

# Machine Learning - HW1

Faizaan Charania  
111463646

We know that  $X_1$  &  $X_2$  are independent and uniformly distributed on  $[0, 1]$

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

Let  $X_3 = 2X_2 \therefore X_3$  is uniformly distributed on  $[0, 2]$

$$\therefore X = \max(X_1, X_3)$$

•  $X$  performs differently in ranges  $0 \leq x \leq 1$  and  $1 \leq x \leq 2$

$\therefore$  for  $\forall x, 0 \leq x \leq 1$

$$C.D.F(X) = F_X(x)$$

$$= P_R(X \leq x)$$

$$= P_R(\max(X_1, X_3) \leq x)$$

$$= P_R(X_1 \leq x, X_3 \leq x)$$

$$= P_R(X_1 \leq x) \cdot P_R(X_3 \leq x) \quad [\because X_1 \perp X_3]$$

$$= F_{X_1}(x) \cdot F_{X_3}(x)$$

$$P.D.F(x) = \frac{d}{dx} F_{x_1}(x) \cdot F_{x_3}(x)$$

$$= F_{x_1}(x) \cdot f_{x_3}(x) + F_{x_3}(x) \cdot f_{x_1}(x)$$

$$= \frac{x-0}{1-0} \cdot \frac{1}{2} + \frac{(x-0)}{(2-0)} \cdot 1$$

$$= \frac{x}{2} + \frac{x}{2} = x //$$

for  $\forall x$ ,  $1 < x \leq 2$

$$\begin{aligned} C.D.F(x) &= P_R(X \leq x) \\ &= P_R(\max(x_1, x_3) \leq x) \\ &= \end{aligned}$$

$\therefore x_1$  is not defined for  $1 < x \leq 2$

$$F_x(x) = P_R($$

$$F_{x_1}(x) = 1 \quad \forall x, 1 < x \leq 2$$

$$\begin{aligned} C.D.F(x) &= P_R(X_1 \leq x) \cdot P_R(X_3 \leq x) \\ &= 1 \cdot \frac{x}{2} \end{aligned}$$

$$= \frac{x}{2}$$

$$\therefore P.D.F(x) = f_x(x) = \frac{d}{dx} F_x(x) = \frac{d}{dx} \frac{x}{2} = \frac{1}{2} //$$

1)  $E(X)$ 

$$\begin{aligned}
 E[X] &= \int_0^2 x f_x(x) dx \\
 &= \int_0^1 x f_x(x) dx + \int_1^2 x f_x(x) dx \\
 &= \int_0^1 x \times x dx + \int_1^2 x \times \frac{1}{2} dx \\
 &= \int_0^1 x^2 dx + \int_1^2 \frac{x}{2} dx \\
 &= \left[ \frac{x^3}{3} \right]_0^1 + \left[ \frac{x^2}{4} \right]_1^2 = \frac{1}{3} + \frac{3}{4}
 \end{aligned}$$

$$= \frac{13}{12} //$$

2)  $\text{Var}(X)$ 

$$\text{Var}(X) = E[X^2] - (E[X])^2$$

$$\begin{aligned}
 E[X^2] &= \int_0^2 x^2 f_x(x) dx \\
 &= \int_0^1 x^2 f_x(x) dx + \int_1^2 x^2 f_x(x) dx
 \end{aligned}$$

$$= \int_0^1 x^3 dx + \int_1^2 \frac{x^2}{2} dx$$

$$= \left[ \frac{x^4}{4} \right]_0^1 + \left[ \frac{x^3}{6} \right]_1^2$$

$$= \frac{1}{4} + \frac{7}{6} = \frac{17}{12}$$

$$\therefore \text{Var}(X) = E[X^2] - (E[X])^2$$

$$= \frac{17}{12} - \left( \frac{13}{12} \right)^2 = \frac{17}{12} - \frac{169}{144}$$

$$= \frac{204 - 169}{144} = \frac{35}{144} //$$

$$\therefore \text{Var}(X) = \frac{35}{144}$$

~~//~~

Q2

2.1.1 Log-likelihood of  $X$  given  $\lambda$

We know that  $X$  given  $\lambda$  is a Poisson distribution

$$\text{Likelihood}(X|\lambda) = P(X|\lambda)$$

$$= P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \lambda)$$

$\because$  All  $X_i$  are independent, we get

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

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

$\therefore$  log-likelihood is

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

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

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

$$\therefore \text{Log-likelihood is } \sum_{i=1}^n (x_i \log \lambda - \log(x_i!)) - n\lambda$$

2) M.L.E for  $\lambda$  in general case

$$\text{M.L.E} = \underset{\lambda}{\operatorname{argmax}} \sum_{i=1}^n (x_i \log \lambda - \log(x_i!)) - n\lambda$$

To find max we differentiate & equate to 0.

$$\therefore \frac{\partial}{\partial \lambda} (\sum (x_i \log \lambda - \log(x_i!)) - n\lambda) = 0$$

$$\sum_{i=1}^n x_i \frac{\partial \log \lambda}{\partial \lambda} - \sum_{i=1}^n \frac{\partial \log(x_i!)}{\partial \lambda} - n \frac{\partial \lambda}{\partial \lambda} = 0$$

$$\therefore \frac{\sum_{i=1}^n x_i}{\lambda} - 0 - n = 0$$

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

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

2.1.3.

To compute M.L.E for  $\lambda$  over observed  $X$   
we can use previous solution 2.1.2

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

$$\therefore \hat{\lambda} = \frac{4+5+3+5+6+9+3}{7}$$

$$= \frac{35}{7} = 5 //$$

9.2  
9.2.1  
9.2.2  
M.A.P

Given  $\lambda \sim \Gamma(\alpha, \beta)$

$\lambda$  is a Gamma distribution with a p.d.f

$$p(\lambda | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \lambda > 0$$

Posterior distribution over  $\lambda$

$$P(\lambda | x) = \frac{P(x|\lambda) \cdot P(\lambda | \alpha, \beta)}{P(x)}$$

$$\therefore P(\lambda | x) \propto P(x|\lambda) \cdot P(\lambda | \alpha, \beta)$$

$$P(x|\lambda) \cdot P(\lambda | \alpha, \beta) = \left( \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right) \times \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \lambda > 0$$

$$= \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod_{i=1}^n (x_i i)!} \times \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot \lambda^{\alpha-1} e^{-\beta\lambda}$$

$$\therefore P(x|\lambda) \cdot P(\lambda|\alpha, \beta)$$

$$= \frac{\sum_{i=1}^n x_i + \alpha - 1}{\prod_{i=1}^n x_i!} \times \frac{\beta^\alpha}{\Gamma(\alpha)} e^{-\lambda(\beta+n)}$$

$$\therefore P(\lambda|x) \propto P(x|\lambda) P(\lambda|\alpha, \beta)$$

$$P(\lambda|x) \propto \frac{\sum_{i=1}^n x_i + \alpha - 1}{\prod_{i=1}^n x_i!} \times \frac{\beta^\alpha}{\Gamma(\alpha)} e^{-\lambda(\beta+n)}$$

Find Maximum a posterior (MAP)

estimate of  $\lambda$

To find M.A.P of  $\lambda$ , we take

derivative of Posterior Distribution of  $\lambda$ .

w.r.t  $\lambda$  and then equate it to zero.

(Take log to make calculation easy)

$$\therefore \underset{\lambda}{\log} (P(\lambda | x)) = 0$$

$$\log(P(\lambda | x)) = \log \left[ \frac{\beta}{\prod_{i=1}^n (x_i!)^{\alpha}} \times \frac{\lambda^{\sum_{i=1}^n x_i + \alpha - 1}}{\Gamma(\alpha)} e^{-\lambda(\beta+n)} \right]$$

$$= \alpha \log \beta + \left( \sum_{i=1}^n x_i + \alpha - 1 \right) \log \lambda - \lambda(\beta+n) -$$

$$\log \left( \prod_{i=1}^n (x_i!)^{\alpha} \right) - \log (\Gamma(\alpha))$$

$$\therefore \frac{\delta}{\delta \lambda} [\log P(\lambda | x)] =$$

$$\left[ \sum_{i=1}^n x_i + \alpha - 1 \right] \times \frac{1}{\lambda} - (\beta + n) = 0$$

$$\therefore \hat{\lambda} = \frac{\left[ \sum_{i=1}^n x_i + \alpha - 1 \right]}{\beta + n}$$

23

## Estimator Bias

Let  $\hat{\eta} = e^{-2x}$

To prove,

$\therefore \text{M.L.E of } \eta = \hat{\eta} = e^{-2x}$

$$\hat{\eta} = \underset{\eta}{\operatorname{argmax}} P(x|\eta)$$

$$\eta = e^{-2\lambda} \quad (\text{given})$$

$$\therefore \log \eta = -2\lambda$$

$$\therefore \lambda = \left( \frac{-\log \eta}{2} \right) \quad \text{---(i)}$$

We know  $X$  is a Poisson s.t

$$X \sim \text{Poisson}(\lambda)$$

$$\therefore P(X|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \{ k=0, 1, 2, \dots \} \quad \text{---(ii)}$$

From (i) & (ii) we get

$$P(X|\eta) = \frac{(\log \eta)^k}{(-2)^k k!} e^{-(\log \eta / 2)} \quad \text{---(iii)}$$

$$\therefore P(X|n) = \frac{(\log n)^k}{(-2)^k k!} n^{1/2} \quad \text{--- (iii)}$$

Then to find out M.L.E (n)

take derivative of (iii) w.r.t n & equate

to zero

$$\frac{\partial}{\partial n} P(X|n) = \frac{k (\log n)^{k-1}}{(-2)^k k!} \left( \frac{1}{n} \right) \times n^{1/2} +$$

$$\frac{(\log n)^k}{(-2)^k k!} \times \frac{1}{2} \times n^{-1/2}$$

$$= \frac{1}{(-2)^k k!} \times (\log n)^{k-1} \left[ \frac{k}{n^{1/2}} + \frac{(\log n)^k}{2 n^{1/2}} \right] = 0$$

$$\therefore \frac{k}{n^{1/2}} + \frac{(\log n)^k}{2 n^{1/2}} = 0$$

$$\therefore (\log n)^k = -k$$

$$\hat{n} = e^{-k}$$

if we only have one observation  
of  $x$

$$\hat{\eta} = e^{-2x}$$

2.3.2 Bias ( $\hat{\eta}$ ) =  $E[\hat{\eta}] - \eta$

$$\eta = e^{-2\lambda}$$

$$E[\hat{\eta}] = \sum_{x \geq 0} \hat{\eta} P(x)$$

$$= \sum_{x \geq 0} e^{-2x} \frac{\lambda^x}{x!} e^{-\lambda}$$

$$= e^{-\lambda} \sum_{x \geq 0} e^{-2x} \frac{\lambda^x}{x!}$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} (\lambda e^{-2})^x \frac{1}{x!}$$

Using the Taylor Series

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

$$\therefore E[\hat{\eta}] = e^{-\lambda} \cdot e^{\lambda e^{-2}}$$

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

$$\therefore \text{bias}(\hat{\eta}) = E[\hat{\eta}] - \eta$$

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

$\therefore \text{bias}(\hat{\eta})$  is just a value, we can also say  $\text{bias}(\hat{\eta}) = e^{-2\lambda} - e^{\lambda(1+e^{-2}-1)}$

We assume  $\hat{\eta} = (-1)^x$

∴ To find out the expected value, we do

$$E[\hat{\eta}] = \sum_{x=0}^{\infty} \hat{\eta} P(x)$$

$$\hat{\eta} = (-1)^x$$

$$\therefore E[\hat{\eta}] = \sum_{x=0}^{\infty} (-1)^x \frac{\lambda^x e^{-\lambda}}{x!}$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(-\lambda)^x}{x!}$$

Using the Taylor series,  $e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$

$$\therefore \sum_{x=0}^{\infty} \frac{(-\lambda)^x}{x!} = e^{-\lambda}$$

$$\therefore E[\hat{\eta}] = e^{-\lambda} * e^{-\lambda} = e^{-2\lambda}$$

$$\begin{aligned}\text{bias}(\hat{\eta}) &= E[\hat{\eta}] - \eta \\ &= e^{-2\lambda} - e^{-2\lambda} \\ &= 0\end{aligned}$$

$\therefore (-1)^X$  is an unbiased estimator

$\hat{\eta} = (-1)^X$  totally depends on ~~the~~  
whether  $X$  is even or odd.

For all even  $X$ ,  $\eta$  is positive &  
 $\eta$  is negative for odd values of  $X$ .

? It is a bad estimator to use.

23

# Regression and M.L.E

$y_i = \omega^T x_i + \epsilon_i$  is generated by  
a linear regression model.

3.1 Assumptions  $\epsilon_1, \dots, \epsilon_n$  are independent  
each with mean 0 & variance  $\sigma_i^2$

$$\therefore \epsilon_i \sim N(0, \sigma_i^2)$$

$$P(Y|w, \sigma) = P(Y=y_1|w, \sigma) \cdot P(Y=y_2|w, \sigma) \cdots$$

$$P(Y=y_n|w, \sigma)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma_i} \exp\left(-\frac{1}{2\sigma_i^2} (y_i - \omega^T x_i)^2\right)$$

∴ log-likelihood is

$$\log(P(Y|w, \sigma)) = \log\left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma_i} \exp\left(-\frac{1}{2\sigma_i^2} (y_i - \omega^T x_i)^2\right)\right)$$

$$= \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}\sigma_i} \right) - \sum_{i=1}^n \left( \frac{1}{2\sigma_i^2} (y_i - w^T x_i)^2 \right)$$

a                    b

$\therefore a$  is independent of  $w$ , to find

$\max P(Y|w, \sigma)$  we ~~maximise~~ find

$\min b$

$\therefore$  M.L.E

$$\hat{w} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \frac{1}{2\sigma_i^2} (y_i - w^T x_i)^2$$

we can put  $\frac{1}{2\sigma_i^2} = \delta_i$

$$\therefore \hat{w} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \delta_i (y_i - w^T x_i)^2$$

$\therefore$  M.L.E can be found by minimising  
the above function.

32

Likelihood is still  $P(Y|w, \epsilon)$

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

$\because \epsilon \sim \text{Laplace}(b)$

Taking log on both sides

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

$$= \sum_{i=1}^n \left[ \log \frac{1}{2b} - \frac{|y_i - w^T x_i|}{b} \right]$$

$$= \underbrace{\sum_{i=1}^n \log \frac{1}{2b}}_a - \underbrace{\frac{1}{b} \sum_{i=1}^n |y_i - w^T x_i|}_b$$

As we have seen in Q. 3.1 a is independent of  $w$ .

$\therefore$  As done before

$$\text{M.L.E of } w = \hat{w} = \underset{w}{\operatorname{argmax}} \frac{1}{b} \sum_{i=1}^n |y_i - w^T x_i|$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n |y_i - w^T x_i|$$

$$\therefore \hat{\omega} = \underset{\omega}{\operatorname{argmin}} |y_i - \omega^T x_i|$$

3.3

Given the model, the residual is

$$r = \sum_{i=1}^n y_i - (\omega^T z_i + \epsilon_i)$$

If we assume  $\epsilon_i$  to be from a Gaussian distribution  $N(0, \sigma^2)$

the residual is depended on the square of parameter  $\sigma$ .

If we assume the noise to be a Laplace distribution, then it is only linearly depended on the value of  $b$ .

∴ In case of noise, the model with Laplace noise will be more robust.

# Programming

Qh

L1

A.1.1

$$\text{Residual } r_i = g_i - (\omega^T x_i + b_i)$$

$\therefore$  In Matrix form

$$R = Y - (W^T X + b)$$

where

$$\begin{aligned}
 R &\rightarrow \mathbb{R}^{n \times 1} \\
 Y &\rightarrow \mathbb{R}^{n \times 1} \\
 W &\rightarrow \mathbb{R}^{d \times 1} \\
 X &\rightarrow \mathbb{R}^{d \times n} \\
 b &\rightarrow \mathbb{R}^{n \times 1}
 \end{aligned}$$

$$\text{Time complexity} = O(n \times d)$$

$\therefore$  Dot product depends on number of non-zero entries in  $X$ .

4.1.2

Update rule for  $b$ 

$$b = \frac{1}{n} [\vec{1} \cdot R] + b$$

where  $\vec{1} \Rightarrow$  a ones-vector of length  $n$ Time complexity =  $O(n)$ 

$\because O(n)$  is required to calculate the vector dot product.

4.1.3

To update  $R$  we need to update the bias term

$$\therefore R = R + b_{\text{old}} - b$$

 $b_{\text{old}} \rightarrow 'b'$  before update $b \rightarrow 'b'$  after updateTime complexity :  $O(n)$ 

$\because$  we need to update  $n$  values with  $O(1)$

Update rule for C<sub>k</sub> using r

$$C_k \leftarrow 2 \sum_{i=1}^n x_{ik} (r_{i-} + w_k x_{ik})$$

∴ In Matrix form,

$$C_k \leftarrow 2 \times (X_k \cdot (r + (X_k \cdot w_k))^T)$$

$X_k \cdot w_k$  takes  $Z_k$  time

but r has 'n' entries, ∴ addition takes  $O(n)$

∴ Time complexity  $O(n)$

$$r \leftarrow r + [w_k \text{ old} - w_k \text{ new}] X_k$$

Time complexity  $O(Z_k)$

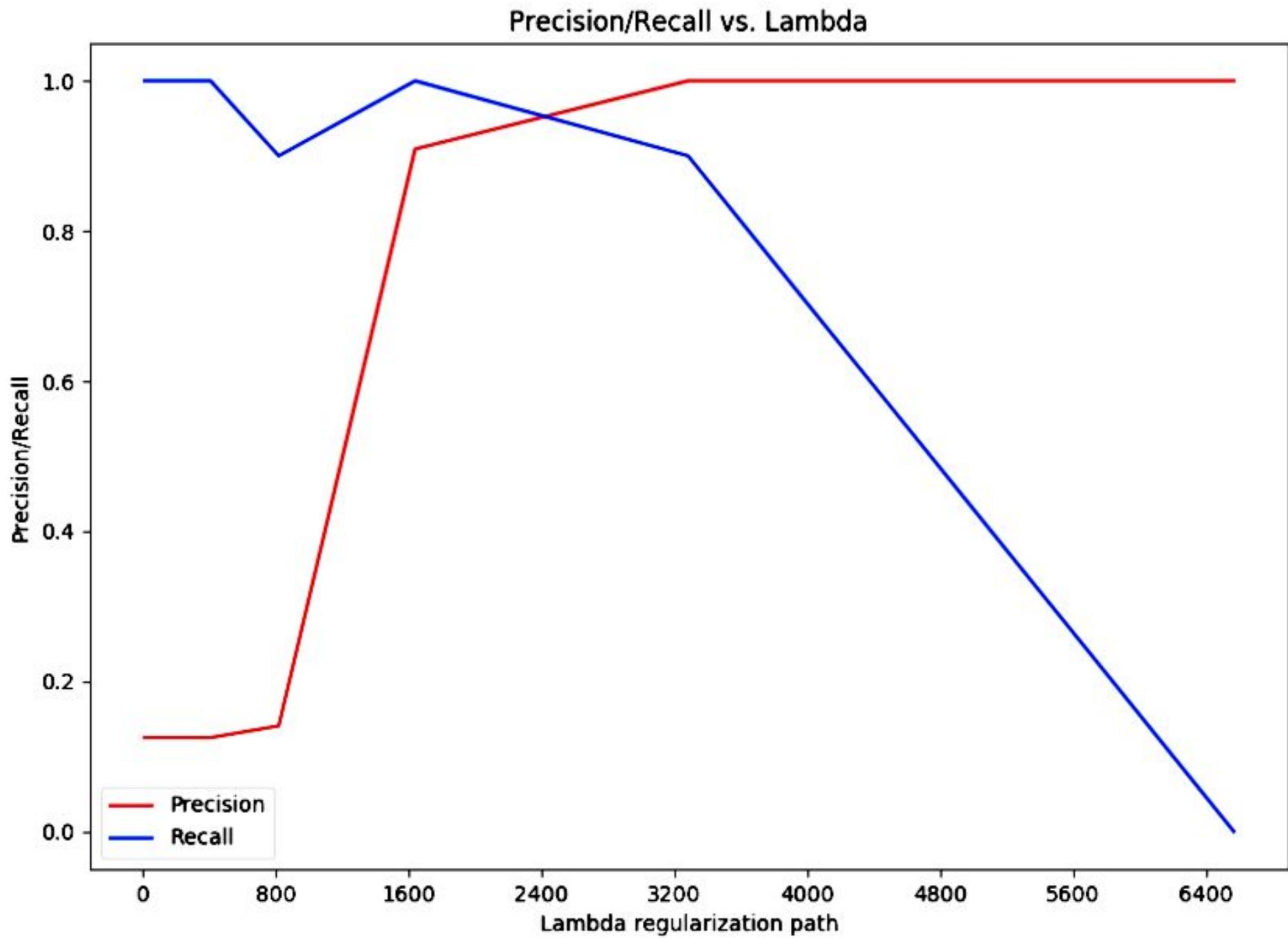
Per iteration complexity

$$O(\max(\max(n, \|X\|_0), nd))$$

Q 4.3.

4.31

The discovery of true non-zeros  
depends on the value of  $\lambda$   
for  $\lambda$  higher precision is high  
while recall is low. As lambda decreases  
precision decreases & recall increases.



4-3-2

We use  $\lambda = 1626.8$ , which gave

precision = 0.909 & recall = 1.0

when used for  $\sigma = 10$ , precision dropped

to 0.1667 while recall stayed at 1

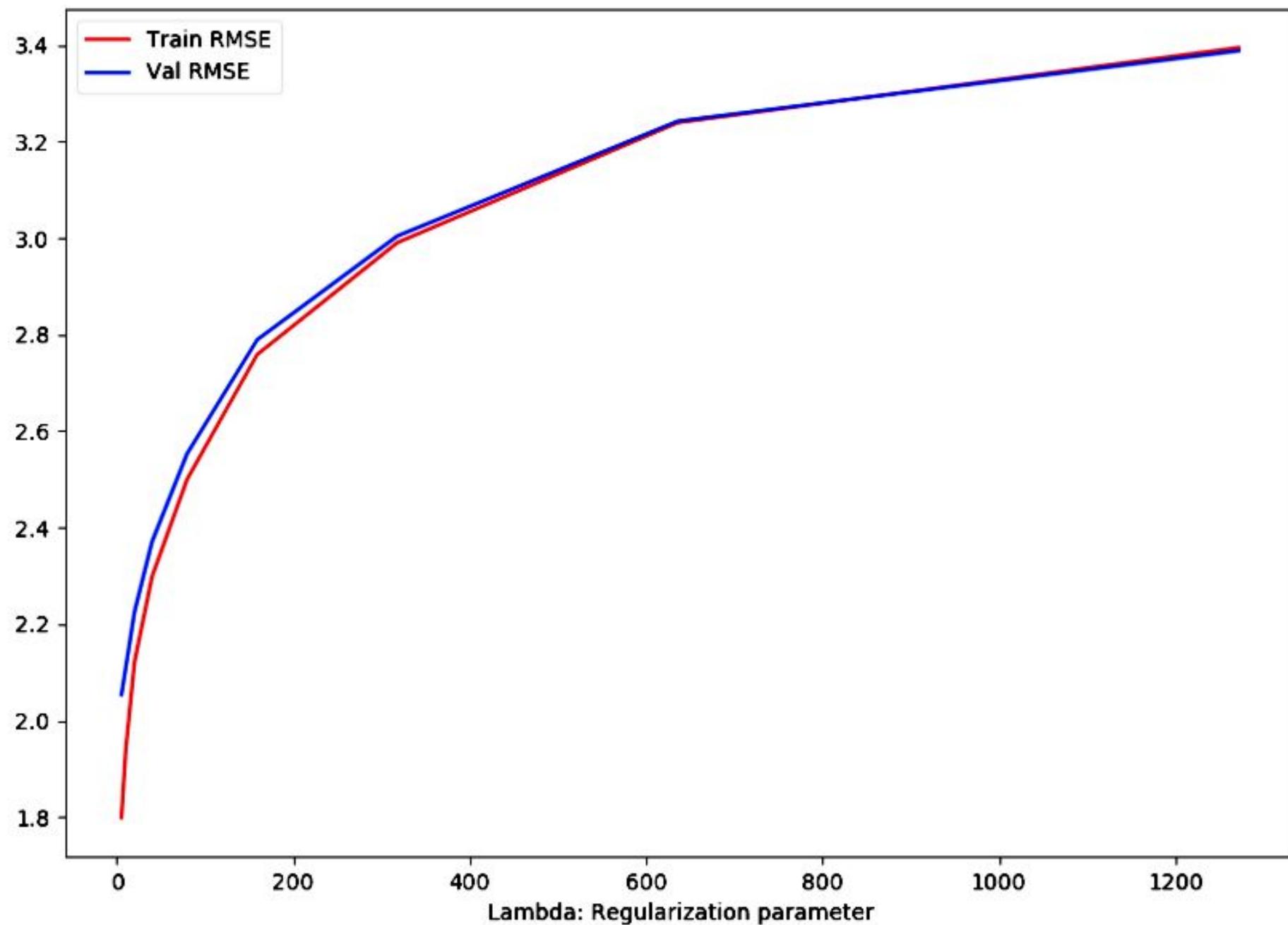
If we increase  $\lambda$ , we will get a higher precision.

6-4.1

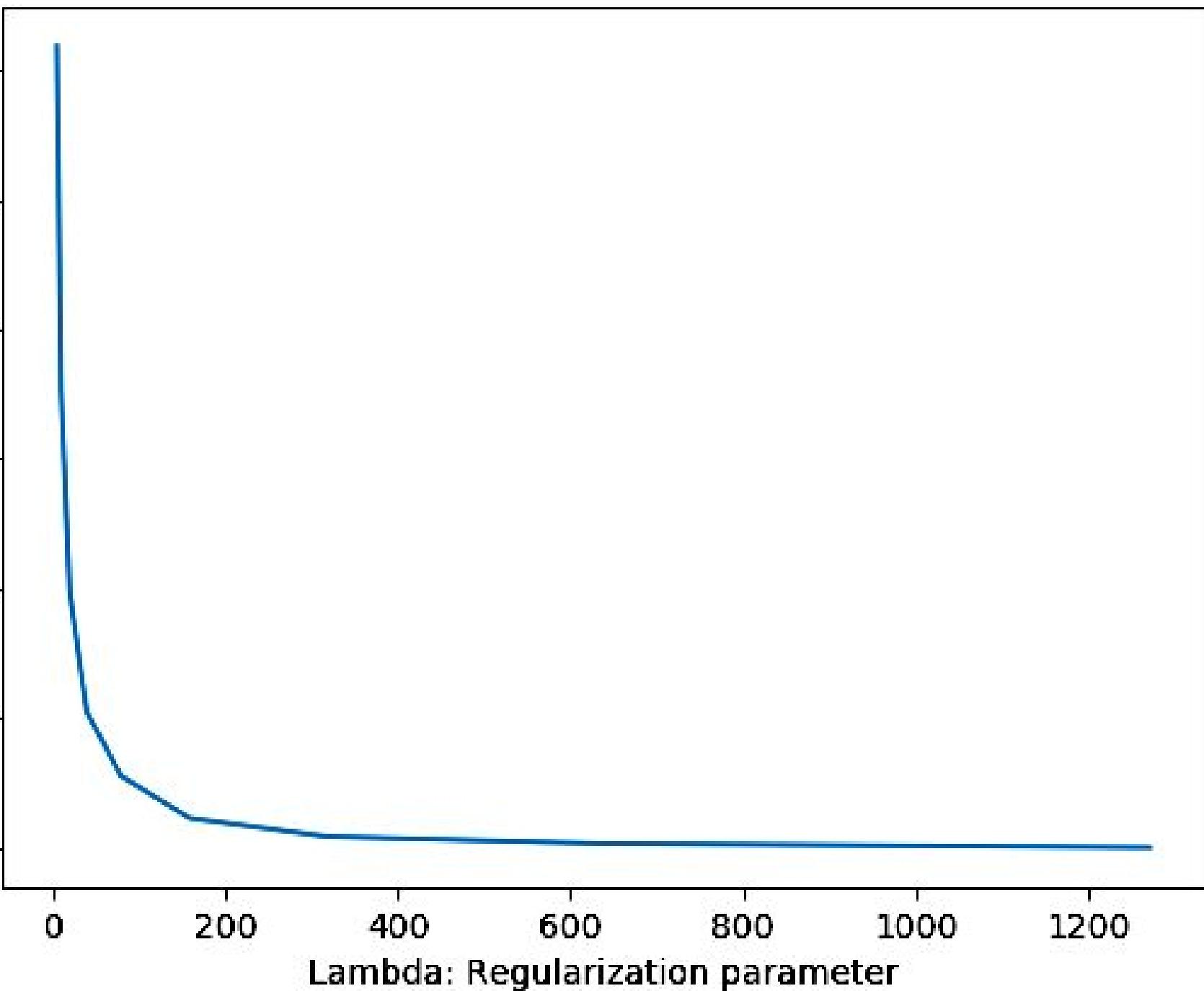
Plots:

(S)

6-4



Number of non zero elements in Weight vector



The  $\lambda$  that achieves the best performance was

$$\lambda = 2.484, \text{ R.M.S.E} = 2.048$$

Top 10 features were

Stars	+	10.93
Spearment	+	10.17
big		9.11
lifesaver		9.09
truly		8.56
sweet black		8.50
nearby		8.27
acidity provides		8.22
lemony		8.02
ageability		7.89

Bottom 10 features

earns	-9.16	sparkler	-6.96
cherry berry	-8.49	brightened	-6.82
black	-8.04	low alcohol	-6.77
soft	-7.49	cuts	-6.71
spices	-7.08	high	-6.54

While some features are intuitive, others might have very high / low weight because they occur with important features.

**submission.csv**

4 hours ago by **FaizaanCharania**

add submission details

**2.05024**

