

P1

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

$$L(\lambda; x_1, \dots, x_n) = \prod_{j=1}^n \frac{e^\lambda}{x_j!} \lambda^{x_j}$$

b)

$$L(\lambda; x_1, \dots, x_n) = -n\lambda - \sum_{j=1}^n (\ln(x_j!)) + (\ln(\lambda)) \sum_{j=1}^n x_j$$

$$\frac{\partial L(\lambda; x_1, \dots, x_n)}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{j=1}^n x_j = 0$$

$$\lambda_{ML} = \frac{1}{n} \sum_{j=1}^n x_j$$

c) Let $\bar{x} = (x_1, \dots, x_n)$ be a vector.

$p(\lambda) = \text{gamma}(a, b)$ as a prior

$$\lambda \mapsto \frac{f(\bar{x}; \lambda) p(\lambda)}{\int_\lambda f(\bar{x}; \lambda) p(\lambda) d\lambda}$$

$$\lambda_{MAP} = \underset{\lambda}{\operatorname{argmax}} \frac{f(\bar{x}; \lambda) p(\lambda)}{\int f(\bar{x}; \lambda) p(\lambda) d\lambda} = \underset{\lambda}{\operatorname{argmax}} f(\bar{x}; \lambda) p(\lambda)$$

$$\lambda_{MAP} = \underset{\lambda}{\operatorname{argmax}} p(\lambda) \prod_{i=1}^n P(X_i | \lambda) = \underset{\lambda}{\operatorname{argmax}} (\log(p(\lambda))) + \sum_{i=1}^n (\log f(x_i | \lambda))$$

$$\log f(x_i | \lambda) = x_i \log(\lambda) - \lambda$$

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

$$\frac{\partial \sum_{i=1}^n \log f(x_i | \lambda)}{\partial \lambda} = \frac{1}{\lambda} \sum_{i=1}^n x_i - n$$

$$\log p(\lambda) = \log \left[\frac{b^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-b\lambda} \right]$$

$$= \log \frac{b^\alpha}{\Gamma(\alpha)} + (\alpha-1) \log(\lambda) - b\lambda$$

$$\frac{\partial \log p(\lambda)}{\partial \lambda} = \frac{\alpha-1}{\lambda} - b$$

$$\frac{\partial (\log(p(\lambda)) + \sum_{i=1}^n \log f(x_i | \lambda))}{\partial \lambda} = \frac{\alpha-1}{\lambda} - b + \frac{1}{\lambda} \sum_{i=1}^n x_i - n = 0$$

$$\frac{\alpha-1 + \sum x_i}{\lambda} = n + b$$

$$\lambda_{MAP} = \frac{\alpha-1 + \sum x_i}{n+b}$$

a)

From Bayes rules,

$$P(\lambda | \bar{X}) = \frac{p(\bar{X}|\lambda) P(\lambda)}{\int_0^\infty p(\bar{X}|\lambda) P(\lambda) d\lambda}$$

$$P(\lambda | \bar{X}) \propto p(\bar{X}|\lambda) P(\lambda)$$

$$P(\bar{X}|\lambda) P(\lambda) = \prod_{i=1}^n e^{-\lambda} \frac{1}{x_i!} \lambda^{x_i} \frac{b^a}{P(a)} \lambda^{a-1} e^{-b\lambda}$$

$$P(\lambda | \bar{X}) \propto \prod_{i=1}^n (x_i! e^{-\lambda}) \lambda^{a-1} e^{-b\lambda} \\ = \lambda^{\sum_{i=1}^n x_i + a - 1} e^{-(n+b)\lambda}$$

which is proportional to pdf of gamma($\sum_{i=1}^n x_i + a, n+b$)
Posterior dist of λ should be gamma($\sum_{i=1}^n x_i + a, n+b$)

e) $\lambda | \bar{X} \sim \text{gamma}(\sum_{i=1}^n x_i + a, n+b)$

$$E(\lambda | \bar{X}) = \frac{\sum_{i=1}^n x_i + a}{n+b} \quad \lambda_{ML} = \frac{\sum x_i}{n} \\ \lambda_{MAP} = \frac{\sum x_i + a - 1}{n+b}$$

$$\text{Var}(\lambda | \bar{X}) = \frac{\sum_{i=1}^n x_i + a}{(n+b)^2}$$

Mean of λ under the Posterior is close to λ_{ML} and λ_{MAP} . In particular, it's only off by $\frac{1}{n+b}$ from the λ_{MAP} . The mean of λ under posterior can also be expressed by a compromise between prior mean $\frac{a}{b}$ and the λ_{ML} .

$$\frac{\sum x_i + a}{n+b} = \lambda_{ML} \left(1 - \frac{b}{n+b}\right) + \frac{a}{b} \left(\frac{b}{n+b}\right)$$

As $n=0$, mean reduces to $\frac{a}{b}$ prior mean

As $n \rightarrow \infty$, mean reduces to λ_{ML}

Problem 2

$$a) \quad y_i \stackrel{iid}{\sim} N(x_i^T w, \sigma^2)$$

$$w_{RR} = (\lambda I + x^T x)^{-1} x^T y$$

$$E(w_{RR}) = E[(\lambda I + x^T x)^{-1} x^T y]$$

$$= (\lambda I + x^T x)^{-1} x^T E(y)$$

$$= (\lambda I + x^T x)^{-1} x^T x_w$$

$$\text{Var}(w_{RR}) = E[(w_{RR} - E(w_{RR})) (w_{RR} - E(w_{RR}))^T]$$

$$= E[w_{RR} w_{RR}^T] - E(w_{RR}) E(w_{RR})^T$$

$$\text{Var}(w_{RR}) = E[(\lambda I + x^T x)^{-1} x^T y y^T x (\lambda I + x^T x)^{-1}] - E(w_{RR}) E(w_{RR})^T$$

$$= (\lambda I + x^T x)^{-1} x^T E(y y^T) x (\lambda I + x^T x)^{-1} - E(w_{RR}) E(w_{RR})^T$$

$$\text{Var}(w_{RR}) = (\lambda I + x^T x)^{-1} x^T (\sigma^2 I + x_w x_w^T x^T) x (\lambda I + x^T x)^{-1} \\ - E(w_{RR}) E(w_{RR})^T$$

$$E(w_{RR}) E(w_{RR})^T = (\lambda I + x^T x)^{-1} x^T x_w w^T x^T x (\lambda I + x^T x)^{-1}$$

$$\text{Var}(w_{RR}) = (\lambda I + x^T x)^{-1} x^T \sigma^2 I x (\lambda I + x^T x)^{-1}$$

$$\text{let } Z = (\lambda (x^T x)^{-1} + I)^{-1}$$

$$(\lambda I + x^T x)^{-1} = (x^T x (\lambda (x^T x)^{-1} + I))^{-1} = (x^T x Z^{-1})^{-1}$$

$$\text{Var}(w_{RR}) = (x^T x Z^{-1})^{-1} x^T \sigma^2 I x (x^T x Z^{-1})^{-1}$$

$$= Z (x^T x)^{-1} x^T \sigma^2 x Z (x^T x)^{-1}$$

$$= \sigma^2 Z (x^T x)^{-1} Z$$

b)

$$w_{RR} = (\lambda I + x^T x)^{-1} x^T y = (\lambda I + x^T x)^{-1} \underbrace{x^T y}_{w_{LS}}$$

$$= [(\lambda I + x^T x)^{-1} + I]^{-1} (x^T x) w_{LS}$$

$$= (\lambda (x^T x)^{-1} + I)^{-1} (x^T x)^{-1} (x^T x) w_{LS}$$

$$= (\lambda (x^T x)^{-1} + I)^{-1} w_{LS}$$

$$\text{let } X = USV^T \quad (X^T X)^{-1} = V S^{-2} V^T$$

$$w_{RR} = (X^T X)^{-1} \tau^{-1} w_{LS}$$

$$= (\lambda V S^{-2} V^T + \tau I)^{-1} w_{LS}$$

$$= V (\lambda S^{-2} + \tau) V^T w_{LS}$$

$$= V M V^T w_{LS}$$

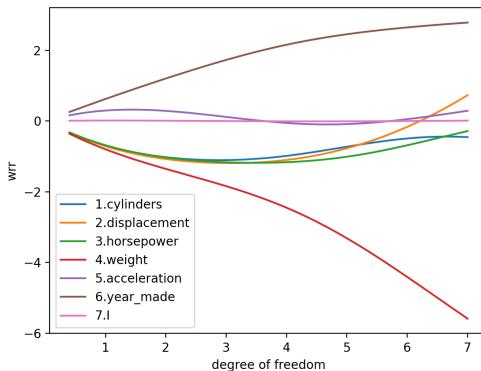
with M to be a diagonal matrix with

$$M_{ii} = \frac{\lambda s_i^2}{\lambda + s_i^2}$$

Problem 3

Part 1

a)



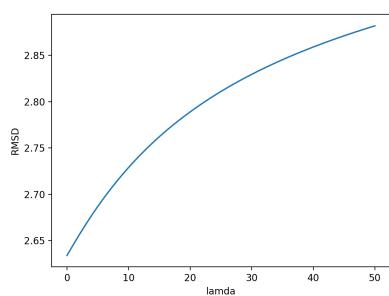
b) dimension 4 and dimension 6 stand out.

As $\lambda \rightarrow \infty$, $df(w) \rightarrow 0$, dimension 6 is always larger than 0, which means it consistently gives positive prediction to the data and its magnitude is greater than 1,2,3,5,7 dims.

It means it's more trustworthy than 1,2,3,5,7 features.

Same as dim 4, which is always negative and its magnitude is larger than 1,2,3,5,7 dims. This means it's gives negative impact to the prediction and it's trustworthy and not very sensitive to λ selection.

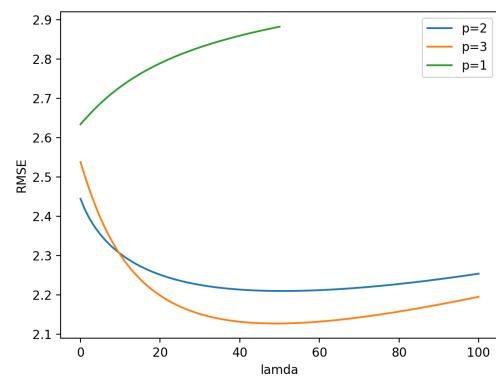
c)



As λ increases, the RMSE increases as well.

This figure tells we should set $\lambda = 0$ in this dataset, which means we should pick least square for this dataset.

P art2
d)



I would choose $p=3$
since it has the lowest
RMSE at $\lambda=40$.
The ideal value of λ
changes to $\lambda=40$ for this
problem.