

HOMEWORK-4

EE559

1. $x_1, x_2, \dots, x_n \rightarrow$ iid Poisson sample

$$\lambda \rightarrow \Gamma(2, 1) \text{ prior distribution. } \Rightarrow P(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} = \lambda^2 e^{-\lambda}$$

To find MAP estimate of λ

solt:

$$= \arg \max_{\lambda} f(\lambda | x_1, x_2, \dots, x_n)$$

$$= \arg \max_{\lambda} \underbrace{f(x | \lambda)}_{f(x_i | \lambda)} \underbrace{f(\lambda)}_{f(\lambda)}$$

$$= f(\lambda) \prod_{i=1}^n f(x_i | \lambda)$$

$$= f(\lambda) \prod_{i=1}^n \frac{x_i!}{\lambda^{x_i} e^{-\lambda}}$$

$$= \lambda^{\sum x_i} \prod_{i=1}^n \frac{e^{-\lambda}}{x_i!} = \lambda^{\sum x_i} e^{-n\lambda} \prod_{i=1}^n \frac{1}{x_i!}$$

taking log

$$L = \arg \max_{\lambda} \left[\sum x_i \log(\lambda) - \underbrace{\lambda - n\lambda}_{-(n+1)\lambda} + \log \left(\prod_{i=1}^n \frac{1}{x_i!} \right) \right]$$

taking derivative & = 0

$$\frac{\sum x_i + 1}{\lambda} - (n+1) = 0$$

$$\boxed{\lambda = \frac{\sum x_i + 1}{n+1}}$$

To ensure this is maximum, we take $\frac{d^2}{d\lambda^2} = -\frac{\sum x_i + 1}{\lambda^2}$

$\frac{\sum x_i + 1}{\lambda^2} > 0$ for $\lambda > 0$, i.e. $\frac{d^2}{d\lambda^2}$ is negative \Rightarrow it is maximum.

$\therefore \lambda_{MAP} = \frac{\sum x_i + 1}{n+1}$ maximizes posterior distribution.

$$P_X(x) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x > 0 \quad \lambda = 0, 1, \dots$$

Likelihood:

$$L(\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

{ Given $x_1, x_2, \dots, x_n \Rightarrow$ iid Poisson}

\therefore log likelihood is

$$l(\lambda) = \sum_{i=1}^n (x_i \log \lambda - \lambda - \log x_i!)$$

$$\text{To find MLE: } \frac{d}{d\lambda} l(\lambda) = 0 \Rightarrow \frac{d}{d\lambda} \left[\sum_{i=1}^n (x_i \log \lambda - \lambda - \log x_i!) \right] \rightarrow \text{goes to 0}$$

$$\Rightarrow \frac{d}{d\lambda} \sum_{i=1}^n \left(\frac{x_i}{\lambda} - 1 \right) = 0 \quad \rightarrow (i)$$

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

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

Given Information

$$I(\lambda) = -E \left[\frac{\partial^2 l(\lambda)}{\partial \lambda^2} \right]$$

from (i)

$$\frac{\partial^2 l(\lambda)}{\partial \lambda^2} = \sum_{i=1}^n \left(-\frac{x_i}{\lambda^2} \right)$$

$$\therefore I(\lambda) = -E \left[\sum_{i=1}^n \left(-\frac{x_i}{\lambda^2} \right) \right]$$

$$I(\lambda) = \frac{1}{\lambda^2} + E \left[\sum_{i=1}^n x_i \right] = \frac{n\bar{x}}{\lambda^2} = \frac{n}{\lambda}$$

\therefore Asymptotic normality:

$$\hat{\lambda}_{MLE} = \bar{x} \approx N \left(\lambda, \frac{1}{\lambda} \right) \Rightarrow \hat{\lambda}_{MLE} \approx N \left(\lambda, \frac{1}{n} \right)$$

Comparison with Central Limit theorem

If X_1, \dots, X_n iid with $E[X] = \mu$
 $\text{Var}(X) = \sigma^2$

then, according to CLT,

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} N(0, 1)$$

for current problem,

$X_i \sim \text{poisson}$

$$E[X_i] = \lambda \quad \text{and} \quad \text{Var}(X_i) = \lambda$$

$$\therefore \frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \xrightarrow{d} N(0, 1)$$

$$\Rightarrow \bar{X} - \lambda \sim N(0, \lambda/n)$$

$$\boxed{\bar{X} \sim N(\lambda, \lambda/n)} \quad \text{for large } n$$

3. $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon \quad \epsilon \sim N(0, \sigma^2)$

to show: MLE = least squares of β vector

MLE:

$$L(\beta_0, \beta_1, \dots, \beta_p, \sigma^2) = \prod_{i=1}^n \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2}{2\sigma^2}\right)$$

$$\begin{aligned} & \text{Exp } \epsilon \\ \therefore \epsilon &= Y - (\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p) \sim N(0, \sigma^2) \\ \& P(\epsilon) = P(Y | X_1, \dots, X_p) \end{aligned}$$

taking log,

$$l(\beta_0, \beta_1, \dots, \beta_p, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2$$

To find MLE, we need to maximize log-likelihood w.r.t $\beta_0, \beta_1, \dots, \beta_p$.

From (i) we notice that we maximize log-likelihood by minimizing

$$\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2 \quad (\text{ii})$$

From least squares, we know,

$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ minimize the RSS

$$\text{i.e. } \hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2 \quad (\text{iii})$$

which gives us normal equation, $\hat{\beta} = (X^T X)^{-1} X^T Y$

∴ From (ii) & (iii) we observe, log-likelihood function is maximized

when RSS is minimized, which is same as

minimizing $\beta_0, \beta_1, \dots, \beta_p$ which minimize the RSS.

∴ MLEs for $\beta_0, \beta_1, \dots, \beta_p$ are same as the least square estimates. ∴ MLEs are also BLUE.

4. Given:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$\hookrightarrow \text{prior distribution of } \beta_i \sim N(0, \sigma^2 / n)$$

MAP estimate:

likelihood function for given β)

$$P(Y|X, \beta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2\right)$$

log likelihood:

$$l(\beta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2$$

Note:
multivariate normal distribution
of linear regression:

$$P(Y) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (Y - \mu)^T \Sigma^{-1} (Y - \mu)\right)$$

$$\Sigma = \sigma^2 I \Rightarrow |\Sigma| = (\sigma^2)^n = \sigma^{2n}$$

Prior distribution:

$$P(\beta) = \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\lambda\beta_i^2}{2\sigma^2}\right)$$

log-prior:

$$= -\frac{p}{2} \log\left(\frac{2\pi\sigma^2}{\lambda}\right) - \frac{\lambda}{2\sigma^2} \sum_{i=1}^p \beta_i^2$$

Posterior distribution:

$$\log p(\beta | Y, X) = l(\beta) + \log P(\beta)$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2 - \frac{\lambda}{2} \sum_{i=1}^p \beta_i^2$$

$$\hat{\beta}_{MAP} = \arg \max_{\beta} \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2 - \frac{\lambda}{2\sigma^2} \sum_{i=1}^p \beta_i^2 \right)$$

should be min

i.e. when $\sum_{i=1}^p \beta_i^2$

$$\therefore \hat{\beta}_{MAP} = \arg \min_{\beta} \left(\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_p X_{ip})^2 + \lambda \sum_{i=1}^p \beta_i^2 \right)$$

Interpretation:

This is equivalent to the Ridge Regression problem
square to shrink

$\lambda \sum_{i=1}^p \beta_i^2$ is added to ordinary least squares to prevent over-fitting.

Coefficients & prevent over-fitting.

$\hat{\beta}_{MAP}$: estimate that maximizes posterior distribution
likelihood & prior beliefs.

5. Given:

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

prior distribution for β_i :

$$\beta_i \sim \text{Laplace}(0, \frac{\sigma^2}{\lambda})$$

MAP Estimate:

likelihood function:

$$P(Y|X, \beta) = \frac{1}{(2\pi\sigma^2)^n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2\right)$$

taking log

$$l(\beta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$$

prior distribution:

$$P(\beta_i) = \frac{\lambda}{2\sigma^2} \exp\left(-\frac{\lambda|\beta_i|}{\sigma^2}\right)$$

for joint,

$$\Rightarrow P(\beta) = \prod_{i=1}^p \frac{\lambda}{2\sigma^2} \exp\left(-\frac{\lambda|\beta_i|}{\sigma^2}\right)$$

log:

$$\log P(\beta) = \sum_{i=1}^p \left(\log \frac{\lambda}{2\sigma^2} - \frac{\lambda|\beta_i|}{\sigma^2} \right)$$

$$= p \log \frac{\lambda}{2\sigma^2} - \frac{\lambda}{\sigma^2} \sum_{i=1}^p |\beta_i|$$

Posterior distribution:

$$= \log p(\beta|y, x) = l(\beta) + \log P(\beta)$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 + \phi$$

$$\text{Maximizing } = \arg \min_{\beta} \left(\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 + \frac{\alpha\lambda}{\sigma^2} \sum_{i=1}^p |\beta_i| \right)$$
$$= \hat{\beta}_{MAP}$$

Interpretation:

- ↳ penalty term, $\frac{\lambda}{\tau^2} \sum_{i=1}^p |\beta_i|$ is added to ordinary least squares to shrink some coefficients to 0.
- ↳ equivalent to 'Lasso Regression Problem'
- ↳ $\hat{\beta}_{MAP}$: point estimate that maximises posterior distribution, considering both data likelihood & prior beliefs.

6.

$$SVD \text{ of } X = U \Sigma V^T$$

$U = n \times n$ orthogonal

let $\Sigma = n \times p$ diagonal with singular values $\tau_1, \tau_2, \dots, \tau_p$

$V = p \times p$ orthogonal matrix.

Ridge regression estimator is given by

$$\hat{\beta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T y \quad (i)$$

$$\text{Substituting } X = U \Sigma V^T \text{ into (i)} = [V \Sigma^T V^T]^{-1} + \lambda I V \Sigma^T V^T y$$

$$\hat{\beta}_{Ridge} = (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T V^T y$$

$$\Rightarrow \text{Let } I = V V^T$$

$$= (\Sigma^T \Sigma + \lambda I) V^T V \Sigma^T V^T y$$

Since V is orthogonal,

$$\hat{\beta}_{Ridge} = V (\Sigma^T \Sigma + \lambda I) \underbrace{V^T V}_{I} \Sigma^T V^T y$$

$$= V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T V^T y \quad (ii)$$

i) predicted values are:

$$\hat{y} = X \hat{\beta}_{Ridge} = U \Sigma V^T \underbrace{V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T}_{I} V^T y$$

$$\Sigma^T \Sigma = \text{diag}(\tau_1^2, \tau_2^2, \dots, \tau_p^2)$$

$$\therefore (\Sigma^T \Sigma + \lambda I)^{-1} = \text{diag}\left(\frac{1}{\tau_1^2 + \lambda}, \dots, \frac{1}{\tau_p^2 + \lambda}\right)$$

$$\therefore \hat{y} = U \cdot \text{diag}\left(\frac{\sigma_1^2}{\sigma_1^2 + \lambda}, \dots, \frac{\sigma_p^2}{\sigma_p^2 + \lambda}\right) U^T y \quad \text{--- (iii)}$$

$$= \sum_{j=1}^p U_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} U_j^T y \quad ; U_j = \text{columns of } U$$

Where, $\frac{\sigma_j^2}{\sigma_j^2 + \lambda}$ acts as a shrinkage factor for corresponding singular vector U_j . \because for fixed λ , shrinkage factor decreases as σ_j decreases i.e if σ_j is smaller, \therefore have greater shrinkage, & less effect on prediction helping reduce overfitting.

(b) From part (a) with SVD of $X = U\Sigma V^T$

we have,

$$(X^T X + \lambda I)^{-1} = V (\Sigma^T \Sigma + \lambda I)^{-1} V^T$$

$$\therefore (X^T X + \lambda I)^{-1} X^T X = V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \Sigma V^T$$

$$\begin{aligned} \therefore X^T X &= V \Sigma^T U^T U \Sigma V^T \\ &= V \Sigma^T \Sigma V^T \end{aligned}$$

$$\therefore \text{tr}[(X^T X + \lambda I)^{-1} X^T X] = \text{tr}[V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \Sigma V^T]$$

$$\text{tr}(AB) = \text{tr}(BA) \Rightarrow V^T = I$$

$$\therefore \text{trace} = \text{tr}[(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T \Sigma] \quad \text{from (iii) part (a)}$$

$$= \text{diag}\left(\frac{\sigma_j^2}{\sigma_j^2 + \lambda}\right)$$

$\therefore \text{trace} = \text{summing of diag elements}$

$$\Rightarrow \text{tr}[(X^T X + \lambda I)^{-1} X^T X] = \sum_{j=1}^p \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$

Conclusion: When $\lambda = 0$, there is no regularization & effective degrees of freedom = p . As $\lambda \uparrow$, degrees of freedom \downarrow reflecting increased shrinkage of co-efficients.