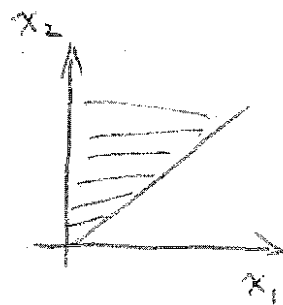


Yuguang Li

1. Probability

$$\begin{aligned}
 1. \quad E(X) &= \int_0^1 \int_0^1 \max(x_1, x_2) P(x_1, x_2) dx_1 dx_2 \\
 &= \int_0^1 \int_0^{x_1} x_2 dx_2 dx_1 + \int_0^1 \int_{x_2}^{x_1} x_1 dx_1 dx_2 \\
 &= \int_0^1 \frac{1}{2} x_2^2 dx_2 + \int_0^1 \frac{1}{2} x_1^2 dx_1 \\
 &= \frac{1}{3}
 \end{aligned}$$



$$\begin{aligned}
 2. \quad \text{Var}(X) &= \int_0^1 \int_0^1 (\max(x_1, x_2) - \frac{1}{3})^2 P(x_1, x_2) dx_1 dx_2 \\
 &= \int_0^1 \int_0^1 (\max(x_1, x_2))^2 - \frac{2}{3} (\max(x_1, x_2)) + \frac{1}{9} dx_1 dx_2 \\
 &= \int_0^1 \int_0^{x_1} x_2^2 dx_2 dx_1 + \int_0^1 \int_{x_2}^{x_1} x_1^2 dx_1 dx_2 - \frac{1}{9} \\
 &= \int_0^1 2 \cdot \frac{1}{3} x_1^3 dx_1 - \frac{1}{9} \\
 &= \frac{1}{6} - \frac{1}{9} = \frac{1}{12}
 \end{aligned}$$

$$\begin{aligned}
 3. \quad \text{Cov}(X, X_1) &= \int_0^1 \int_0^1 (\max(x_1, x_2) - x_1) (\max(x_1, x_2) - x_2) dx_1 dx_2 \\
 &= \int_0^1 \int_0^1 (\max(x_1, x_2))^2 - (x_1 + x_2) \max(x_1, x_2) + x_1 x_2 dx_1 dx_2 \\
 &= \frac{1}{6} - \int_0^1 \int_0^1 (x_1 + x_2) \max(x_1, x_2) + x_1 x_2 dx_1 dx_2 \\
 &= \frac{1}{6} - \left(\int_0^1 \int_0^{x_1} (x_1 + x_2) x_2 + x_1 x_2 dx_2 dx_1 + \int_0^1 \int_{x_2}^{x_1} (x_1 + x_2) \cdot x_1 + x_1 x_2 dx_1 dx_2 \right) \\
 &= \frac{1}{6} - \int_0^1 \frac{8}{3} x_1^3 dx_1 \\
 &= -\frac{1}{2}
 \end{aligned}$$

$$2. \quad 1. \quad P(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$$

$$P(G|\lambda) = \frac{\lambda^{\sum G_i}}{\prod G_i!} e^{-n\lambda}$$

$$\ln(P(G|\lambda)) = \ln\left(\frac{\lambda^{\sum G_i}}{G_i!} e^{-n\lambda}\right)$$

$$= \sum G_i \ln \lambda - \ln(G_i!) - n\lambda$$

$$= \sum (G_i \ln \lambda - \ln G_i!) - n\lambda$$

$$2. \quad \frac{d \ln(P(G|\lambda))}{d\lambda} = \frac{\sum G_i}{\lambda} - n = 0 \quad \lambda = \frac{n}{\sum G_i}$$

$$3. \quad \frac{4+1+\dots+8}{n} - 8 = \frac{25}{n} - 8$$

$$\hat{\lambda}_{MLE} = 8/25 = 0.32$$

7/10

3. 1. ~~$O(n^2d)$~~ Complexity for $X^T X$

$$O(n^2d)$$

Complexity for $(X^T X)^{-1}$

$$O(n^2d + d^3)$$

$$\begin{aligned} 2. \hat{y}_i &= x_i \hat{w} = x_i (X^T X)^{-1} X^T Y \\ &= h_i Y \end{aligned}$$

$$3. SSE_Z = (Y - Z)^T (Y - Z)$$

$$= \sum_{j=1}^n \sum_{i=1, i \neq j}^n (y_i - z_j)^2 + \sum_{j=1}^n (y_i - z_j)^2$$

$$= \sum_{j=1}^n \sum_{i=1, i \neq j}^n (y_i - y_i)^2 + \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2 = LOOCV$$

$$\begin{aligned} 4. \hat{y}_i^{(-i)} &= x_i \hat{w}^{(-i)} = x_i \cdot (X^T X)^{-1} X^T Z \\ &= h_i \cdot Z \end{aligned}$$

$$5. \hat{y}_i - \hat{y}_i^{(-i)} = \hat{y}_i - h_i Z$$

$$= \sum_i h_i y_i - h_i \hat{y}_i^{(-i)} - \sum_{j \neq i} h_{ij} y_j$$

$$= h_{ii} y_i - h_{ii} \hat{y}_i^{(-i)}$$

$$6. O(3d + 3n)$$

4.

1. (a)

ridge

LASSO

↓

↓

~~with~~

~~with~~

~~with~~

~~with~~

~~with~~

~~with~~

(b)

↗

↗

• Less Constraints on weights means more likely overfitting on training dataset.

• Overfitting on training dataset will make the model less accurate. So we will observe higher error on testing dataset

(c)

↗

↗

• Less Constraints on weights will leave weights with high values.

(d)

↗

↗

• Also more non-zeros elements of $\hat{\beta}$

• Because the regularization of two regression methods are similar. (L1 & L2 norm). So we will observe similar effects.

2. (a) Both ↗, because this will ~~be~~ constrain the weights to be small.

The model will be constrained to be simple. So the training error should ~~be~~ increase due to the simplicity of the model

(b) Both could be either increase or decrease. The over simplified model might be not good enough to describe the relation, which will result in the increase of testing error. On the other hand, a large λ might prevent overfitting. So the testing error might decrease.

(c) Both decrease. & A larger λ will control $\hat{\beta}$ to be small

(d) Both decrease same as (c)

5. 5.1 1. $L(f) = E_{X,Y} ((Y - f(x))^2)$

$$= E_{X,Y} ((Y - E(Y|X) + E(Y|X) - f(x))^2)$$

where $E(Y|X)$ is the conditional expectation of Y given X

$$L(f) = E_{X,Y} (Y^2 - 2E(Y|X)Y + E(Y|X)^2) + E_{X,Y} (2E(Y|X)Y + 2E(Y|X)f(x) - 2E(Y|X)^2 - 2f(x)Y) + E_{X,Y} (E(Y|X)^2 - 2E(Y|X)f(x) + f(x)^2)$$

$$= E_{X,Y} ((Y - E(Y|X))^2) + E_{X,Y} ((E(Y|X) - f(x))^2)$$

$$+ 2E_{X,Y} ((E(Y|X) - f(x))(Y - E(Y|X)))$$

Because $E(Y|X) - f(x)$ and $Y - E(Y|X)$ are independent random variables.

$$E_{X,Y} ((E(Y|X) - f(x))(Y - E(Y|X))) = 0$$

$$\text{So } L(f) = E_{X,Y} ((Y - E(Y|X))^2) + E_X ((E(Y|X) - f(x))^2)$$

2. $L(f) = E_X ((E(Y|X) - f(x))^2) + E_{X,Y} ((Y - E(Y|X))^2)$

In order to minimize $L(f)$ we can have $L(f_{\text{Bayes}}) = E(Y|X)$

while we can't do anything to the second term

5.2 $E_L L(\hat{f}) = E_L E_{X,Y} ((Y - E(Y|X))^2) + E_L E_X ((E(Y|X) - \hat{f}(x))^2)$

$$= E_L E_{X,Y} ((Y - E(Y|X))^2) + E_L E_X ((E(Y|X)^2 + \hat{f}(x)^2 - 2E(Y|X)\hat{f}(x)))$$

$$+ (E \hat{f}(x)^2 - 2 \bar{f}(x) \hat{f}(x) + \bar{f}(x)^2) + 2(\bar{f}(x) \hat{f}(x) + E(Y|X) \hat{f}(x) - \bar{f}(x)^2 - 2E(Y|X) \hat{f}(x))$$

$$= E_{X,Y} ((Y - E(Y|X))^2) + E_X ((E(Y|X) - \bar{f}(x))^2)$$

$$+ E_L E_X ((\hat{f}(x) - \bar{f}(x))^2) + 2E_L E_X ((\bar{f}(x) - E(Y|X))(\hat{f}(x) - \bar{f}(x)))$$

Because $\hat{f}(x) - \bar{f}(x)$ and $\bar{f}(x) - E(Y|x)$ are two independent random variable.

$$E_z E_x ((\bar{f}(x) - E(Y|x))(\hat{f}(x) - \bar{f}(x))) = 0.$$

$$\text{So } E_z L(\hat{f}) = E_{x,Y} (LY - E(Y|x))^2 + E_x ((E(Y|x) - \bar{f}(x))^2) + E_z E_x ((\hat{f}(x) - \bar{f}(x))^2)$$

2.	$E_{x,Y} ((Y - E(Y x))^2)$	noise
	$E_x ((E(Y x) - \bar{f}(x))^2)$	variance
	$E_z E_x ((\hat{f}(x) - \bar{f}(x))^2)$	bias

3. No f^* is the best estimator of $\bar{f}(x)$ it might be biased

4. With a large sample size N

- $\bar{f}(x)$ is the theoretical Y value at point x
- $E(Y|x)$ is the observed mean of Y given x
- $\hat{f}(x)$ is the estimated Y value using our model at point x . Our model is biased for most of the case (or is just optimized, but not the real one)

674

1. threshold 0.5

0/1 loss training: 0.0755

square loss training: 0.0505

2. 0/1 loss test: 0.0764

square loss test: 0.0515

3. Because the relation between features & labels might not be linear. Plus there are so many points lying around 0.5, which does not represent probability.

4. $\hat{y}_i = w^T x_i$ ~~is~~ ← Bayes optimal predictor.

$$\text{hard predictor } \hat{y}_i = \begin{cases} 0 & w^T x_i \leq 0.5 \\ 1 & w^T x_i > 0.5 \end{cases}$$

$$y_i = w^T x_i + \epsilon_i \quad \epsilon_i \text{ is a Gaussian random Variable with 0 mean.}$$

7.

7.1

1. Initialize $w_0^0 = 0$

while not converged do

$$w_0^{(t+1)} \leftarrow \sum_{i=1}^N (y_i - \hat{y}_i^{(t)}) / N + w_0^{(t)}$$

for $k \in \{1, 2, \dots, d\}$ do

same as Algorithm 1.

end

end

2. $O(\|X\|_0)$

3. $O(\|X\|_0 + 2)$, $O(N + 2)$

4. $O(z_1)$

5. $\hat{y}_i^{(t+1)} = x_i w^{(t)} + w_0^{(t+1)}$

$$= \hat{y}_i^{(t)} - w_0^{(t)} + w_0^{(t+1)}$$

Complexity $O(N)$ or $O(z_1)$

6. $C_k^{(t+1)} = 2 \sum_{i=1}^N x_{ik} (y_i - \hat{y}_i^{(t)} + w_k^{(t)} x_{ik})$

$$a_k^{(t+1)} = 2 \sum_{i=1}^N x_{ik}^2$$

$$w_k^{(t+1)} = \begin{cases} (C_k^{(t+1)} + \bar{r}) / a_k^{(t+1)} & C_k < -\bar{r} \\ 0 & C_k \in [-\bar{r}, \bar{r}] \\ (C_k^{(t+1)} - \bar{r}) / a_k^{(t+1)} & C_k > +\bar{r} \end{cases}$$

$$\hat{y}_i^{(t+1)} = \hat{y}_i^{(t)} + x_{ik} (w_k^{(t+1)} - w_k^{(t)})$$

Complexity $O(N)$

7

$O(Nd)$

7 Programming Question 2: Lasso

7.3 Try out your work on synthetic data

Question 1

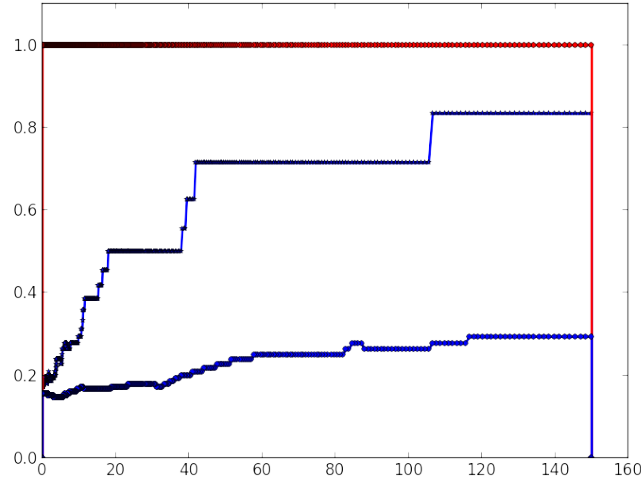


Figure 1: Recall(red) and precision(blue) plots of σ 1(star) and 10(circle) values against λ .

A large λ helps to more accurately discover non-zeros, because more false positive are discovered. As the λ goes lower, the precision goes down. But recalls are always 100%.

Question 2

Plots are found above with circle markers.

When I changed σ from 1.0 to 10.0, the precision become lower, because more false positive are discovered. This is caused by the noise introduced by the larger σ . But the recalls stay at 100%.

In this case, I'd like to start with a larger λ and decreased from there.

7.4 Become a data scientist at Yelp

Question 1

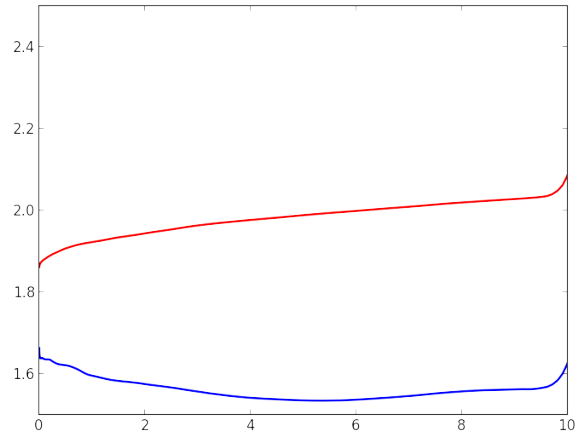


Figure 2: RMSE values against λ from training (red line) and testing (blue line) datasets.

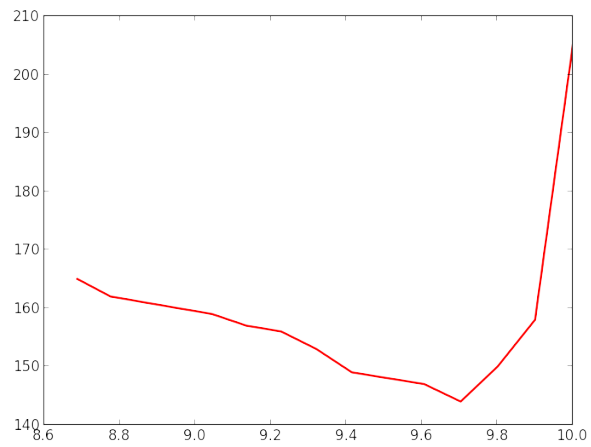


Figure 3: Number of nonzero weights against λ from training (red line) and testing (blue line) datasets.

Question 2

Using the best λ value 4.9983 from the validation performance, I got the RMSE value 1.1829717559 from the testing dataset.

Question 3

The best 10 features I discovered are :

```
sqrt(ReviewNumCharacters*UserCoolVotes) : [ 67.64008046]
sqrt(UserCoolVotes*BusinessNumStars) : [ 33.55206007]
sqrt(ReviewNumCharacters*UserFunnyVotes) : [ 14.92750312]
sqrt(UserFunnyVotes*InPhoenix) : [ 61.34628919]
BusinessNumReviews*InGlendale : [ 17.02279332]
ReviewInFall*InGlendale : [ 16.25769229]
UserUsefulVotes*InScottsdale : [ 29.18772702]
log(ReviewNumCharacters*UserUsefulVotes) : [ 35.83890522]
sq(UserUsefulVotes*IsRestaurant) : [ 42.39206274]
sqrt(ReviewNumWords*UserUsefulVotes) : [ 56.53814474]
```

Question 4

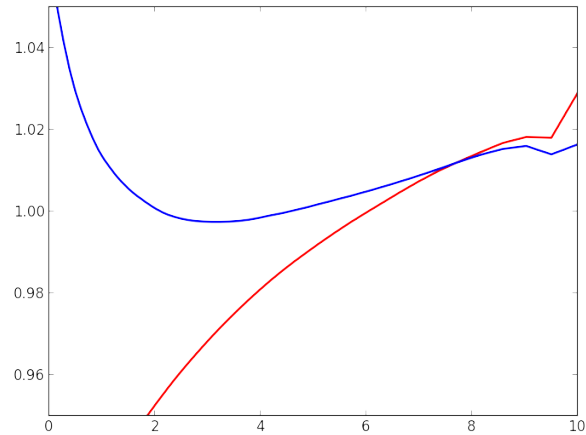


Figure 4: RMSE values against λ from training (red line) and testing (blue line) datasets.

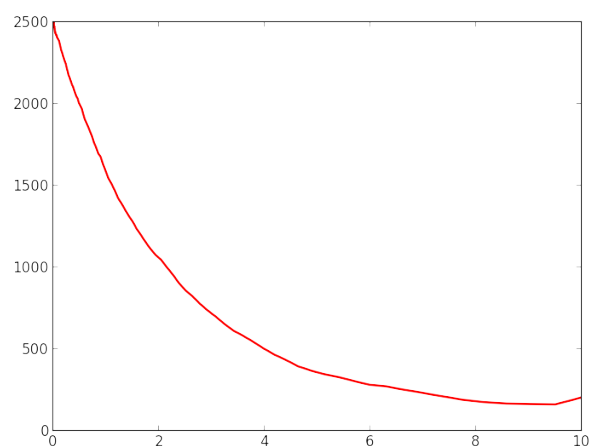


Figure 5: Number of nonzero weights against λ from training (red line) and testing (blue line) datasets.

Using the best λ value 3.0 from the validation performance, I got the RMSE value 0.5608 from the testing dataset.

The best 10 features I discovered are :

great : [21.10928191]

best : [18.00001071]

amazing : [14.74667601]

love : [11.76668926]

delicious : [11.09071936]

awesome : [12.56298077]

perfect : [9.47380946]

excellent : [8.90874206]

wonderful : [8.44758452]

friendly : [8.64004513]