

CSE 446: Homework 1  
ayush29f@cs.washington.edu  
Ayush Saraf

---

**Problem 1**

---

Let  $T = \text{Test is positive}$ ,  $D = \text{Person has the disease}$

$$P(D) = 10^{-4}$$

$$P(T|D) = 0.99$$

Now we have to calculate  $P(D|T)$  using Bayes' Rule;

$$P(D|T) = \frac{P(T|D)P(D)}{P(T)}$$

$$P(D|T) = \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|D')P(D')}$$

$$P(D|T) = \frac{0.99 * 10^{-4}}{0.99 * 10^{-4} + (1 - 0.99) * (1 - 10^{-4})}$$

$$P(D|T) = 0.0098$$

---

**Problem 2.1**

---

In order to solve this question we will maximize the log likelihood ( $L$ ) of the this distribution to be part of Poisson Distribution.

$$L(\lambda|\mathbf{x}) = \prod_{i=1}^6 \text{Poi}(x_i|\lambda)$$

$$L(\lambda|\mathbf{x}) = \prod_{i=1}^6 e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}$$

Taking Log on both sides,

$$LL(\lambda|\mathbf{x}) = \sum_{i=1}^6 -\lambda + x_i \ln(\lambda) - \ln(x_i!)$$

Setting the derivative to 0,

$$\frac{\partial LL(\lambda|\mathbf{x})}{\partial \lambda} = \sum_{i=1}^6 \frac{x_i}{\lambda} - 1 = 0$$

$$\lambda = \sum_{i=1}^6 \frac{x_i}{6}$$

$$\lambda = \frac{13}{6}$$

Checking 2nd derivative for maxima or minima,

$$\frac{\partial^2 LL(\lambda|\mathbf{x})}{\partial \lambda^2} = \sum_{i=1}^6 \frac{-x_i}{\lambda^2} < 0$$

plugging  $\lambda = 13/6$ , and since  $x_i > 0$  this is a maxima and therefore maximizing Likelihood at  $\lambda = 13/6$

## Problem 2.2

Given,  $n_{clear} = 11$ ,  $n_{cloudy} = 24$ ,  $n_{rain} = 75$

Calculate  $p_{clear}$ ,  $p_{cloudy}$ ,  $p_{rain}$  We know,

$$p_{clear} + p_{cloudy} + p_{rain} = 1$$

$$L(p_{clear}, p_{cloudy} | n_{clear}, n_{cloudy}, n_{rain}) = \frac{110!}{11! * 24! * 75!} p_{clear}^{11} p_{cloudy}^{24} (1 - p_{clear} - p_{cloudy})^{75}$$

Taking Log on both sides,

$$LL(p_{clear}, p_{cloudy} | \mathbf{n}) = \log\left(\frac{110!}{11! * 24! * 75!}\right) + 11\log(p_{clear}) + 24\log(p_{cloudy}) + 75\log(1 - p_{clear} - p_{cloudy})$$

Take partial derivatives and set them to 0

$$\frac{\partial LL(p_{clear}, p_{cloudy} | \mathbf{n})}{\partial p_{clear}} = \frac{11}{p_{clear}} - \frac{75}{(1 - p_{clear} - p_{cloudy})} = 0$$

$$\frac{\partial LL(p_{clear}, p_{cloudy} | \mathbf{n})}{\partial p_{cloudy}} = \frac{24}{p_{cloudy}} - \frac{75}{(1 - p_{clear} - p_{cloudy})} = 0$$

Calculating for 2 variables and 2 equations,

$$p_{clear} = \frac{11}{11 + 24 + 75} = \frac{11}{110}$$

$$p_{cloudy} = \frac{24}{11 + 24 + 75} = \frac{24}{110}$$

$$p_{rain} = 1 - p_{clear} - p_{cloudy} = \frac{75}{110}$$

Checking 2nd derivative for confirming it is maxima for these values of  $p$ .

---

**Problem 3.1**

---

In order to estimate we only need 2 statistics from the ones shown above:  $C_{xx}^{(n)}, C_{xy}^{(n)}$

---

**Problem 3.2**

---

In order to estimate we only need 4 statistics from the ones shown above:  $\bar{x}^{(n)}, \bar{y}^{(n)}, C_{xx}^{(n)}, C_{xy}^{(n)}$

---

**Problem 3.3**

---

According to the definition of mean,

$$\begin{aligned}\bar{x}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} x_i \\ (n+1)\bar{x}^{(n+1)} &= \sum_{i=1}^{n+1} x_i \\ \frac{(n+1)}{n} \bar{x}^{(n+1)} &= \frac{1}{n} \sum_{i=1}^{n+1} x_i \\ \frac{(n+1)}{n} \bar{x}^{(n+1)} &= \frac{1}{n} \sum_{i=1}^n x_i + \frac{x_{n+1}}{n} \\ \frac{(n+1)}{n} \bar{x}^{(n+1)} &= \bar{x}^{(n)} + \frac{x_{n+1}}{n} \\ \frac{(n+1)}{n} \bar{x}^{(n+1)} &= \bar{x}^{(n)} + \frac{x_{n+1}}{n} \\ \bar{x}^{(n+1)} &= \frac{n}{(n+1)} \left( \bar{x}^{(n)} + \frac{x_{n+1}}{n} \right) \\ \bar{x}^{(n+1)} &= \frac{n\bar{x}^{(n)}}{(n+1)} + \frac{x_{n+1}}{n+1} \\ \bar{x}^{(n+1)} &= \bar{x}^{(n)} - \frac{\bar{x}^{(n)}}{(n+1)} + \frac{x_{n+1}}{n+1} \\ \bar{x}^{(n+1)} &= \bar{x}^{(n)} - \frac{\bar{x}^{(n)}}{(n+1)} + \frac{x_{n+1}}{n+1}\end{aligned}$$

Hence Proved,

$$\bar{x}^{(n+1)} = \bar{x}^{(n)} + \frac{1}{n+1} \left( x_{n+1} - \bar{x}^{(n)} \right)$$

Similarly for  $\bar{y}^{(n)}$

$$\begin{aligned}
\bar{y}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} y_i \\
(n+1)\bar{y}^{(n+1)} &= \sum_{i=1}^{n+1} y_i \\
\frac{(n+1)}{n} \bar{y}^{(n+1)} &= \frac{1}{n} \sum_{i=1}^{n+1} y_i \\
\frac{(n+1)}{n} \bar{y}^{(n+1)} &= \frac{1}{n} \sum_{i=1}^n y_i + \frac{y_{n+1}}{n} \\
\frac{(n+1)}{n} \bar{y}^{(n+1)} &= \bar{y}^{(n)} + \frac{y_{n+1}}{n} \\
\frac{(n+1)}{n} \bar{y}^{(n+1)} &= \bar{y}^{(n)} + \frac{y_{n+1}}{n} \\
\bar{y}^{(n+1)} &= \frac{n}{(n+1)} \left( \bar{y}^{(n)} + \frac{y_{n+1}}{n} \right) \\
\bar{y}^{(n+1)} &= \frac{n\bar{y}^{(n)}}{(n+1)} + \frac{y_{n+1}}{n+1} \\
\bar{y}^{(n+1)} &= \bar{y}^{(n)} - \frac{\bar{y}^{(n)}}{(n+1)} + \frac{y_{n+1}}{n+1} \\
\bar{y}^{(n+1)} &= \bar{y}^{(n)} - \frac{\bar{y}^{(n)}}{(n+1)} + \frac{y_{n+1}}{n+1}
\end{aligned}$$

Hence Proved,

$$\bar{y}^{(n+1)} = \bar{y}^{(n)} + \frac{1}{n+1} \left( y_{n+1} - \bar{y}^{(n)} \right)$$

---

### Problem 3.4

---

According to the definition of  $C_{xy}$ ,

$$\begin{aligned}
C_{xy}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} (x_i - \bar{x}^{(n+1)})(y_i - \bar{y}^{(n+1)}) \\
C_{xy}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} (x_i y_i - x_i \bar{y}^{(n+1)} - y_i \bar{x}^{(n+1)} + \bar{x}^{(n+1)} \bar{y}^{(n+1)}) \\
C_{xy}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} \left( x_i y_i - x_i \bar{y}^{(n+1)} - y_i \bar{x}^{(n+1)} \right) + \bar{x}^{(n+1)} \bar{y}^{(n+1)}
\end{aligned}$$

$$C_{xy}^{(n+1)} = \frac{1}{n+1} \sum_{i=1}^{n+1} x_i y_i - \frac{1}{n+1} \sum_{i=1}^{n+1} x_i \bar{y}^{(n+1)} - \frac{1}{n+1} \sum_{i=1}^{n+1} y_i \bar{x}^{(n+1)} + \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{1}{n+1} \left( \sum_{i=1}^{n+1} x_i y_i \right) - 2\bar{x}^{(n+1)} \bar{y}^{(n+1)} + \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{1}{n+1} \left( \sum_{i=1}^{n+1} x_i y_i \right) - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{x_{n+1} y_{n+1}}{n+1} + \frac{1}{n+1} \left( \sum_{i=1}^n x_i y_i \right) - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{x_{n+1} y_{n+1}}{n+1} + \frac{1}{n+1} \left( \sum_{i=1}^n x_i y_i - x_i \bar{y}^{(n)} - y_i \bar{x}^{(n)} + \bar{x}^{(n)} \bar{y}^{(n)} + x_i \bar{y}^{(n)} + y_i \bar{x}^{(n)} - \bar{x}^{(n)} \bar{y}^{(n)} \right) - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{x_{n+1} y_{n+1}}{n+1} + \frac{1}{n+1} \left( \sum_{i=1}^n x_i y_i - x_i \bar{y}^{(n)} - y_i \bar{x}^{(n)} + \bar{x}^{(n)} \bar{y}^{(n)} \right) + \frac{1}{n+1} \sum_{i=1}^n x_i \bar{y}^{(n)} +$$

$$\frac{1}{n+1} \sum_{i=1}^n y_i \bar{x}^{(n)} - \frac{1}{n+1} \sum_{i=1}^n \bar{x}^{(n)} \bar{y}^{(n)} - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{x_{n+1} y_{n+1}}{n+1} + \frac{n C_{xy}^{(n)}}{n+1} + \frac{1}{n+1} \sum_{i=1}^n x_i \bar{y}^{(n)} + \frac{1}{n+1} \sum_{i=1}^n y_i \bar{x}^{(n)} - \frac{1}{n+1} \sum_{i=1}^n \bar{x}^{(n)} \bar{y}^{(n)} - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{x_{n+1} y_{n+1}}{n+1} + \frac{n C_{xy}^{(n)}}{n+1} + \frac{n \bar{y}^{(n)} \bar{x}^{(n)}}{n+1} - \bar{x}^{(n+1)} \bar{y}^{(n+1)}$$

$$C_{xy}^{(n+1)} = \frac{1}{n+1} \left( x_{n+1} y_{n+1} + n C_{xy}^{(n)} + n \bar{y}^{(n)} \bar{x}^{(n)} - (n+1) \bar{x}^{(n+1)} \bar{y}^{(n+1)} \right)$$

According to the definition of  $C_{xx}$ ,

$$\begin{aligned}
C_{xx}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} (x_i - \bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{1}{n+1} \sum_{i=1}^{n+1} (x_i^2 + (\bar{x}^{(n+1)})^2 - 2x_i\bar{x}^{(n+1)}) \\
C_{xx}^{(n+1)} &= \frac{1}{n+1} \left( \sum_{i=1}^{n+1} x_i^2 \right) - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{n}{n(n+1)} \left( \sum_{i=1}^n x_i^2 \right) + \frac{x_{n+1}^2}{n+1} - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{n}{n(n+1)} \left( \sum_{i=1}^n x_i^2 + (\bar{x}^{(n)})^2 - (\bar{x}^{(n)})^2 - 2x_i\bar{x}^{(n)} + 2x_i\bar{x}^{(n)} \right) + \frac{x_{n+1}^2}{n+1} - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{n}{n(n+1)} \left( \sum_{i=1}^n (x_i - \bar{x}^{(n)})^2 - (\bar{x}^{(n)})^2 + 2x_i\bar{x}^{(n)} \right) + \frac{x_{n+1}^2}{n+1} - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{n}{n(n+1)} \sum_{i=1}^n (x_i - \bar{x}^{(n)})^2 - \frac{n}{n(n+1)} \sum_{i=1}^n (\bar{x}^{(n)})^2 + \frac{n}{n(n+1)} \sum_{i=1}^n 2x_i\bar{x}^{(n)} + \frac{x_{n+1}^2}{n+1} - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{n}{(n+1)} C_{xx}^{(n)} + \frac{n(\bar{x}^{(n)})^2}{(n+1)} + \frac{x_{n+1}^2}{n+1} - (\bar{x}^{(n+1)})^2 \\
C_{xx}^{(n+1)} &= \frac{1}{n+1} \left( nC_{xx}^{(n)} + n(\bar{x}^{(n)})^2 + x_{n+1}^2 - (n+1)(\bar{x}^{(n+1)})^2 \right)
\end{aligned}$$

---

### Problem 3.5

---

The two examples could be considered for online regression over batch regression:

1. Flight Price Prediction: Since there are so many flights with multiple incoming data points every minute across the world, it makes much more sense to have an online model that holds the statistics instead of recomputing the parameters
2. Twitter Retweet Prediction: Since there are so tweets and retweets with multiple incoming data points every millisecond across the world, it makes much more sense to have an online model that holds the statistics instead of recomputing the parameters

---

**Problem 4.1.1**

---

- (a) The error on the training set will decrease compared to where  $\lambda > 0$  because it will over fit to training set
- (b) The error on the training set will increase compared to where  $\lambda > 0$  because it will over fit to training set and not test set
- (c) the magnitude of the values in  $\hat{w}$  will be larger because we are not penalizing the  $L1(\hat{w})$  in the loss function
- (d) the number of non-zero elements in  $\hat{w}$  will be more because there will be no feature selection

---

**Problem 4.1.2**

---

- (a) The training error on the training set will increase because now a lot of features are eliminated and it will under fit to the training set
- (b) The testing error will also increase for similar reasons given that we over estimated the  $\lambda$  because the model has eliminated too many features
- (c) the magnitude of the values in  $\hat{w}$  will be really small because we are penalizing the  $L1(\hat{w})$  too much in the loss function
- (d) the number of non-zero elements in  $\hat{w}$  will be lower because there will be a lot of feature selection

---

**Problem 4.2**

---

1. Taking the derivative w.r.t. to  $\hat{w}_i$

$$f(\hat{w}) = \lambda \|\hat{w}\|_1$$

$$f(\hat{w}) = \lambda \sum_i |\hat{w}_i|$$

$$\frac{\partial f(\hat{w})}{\partial \hat{w}_i} = \lambda \quad \text{if } \hat{w}_i > 0$$

$$\frac{\partial f(\hat{w})}{\partial \hat{w}_i} = -\lambda \quad \text{if } \hat{w}_i < 0$$

2. Taking the derivative w.r.t. to  $\hat{w}_i$

$$f(\hat{w}) = \lambda \|\hat{w}\|_2$$

$$f(\hat{w}) = \lambda \sum_i \hat{w}_i^2$$

$$\frac{\partial f(\hat{w})}{\partial \hat{w}_i} = 2\lambda \hat{w}_i$$

- 3.
- In the case of  $\hat{w}_i < 1/2$ , lasso will push more strongly again away from SSE solution because the differential has a magnitude of  $\lambda$  regardless the value of  $\hat{w}_i$  unlike ridge where the magnitude of  $2\lambda \hat{w}_i$  and if  $\hat{w}_i < 1/2$  then the differential will be lower and will have lower impact.
  - Similarly, if  $\hat{w}_i > 1/2$  ridge will have higher impact because now the  $2\lambda \hat{w}_i > \lambda$ .
  - If we have a really large  $\lambda$  then there will be a lot of regularization on the features in both cases. However, when we have a co related features ridge will end up pushing all the co-relating features to lower values and hence having a worse model. While in ridge regression since we use co-ordinate decent we will have one of the co-relating feature which is higher  $\hat{w}_i$  value and others will be pushed to 0 and hence eliminating the relative features.