

Homework 1

Mingfeng Li
cmu4209)

Problem 1 (written) – 25 points

Imagine you have a sequence of N observations (x_1, \dots, x_N) , where each $x_i \in \{0, 1, 2, \dots, \infty\}$. You model this sequence as i.i.d. from a Poisson distribution with unknown parameter $\lambda \in \mathbb{R}_+$, where

$$p(X|\lambda) = \frac{\lambda^X}{X!} e^{-\lambda}$$

- (a) What is the joint likelihood of the data (x_1, \dots, x_N) ?
- (b) Derive the maximum likelihood estimate λ_{ML} for λ .

To help learn λ , you use a prior distribution. You select the distribution $p(\lambda) = \text{gamma}(a, b)$.

- (c) Derive the maximum a posteriori (MAP) estimate λ_{MAP} for λ ?
- (d) Use Bayes rule to derive the posterior distribution of λ and identify the name of this distribution.
- (e) What is the mean and variance of λ under this posterior? Discuss how it relates to λ_{ML} and λ_{MAP} .

$$a) \quad p(x_1, \dots, x_N | \lambda) = \prod_{i=1}^N p(x_i | \lambda) = \prod_{i=1}^N \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$$

$$b) \quad \ln \prod_{i=1}^N p(x_i | \lambda) = \sum_{i=1}^N \ln \left(\frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = \sum_{i=1}^N (x_i \ln \lambda - \ln x_i! - \lambda)$$

$$\nabla_{\lambda} \ln \prod_{i=1}^N p(x_i | \lambda) = 0$$

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

$$\frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

$$\frac{1}{\lambda} \sum_{i=1}^N x_i = N$$

$$\bar{x} = \lambda_{ML}$$

$$c) \quad \lambda_{MAP} = \arg \max_{\lambda} \ln p(\lambda | x_1, \dots, x_N)$$

$$= \arg \max_{\lambda} \ln \frac{p(x_1, \dots, x_N | \lambda) p(\lambda)}{p(x_1, \dots, x_N)}$$

$$= \arg \max_{\lambda} \ln p(x_1, \dots, x_N | \lambda) + \ln p(\lambda)$$

$$\text{Let } L = \ln p(x_1, \dots, x_N | \lambda) + \ln p(\lambda)$$

$$= \sum_{i=1}^N (x_i \ln \lambda - \ln x_i! - \lambda) + \ln \frac{b^a}{\Gamma(a)} + (a-1) \ln \lambda - b\lambda$$

$$= \ln \lambda \sum_{i=1}^N x_i - \sum_{i=1}^N \ln x_i! - n\lambda + \ln \frac{b^a}{\Gamma(a)} + (a-1) \ln \lambda - b\lambda$$

=

$$\tau_{\lambda} \mid = \frac{\sum_{i=1}^n \lambda_i}{\lambda} - n + \frac{a-1}{\lambda} - b = 0$$

$$\frac{\sum_{i=1}^n \lambda_i + a - 1}{\lambda} = n + b$$

$$\lambda_{\text{MAP}} = \frac{a + \sum_{i=1}^n \lambda_i - 1}{n + b}$$

$$\begin{aligned} a) \quad P(\lambda | x_1, \dots, x_n) &= \frac{P(x_1, \dots, x_n | \lambda) P(\lambda)}{P(x_1, \dots, x_n)} \\ &\propto P(x_1, \dots, x_n | \lambda) P(\lambda) \\ &\propto \lambda^{\sum_{i=1}^n x_i + a - 1} e^{-(n+b)\lambda} \end{aligned}$$

$$\Rightarrow P(\lambda | x_1, \dots, x_n) \text{ is gamma} \left(\sum_{i=1}^n x_i + a, n + b \right)$$

$$e) \quad \hat{\lambda}_{\text{post}} = \frac{\sum_{i=1}^n x_i + a}{n + b} \quad \text{Var}(\lambda_{\text{post}}) = \frac{\sum_{i=1}^n x_i + a}{(n + b)^2}$$

$$\begin{aligned} \hat{\lambda}_{\text{post}} &= \frac{\sum_{i=1}^n x_i + a}{n + b} = \left(\frac{n}{n + b} \right) \frac{\sum_{i=1}^n x_i}{n} + \frac{a}{n + b} \\ &= \frac{n}{n + b} \lambda_{\text{ML}} + \frac{a}{n + b} \end{aligned}$$

$$\begin{aligned} \hat{\lambda}_{\text{post}} &= \frac{\sum_{i=1}^n x_i + a}{n + b} = \frac{\sum_{i=1}^n x_i + a - 1}{n + b} + \frac{1}{n + b} \\ &= \lambda_{\text{MAP}} + \frac{1}{n + b} \end{aligned}$$

Problem 2 (written) – 15 points

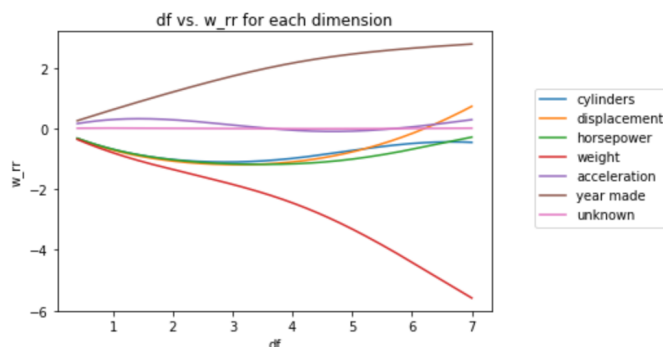
You have data (x_i, y_i) for $i = 1, \dots, n$, where $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$. You model this as $y_i \stackrel{iid}{\sim} N(x_i^T w, \sigma^2)$. You use the data you have to approximate w with $w_{RR} = (\lambda I + X^T X)^{-1} X^T y$, where X and y are defined as in the lectures. Derive the results for $\mathbb{E}[w_{RR}]$ and $\mathbb{V}[w_{RR}]$ given in the slides.

$$\begin{aligned}\mathbb{E}[w_{RR}] &= \mathbb{E}[(\lambda I + X^T X)^{-1} X^T y] \\ &= (\lambda I + X^T X)^{-1} X^T \mathbb{E}[y] \\ &= (\lambda I + X^T X)^{-1} X^T X w\end{aligned}$$

$$\begin{aligned}\text{Var}[w_{RR}] &= \mathbb{E}[(w_{RR} - \mathbb{E}[w_{RR}])(w_{RR} - \mathbb{E}[w_{RR}])^T] \\ &= \mathbb{E}[w_{RR} w_{RR}^T] - \mathbb{E}[w_{RR}] \mathbb{E}[w_{RR}]^T \\ &= \mathbb{E}[(\lambda I + X^T X)^{-1} X^T y y^T X (\lambda I + X^T X)^{-1}] - (\lambda I + X^T X)^{-1} X^T X w w^T X^T X (\lambda I + X^T X)^{-1} \\ &= (\lambda I + X^T X)^{-1} X^T (\sigma^2 I + \underbrace{X w w^T X^T}) X (\lambda I + X^T X)^{-1} - (\lambda I + X^T X)^{-1} X^T X w w^T X^T X (\lambda I + X^T X)^{-1} \\ &= (\lambda I + X^T X)^{-1} X^T \sigma^2 I X (\lambda I + X^T X)^{-1} \\ &= (\lambda I + X^T X)^{-1} X^T \sigma^2 I X (\lambda I + X^T X)^{-1}\end{aligned}$$

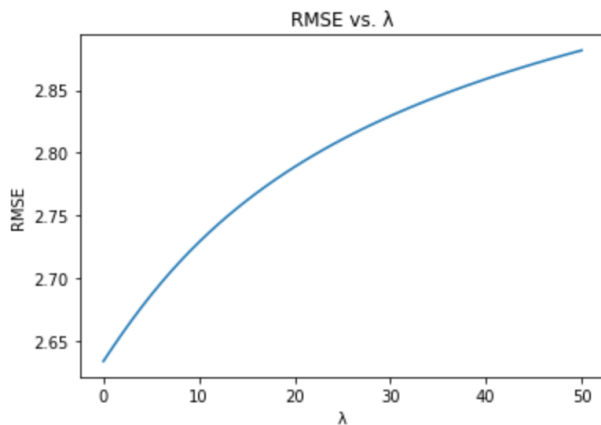
Problem 3.

Part I. a)



- b) Year_Made and Weight clearly stands out from others. This suggests that as $df \nearrow$, $\lambda \rightarrow 0$, The w_{RR} getting closer to w_{LS} . Hence, Year_Made and Weight are two dimensions with significant w_{LS} value.

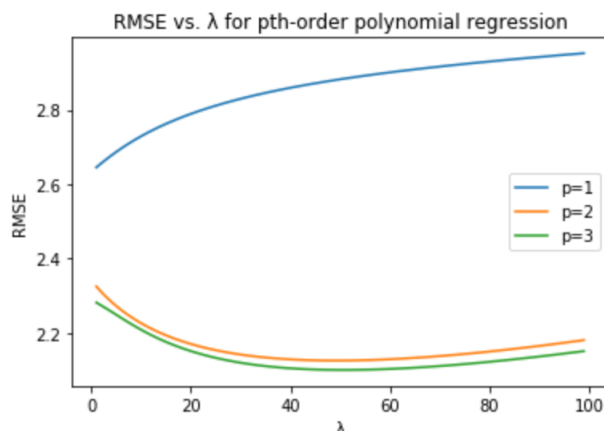
c)



The RMSE increases monotonically as λ increases. This indicates that we would prefer a smaller λ value and we would prefer least squares to ridge regression.

Part II.

d)



According to the graph, I would choose value of $p = 2$, as 2nd-order polynomial model have significant better (smaller RMSE) than 1st-Order and pretty close to 3-rd Order RMSE.

The RMSE decreases monotonically from 0 to around 40. However, that large λ value would push w -rr closer to 0. A plot of w -rr with df for 2nd-Order Model would help here.