$P_1$ (a) $P(x_1, x_2 \cdots x_N | \lambda) = \prod_{i=1}^{N} P(x_i | \lambda) = \prod_{i=1}^{N} \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = \frac{\lambda^{\sum_{i=1}^{N} x_i} e^{-N\lambda}}{\prod_{i=1}^{N} x_i!}$

(b) $\lambda_{ML} = \arg\max_{\lambda} P(x_1, x_2 \cdots x_N) = \arg\max_{\lambda} \prod_{i=1}^{N} \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right)$

$\nabla_\lambda \prod_{i=1}^{N} \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = 0$

$\because \ln \left( \prod_{i}^{N} \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = \sum_{i}^{N} \ln \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right)$

$\therefore \nabla_\lambda \sum_{i=1}^{N} \ln \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = \sum_{i=1}^{N} \nabla_\lambda \ln \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) = \sum_{i=1}^{N} \nabla_\lambda \left( \ln \frac{\lambda^{x_i}}{x_i!} - \lambda \right) = \sum_{i=1}^{N} \nabla_\lambda \left( \ln \lambda^{x_i} - \ln x_i! - \lambda \right)$

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

$\therefore \lambda_{ML} = \frac{1}{N} \sum_{i=1}^{N} x_i$

(c) $\lambda_{MAP} = \arg\max_{\lambda} \ln P(\lambda | x)$

$= \arg\max_{\lambda} \ln \frac{P(x|\lambda) P(\lambda)}{P(x)}$

$= \arg\max_{\lambda} \left( \ln P(x|\lambda) + \ln P(\lambda) - \ln P(x) \right)$ $\qquad \Gamma(a) = (a-1)!$

$\lambda_{MAP} = \arg\max_{\lambda} \left( \sum_{i=1}^{N} \ln \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) + \ln \frac{b^a \lambda^{a-1} e^{-b\lambda}}{\Gamma(a)} \right)$

$= \arg\max_{\lambda} \left[ \sum_{i=1}^{N} \left( \ln \lambda^{x_i} - \ln x_i! - \lambda \right) + \ln \beta^a + \ln \lambda^{a-1} - b\lambda - \ln \Gamma(a) \right]$

$\nabla_\lambda \left[ \sum_{i=1}^{N} \left( \ln \lambda^{x_i} - \ln x_i! - \lambda \right) + \ln \lambda^{a-1} - b\lambda + Const \right]$

$= \sum_{i=1}^{N} \left( \frac{x_i \lambda^{x_i-1}}{\lambda^{x_i}} - 1 \right) + \frac{(a-1)\lambda^{a-2}}{\lambda^{a-1}} - b$

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

$\frac{\sum_{i=1}^{N} x_i - \lambda N}{\lambda} + \frac{a-1}{\lambda} - b = 0$

$\sum_{i=1}^{N} x_i - \lambda N + a - 1 - \lambda = 0$

$\lambda_{MAP} = \frac{\sum_{i=1}^{N} x_i + a - 1}{N + b}$

(d) $P(\lambda | x) = \frac{P(x|\lambda) P(\lambda)}{P(x)}$

$\propto P(x|\lambda) P(\lambda)$

$\propto \prod_{i=1}^{N} \left( \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) \cdot \frac{\lambda^{a-1} b^a e^{-b\lambda}}{(a-1)!}$

$\propto \prod_{i=1}^{N} \left( \frac{\lambda^{x_i}}{x_i!} \right) e^{-(N+b)\lambda} \cdot \frac{\lambda^{a-1} b^a}{(a-1)!}$

$$\propto \frac{\lambda^{\sum_{i=1}^{N} x_i + (a-1)} b^a e^{-(N+b)\lambda}}{\prod_{i=1}^{N} x_i! \, (a-1)!}$$

$$\propto \lambda^{\sum_{i=1}^{N} x_i + a - 1} e^{-(N+b)\lambda}$$

$$\therefore \, P(\lambda|x) \sim Gamma\left(\sum_{i=1}^{N} x_i + a, \, N+b\right)$$

(e)
$$mean = \frac{\sum_{i=1}^{N} x_i + a}{N+b}$$

$$var = \frac{\sum_{i=1}^{N} x_i + a}{(N+b)^2}$$

Discuss: $\lambda_{MAP}$ is the mode of the mean of this posteriori distribution. it's the maximum value of this ~~mean~~ mean except for $a, b$. $\lambda_{ML}$ is similar to the $\lambda$. When dataset is large $\lambda_{ML}$ equals $\lambda_{MAP}$.    $\lambda_{MAP}$ considers prior distribution compared to $\lambda_{ML}$.

P2  $E[W_{RR}] = E[(\lambda I + X^T X)^{-1} X^T y]$

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

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

$Var[W_{RR}] = E[W_{RR} W_{RR}^T] - E[W_{RR}] E[W_{RR}]^T$     Let $A = (\lambda I + X^T X)^{-1}$

$\because E[yy^T]$
$= \sigma^2 I + X w w^T X^T$

$$= A X^T E[yy^T] X A^T - A X^T X w w^T X^T X A^T$$

$$= A X^T \sigma^2 I X A^T + A X^T X w w^T X^T X A^T - A X^T X w w^T X^T X A^T$$

$$= A X^T \sigma^2 I X A^T$$

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

Let $z = (I + \lambda(X^T X)^{-1})^{-1}$

$\therefore A = z(X^T X)^{-1}$      $((X^T X)^{-1})^T = ((X^T X)^T)^{-1} = (X^T X)^{-1}$

$\therefore Var[W_{RR}] = z(X^T X)^{-1} X^T \sigma^2 X (X^T X)^{-1} z^T = z(X^T X) z^T \sigma^2$

P3.  Part 1.

(a)

df(λ)-Wrr figure
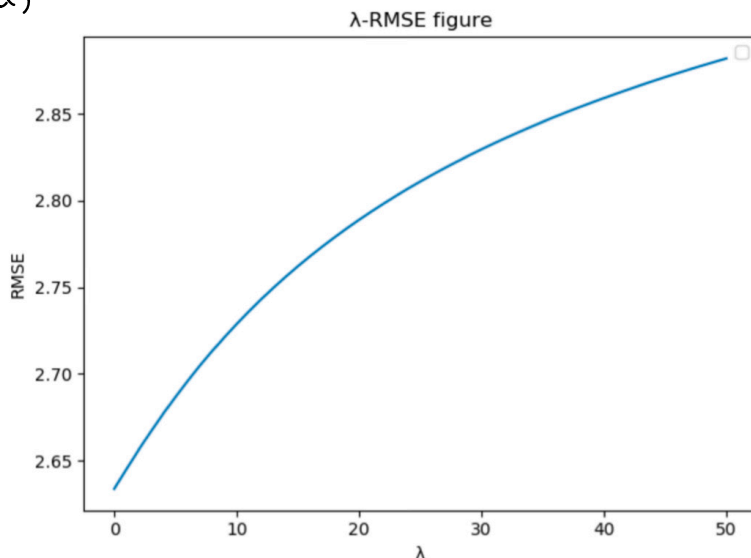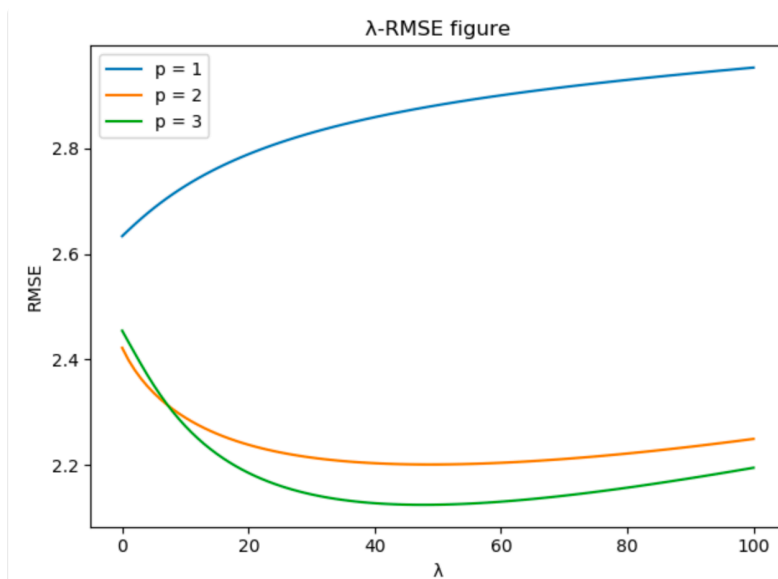
(b) The features "Year made" and "weight" clearly stand out.
From the figure, we can see that with the increase of df(λ). "Year made"
feature's Wrr increase a lot, which indicating that with the decrease of λ.
increase of df(λ), this "Year made" feature Wrr weights more and becomes
more important. As for "weight", with the increase of df(λ), decrease of λ.
this "weight" feature's Wrr decreases a lot, indicating that this feature
weights less and becomes less important.

(c)



λ-RMSE figure

From the figure, we can see that with the increase of $\lambda$, the "RMSE" becomes larger and larger, which is not idea. For this problem, the figure shows that it's better to choose small $\lambda$ which can get smaller RMSE. $\lambda = 0$, which is the least square solution, meaning that least square solution is better for this problem.

(d)



Base on the plot, when $\lambda$ is less than (about) 10, I'll choose p=2 since it'll lead to the smallest RMSE. When $\lambda$ is greater than (about) 10. I'll choose p=3 since it'll lead to the smallest RMSE.

When p=1. the plot shows that with the increase of $\lambda$. RMSE increases. So the idea $\lambda$ is 0 for p=1.

When p=2 and 3, the plot shows that RMSE first decreases as $\lambda$ increasing. When $\lambda$ is about 40, RMSE starts increasing. So the idea $\lambda$ for p=2 and 3 is 40