2.

② 1. Given $\quad y = w^T \phi(x).$

Considering gaussian distribution for each dimension of $y$.

Let $\quad y_i = [ y_i^{(1)} \; y_i^{(2)} \ldots y_i^{(k)} ] \quad 1 \leq i \leq N.$

$\phi(x) = [a_{ij}]_{M \times N} \quad$ where $a_{ij} = j^{th}$ dimension value of $x_j$

Let $\quad y_p' = [ y_1^{(p)} \; y_2^{(p)} \ldots y_N^{(p)} ]^T$

$\Rightarrow \quad y_p' | x \sim \mathcal{N}(0, \sigma^2) + \phi(x)^T w_p$

$$L = \exp\left( -\frac{1}{2\sigma^2} \left( \phi(x)^T w_p - y_p' \right)^2 \right)$$

$$\log L = \frac{-1}{2\sigma^2} \left( \phi(x)^T w_p - y_p' \right)^2$$

find optimal value for $w_i$

$\Rightarrow \quad \textcircled{0} \quad \dfrac{\partial \log(L)}{\partial w_i} = 0$

$\Rightarrow$

$\quad \phi(x) y_p' - \phi(x) \phi(x)^T w_p = 0.$

$\Rightarrow$

$\quad w_p = \left( \phi(x) \phi(x)^T \right)^{-1} \phi(x) y_p'$

$\Rightarrow$

$\quad w = \left( \phi(x) \phi(x)^T \right)^{-1} \phi(x) y'$

∴ MLE of $w$

$$w = (\phi(x)\phi(x)^T)^{-1}\phi(x)y.$$

## MAP estimate

Assuming a gaussian prior for $w_i$ with parameter as $\lambda$.

$$p(w_i|\lambda) = \left(\frac{\lambda}{2\pi}\right)^{M/2} \exp\left(-\frac{\lambda}{2} w_i^T w\right)$$

from Baye's rule, we have

$$p(w/x,y,\lambda) \propto \underbrace{p(x,y/w)}_{\downarrow} \cdot \underbrace{p(w/\lambda)}_{\downarrow \text{prior}}$$

$$\underbrace{\phantom{p(w/x,y,\lambda)}}_{\downarrow \text{Posterior}}$$

$$\Rightarrow p(w/x,y,\lambda) \propto p(x,y/w) \cdot \prod_{i=1}^{K} p(w_i/\lambda).$$

$$\alpha \ \exp\left( \frac{-1}{2} \sum_{i=1}^{N} \sum_{j=1}^{k} \frac{y_i^{(j)} - \omega_j^T \phi(x_i)}{\sigma_j^2} \right) \exp\left( -\sum_{j=1}^{K} \frac{\lambda}{2} \omega_j^T \omega_j \right)$$

$$\Rightarrow \log\left( P(\omega \mid x, y, \lambda) \right) = \frac{-1}{2} \sum_{i=1}^{N} \sum_{j=1}^{k} \frac{y_i^{(j)} - \omega_j^T \phi(x_i)}{\sigma_j^2}$$

$$- \sum_{j=1}^{k} \frac{\lambda}{2} \omega_j^T \omega_j + \text{Constant}$$

maximizing log-MAP to find the estimate for $\omega_j$.

$$\frac{\partial \log\left( P(\omega \mid x, y, \lambda) \right)}{\partial \omega_j} = 0$$

$$\Rightarrow \quad -\phi(x) y_j' + \phi(x)\phi(x)^T \omega_j + \lambda \omega_j = 0.$$

$$\Rightarrow \quad \omega_j = \left( \phi(x)\phi(x)^T + \lambda I \right)^{-1} \phi(x) y_j'$$

MAP estimate of $W$

$$\Rightarrow \quad W = \left( \phi(x)\phi(x)^T + \lambda I \right)^{-1} \phi(x \mid y.$$

which is same as solution for ridge regression.

Q4 Considering $\phi(0) = (1, 0)^T$, $\phi(1) = (0, 1)^T$

$$\phi(x) = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \qquad y = w^T \phi(x)$$

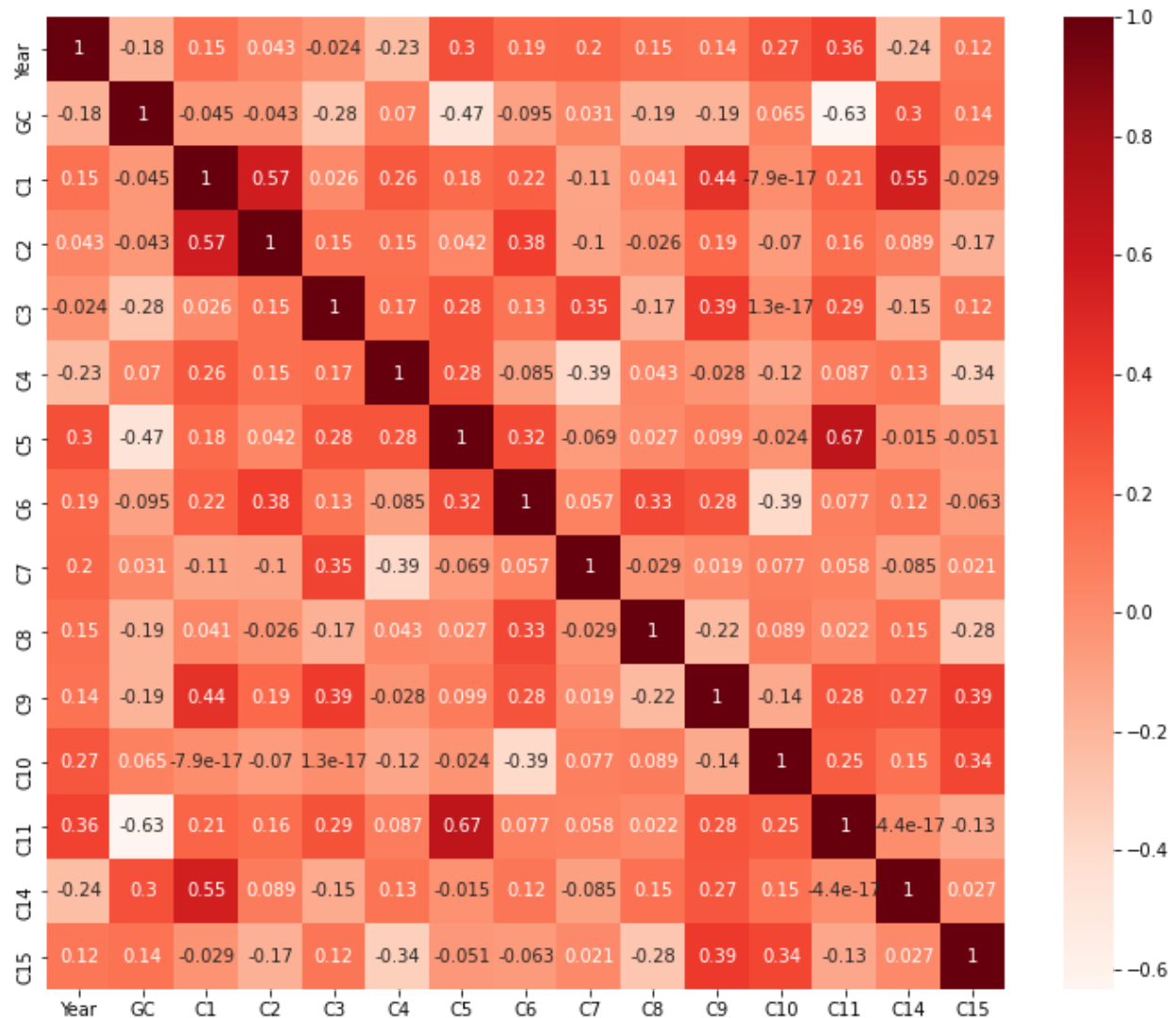$$y = \begin{bmatrix} -1 & -1 \\ -1 & -2 \\ -2 & -1 \\ 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}$$

MLE of $w = (\phi(x) \, \phi^T(x))^{-1} \phi(x) y$

$$= \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ -1 & -2 \\ -2 & -1 \\ 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & \frac{1}{3} \end{bmatrix} \begin{bmatrix} -4 & -4 \\ 4 & 4 \end{bmatrix} = \begin{bmatrix} -\frac{4}{3} & -\frac{4}{3} \\ \frac{4}{3} & \frac{4}{3} \end{bmatrix}.$$

$\therefore \quad w = \begin{bmatrix} -\frac{4}{3} & -\frac{4}{3} \\ \frac{4}{3} & \frac{4}{3} \end{bmatrix}$.

3. Since, the number of deaths due to horse kick is not correlated with year number and the input has no other features, we will take expectation of poisson distribution as the predicted value for any input.

| | Year | GC | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 | C14 | C15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Year | 1 | -0.18 | 0.15 | 0.043 | -0.024 | -0.23 | 0.3 | 0.19 | 0.2 | 0.15 | 0.14 | 0.27 | 0.36 | -0.24 | 0.12 |
| GC | -0.18 | 1 | -0.045 | -0.043 | -0.28 | 0.07 | -0.47 | -0.095 | 0.031 | -0.19 | -0.19 | 0.065 | -0.63 | 0.3 | 0.14 |
| C1 | 0.15 | -0.045 | 1 | 0.57 | 0.026 | 0.26 | 0.18 | 0.22 | -0.11 | 0.041 | 0.44 | -7.9e-17 | 0.21 | 0.55 | -0.029 |
| C2 | 0.043 | -0.043 | 0.57 | 1 | 0.15 | 0.15 | 0.042 | 0.38 | -0.1 | -0.026 | 0.19 | -0.07 | 0.16 | 0.089 | -0.17 |
| C3 | -0.024 | -0.28 | 0.026 | 0.15 | 1 | 0.17 | 0.28 | 0.13 | 0.35 | -0.17 | 0.39 | 1.3e-17 | 0.29 | -0.15 | 0.12 |
| C4 | -0.23 | 0.07 | 0.26 | 0.15 | 0.17 | 1 | 0.28 | -0.085 | -0.39 | 0.043 | -0.028 | -0.12 | 0.087 | 0.13 | -0.34 |
| C5 | 0.3 | -0.47 | 0.18 | 0.042 | 0.28 | 0.28 | 1 | 0.32 | -0.069 | 0.027 | 0.099 | -0.024 | 0.67 | -0.015 | -0.051 |
| C6 | 0.19 | -0.095 | 0.22 | 0.38 | 0.13 | -0.085 | 0.32 | 1 | 0.057 | 0.33 | 0.28 | -0.39 | 0.077 | 0.12 | -0.063 |
| C7 | 0.2 | 0.031 | -0.11 | -0.1 | 0.35 | -0.39 | -0.069 | 0.057 | 1 | -0.029 | 0.019 | 0.077 | 0.058 | -0.085 | 0.021 |
| C8 | 0.15 | -0.19 | 0.041 | -0.026 | -0.17 | 0.043 | 0.027 | 0.33 | -0.029 | 1 | -0.22 | 0.089 | 0.022 | 0.15 | -0.28 |
| C9 | 0.14 | -0.19 | 0.44 | 0.19 | 0.39 | -0.028 | 0.099 | 0.28 | 0.019 | -0.22 | 1 | -0.14 | 0.28 | 0.27 | 0.39 |
| C10 | 0.27 | 0.065 | -7.9e-17 | -0.07 | 1.3e-17 | -0.12 | -0.024 | -0.39 | 0.077 | 0.089 | -0.14 | 1 | 0.25 | 0.15 | 0.34 |
| C11 | 0.36 | -0.63 | 0.21 | 0.16 | 0.29 | 0.087 | 0.67 | 0.077 | 0.058 | 0.022 | 0.28 | 0.25 | 1 | 4.4e-17 | -0.13 |
| C14 | -0.24 | 0.3 | 0.55 | 0.089 | -0.15 | 0.13 | -0.015 | 0.12 | -0.085 | 0.15 | 0.27 | 0.15 | -4.4e-17 | 1 | 0.027 |
| C15 | 0.12 | 0.14 | -0.029 | -0.17 | 0.12 | -0.34 | -0.051 | -0.063 | 0.021 | -0.28 | 0.39 | 0.34 | -0.13 | 0.027 | 1 |

For proving that the feature year is not correlated we will first plot the Pearson correlation heatmap and see there isn't much correlation of corps with year. Hence, we use expectation as the predicted value.

**For maximum likelihood estimation:**

Given observations $y_1, y_2, \ldots y_n$ for input $x_1, x_2 \ldots x_n$.

$$L(\theta) = f\left(x_1, x_2, \ldots, x_n \mid \theta\right)$$ if $\theta$ is true val of param. the probability that we observe $x_1, x_2 \ldots x_n$.

for MLE, we maximize $L(\theta)$.

As our data points are iid,

maximize $L(\theta) = f\left(x_1 \mid \theta\right) \cdot f\left(x_2 \mid \theta\right) \ldots \cdot f\left(x_n \mid \theta\right)$

For likelihood on poisson distribution,

$$P(Y \mid \lambda) = f\left(y_1 \mid \lambda\right) \cdot f\left(y_2 \mid \lambda\right) \ldots \cdot f\left(y_n \mid \lambda\right)$$

$$\Rightarrow P(Y \mid \lambda) = \frac{e^{-n\lambda} \; \lambda^{\sum_{i=1}^{n} y_i}}{\prod_{i=1}^{n} y_i!} \quad ; \quad \left(\because f\left(y_1 \mid \lambda\right) = \frac{e^{-\lambda} \lambda^{y_i}}{y_i!}\right)$$

$$\& \; Y = y_1, y_2, \ldots y_n.$$

$$\lambda_{ML} = \arg\max_{\lambda} \log\left(P(Y \mid \lambda)\right)$$

$$\Rightarrow \frac{d}{d\lambda}\left(-n\lambda + \sum_{i=1}^{n} y_i \, \log\lambda + (-1) \log\left(\prod_{i=1}^{n} y_i!\right)\right) = 0$$

$$\Rightarrow \lambda_{ML} = \frac{\sum_{i=1}^{n} y_i}{n}.$$

Refer **Table 1** below for the poisson parameters (ML) and rmse values.

**For maximum aposteriori estimation:**

The gamma distribution is the conjugate prior for the likelihood function - poisson. Hence, we choose gamma distribution for prior distribution. It has

two hyper parameters: (alpha. beta) which can be found via grid search over training set.

*gamma distribution with parameters $\rightarrow (\alpha, \beta)$:*

$$P\left(\lambda \mid \alpha, \beta\right) = \frac{\beta^{\alpha} \lambda^{\alpha-1} e^{-\beta\lambda}}{\Gamma(\alpha)}$$

$$P\left(\lambda \mid y, \alpha, \beta\right) \propto P(Y \mid \lambda) \cdot P\left(\lambda \mid \alpha, \beta\right)$$

$$\propto \frac{e^{-(\beta+n)\lambda} \; \lambda^{\left(\sum_{i=1}^{n} y_i + \alpha - 1\right)}}{\Gamma(\alpha) \prod_{i=1}^{n} y_i!}$$

$$\lambda_{MAP} = \underset{\lambda}{\text{argmax}} \; \log P\left(\lambda \mid y, \alpha, \beta\right)$$

$$\Rightarrow \frac{d}{d\lambda}\left(-(\beta+n)\lambda + \left(\sum_{i=1}^{n} y_i + \alpha - 1\right) \log \lambda - \log\left(\Gamma(\alpha) \prod_{i=1}^{n} y_i!\right)\right) = 0$$

$$\Rightarrow -(\beta+n) + \frac{\sum_{i=1}^{n} y_i + \alpha - 1}{\lambda} = 0 \Rightarrow \lambda = \frac{\sum_{i=1}^{n} y_i + \alpha - 1}{n + \beta}$$

After doing a grid search over possible values of (alpha, beta) from {1..10} and taking Σ(rmse) as the cost metric over all corps on the training set, we find that alpha = 2 and beta = 1 is the best fit for it.
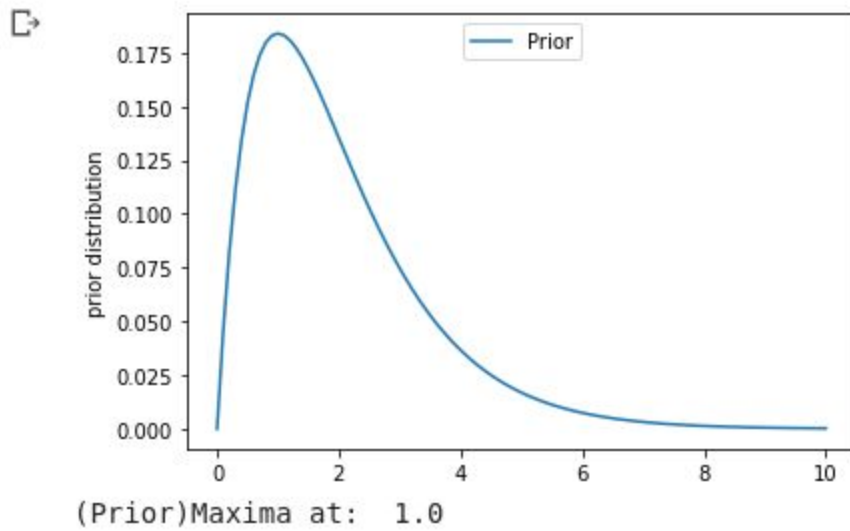
Finally **Table 1** below shows the MAP parameters and rmse values.

| Corp | ML | | MAP | |
| --- | --- | --- | --- | --- |
| | λ | RMSE | λ | RMSE |
| G | 1 | 0.755929 | 1 | 0.755929 |
| I | 0.692308 | 1.11244 | 0.714286 | 1.10657 |
| II | 0.615385 | 0.729756 | 0.642857 | 0.731925 |
| III | 0.615385 | 0.729756 | 0.642857 | 0.731925 |
| IV | 0.461538 | 0.484764 | 0.5 | 0.5 |
| V | 0.384615 | 0.587989 | 0.428571 | 0.553283 |
| VI | 0.846154 | 0.989804 | 0.857143 | 0.989743 |
| VII | 0.538462 | 0.898011 | 0.571429 | 0.892143 |
| VIII | 0.307692 | 0.509421 | 0.357143 | 0.5 |
| IX | 0.692308 | 0.738393 | 0.714286 | 0.742307 |
| X | 0.538462 | 1.15969 | 0.571429 | 1.14286 |
| XI | 1 | 1.13389 | 1 | 1.13389 |
| XIV | 1.46154 | 1.02381 | 1.42857 | 1 |
| XV | 0.307692 | 0.941214 | 0.357143 | 0.928571 |

**Table 1**

# Graphs

## Prior



(Prior)Maxima at:   1.0

# Corp 2



(Likelihood)Maxima at:   0.6
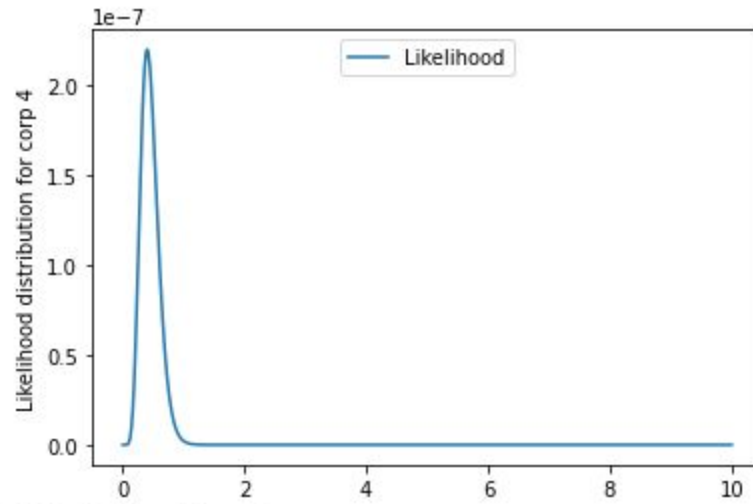


(Posterior)Maxima at:   0.62
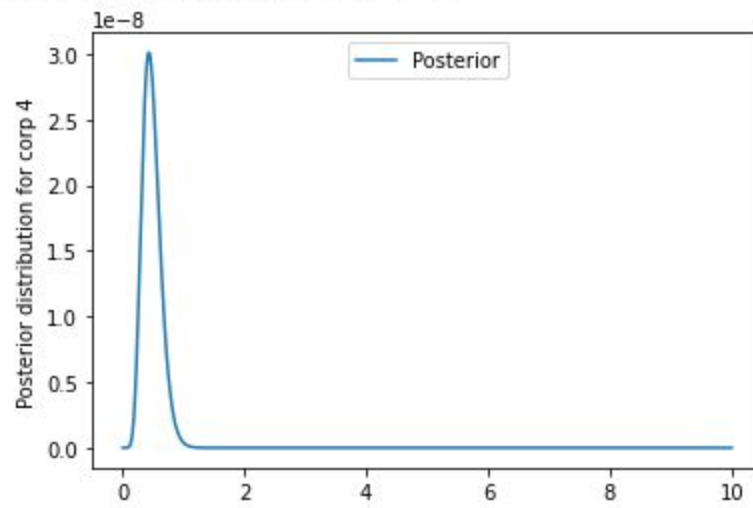
# Corp 4



(Likelihood)Maxima at:  0.4



(Posterior)Maxima at:  0.43

# Corp 6
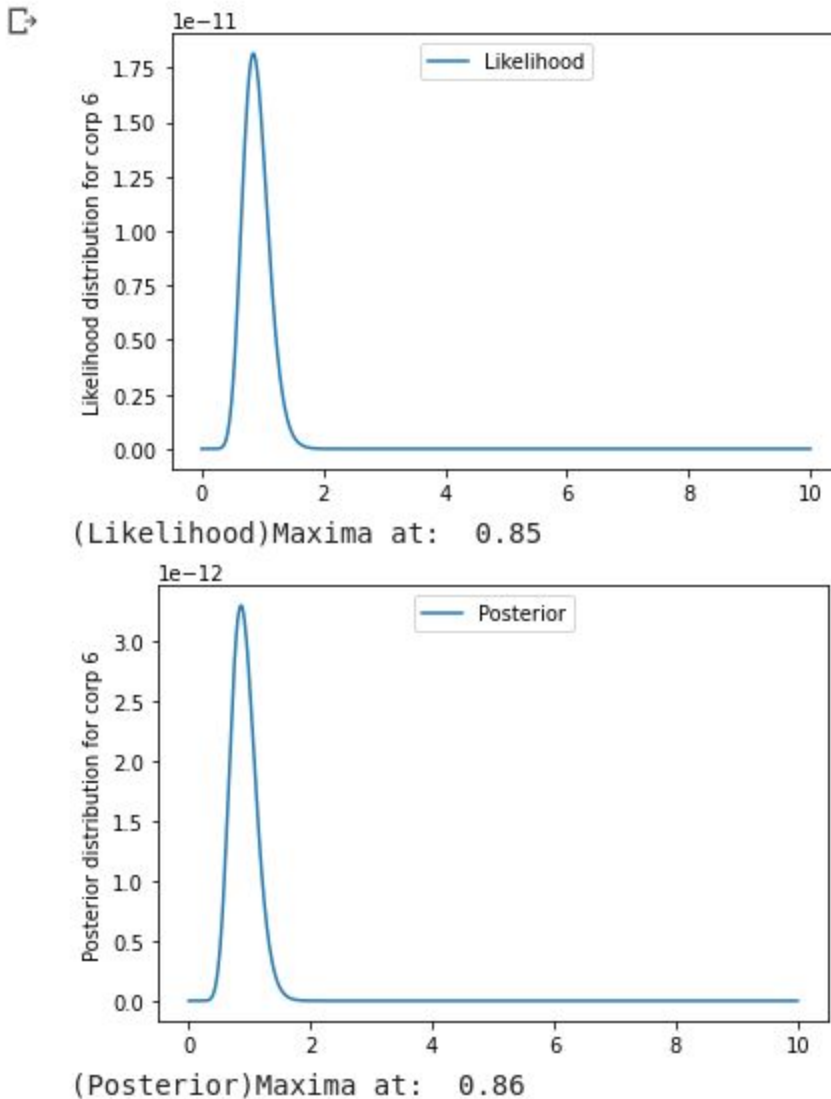


(Likelihood)Maxima at:  0.85

(Posterior)Maxima at:  0.86

The above maximas are obtained after considering the whole dataset for Likelihood and Prior graphs. The graph has been plotted against lambda values with precision 0.1, Maximas for which have been written below the corresponding ones.