,strong})}{P(\text{yes}|\text{over,hot,high,strong})} = \frac{P(\text{over}|\text{no})}{P(\text{over}|\text{yes})} \cdot \frac{P(\text{hot}|\text{no})}{P(\text{hot}|\text{yes})} \cdot \frac{P(\text{high}|\text{no})}{P(\text{high}|\text{yes})} \cdot \frac{P(\text{strong}|\text{no})}{P(\text{strong}|\text{yes})} \cdot \frac{P(\text{no})}{P(\text{yes})}$$

When $\beta = 0$,

$$\frac{P(\text{no}|\text{over,hot,high,strong})}{P(\text{yes}|\text{over,hot,high,strong})} = 0 < 1$$

So the classifier outputs yes.

When $\beta = 1$,

$$\frac{P(\text{no}|\text{over,hot,high,strong})}{P(\text{yes}|\text{over,hot,high,strong})} = \frac{\frac{1}{8} \cdot \frac{3}{8} \cdot \frac{5}{7} \cdot \frac{4}{7} \cdot \frac{5}{14}}{\frac{5}{12} \cdot \frac{1}{4} \cdot \frac{4}{11} \cdot \frac{4}{11} \cdot \frac{9}{14}} < 1$$

So the classifier outputs yes.

When $\beta = \infty$

$$P(\text{over}|\text{no}) = P(\text{over}|\text{yes}) = \frac{1}{3}$$

$$P(\text{hot}|\text{no}) = P(\text{hot}|\text{yes}) = \frac{1}{3}$$

$$P(\text{high}|\text{no}) = P(\text{high}|\text{yes}) = \frac{1}{2}$$

$$P(\text{strong}|\text{no}) = P(\text{strong}|\text{yes}) = \frac{1}{2}$$

$$\frac{P(\text{no}|\text{over,hot,high,strong})}{P(\text{yes}|\text{over,hot,high,strong})} = \frac{\frac{5}{14}}{\frac{9}{14}} < 1$$

So the classifier outputs yes.

4. The prediction for (*Overcast, Hot, High, Strong*) is *Yes* and for (*Rain, Cool, High, Strong*) is *No*.

Problem 4 (15 points). Recall the Gaussian distribution with mean μ , and variance σ^2 . The density is given by:

$$p_{\mu,\sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Given n independent samples X_1, \dots, X_n from a Gaussian distribution with unknown mean, and variance, let μ_{ML} , and σ_{ML}^2 denote the maximum likelihood estimates of mean and variance.

1. Show that

$$\mu_{ML} = \frac{\sum_{i=1}^n X_i}{n}.$$

2. Show that

$$\sigma_{ML}^2 = \frac{1}{n} \left(\sum_{i=1}^n (X_i - \mu_{ML})^2 \right).$$

3. What is the expectation and variance of μ_{ML} ?

Solution:

1 and 2.

The likelihood can be written as follows

$$\begin{aligned}
 p(X_1, X_2 \cdots X_n | \mu, \sigma) &= \prod_{i=1}^n p(X_i | \mu, \sigma) \\
 &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right) \\
 &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2}\right)
 \end{aligned} \tag{1}$$

By monotonicity of logarithms it suffices to maximize the log-likelihood:

$$\log p(X_1, X_2 \cdots X_n | \mu, \sigma) = n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2} = -\frac{1}{2}n \log 2\pi\sigma^2 - \frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2} \tag{2}$$

$$\hat{\mu}_{ML} = \arg \max_{\mu} \log(p(X_1, X_2 \cdots X_n | \mu, \sigma)).$$

$$\hat{\sigma}_{ML} = \arg \max_{\sigma} \log(p(X_1, X_2 \cdots X_n | \mu, \sigma)).$$

Taking the partial derivative of (2) w.r.t. μ , and noting that the first term is independent of μ , we obtain:

$$\left. \frac{\partial \log p(X_1, X_2 \cdots X_n | \mu, \sigma)}{\partial \mu} \right|_{\mu=\hat{\mu}_{ML}, \sigma=\hat{\sigma}_{ML}} = \frac{\sum_{i=1}^n (X_i - \hat{\mu}_{ML})}{\sigma^2} = 0.$$

Therefore,

$$\hat{\mu}_{ML} = \frac{\sum_{i=1}^n X_i}{n}$$

$$\left. \frac{\partial \log p(X_1, X_2 \cdots X_n | \mu, \sigma)}{\partial \sigma} \right|_{\mu=\hat{\mu}_{ML}, \sigma=\hat{\sigma}_{ML}} = -\frac{n}{\hat{\sigma}_{ML}} + \frac{\sum_{i=1}^n (X_i - \hat{\mu}_{ML})^2}{\hat{\sigma}_{ML}^3} = 0.$$

$$\left. \frac{\partial \log p(X_1, X_2 \cdots X_n | \mu, \sigma)}{\partial \sigma^2} \right|_{\mu=\hat{\mu}_{ML}, \sigma=\hat{\sigma}_{ML}} = -\frac{n}{2\hat{\sigma}_{ML}^2} + \frac{\sum_{i=1}^n (X_i - \hat{\mu}_{ML})^2}{2\hat{\sigma}_{ML}^4} = 0.$$

Therefore,

$$\hat{\sigma}_{ML}^2 = \frac{\sum_{i=1}^n (X_i - \hat{\mu}_{ML})^2}{n}.$$

3.

$$\mathbb{E}[\hat{\mu}_{ML}] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \mu.$$

$$\text{Var}(\hat{\mu}_{ML}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n}$$