

## DS5220 Homework - 03.

KAVANA VENKATESH.

1. Given  $y^i \sim N(\exp(\omega x^i), 1)$ 

$$\text{wkt, Normal distribution} = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(y^i - \mu)^2}{2\sigma^2}} \quad \text{--- ①}$$

Consider  $\sum_{i=1}^N x^i \exp(\omega x^i)$ ,

$$\therefore \text{wkt, } \mu = e^{\omega x^i} \text{ and } \sigma^2 = 1.$$

Assuming random variables are IID,  
we've

$$\mathcal{L}(\mu, \sigma^2 | x) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{(y^i - e^{\omega x^i})^2}{2}}.$$

Taking  $\ln$  on both sides,

$$\ln(\mathcal{L}(\mu, \sigma^2 | x)) = \ln\left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{(y^i - e^{\omega x^i})^2}{2}}\right)$$

$$= \ln\left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi}}\right) + \ln \prod_{i=1}^N e^{-\frac{(y^i - e^{\omega x^i})^2}{2}}$$

$$= N \cdot \ln \frac{1}{\sqrt{2\pi}} + \sum_{i=1}^N \frac{1}{2} \cdot -(y^i - e^{\omega x^i})^2.$$

$$= \frac{-N}{2} \ln 2\pi - \frac{1}{2} \sum_{i=1}^N (y^i - e^{\omega x^i})^2.$$

Differentiating the above eq. and  
setting it to 0.



$$\frac{\partial}{\partial \omega} \ln \left( \mathcal{L}(e^{\omega x^i}, 1/x) \right) = \frac{\partial}{\partial \omega} \left[ -\frac{N}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^N (y^i - e^{\omega x^i})^2 \right]$$

$$\Rightarrow \frac{\partial}{\partial \omega} \left( -\frac{N}{2} \ln(2\pi) \right) - \frac{1}{2} \frac{\partial}{\partial \omega} \sum_{i=1}^N (y^i - e^{\omega x^i})^2$$

$$\Rightarrow 0 - \frac{1}{2} \cdot 2 \cdot \sum_{i=1}^N (y^i - e^{\omega x^i}) (0 - e^{\omega x^i} \cdot x^i)$$

$$\Rightarrow 0 + \sum_{i=1}^N (y^i - e^{\omega x^i}) (e^{\omega x^i} \cdot x^i) = 0$$

$$\Rightarrow \sum_{i=1}^N y^i x^i \cdot e^{\omega x^i} - \sum_{i=1}^N e^{\omega x^i} \cdot e^{\omega x^i} \cdot x^i = 0$$

$$\Rightarrow \sum_{i=1}^N x^i \cdot e^{2\omega x^i} = \sum_{i=1}^N y^i x^i e^{\omega x^i}$$

$$\Rightarrow \boxed{\sum x^i \exp(2\omega x^i) = \sum x^i y^i \exp(\omega x^i)}$$

$\therefore E$  is the solution.

## 2. MAP Estimate:-

Given total no. of trials,  $N = N_0 + N_1$

Bernoulli random variable  $x$ ,  $p(x=1) = \theta$

Given dataset,  $\mathcal{D} = \{x_1, \dots, x_N\}$

$N_0 \rightarrow$  No. of trials when  $x_i = 0$ .

$N_1 \rightarrow$  No. of trials when  $x_i = 1$ .



$$P(\theta) = \begin{cases} 0.2 & \text{if } \theta = 0.6 \\ 0.8 & \text{if } \theta = 0.8 \\ 0 & \text{otherwise} \end{cases}$$

2.1) write down the log likelihood function  $P(\theta/\theta)$ . what's the max likelihood solution for  $\theta$ ?

likelihood function for  $P(\theta/\theta)$ .

$$\log(P(\theta/\theta)) = \log \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i}$$

$$= \sum_{i=1}^N \left[ \log \theta^{x_i} + \log (1-\theta)^{1-x_i} \right]$$

$$= \sum_{i=1}^N \left[ x_i \log \theta + (1-x_i) \log (1-\theta) \right]$$

Maximum likelihood solution for  $\theta$

$$\frac{d}{d\theta} \log P(\theta/\theta) = \frac{d}{d\theta} \left( \sum_{i=1}^N \left[ x_i \log \theta + (1-x_i) \log (1-\theta) \right] \right) = 0$$

$$\Rightarrow \sum_{i=1}^N \left[ x_i \cdot \frac{d}{d\theta} \log \theta + (1-x_i) \frac{d}{d\theta} \log (1-\theta) \right] = 0$$

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



$$= \sum_{i=1}^N \frac{x_i}{\theta} = \sum_{i=1}^N \frac{(1-x_i)}{1-\theta}$$

$$\Rightarrow \frac{1}{\theta} \sum_{i=1}^N x_i = \frac{1}{1-\theta} \cdot \sum_{i=1}^N (1-x_i)$$

$$\Rightarrow (1-\theta) \sum_{i=1}^N x_i = \theta \cdot \sum_{i=1}^N (1-x_i)$$

$$\Rightarrow \sum_{i=1}^N x_i - \theta \cdot \sum_{i=1}^N x_i = \theta \cdot \sum_{i=1}^N (1-x_i)$$

$$\Rightarrow \sum_{i=1}^N x_i - \theta \sum_{i=1}^N x_i = \theta \cdot \sum_{i=1}^N 1 - \theta \cdot \sum_{i=1}^N x_i$$

$$\Rightarrow \sum_{i=1}^N x_i = \theta \cdot N$$

$$\therefore \hat{\theta}_{MLE} = \frac{1}{N} \cdot \sum_{i=1}^N x_i = \frac{N_1}{N}$$

- 2.2) Consider maximizing the posterior distribution,  $p(D|\theta) \cdot p(\theta)$  that takes advantage of the prior. what is MAP estimation?

posterior =  $p(D|\theta) \cdot p(\theta)$  where  
 $p(D|\theta)$  = likelihood,  
 $p(\theta)$  = prior.

$$\therefore \hat{\theta}_{map} = \operatorname{argmax} \log \text{likelihood} + \text{prior}$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ \sum_{i=1}^N (x_i \log \theta + (1-x_i) \log(1-\theta)) \right] + \log p(\theta).$$

$$\text{given } p(\theta) = \begin{cases} 0.2 & ; \theta = 0.6 \\ 0.8 & \theta = 0.8 \\ 0 & \text{otherwise.} \end{cases}$$

$\therefore \log p(\theta)$  will always be constant.

$\therefore \boxed{\hat{\theta}_{\text{map}} = \hat{\theta}_{\text{MLE}}}$  because derivation of constant is 0.

we will end up with maximum likelihood equation.



### 3. Constrained optimization.

3.1) Given regression;  $\min_{\theta} \frac{1}{2} \|y - X\theta\|_2^2$  s.t.

$w^T \theta = b$ . Find the closed form solution.

$$\begin{aligned} \mathcal{L}(\theta, \alpha) &= \frac{1}{2} \|y - X\theta\|_2^2 + \sum \alpha_i (w^T \theta - b) \\ &= \frac{1}{2} (y - X\theta)^T (y - X\theta) + \sum \alpha_i (w^T \theta - b) \end{aligned}$$

Taking derivative wrt  $\theta$ ,

$$\begin{aligned} \frac{\partial}{\partial \theta} \mathcal{L}(\theta, \alpha) &= \frac{\partial}{\partial \theta} \left( \frac{1}{2} \right) \left[ y^T y - y^T X\theta - X^T \theta^T y \right. \\ &\quad \left. + X^T \theta^T X \theta \right] + \left[ \sum \alpha_i w^T \right]^T \end{aligned}$$

$$= \frac{1}{2} \left[ -[y^T X]^T - X^T y + 2X^T X \theta \right] + \alpha^{*T} w \quad \text{where } \sum \alpha_i = \alpha^*$$

$$= \frac{1}{2} \left[ -2X^T y + 2\theta \right] + \alpha^{*T} w$$

( $\because X^T X = I$ )

Equating to 0, we have

$$-X^T Y + \Theta = -\mathcal{L}^{*T} \omega.$$

$$\Theta^* = X^T Y - \mathcal{L}^{*T} \omega \quad \text{--- ①.}$$

or

$$\Theta^* = X^T Y - \omega^T \mathcal{L}.$$



$$3.2) \text{ Now, } \frac{dL}{d\alpha}(\theta^*, \alpha^*) = \frac{\partial}{\partial \alpha} \left[ \frac{1}{2} \left[ y^T y - y^T x \theta - \theta^T x^T y - \|\theta\|_2^2 \right] + \sum_{p=1}^N \alpha_p (\omega^T \theta^* - b) \right] = 0.$$

$$= \omega^T \theta^* - b = 0.$$

$$\Rightarrow \omega^T (x^T y - \sum_{p=1}^N \alpha_p^T \omega) - b = 0.$$

$$\omega^T x^T y - \omega^* \alpha^* \omega - b = 0.$$

$$\& \sum_{p=1}^N \alpha_p^T = \alpha^*$$

$$\Rightarrow \omega^T x^T y - b = \omega^T \alpha^* \omega.$$

$$\boxed{\alpha^* = \frac{(\omega^T x^T y - b)}{\|\omega\|_2^2}} \rightarrow \textcircled{2}.$$

Substitute  $\textcircled{2}$  in  $\textcircled{1}$ ,



$$\theta^* = X^T Y - \alpha^* \omega = X^T Y - \left( \frac{\omega^T X^T Y}{\|\omega\|_2^2} \right) \omega$$

$$\begin{aligned} X^T Y - \frac{(\omega^T X^T Y - b) \omega}{\omega^T \omega} &= X^T Y - \frac{\omega^T X^T Y}{\omega^T} + \frac{b}{\omega^T} \\ &= X^T Y - X^T Y + \frac{b}{\omega^T} \end{aligned}$$

$$\Rightarrow \theta^* = \frac{b}{\omega^T}$$

$$\Rightarrow \boxed{\omega^T \theta^* = b}$$

Thus  
Verified.

3.3) without the constraint, we have

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{2} \|Y - X\theta\|_2^2 = \frac{1}{2} \left[ Y^T Y - Y^T X \theta - \right. \\ &\quad \left. X^T \theta^T Y + X^T \theta^T X \theta \right] \end{aligned}$$

Taking the derivative wrt  $\theta$ ,

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \frac{1}{2} \left[ -2X^T Y + 2\theta \right]$$

Equating to zero, we have

$$X^T Y = \theta \Rightarrow \boxed{\theta = X^T Y}$$

#### 4. Convexity of logistic regression.

Log likelihood of logistic regression is,

$$\ell(\omega) = \sum_{i=1}^n \left[ y^i \phi(x_i)^T \omega - \log(1 + e^{\omega^T \phi(x_i)}) \right]$$

→ To prove log likelihood of logistic regression is convex, we prove Hessian is positive semi-definite.

$$\nabla_{\vec{\omega}}^2 (-\ell(\omega)) \geq 0 \quad [\because \text{+ve semi-definite}]$$

$$-\ell(\omega) = - \sum_{i=1}^n \left[ y^i \phi(x_i)^T \omega - \log(1 + e^{\omega^T \phi(x_i)}) \right]$$

$$\Rightarrow - \left[ y \phi^T \omega - \log(1 + e^{\phi^T \omega}) \right]$$

$$\rightarrow \text{Gradient, } \nabla_{\vec{\omega}} = - \left( y \phi^T \right) + \frac{1}{1 + e^{\phi^T \omega}} \cdot (0 + e^{\phi^T \omega}) \cdot \phi$$

$$\text{Hessian, } \nabla_{\vec{\omega}}^2 = 0 + \frac{\partial}{\partial \omega} \cdot \left( \frac{\phi e^{\phi^T \omega}}{1 + e^{\phi^T \omega}} \right)$$

$$\Rightarrow \frac{\partial}{\partial \omega} \cdot \left( \frac{\phi e^{\phi^T \omega}}{1 + e^{\phi^T \omega}} \right)$$



$$\Rightarrow \phi \cdot \left( \frac{\partial}{\partial \omega} \cdot \left( \frac{e^{\phi^T \omega}}{1 + e^{\phi^T \omega}} \right) \right)$$

$$\omega K T, \quad \frac{\partial}{\partial x} \left( \frac{u}{V} \right) = \left( \frac{V \frac{du}{dx} - u \frac{dV}{dx}}{V^2} \right)$$

$$u = e^{\phi^T \omega}, \quad V = 1 + e^{\phi^T \omega}.$$

$$\Rightarrow \phi \cdot \left( \frac{(1 + e^{\phi^T \omega}) \left( \frac{d}{d\omega} e^{\phi^T \omega} \right) - e^{\phi^T \omega} \left( \frac{d}{d\omega} (1 + e^{\phi^T \omega}) \right)}{(1 + e^{\phi^T \omega})^2} \right)$$

$$= \phi \left( \frac{e^{\phi^T \omega} \cdot (\phi^T)^T (1 + e^{\phi^T \omega}) - e^{\phi^T \omega} (e^{\phi^T \omega} \cdot \phi^T)}{(1 + e^{\phi^T \omega})^2} \right)$$

$$= \phi \left( \frac{(\phi^T)^T e^{\phi^T \omega} (1 + e^{\phi^T \omega}) - \phi^T e^{2\phi^T \omega}}{(1 + e^{\phi^T \omega})^2} \right)$$

$$\nabla_{\omega}^2 = \frac{\phi \phi^T e^{\phi^T \omega}}{(1 + e^{\phi^T \omega})^2} \geq 0. \quad \left( \because \text{+ve Value} \geq 0 \right. \\ \left. \text{Squared Value} \right)$$

$\therefore$  This is a convex function.

5) Let all the points using  $\phi(x)$  be obtained as,

5.1)  $\phi(x) = [1, x_1, x_2, x_1 x_2]$ .

class.	$x(x_1, x_2)$	$\phi(x) [1, x_1, x_2, x_1 x_2]$
+	(1, 1)	(1, 1, 1, 1)
+	(-1, -1)	(1, -1, -1, 1)
-	(-1, 1)	(1, -1, 1, -1)
-	(1, -1)	(1, 1, -1, -1)

we can see that we can ignore the first coordinate as it's the same for all the points. From simple intuition, we can see that the points can be linearly separated by boundary  $C_4 = 0$  and hence,  $w = [0, 0, 0, 1]$  and this is the hyperplane that linearly separates +ve & -ve points.

Note:- Both classes are not linearly separable in original space but can be linearly separated if we transform the feature space.



$$5.2) \quad K(X, Z) = \langle \phi(X), \phi(Z) \rangle$$

$$= [1, X_1, X_2, X_1 X_2]^T [1, Z_1, Z_2, Z_1 Z_2]$$

$$= 1 + X_1 Z_1 + X_2 Z_2 + X_1 X_2 Z_1 Z_2$$

$$= 1 + X^T Z + X_1 X_2 Z_1 Z_2$$

### 6.1) Constructing Kernels:

$$\text{given } K_1(x, z) : \mathbb{R}^d * \mathbb{R}^d \rightarrow \mathbb{R}$$

$$K_2(x, z) : \mathbb{R}^d * \mathbb{R}^d \rightarrow \mathbb{R}$$

To prove:  $K(x, z) \Rightarrow c_1 K_1(x, z) + c_2 K_2(x, z)$  for  $c_1, c_2 \geq 0$

$K_1$  has its feature map  $\phi_1$  and inner product  $\langle \cdot \rangle_{HK_1}$  and  $K_2$  has its feature map  $\phi_2$  and inner product  $\langle \cdot \rangle_{HK_2}$ .

$$\text{So, } K(x, z) = \langle \overline{Jc_1} \phi_1(x), \overline{Jc_1} \phi_1(z) \rangle_{HK_1} +$$

$$\langle \overline{Jc_2} \phi_2(x), \overline{Jc_2} \phi_2(z) \rangle_{HK_2}$$

$$= \langle [\overline{Jc_1} \phi_1(x), \overline{Jc_2} \phi_2(z)],$$

$$\overline{Jc_1} \phi_1(z), \overline{Jc_2} \phi_2(z) \rangle_{HK_{\text{new}}}$$

---



$$6.2) \quad K(x, z) = K_1(x, z) \cdot K_2(x, z)$$

$$= \phi_1(x)^T \phi_1(z) \phi_2(x)^T \phi_2(z)$$

$$= \sum_i \phi_1(x)_i \phi_1(z)_i \sum_j \phi_2(x)_j \phi_2(z)_j$$

$$= \sum_{i,j} \phi_1(x)_i \phi_1(z)_i \phi_2(x)_j \phi_2(z)_j$$

$$= \sum_{i,j} \phi_1(x)_i \phi_2(x)_j \phi_1(z)_i \phi_2(z)_j$$

$$= \underline{\underline{\phi(x)^T \cdot \phi(z)}} \quad \text{where}$$

$$\underline{\underline{\phi(x) = \phi_1(x) \otimes \phi_2(x)}} \quad (\text{Kronecker product}).$$