

## 1 Maximizing Likelihood &amp; Minimizing Cost

Maximum Likelihood Estimation (MLE) is a method for estimating the parameters of a statistical model given observations.

**Data** Suppose we obtain  $n$  discrete *observations* belonging to  $B := \{1, 2, 3, 4\}$ . Our dataset looks something like the following.

$$\begin{aligned} r_1 &= 1 \\ r_2 &= 1 \\ r_3 &= 3 \\ &\vdots \\ r_n &= 1 \end{aligned}$$

**Assumptions** Suppose we aim to estimate the occurrence probabilities of each class in  $B$  based on the observed data. We additionally assume that observations are independent and identically distributed (i.i.d.). In particular, this assumption implies that the order of the data does not matter.

**Model** Based on these assumptions, a natural model for our data is the multinomial distribution. In a multinomial distribution, the order of the data does not matter, and we can equivalently represent our dataset as  $(y, c_y)_{y \in B}$ , where  $c_y$  is the number of items of class  $y$ .

The probability mass function (PMF) of the multinomial distribution—this is, the probability in  $n$  trials of obtaining each class  $i$   $x_i$  times—is

$$P(x_1, \dots, x_k) = n! \prod_{i=1}^k \frac{p_i^{x_i}}{x_i!}.$$

- (a) Derive an expression for the likelihood for this problem. What are the observations? What are the parameters? What parameters are we trying to estimate with MLE?

observations :  $x_1, \dots, x_k$

parameters :  $p_1, \dots, p_k, n, k$

MLE :  $p_1, \dots, p_k$

- (b) Typically, the log-likelihood  $\ell(\theta) = \log L(\theta)$  is used instead of  $L(\theta)$ . Write down the expression for  $\ell(\theta)$ . Why might this be a good idea?

$$\begin{aligned} \ell(\theta) = \log L(\theta) &= \log \left( n! \prod_{i=1}^k \frac{p_i^{x_i}}{x_i!} \right) = \left( \sum_{i=1}^n \log i \right) + \sum_{i=1}^k \log \left( \frac{p_i^{x_i}}{x_i!} \right) \\ &= \left( \sum_{i=1}^n \log i \right) + \sum_{i=1}^k \left( x_i \log p_i - \sum_{j=1}^{x_i} \log j \right) \end{aligned}$$

- (c) Another idea might be to minimize the cross-entropy based on raw observations, corresponding to the following program

$$\underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmin}} - \sum_{i=1}^n \sum_{y \in B} \delta_{r_i y} \log p_y$$

where  $p$  is the vector of probabilities per class  $[p_1 \ p_2 \ p_3 \ p_4]^T$ , and  $\delta_{r_i y}$  is the Kronecker delta that outputs 1 if  $r_i = y$  and 0 otherwise.

Show that this program is equivalent to the MLE program.

$$\begin{aligned} \underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmax}} \mathcal{L}(p; x_1, x_2, x_3, x_4) &= \underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmax}} \left( \log n! + \sum_{i=1}^k x_i \log p_i - \sum_{i=1}^k \log x_i! \right) \\ &= \underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmax}} \sum_{i=1}^k x_i \log p_i \\ &= \underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmax}} \sum_{i=1}^n \sum_{y \in B} \delta_{r_i y} \log p_y \\ &= \underset{\substack{p \in \mathbb{R}_+^4 \\ \|p\|_1 = 1}}{\operatorname{argmin}} - \sum_{i=1}^n \sum_{y \in B} \delta_{r_i y} \log p_y \end{aligned}$$

## 2 Independence and Multivariate Gaussians

As described in lecture, a covariance matrix  $\Sigma \in \mathbb{R}^{N \times N}$  for a random variable  $X \in \mathbb{R}^N$  with the following values, where  $\text{cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$  is the covariance between the  $i$ -th and  $j$ -th elements of the random vector  $X$ :

$$\Sigma = \begin{bmatrix} \text{cov}(X_1, X_1) & \dots & \text{cov}(X_1, X_n) \\ \dots & & \dots \\ \text{cov}(X_n, X_1) & \dots & \text{cov}(X_n, X_n) \end{bmatrix}. \quad (1)$$

Recall that the density of an  $N$  dimensional Multivariate Gaussian Distribution  $\mathcal{N}(\mu, \Sigma)$  is defined as follows when  $\Sigma$  is positive definite:

$$f(x) = \frac{1}{\sqrt{(2\pi)^N |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}. \quad (2)$$

Here,  $|\Sigma|$  denotes the determinant of the matrix  $\Sigma$ .

(a) Consider the random variables  $X$  and  $Y$  in  $\mathbb{R}$  with the following conditions.

- (i)  $X$  and  $Y$  can take values  $\{-1, 0, 1\}$ .
- (ii) When  $X$  is 0,  $Y$  takes values 1 and -1 with equal probability ( $\frac{1}{2}$ ). When  $Y$  is 0,  $X$  takes values 1 and -1 with equal probability ( $\frac{1}{2}$ ).
- (iii) Either  $X$  is 0 with probability ( $\frac{1}{2}$ ), or  $Y$  is 0 with probability ( $\frac{1}{2}$ ).

**Are  $X$  and  $Y$  uncorrelated? Are  $X$  and  $Y$  independent? Prove your assertions.** *Hint:* Write down the joint probability of  $(X, Y)$  for each possible pair of values they can take.

$$\begin{aligned} P(X=0, Y=1) &= \frac{1}{4} & P(X=0, Y=-1) &= \frac{1}{4} \\ P(X=1, Y=0) &= \frac{1}{4} & P(X=-1, Y=0) &= \frac{1}{4} \\ \text{cov}(X, Y) &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \\ &= \mathbb{E}(XY) = 0 \\ X \text{ and } Y &\text{ are uncorrelated.} \\ P(X=0, Y=0) &= 0 \neq P(X=0) \cdot P(Y=0) \\ X \text{ and } Y &\text{ are not independent.} \end{aligned}$$

- (b) For  $X = [X_1, \dots, X_n]^T \sim \mathcal{N}(\mu, \Sigma)$ , verify that if  $X_i, X_j$  are independent (for all  $i \neq j$ ), then  $\Sigma$  must be diagonal, i.e.,  $X_i, X_j$  are uncorrelated.

$$X_i, X_j \text{ are independent} \Rightarrow X_i, X_j \text{ are uncorrelated} \\ \Rightarrow \text{cov}(X_i, X_j) = 0 \quad (i \neq j)$$

so  $\Sigma$  must be diagonal

- (c) Let  $N = 2$ ,  $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ , and  $\Sigma = \begin{pmatrix} \alpha & \beta \\ \beta & \gamma \end{pmatrix}$ . Suppose  $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N}(\mu, \Sigma)$ . Show that  $X_1, X_2$  are independent if  $\beta = 0$ . Recall that two continuous random variables  $W, Y$  with joint density  $f_{W,Y}$  and marginal densities  $f_W, f_Y$  are independent if  $f_{W,Y}(w, y) = f_W(w)f_Y(y)$ .

$$\begin{aligned} \mathcal{N}(\mu, \Sigma) &= \int_{\mathbb{R}^2} \frac{1}{2\pi\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right) dx \\ \text{if } \beta=0, \mathcal{N}(\mu, \Sigma) &= \int_{\mathbb{R}^2} \frac{1}{2\pi\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}\left(\alpha^{-1}(x_1-\mu_1)^2 + \gamma^{-1}(x_2-\mu_2)^2\right)\right) dx \\ &= \int_{\mathbb{R}^2} \frac{1}{2\pi\sqrt{|\Sigma|}} \exp\left(-\frac{(x_1-\mu_1)^2}{2\alpha}\right) \exp\left(-\frac{(x_2-\mu_2)^2}{2\gamma}\right) dx \\ &= \left(\int_{\mathbb{R}} \frac{1}{\sqrt{2\pi\alpha}} e^{-\frac{(x_1-\mu_1)^2}{2\alpha}} dx_1\right) \left(\int_{\mathbb{R}} \frac{1}{\sqrt{2\pi\gamma}} e^{-\frac{(x_2-\mu_2)^2}{2\gamma}} dx_2\right) = \mathcal{N}(\mu_1, \alpha) \cdot \mathcal{N}(\mu_2, \gamma) \end{aligned}$$

- (d) Consider a data point  $x$  drawn from an  $N$ -dimensional zero mean Multivariate Gaussian distribution  $\mathcal{N}(0, \Sigma)$ , as shown above. Assume that  $\Sigma^{-1}$  exists. Prove that there exists a matrix  $A \in \mathbb{R}^{N \times N}$  such that  $x^T \Sigma^{-1} x = \|Ax\|_2^2$  for all vectors  $x$ . What is the matrix  $A$ ?

$$x^T \Sigma^{-1} x = (Ax)^T (Ax) = x^T A^T A x$$

$$\Sigma^{-1} = A^T A$$

$$\Sigma = Q \Lambda Q^{-1}$$

$$\Sigma^{-1} = Q \Lambda^{-1} Q^{-1}$$

$$A = Q \Lambda^{-\frac{1}{2}} Q$$

$$\lambda_A = (\lambda_\Sigma)^{-\frac{1}{2}} = \sqrt{\frac{1}{\lambda_\Sigma}}$$

### 3 Least Squares (using vector calculus)

In ordinary least-squares linear regression, we typically have  $n > d$  so that there is no  $\mathbf{w}$  such that  $\mathbf{X}\mathbf{w} = \mathbf{y}$  (these are typically overdetermined systems — too many equations given the number of unknowns). Hence, we need to find an approximate solution to this problem. The residual vector will be  $\mathbf{r} = \mathbf{X}\mathbf{w} - \mathbf{y}$  and we want to make it as small as possible. The most common case is to measure the residual error with the standard Euclidean  $\ell^2$ -norm. So the problem becomes:

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2,$$

where  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{w} \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^n$ .

Assume that  $\mathbf{X}$  is full rank.

(a) How do we know that  $\mathbf{X}^T \mathbf{X}$  is invertible?

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \mathbf{v} &= \mathbf{0} & \left\{ \begin{array}{l} \mathbf{v} \text{ must be zero vector} \\ \text{so } \mathbf{X}^T \mathbf{X} \text{ is full rank} \\ \mathbf{X}^T \mathbf{X} \text{ is invertible.} \end{array} \right. \\ \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} &= 0 \\ \|\mathbf{X} \mathbf{v}\|_2^2 &= 0 \\ \mathbf{X} \mathbf{v} &= \mathbf{0} \end{aligned}$$

(b) Derive using vector calculus an expression for an optimal estimate for  $\mathbf{w}$  for this problem.

$$\begin{aligned} \frac{\partial \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2}{\partial \mathbf{w}} &= \frac{\partial \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2}{\partial (\mathbf{X}\mathbf{w} - \mathbf{y})} \cdot \frac{\partial (\mathbf{X}\mathbf{w} - \mathbf{y})}{\partial \mathbf{w}} \\ &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \cdot \mathbf{X} \\ \nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 &= \left( \frac{\partial \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2}{\partial \mathbf{w}} \right)^T = 2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0} \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

(c) What should we do if  $\mathbf{X}$  is not full rank?

~~gradient descent~~