

HW4

1. The pdf of a $\Gamma(2, 1)$ random variable is $p(z) = z \exp(-z)$, $z > 0$, and the pmf of a Poisson random variable X is $p_X(x) = \lambda^x e^{-\lambda} / x!$, $\lambda > 0$, $x = 0, 1, \dots$. Assuming that X_1, X_2, \dots, X_n is an i.i.d Poisson sample given that λ has a $\Gamma(2, 1)$ prior distribution, find the MAP estimate of λ and prove that what you find is actually a value that maximizes the posterior. (10 pts)

\because iid

$$\therefore p(X_1, X_2, \dots, X_n | \lambda) = \prod_{i=1}^n \frac{\lambda^{X_i} e^{-\lambda}}{X_i!} = \frac{\lambda^{\sum_{i=1}^n X_i} e^{-n\lambda}}{\prod_{i=1}^n X_i!}$$

$$\therefore P(\lambda | X_1, X_2, \dots, X_n) \propto p(X_1, X_2, \dots, X_n | \lambda) P(\lambda) = \frac{\lambda^{\sum_{i=1}^n X_i} e^{-n\lambda}}{\prod_{i=1}^n X_i!} \cdot \lambda e^{-\lambda}$$

$$\therefore P(\lambda | X_1, X_2, \dots, X_n) \propto \lambda^{1 + \sum_{i=1}^n X_i - (n+1)\lambda} = \lambda^{(2 + \sum_{i=1}^n X_i) - 1} e^{-(n+1)\lambda}$$

\therefore posterior is a $\text{Gamma}(2 + \sum_{i=1}^n X_i, n+1)$ distribution

$$\therefore \hat{\lambda}^{\text{MAP}} = \underset{\lambda}{\operatorname{argmax}} P(\lambda | X_1, X_2, \dots, X_n) = \text{mode of } \text{Gamma}(2 + \sum_{i=1}^n X_i, n+1)$$

$$\therefore \hat{\lambda}^{\text{MAP}} = \frac{\alpha - 1}{\beta} = \frac{1 + \sum_{i=1}^n X_i}{n+1}$$

proof: let $g(\lambda) = \lambda^{(2 + \sum_{i=1}^n X_i) - 1} e^{-(n+1)\lambda}$

$$\therefore \ln g(\lambda) = (1 + \sum_{i=1}^n X_i) \ln \lambda - (n+1)\lambda$$

$$\therefore \text{let } \frac{d}{d\lambda} g(\lambda) = \frac{1 + \sum_{i=1}^n X_i}{\lambda} - (n+1) = 0$$

$$\therefore \lambda = \frac{1 + \sum_{i=1}^n X_i}{n+1} = \hat{\lambda}^{\text{MAP}}$$

$$\therefore \frac{d^2}{d\lambda^2} g(\lambda) = \frac{-(1 + \sum_{i=1}^n X_i)}{\lambda^2} < 0, \quad \forall \lambda > 0$$

$\therefore \hat{\lambda}^{\text{MAP}}$ maximizes the posterior

2. Assume that you have an i.i.d sample from a population with Poisson pmf, i.e. $p_X(x) = \lambda^x e^{-\lambda} / x!$, $\lambda > 0$, $x = 0, 1, \dots$. Calculate the MLE of λ and its asymptotic distribution by calculating Fisher information and compare the results with those of the Central Limit Theorem. (10 pts)

$$\therefore L(\lambda) = L(\lambda | D) = \prod p(x_i | \lambda) = \prod \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod x_i!}$$

$$\therefore l(\lambda) = (\sum x_i) \ln \lambda - n\lambda - \sum \ln x_i$$

$$\therefore \text{let } \frac{d}{d\lambda} l(\lambda) = \frac{\sum x_i}{\lambda} - n = 0$$

$$\therefore \widehat{\lambda}_{MLE} = \frac{1}{n} \sum x_i = \bar{x}$$

$$\therefore \frac{d^2}{d\lambda^2} l(\lambda) = -\frac{\sum x_i}{\lambda^2} \quad \text{and} \quad E[x_i] = \lambda$$

$$\therefore I(\lambda) = -E\left[\frac{d^2}{d\lambda^2} l(\lambda)\right] = -E\left[-\frac{\sum x_i}{\lambda^2}\right] = \frac{n\lambda}{\lambda^2} = \frac{n}{\lambda}$$

$$\therefore \widehat{\lambda}_{MLE} \sim N\left(\lambda, \frac{1}{I(\lambda)}\right) = N\left(\lambda, \frac{\lambda}{n}\right)$$

$$\therefore \text{in CLT } \bar{x} \sim N(E[x], \frac{V[x]}{n}) = N\left(\lambda, \frac{\lambda}{n}\right)$$

\therefore same

3. Assume that $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Show that the MLE and least squares estimates of the β vector are the same, which means MLE is also BLUE according to Gauss-Markov. Remember that the log-likelihood function $l(\beta_0, \beta_1, \dots, \beta_p)$ here is based on the conditional density $p(Y|X_1, \dots, X_p)$. (10 pts)

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

$$\therefore Y = X\beta + \epsilon \sim N(X\beta, \sigma^2 I)$$

$$\begin{aligned} \therefore L(\beta) &= p(Y|X, \beta) = \frac{1}{(2\pi)^{\frac{n}{2}} |\det \sigma^2 I|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (Y - X\beta)^T (\sigma^2 I)^{-1} (Y - X\beta)\right) \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \|Y - X\beta\|^2\right) \end{aligned}$$

$$\therefore l(\beta) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|Y - X\beta\|^2$$

$$\therefore \hat{\beta}^{MLE} = \underset{\beta}{\operatorname{argmax}} l(\beta) = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2 = \hat{\beta}^{LS}$$

$$\therefore \hat{\beta}^{MLE} = \hat{\beta}^{LS} = (X^T X)^{-1} X^T Y$$

$\therefore \hat{\beta}^{MLE}$ is BLUE

4. Find the MAP estimate of β under the assumption that $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and that the prior distribution of (independent) β_i , $i = 1, 2, \dots, p$ is $\mathcal{N}(0, \sigma^2/\lambda)$. Interpret your results. (15 pts)

according to Q3 $p(Y|X, \beta) \propto \exp(-\frac{1}{2\sigma^2} \|Y - X\beta\|^2)$

$\because \beta_i \sim \mathcal{N}(0, \sigma^2/\lambda)$, $i = 1, 2, \dots, p$

$\therefore p(\beta) \propto \exp(-\frac{1}{2\sigma^2} \beta^T \Lambda \beta)$, $\Lambda = \text{diag}(0, \lambda, \dots, \lambda)$

$\therefore p(\beta|X, Y) \propto p(Y|X, \beta) p(\beta) \propto \exp(-\frac{1}{2\sigma^2} \|Y - X\beta\|^2 - \frac{1}{2\sigma^2} \beta^T \Lambda \beta)$

$\therefore \ln p(\beta|X, Y) \propto -(\|Y - X\beta\|^2 + \beta^T \Lambda \beta)$ similar to ridge regression

$$\therefore \widehat{\beta}^{\text{MAP}} = (X^T X + \Lambda)^{-1} X^T Y$$

$\therefore \widehat{\beta}^{\text{MAP}}$ corresponds to ridge regression.

use Gaussian prior to achieve l_2 -penalty

but only regularize non-intercept coefficients

5. Find the MAP estimate of β under the assumption that $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and that the prior distribution of (independent) β_i , $i = 1, 2, \dots, p$ is $\text{Lap}(0, \sigma^2/\lambda)$. Interpret your results. (15 pts)

Similarly, $p(Y|X, \beta) \propto \exp(-\frac{1}{2\sigma^2} \|Y - X\beta\|^2)$

$\because \beta_i \sim \text{Lap}(0, \sigma^2/\lambda)$, $i = 1, 2, \dots, p$

$$\therefore p(\beta_i) = \frac{1}{2\sigma^2/\lambda} \exp\left(-\frac{|\beta_i|}{\sigma^2/\lambda}\right) = \frac{\lambda}{2\sigma^2} \exp\left(-\frac{\lambda}{\sigma^2} |\beta_i|\right) \quad i = 1, 2, \dots, p$$

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

$$\therefore p(\beta|X, Y) \propto p(Y|X, \beta) p(\beta) \propto \exp\left(-\frac{1}{2\sigma^2} \|Y - X\beta\|^2 - \frac{\lambda}{\sigma^2} \sum_{i=1}^p |\beta_i|\right)$$

$$\therefore \ln p(\beta|X, Y) \propto -\frac{1}{2} \|Y - X\beta\|^2 - \lambda \sum_{i=1}^p |\beta_i|$$

$$\therefore \widehat{\beta}^{\text{MAP}} = \underset{\beta}{\operatorname{argmax}} \ln p(\beta|X, Y) = \underset{\beta}{\operatorname{argmin}} (\|Y - X\beta\|^2 + 2\lambda \sum_{i=1}^p |\beta_i|)$$

$\therefore \widehat{\beta}^{\text{MAP}}$ corresponds to Lasso regression.

use Laplace prior to achieve l_1 -penalty

6. In the regularized least squares problem, assume that the singular value decomposition of \mathbf{X} is $\mathbf{U}\Sigma\mathbf{V}^T$.

(a) Show that the vector of predicted values is: (10 pts)

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}^{\text{Ridge}} = \sum_{j=1}^p \mathbf{u}_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y}$$

where \mathbf{u}_j are the columns of \mathbf{U} . Conclude that greater amount of shrinkage is applied to basis vectors \mathbf{u}_j that have smaller singular values σ_j , for a fixed $\lambda \geq 0$.

$$\begin{aligned}
 (a) \because \widehat{\beta}^{\text{Ridge}} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= ((\mathbf{U}\Sigma\mathbf{V}^T)^T (\mathbf{U}\Sigma\mathbf{V}^T) + \lambda \mathbf{I})^{-1} (\mathbf{U}\Sigma\mathbf{V}^T)^T \mathbf{y} \\
 &= (\mathbf{V}^T \Sigma \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T + \lambda \mathbf{I})^{-1} \mathbf{V}^T \Sigma \mathbf{U}^T \mathbf{y} \quad (\because \mathbf{U}^{-1} = \mathbf{U}^T, \mathbf{V}^{-1} = \mathbf{V}^T) \\
 &= (\mathbf{V}^T (\Sigma^2 + \lambda \mathbf{I}))^{-1} \mathbf{V}^T \Sigma \mathbf{U}^T \mathbf{y} \\
 &= \mathbf{V}^T (\Sigma^2 + \lambda \mathbf{I})^{-1} \Sigma \mathbf{U}^T \mathbf{y} \\
 \therefore \hat{\mathbf{y}} &= \mathbf{X} \widehat{\beta}^{\text{Ridge}} = (\mathbf{U}\Sigma\mathbf{V}^T) \mathbf{V} (\Sigma^2 + \lambda \mathbf{I})^{-1} \Sigma \mathbf{U}^T \mathbf{y} \\
 &= \mathbf{U} \Sigma (\Sigma^2 + \lambda \mathbf{I})^{-1} \Sigma \mathbf{U}^T \mathbf{y} \\
 &= \mathbf{U} \frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}} \mathbf{U}^T \mathbf{y} \\
 \therefore \frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}} &= \text{diag} \left(\frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) \\
 \therefore \hat{\mathbf{y}} &= \sum_{j=1}^p \mathbf{u}_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y}
 \end{aligned}$$

(b) Use SVD to show that (10 pts)

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

This quantity is equal to the degrees of freedom p when $\lambda = 0$ and is called the effective degrees of freedom for the Ridge-regularized model.

$$\begin{aligned} (b) \because \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T &= (\mathbf{U} \Sigma \mathbf{V}^T)(\Sigma^2 + \lambda \mathbf{I})^{-1}(\mathbf{U} \Sigma \mathbf{V}^T)^T \\ &= \mathbf{U} \Sigma \mathbf{V}^T \mathbf{V} (\Sigma^2 + \lambda \mathbf{I})^{-1} \mathbf{V}^T \mathbf{V} \Sigma \mathbf{V}^T \\ &= \mathbf{U} \Sigma (\Sigma^2 + \lambda \mathbf{I})^{-1} \Sigma \mathbf{V}^T \\ &= \mathbf{U} \frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}} \mathbf{U}^T \end{aligned}$$

$$\begin{aligned} \therefore \text{tr}[\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T] &= \text{tr}\left[\mathbf{U} \frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}} \mathbf{U}^T\right] \\ &= \text{tr}\left[\frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}} \mathbf{U}^T \mathbf{U}\right] \quad (\because \text{tr}(ABC) = \text{tr}(BCA)) \\ &= \text{tr}\left[\frac{\Sigma^2}{\Sigma^2 + \lambda \mathbf{I}}\right] \end{aligned}$$

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

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