Statistical Learning

Homework 1

.1  Prediction Problem

Given

$$K_1 = 96$$
$$n = 20$$

(1) So a High-Dimensional data Set-up... It is obvious we can't use OLS, we have to Include some penalty ...

(2) We also checked Multi-collinearity and turns out it's not a genuine issue (very low correlation) ——

So, the two most obvious choices are

.1 Ridge

.2 Lasso

furthermore, we know we can use all 96 variables ....

So, we used both Ridge and lasso, on average Ridge was scoring better (using Cross validations)

finally, Even though Ridge was performing slightly better than lasso, we decided to predict using lasso as it was a random choice -- (very low error difference)...

The final equation then can be rewritten as

$$RSS + \lambda \sum_{J=1}^{P} |\beta_J| \quad \text{---(A)}$$

Residual sum of square

The code can explain further ...

Problem: $\hat{\beta}(z, m) = \underset{\beta \in R^n}{\text{argmin}} \sum_{i=1}^{n} \frac{1}{M} \sum_{m=1}^{M} -\ln p_B(y_i \mid \tilde{x}_i^{(m)})$

$= \underset{\beta}{\text{argmin}} \sum_{i=1}^{n} \frac{1}{M} \sum_{m=1}^{M} (y_i - x_i^T \beta)^2$

1) Algorithm: for each $z$

$z \rightarrow$ A vector of probabilities

$M \rightarrow$ A vector of Integers

delta-run ($z$, M, features, response):

for each $z$:

for each M:

Create Bernoulli-Random-Vector

Create New features (add Noise)

Perform OLS

Calculate error

Save error

$\rightarrow$ Finally Select $z$, M. for minimum error.

2) We applied Ridge Regression & Lasso Regression to the Supernova Dataset & finally compared errors among all three models.

③ Supernova Dataset was split into Train-1, Train 2, Train 3, Train 4, and Test-1, Test-2, Test-3, Test 4 and Dropout technique was performed on all these splitted Datasets along with Ridge & Lasso.

4)

| | Average Error | Minimum Error | Maximum Error |
|---|---|---|---|
| Dropout Technique: | 167 | 0.1 | 5 |
| Ridge Regression : | 2 | 3.6 | 4 |
| Lasso Regression : | 4 | 1.2 | 4 |

Average Error: Average Mean square Error over all Test Datasets.

Minimum Error: Minimum Mean square Error among all Test Datasets

Maximum Error: Maximum Mean square Error among all Test Datasets.

**Remarks :** Generating The Noisy Model :

Bernoulli Random Vector is created.

$$\xi \sim Bernoulli\ (1, \beta, m) \qquad \{ m \text{ is no of rows of our dataset} \}$$

$$\xi_e = \frac{Bernoulli\ (1, \beta, m)}{(1-\beta)}$$

So $\xi$ is a vector of $\left(\beta s, \frac{1}{(1-\beta)}\right)$ values.

For every column in the original dataset, we multiply the

→ $\xi$ vector elementwise with every column.

# PART-II

Simulation Study:

1) a) Generating feature set:

We generate 100 features from uniform distribution with observations between 1 & 100.

$$X_i \sim Unif(1, 100, 100)$$

b) Generating Response

We generate $Y_i$ vector from a normal distribution having mean as the overall mean of $X_i$'s scaled by a factor of 100, and variance as 10.

$$Y_i \sim N(mean(X_i)*100, 10, 75)$$

2) After Generating the feature set & $Y_i$, we repeat the same procedure as done in part-I on this dataset. Dropout techniques Ridge & Lasso were used to fit the model.

3) The Generated Dataset was split into 5 samples, Train-1, Train-2, Train-3, Train-4, Train 5 & <u>Test-1</u>, <u>Test-2</u> <u>Test-3</u>, <u>Test-4</u>, <u>Test-5</u> and the average, maximum & minimum errors were computed over all these datasets.

| | Average Error | Minimum Error | Max Error |
|---|---|---|---|
| Dropout techniques: | 34.1 | 2.2e-26 | 100 |
| Ridge Regression: | 154 | 56 | 178 |
| LASSO Regression: | 111 | 100 | 214 |