Q.1. Probability

$X_1, X_2$ ind. continuous random var. uniformly distr. on $[0,1]$.

$X = \max(X_1, 2X_2)$.

(1) $E[X]$.

cdf: $F_X(x) = P(X \leq x) = Pr(\max(X_1, 2X_2))$

$\qquad = \underset{F_{X_1}(x)}{Pr(X_1 \leq x)} \cdot \underset{F_{X_2}(x)}{Pr(2X_2 \leq x)}$

pdf: $f_X(x) = \dfrac{d}{dx} F_X(x)$

$\qquad = \dfrac{d}{dx} F_{X_1}(x) \cdot F_{X_2}(x)$

$\qquad = F_{X_1}(x) f_{X_2}(x) + f_{X_1}(x) \cdot F_{X_2}(x)$.

$f_{X_1}(x) = $ pdf of $X_1 = 1$.

$f_{X_2}(x) = $ pdf of $2X_2 = \dfrac{1}{2}$.

$\therefore f_X(x) = \displaystyle\int_0^1 x\, dx + \int_1^2 \dfrac{x}{2}\, dx$

$\qquad = \left[\dfrac{x^2}{2}\right]_0^1 + \left[\dfrac{x^3}{6}\right]_1^2 = \dfrac{13}{12}$.

(2) $Var(X) = E[X^2] - (E[X])^2$

$E[X^2] = \displaystyle\int_0^1 x^2 \cdot x\, dx + \int_1^2 \dfrac{x^2}{2}\, dx$

$\qquad = \left[\dfrac{x^4}{4}\right]_0^1 + \left[\dfrac{x^3}{6}\right]_1^2 = \dfrac{34}{24} = \dfrac{17}{12}$

Q.2. Parameter Estimation.

$$P(X = k | \lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \qquad k \in \{0, 1, 2, \ldots\}$$

(2.1)   Give log-likelihood fn. of X given $\lambda$

① $L(X_1 \ldots X_n | \lambda) = P(X_1|\lambda) \cdot P(X_2|\lambda) \cdots P(X_n|\lambda)$

$$= \prod_{i=1}^{n} P(X_i|\lambda)$$

Taking log likelihood

$$= \log\left(\prod_{i=1}^{n} P(X_i|\lambda)\right) = \log\left(\prod_{i}^{n} \frac{\lambda^{X_i} e^{-\lambda}}{X_i!}\right)$$

$$= \sum_{i=1}^{n} \log\left(\frac{\lambda^{X_i} e^{-\lambda}}{X_i!}\right)$$

$$= \sum_{i=1}^{n} \left[\log(\lambda^{X_i} e^{-\lambda}) - \log X_i!\right]$$

$$\boxed{L(\lambda) = \log \lambda \cdot \sum_{i=1}^{n} X_i - n\lambda - \sum_{i=1}^{n} \log(X_i!)}$$

② Compute MLE of $\lambda$ in general case.

~~Take derivative~~ $\underset{\substack{arg \\ max \\ \lambda}}{L(\lambda)} = \log \lambda \cdot \sum_{i=1}^{n} X_i - n\lambda - \sum_{i=1}^{n} \log(X_i!)$

Take derivative wrt. $\lambda$ and set it to 0.

i.e $\frac{\partial}{\partial \lambda} L(\lambda) = 0$. $\Rightarrow \frac{\partial}{\partial \lambda}\left(\log \lambda \cdot \sum_{i=1}^{n} X_i - n\lambda - \sum_{i=1}^{n} \log(X_i!)\right) = 0$.

$\therefore \quad \frac{1}{\lambda} \sum_{i=1}^{n} X_i - n - 0 = 0$.

$\therefore \quad \boxed{\lambda = \frac{1}{n} \sum_{i=1}^{n} X_i}$

③ Compute MLE for $\lambda$ using the observed X.

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

$$\therefore \lambda = \frac{1}{7} \times [4+5+3+5+6+9+3] = \frac{35}{7} = 5.$$

(2.2) **MAP.**
$$P(\lambda|\alpha,\beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot \lambda^{\alpha-1} e^{-\beta\lambda}, \quad \lambda > 0$$

① Compute posterior distribution over $\lambda$.

$$P(\lambda|x_1 \ldots x_n) = P(x_1 \ldots x_n|\lambda) \cdot P(\lambda) \qquad --- ①$$

$$P(x_1 \ldots x_n|\lambda) = \prod_{i=1}^{n} P(x_i|\lambda)$$

we already have calculate $\prod_{i=1}^{n} P(x_i|\lambda)$ in (2.1.1)

$$\therefore P(x_1 \ldots x_n|\lambda) = \log\lambda \sum_{i=1}^{n} x_i - n\lambda - \sum_{i=1}^{n} \log(x_i!)$$

$$= \frac{\lambda^{\sum_{i=1}^{n} x_i} e^{-n\lambda}}{\prod_{i=1}^{n}(x_i!)}$$

Put in equation ① .

$$P(\lambda|x_1 \ldots x_n) = \frac{\lambda^{\sum_{i=1}^{n} x_i} e^{-n\lambda}}{\prod_{i=1}^{n} x_i!} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot \lambda^{\alpha-1} \cdot e^{-\beta\lambda}$$

$$\boxed{P(\lambda|x_1 \ldots x_n) = e^{-\lambda(n+\beta)} \cdot \lambda^{\sum x_i + \alpha - 1} \cdot \frac{\beta^\alpha}{\prod_{i=1}^{n} x_i! \cdot \Gamma(\alpha)}}$$

(2-2)

② Derive an analytic expression for the maximum a posterior (MAP) estimate of $\lambda$.

$$\lambda_{MAP} = \underset{\lambda}{argmax} \; Log\left( e^{-\lambda(n+\beta)} \cdot \lambda^{\sum_{i=1}^{n} X_i + \alpha - 1} \cdot \frac{\beta^{\alpha}}{\prod_{i=1}^{n} X_i! \, \Gamma(\alpha)} \right)$$

$$\cancel{\text{Take derivative}} = \underset{\lambda}{argmax}\left[ -\lambda(n+\beta) + \left(\sum_{i=1}^{n} X_i + \alpha - 1\right) Log\,\lambda \right.$$

$$\left. + \; Log\left( \frac{\beta^{\alpha}}{\prod_{i=1}^{n} X_i! \, \Gamma(\alpha)} \right) \right]$$

Take derivative wrt $\lambda$ and set it to 0.

$$\therefore \quad 0 = -(n+\beta) + \frac{1}{\lambda} \sum_{i=1}^{n} X_i + \alpha - 1.$$

$$\therefore \quad \boxed{\lambda = \frac{\sum_{i=1}^{n} X_i + \alpha - 1}{n+\beta}}$$

(2.3) $\qquad X \sim poisson(\lambda), \quad n = e^{-2\lambda}$

① $\hat{n} = e^{-2X}$. Show that $\hat{n}$ is the MLE of $n$.

We have, $\quad P(X=k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \qquad \{poisson's \; distr.\}$

$\cancel{\text{bt}}$ Also, $\quad n = e^{-2\lambda} \Rightarrow \lambda = \frac{-Log\,n}{2}$.

Subs. value of $\lambda$.

$$P(X=k|\lambda) = \frac{\left(\frac{-Log\,n}{2}\right)^k e^{-\left(\frac{-Log\,n}{2}\right)}}{k!}$$

For MLE of $n$, take derivative wrt $n$ & set to 0.

i.e $\frac{d}{d\eta} = 0$ . $\Rightarrow \eta_{MLE} = \frac{df}{d\eta}$

Take log likelihood : $L(\eta) = k \log\left(-\frac{\log\eta}{2}\right) + \frac{\log\eta}{2} - \log k!$

Take derivative and set to zero.

$\therefore \quad 0 = \frac{d}{d\eta}\left(k\log\left(-\frac{\log\eta}{2}\right) + \frac{\log\eta}{2} - \log k!\right)$

$$\frac{k}{\eta\log\eta} + \frac{1}{2\eta} = 0$$

$\therefore \quad \boxed{\eta_{MLE} = e^{-2k}}$ $\quad \therefore \hat{\eta}$ is the MLE of $\eta$.

② Show that: bias of $\hat{\eta}$ is $e^{-2\lambda} - e^{\lambda\left(\frac{1}{e^2}-1\right)}$

We know that : bias $(\hat{\eta}) = E[\hat{\eta}] - \eta$ . —①

$E[\hat{\eta}] = E[e^{-2X}] = \sum_{k=0}^{\infty} e^{-2k} P(X=k)$ $\quad \{\sum k \cdot P(X=k)\}$

$P(X=k) = \frac{\lambda^k \cdot e^{-\lambda}}{k!}$ $\quad$ {poisson's distr}

$\therefore E[e^{-2X}] = \sum_{k=0}^{\infty} e^{-2k} \cdot \frac{\lambda^k \cdot e^{-\lambda}}{k!}$

$= e^{-\lambda} + \frac{e^{-2}\lambda e^{-\lambda}}{1!} + \frac{e^{-4}\lambda^2 e^{-\lambda}}{2!} + \cdots$

$= e^{-\lambda}\left[1 + \frac{e^{-2}\lambda}{1!} + \frac{e^{-4}\lambda^2}{2!} + \cdots\right]$

Let's put $e^{-2\lambda} = Z$

∴ ~~BEX~~ $E[e^{-2X}] = e^{-\lambda}\left[1 + \dfrac{Z}{1!} + \dfrac{Z^2}{2!} + \cdots\right]$

Now, $1 + Z + \dfrac{Z^2}{2!} + \cdots$ is Taylor expansion of $e^Z$.

∴ $E[e^{-2X}] = \left(\dfrac{Z^k}{k!}\right) e^{-\lambda} \cdot = e^Z \cdot e^{-\lambda}$.

Reput $Z = e^{-2\lambda}$.

$E[e^{-2X}] = e^{e^{-2}\cdot\lambda}e^{-\lambda}$

$= e^{\lambda(-1 + e^{-2})}$

$= e^{\lambda\left(\frac{1}{e^2} - 1\right)}$

Subs. in ①.

② bias$(\hat{n}) = E[\hat{n}] - \hat{n}$

$\hat{n} = e^{-2\lambda}$.

∴ bias$(\hat{n}) = e^{\lambda\left(\frac{1}{e^2} - 1\right)} - e^{-2\lambda}$.

$\boxed{\text{bias}(\hat{n}) = -e^{-2\lambda} + e^{\lambda(1/e^2 - 1)}}$.

(2.3)

③ we know: bias$(\hat{n}) = E[\hat{n}] - n$ —— (i)

Here, $\hat{n} = (-1)^X$.

∴ $E[\hat{n}] = E[(-1)^X] = \displaystyle\sum_{k=0}^{\infty} (-1)^k \cdot \dfrac{\lambda^k e^{-\lambda}}{k!}$ $\left\{\begin{array}{l} E(X) = \\ \sum k \cdot P(X=k) \end{array}\right\}$

$= e^{-\lambda} - \dfrac{\lambda e^{-\lambda}}{1!} + \dfrac{\lambda^2 e^{-\lambda}}{2!} - \cdots$

$$E[(-1)^x] = e^{-\lambda}\left(1 - \frac{\lambda}{1!} + \frac{\lambda^2}{2!} - \frac{\lambda^3}{3!} + \cdots\right)$$

$1 - \lambda + \frac{\lambda^2}{2!} - \cdots$    is Taylor exp$^n$ of $e^{-\lambda}$.

$$\therefore \quad E[(-1)^x] = e^{-\lambda} \cdot e^{-\lambda} = e^{-2\lambda}.$$

Put in ①.

$$\text{bias}(\hat{n}) = E[\hat{n}] - n.$$
$$= e^{-2\lambda} - e^{-2\lambda}$$
$$= 0.$$

When bias $= 0$, $\Rightarrow$ estimator is unbiased.

Here, we considered only one observation, which came out to be really close to real value. This means that, though it worked well (unbiased) for 1 sample, doesn't generalize when sample space increases.


(3) <u>Regression and MLE.</u>

(3.1)    $y_i = w^T x_i + \epsilon_i$,    $\epsilon_i$ are $N(0, \sigma_i^2)$

Taking conditional likelihood.

$$P(y \mid x_i, w, \sigma_i) = \prod_{i=1}^{n} P(y_i \mid x_i, w, \sigma_i)$$
$$= \prod_{i=1}^{n} \left(\frac{1}{\sqrt{2\pi}\,\sigma_i}\right) \exp\left[\frac{-1}{2\sigma_i^2}(y_i - w^T x_i)^2\right]$$

Taking log likelihood.
$$= \left(\frac{1}{\sqrt{2\pi}}\right)^n \left[\log\left(\frac{1}{\sigma_i}\right) + \log\left(\exp\left(\sum_{i=1}^{n} \frac{-1}{2\sigma_i^2}(y_i - w^T x_i)^2\right)\right)\right]$$

calculating MLE.

$$L(w) = \underset{w}{\arg\max} \quad \left(\frac{1}{\sqrt{2\pi}}\right)^n \left[ \log\left(\frac{1}{\sigma_i}\right) + \log\left(\exp\left(\sum_i^n \frac{-1}{2\sigma_i^2}(y_i - w^T x_i)^2\right)\right)\right]$$

Take partial derivative wrt. w & set to 0.

First, we can simplify $L(w)$.

consider $\quad s_i = \frac{1}{2\sigma_i^2 \pi}$.

$$\therefore \quad L(w) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \left[\log\left(\frac{1}{\sigma_i}\right) + \sum_i^n s_i (y_i - w^T x_i)^2\right]$$

$$0 = \frac{d}{dw}\left(-\sum_i^n s_i \left(y_i^2 - 2 y_i w^T x_i + w^T \cdot w x_i^2\right)\right)$$

$$0 = \frac{d}{dw}\left(- s_i \| y - x^T w \|^2\right)$$

ie $\quad -2 x s^T (y - x^T w) = 0$.

$$x s^T y = x s^T x^T w$$

$$\therefore \quad w = (x s^T x^T)^{-1} x s^T y$$

(3.2) $\quad p(\epsilon_i) = \frac{1}{2b} \exp\left(-\frac{|\epsilon_i|}{b}\right)$

$$L(w) = \underset{w}{\arg\min} \quad - \log \prod_i^n P(y_i | x_i, w, b)$$

$$= \underset{w}{\arg\min} \quad - \prod_{i=1}^n \frac{1}{2b} \exp\left(-\frac{|\epsilon_i|}{b}\right)$$

$$= \quad - \prod_{i=1}^n \frac{1}{2b} \exp\left(-\frac{|y_i - w^T x_i|}{b}\right)$$

$$= \quad - \left(\frac{1}{2b}\right)^n \exp\left[\frac{1}{b} \sum_i^n -|y_i - w^T x_i|\right]$$

**Q.4.3 Try out your work on synthetic data**

1.

Q. Compare Sparsity Pattern of your Lasso solution to the true model of parameters of w*

Sparsity pattern of the Lasso Solution: 57

Sparsity Pattern of the true model w* = 10

Q. Precision: 0.175438596491

Q. Recall: 1

Q. How well are you able to discover the true nonzeros?

The true non-zeros are the first k elements of lasso solution which are non-zero compared with the first k elements of w* which are non-zero. In my case, it is exactly same (k non-zero). So, the ability to discover true non-zeros is 100%

Q. Comment on how Lambda affects these results

For Precision: The value of precision is 1 and then becomes zero as Lambda reaches 0

For Recall: The value of recall increases linearly with lambda until at a point it reaches 1 and then flattens out

2. Q. How are precision and recall affected when sigma=10?

After setting Lambda = 3, the precision and recall values change as follows:
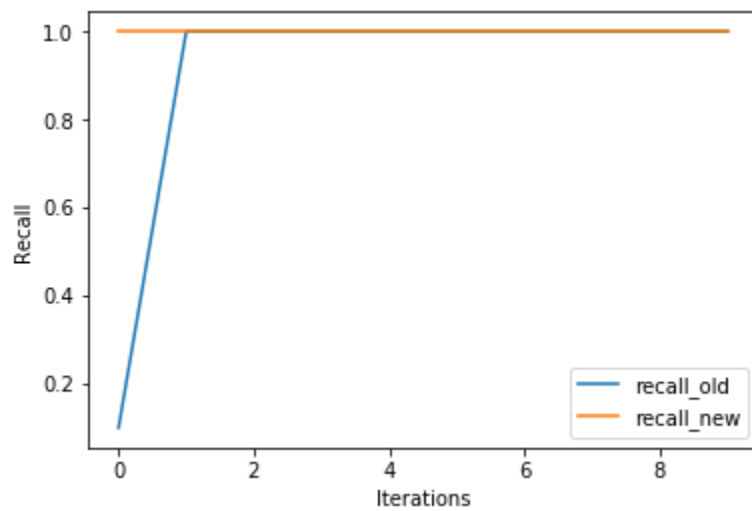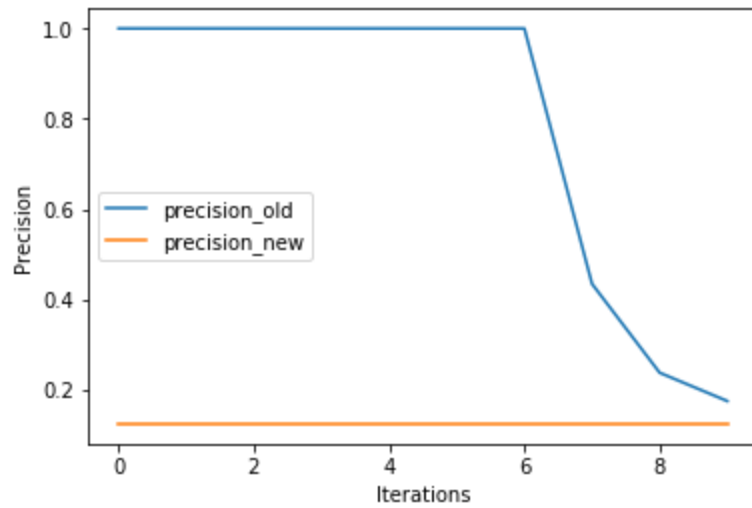
Precision:

[0.125, 0.125, 0.125, 0.125, 0.125, 0.125, 0.125, 0.125, 0.125, 0.125]

Recall:

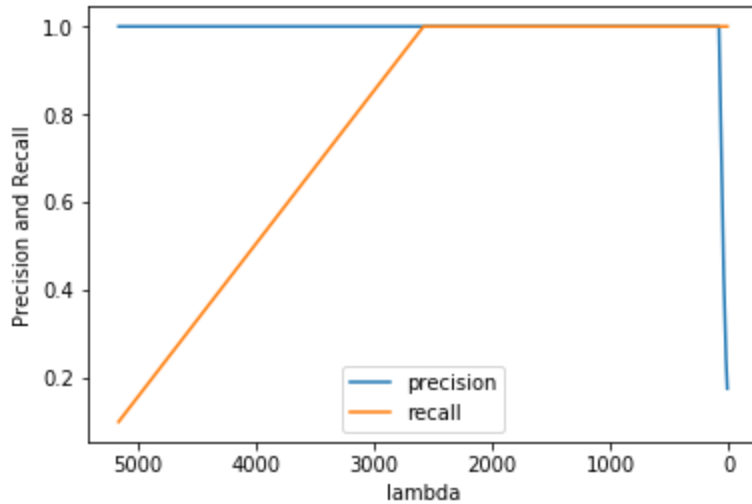[1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]

Once we get an optimal Lambda value, then the precision remains the same (close to) as obtained for sigma = 1. Also, the recall is 1 for all iterations. This means that changing the standard deviation of the noise does not affect the precision or the recall

How might you change the Lambda to achieve better precision or recall?

Changing of lambda values does not change the Precision and Recall

Plots:

**Q.4.4 Predicting goodness points of a wine given its review**

1.

<u>RMSE Training</u>:

[3.6659156152848538, 3.2726658154603703, 2.9817760693142046, 2.7292553433116624, 2.477646256792003, 2.286892526860811, 2.1212885463916873, 1.9524599775424407, 1.8111901053647608, 1.710865888501754, 1.6477442630087837]
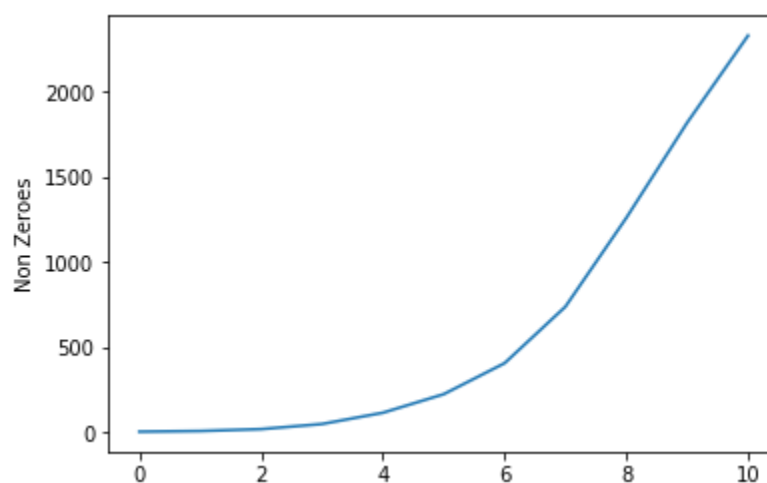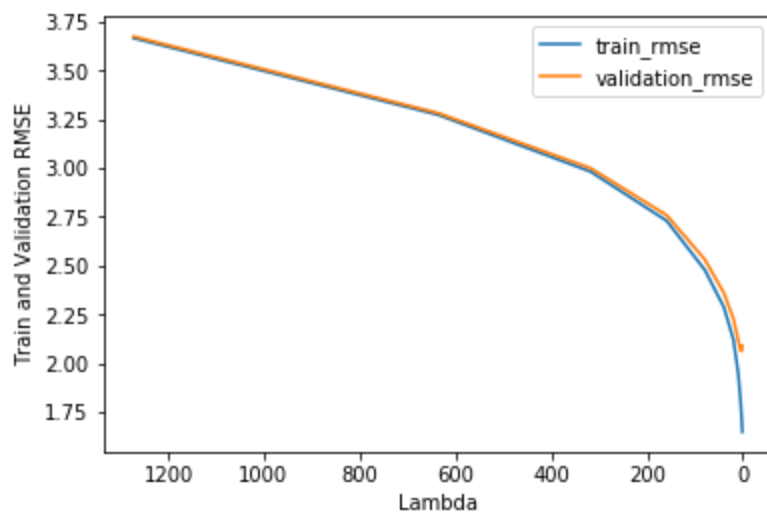
<u>RMSE Validation</u>:

[3.674432595136281, 3.281769756113443, 2.999308942322199, 2.7580247829868267, 2.531297191857193, 2.3604884258078784, 2.2285187942171194, 2.118712422561662, 2.064567203663436, 2.062058813848828, 2.0888203249606496]

<u>Number of Non-Zeros</u>:

[6, 10, 21, 51, 117, 226, 408, 740, 1261, 1820, 2331]

<u>Graphs</u>:

2. Lambda = 2.4841204121093856

   Top 10 Features with highest weight:

   ['acidity provides', 'truly', 'nearly', 'sweet black', 'lemony', 'ageability', 'lifesaver', 'big', 'stars', 'spearmint']

   Top 10 Features with lowest weight:

    ['earns', 'high', 'cherry berry', 'soft', 'sparkler', 'liqueur', 'cuts', 'semillon', 'brightened', 'banana']

3. RMSE from Kaggle: 1.94250