

# Homework 1: k-Nearest Neighbor classification and Statistical Estimation

## #1 Statistical Estimation

1)  $D = \{x_1, x_2, \dots, x_n\}$  from  $N$

Pois ( $x = x_i | \lambda$ )

$$p(x=x) = \frac{e^{-\lambda} \lambda^x}{x!}, n=0,1,2,\dots$$

$$P(D|\lambda) = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$$

$$\log P(D|\lambda) = -n\lambda + \sum_{i=1}^n x_i \log \lambda - \log \left[ \prod_{i=1}^n x_i! \right]$$

$$\frac{\partial}{\partial \lambda} \log P(D|\lambda) = -n + \frac{\sum_{i=1}^n x_i}{\lambda} = 0$$

$$\Rightarrow \frac{\sum_{i=1}^n x_i}{\lambda} = n \Rightarrow \hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\Rightarrow \hat{\lambda} = \bar{x} \rightarrow \lambda \text{ is the average \# of occurrence}$$

2)  $\odot$  SO,  $f(x_i|\lambda) = e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}$ ,  $x_i = 0, 1, 2, \dots$

$$\Rightarrow f(x) = \frac{e^{-\lambda} \lambda^{x-1}}{(x-1)!} \beta^{-1}, \lambda > 0$$

$$\text{SO, } f(x_i, \lambda) = \frac{N}{\prod_{i=1}^N} f(x_i | \lambda) f(\lambda)$$

$$= e^{-N\lambda} \frac{\lambda^{\sum_{i=1}^N x_i}}{\prod_{i=1}^N x_i!} = \frac{e^{-N\lambda} \lambda^{N-1}}{(N-1)!} \beta^{-1}$$

$$\Rightarrow f(x) = \int_0^\infty f(x, \lambda) d\lambda$$

$$= \int_0^\infty \frac{e^{-N\lambda}}{\prod_{i=1}^N x_i!} d\lambda = e^{-N\lambda} e^{-\lambda} \lambda^{\sum_{i=1}^N x_i} \lambda^{-1} d\lambda$$

$$= \frac{e^{-N\lambda}}{\prod_{i=1}^N x_i!} \int_0^\infty e^{-\lambda(N+1)} \lambda^{\sum_{i=1}^N x_i + \alpha - 1} d\lambda$$

$$\text{Here, } \Rightarrow P(X|D) = \frac{f(x)}{f(\lambda)}$$

$$\Rightarrow \frac{e^{-N\lambda}}{\prod_{i=1}^N x_i!} \frac{e^{-\lambda(N+1)} \lambda^{\sum_{i=1}^N x_i + \alpha - 1}}{e^{-N\lambda} \frac{\lambda^{\sum_{i=1}^N x_i}}{\prod_{i=1}^N x_i!}}$$

$$\Rightarrow \frac{1}{\Gamma(\sum_{i=1}^N x_i + \alpha)} e^{-\lambda(N+1)} \lambda^{\sum_{i=1}^N x_i + \alpha - 1}$$

$$\Rightarrow \frac{1}{\Gamma(\sum_{i=1}^N x_i + \alpha)} e^{-\lambda(N+1)} \lambda^{\sum_{i=1}^N x_i + \alpha - 1}$$

$$\ln P(X|D) \Rightarrow -\ln \left[ \Gamma\left(\sum_{i=1}^N x_i + \alpha\right) \right] - \lambda(N+1) + \left(\sum_{i=1}^N x_i + \alpha - 1\right) \ln \lambda + \left(\sum_{i=1}^N x_i + \alpha\right) \ln(N+1)$$

$$= C - \lambda(N+1) + \ln \lambda \left[ \sum_{i=1}^N x_i + \alpha - 1 \right]$$

$$= \frac{d}{d\lambda} \cdot P(X|D) = 0 + \frac{d}{d\lambda} \left[ -\lambda(N+1) \right] + \frac{1}{d\lambda} \left[ \ln \lambda \left[ \sum_{i=1}^N x_i + \alpha - 1 \right] \right]$$

$$= -(N+1) + \frac{\sum_{i=1}^N x_i + \alpha - 1}{\lambda}$$

$$\text{SO, } \frac{d}{d\lambda} \ln P(X|D) = 0$$

$$\Rightarrow - (N+1) + \frac{\left(\sum_{i=1}^N x_i + \alpha - 1\right)}{\lambda} = 0$$

$$\Rightarrow \left(\sum_{i=1}^N x_i + \alpha - 1\right) = \lambda(N+1)$$

$$\Rightarrow \lambda = \frac{\sum_{i=1}^N x_i + \alpha - 1}{N+1} = \hat{\lambda}$$

3) MAP estimator means maximum of a posterior estimator.

for example, here in this problem the maximum value of a

parameter  $\lambda$  which maximizes the posterior probability

distribution of  $P(\lambda|D)$ . Therefore,  $\lambda = \frac{\sum_{i=1}^N x_i + \alpha - 1}{N+1}$  is a MAP

estimator of  $\lambda$

3) Let  $D|\lambda = y_1, y_2, \dots, y_n | \lambda \stackrel{iid}{\sim} \text{Poisson}(\lambda)$  and  $\lambda \sim \text{Gamma}(\alpha, \beta)$

The Conditional joint density of  $y_1, y_2, \dots, y_n | \lambda$  is  $f(y_1, y_2, \dots, y_n | \lambda) = \prod_{i=1}^n \left( \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} \right) = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n y_i}}{\prod_{i=1}^n y_i!}$

The Prior is  $\pi(\lambda) = \frac{e^{-\beta\lambda} \lambda^{\alpha-1} \beta^\alpha}{\Gamma(\alpha)} \quad \lambda > 0$

The posterior distribution  $\Rightarrow \pi(\lambda | D) = \pi(\lambda | y_1, \dots, y_n) = \frac{f(y_1, y_2, \dots, y_n | \lambda) \pi(\lambda)}{\int_0^\infty f(y_1, \dots, y_n | \lambda) \pi(\lambda) d\lambda}$  This is independent of  $\lambda$ , meaning

$$\lambda | D \sim \text{Gamma}(\alpha_p, \beta_p) \text{ where } \alpha_p = \sum_{i=1}^n y_i + \alpha - 1, \quad \beta_p = \beta + n$$

$$\pi(\lambda | D) \propto f(y_1, y_2, \dots, y_n | \lambda) \pi(\lambda) \Rightarrow \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n y_i}}{\left( \prod_{i=1}^n y_i! \right)} \times \frac{e^{-\beta\lambda} \lambda^{\alpha-1} \beta^\alpha}{\Gamma(\alpha)}$$

Hence the gamma distribution is a conjugate prior to the poisson as seen.

$$\propto e^{-\lambda(\beta+n)} \cdot \lambda^{\left(\sum_{i=1}^n y_i + \alpha - 1\right)}$$

## #2 K-Nearest Neighbor (KNN)

4) One-hot encoding technique would make it hard when calculating euclidean distance. The values are all 1s and 0s. Since it is all zeros and ones this can result in a zero value for each label. (This can be handled by setting Unknown to ignore)

↳ example:  $\begin{bmatrix} 1, 0, 0 \end{bmatrix}, \begin{bmatrix} 0, 1, 0 \end{bmatrix}, \begin{bmatrix} 0, 0, 1 \end{bmatrix}$

Categorical variables as ordinals would potential provide smaller matrixes when solving for Euclidean distance in knn.

↳ example:  $\begin{bmatrix} 1 \end{bmatrix}, \begin{bmatrix} 2 \end{bmatrix}, \begin{bmatrix} 3 \end{bmatrix}$

5) 100% of the training data has an income > 50k. This can affect the model because we only have data from higher income individuals. This means if it was trying to make a prediction for someone who makes <= 50k it would be very lost because it has no data to go off of. 70% accuracy is good depending on what the model is testing for, whether it is good or not will vary b/w situations. This model has 12 dimensions

6) Let  $x = (x_1, x_2, \dots, x_d)^T$  be a vector. Zero vector  $\Rightarrow x = (0, 0, \dots, 0)^T$  so, The Euclidean distance b/w  $x$  and the zero vector  $\Rightarrow \sqrt{(x_1-0)^2 + (x_2-0)^2 + \dots + (x_d-0)^2} \Rightarrow \|x\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$

Let  $z = (z_1, z_2, \dots, z_d)^T$  The Euclidean distance b/w  $x$  and  $z$  vector  $\Rightarrow \sqrt{(x_1-z_1)^2 + (x_2-z_2)^2 + \dots + (x_d-z_d)^2}$

$\Rightarrow \|x-z\|_2 = \sqrt{\sum_{i=1}^d (x_i-z_i)^2} \rightarrow$  implies that, distance b/w  $x$  and  $z$  can be written as an  $L_2$  norm

7) coding

8) coding

9) The best  $k$  observed is 9 or 7 because overall they had the highest training accuracy and validation accuracy. When  $k=1$  the training error is close to 0% when you choose a test sample closest to your training set.  $k=1$  will not be zero when the test sample is not close to your training set. A trend you may observe is the higher  $k$  is the lower the training accuracy but the higher the validation accuracy is.

Underfitting occurs when the training data does not quite fit the trends of real data, meaning our model does not fit well enough. This can kinda be seen by lower  $k$ 's like  $k$  of 1 because validation is 78.72% but it does well with training, getting a 98.74, so it's not the best example.

An example of overfitting can be seen when  $k$  is too large and it causes training to go awry. This can be seen with  $k=8000$  and train acc. = 75.41%.

10)

### #3 Debriefing

- 1) 3.5 days (yes I barely started it the day before it was due and no that will not be happening again)
- 2) difficult, I was literally brought to tears
- 3) Mostly alone, was gonna ask Tas but I realized there are no office hours on Thursdays and I missed Wednesday but Fridays was helpful
- 4) I think at first 30% but after struggling through maybe like 85% now
- 5) sorry it's so slow, I did not want to code anymore