

Homework 1

Guozhen Li

STA 208 - Statistical Machine Learning

April 18, 2018

1 Exercise 1

1.1 Predictor minimizing true risk

The true risk is

$$R(g) = \mathbb{E}[\ell(Y, g(X))]$$

By conditioning on X , we can write $R(g)$ as

$$R(g) = \mathbb{E}_X \ell(Y, g(X)) \cdot \Pr(Y|X)$$

Minimize $R(g)$ pointwise:

$$\hat{g}(x) = \arg \min_{h \in \{0,1\}} \ell(Y, h) \Pr(Y|X = x)$$

With $\ell()$ being the Hamming loss function this simplifies to:

$$\hat{g}(x) = \arg \max_{h \in \{0,1\}} \Pr(Y = h|X = x)$$

Because Y only takes value of 0 or 1, $\hat{g}(x)$ in practical can be written as

$$\hat{g}(x) = 1\{\Pr(Y = 1|X = x) > \frac{1}{2}\}$$

Reference: ESL pg. 20

1.2 True risk of Bayes classifier

The true risk is

$$\begin{aligned} R(g) &= \mathbb{E}[\ell(Y, g(X))] \\ &= \mathbb{E}_X \ell(Y, g(X)) \cdot \Pr(Y|X) \\ &= \ell(Y, g(X)) \cdot \Pr(Y|X) \end{aligned}$$

1.3 Classify with a single x_j

To fit a model like $h(x) = 1\{x_j > 0\}$, we can try all possible j 's and find the one that gives smallest empirical risk. Pseudo code for this algorithm as follows:

```
best_R = Inf
best_j = None
for j in [1, 2, ..., p]:
    pred = [1 if x[i, j] > 0 else 0 for i in 1:N]
    loss = xor(pred, y)
    risk = sum(loss)/N
    if risk < best_R:
        best_R = risk
        best_j = j
return best_j
```

This algorithm returns the best j .

When making a prediction based on a x_{new} , simply check the value of its j th component, and if $x_{new,j} > 0$, predict 1, otherwise predict 0.

1.4 Number of samples needed

$$\mathbb{P}\{R(\hat{g}) < R(h) + 0.1\} \geq 0.95$$

would be the same as

$$\mathbb{P}\{R(\hat{g}) > R(h) + 0.1\} \leq 0.05$$

meanwhile,

$$\begin{aligned} \mathbb{P}\{R(\hat{g}) > R(h) + 0.1\} &= \mathbb{P}\{R(\hat{g}) - R(h) > 0.1\} \\ &= \mathbb{P}\left\{\bigcup_{i=1}^n [R_n(h_i) - R(h) > 0.1]\right\} \\ &\leq \sum_{i=1}^n \mathbb{P}\{R_n(h_i) - R(h) > 0.1\} \\ &\leq 2n \exp(-2n \times 0.1^2) \end{aligned}$$

To make that ≤ 0.05 , we can make

$$\begin{aligned} 2n \exp(-2n \times 0.1^2) &\leq 0.05 \\ 2n \exp(-0.02n) &\leq 0.05 \end{aligned}$$

Solve this to get $n \geq 495$. At least 495 samples needed.

2 Exercise 2

2.1 $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$

In linear regression, we already know that

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})_{-1}\mathbf{X}^T\mathbf{y}$$

Here if we make $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})_{-1}\mathbf{X}^T$, then $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$

In the case of kNN model, we can make a matrix \mathbf{H} such that:

$$H_{i,j} = \begin{cases} \frac{1}{k}, & \text{if } x_j \in N_k(x_i) \\ 0, & \text{otherwise} \end{cases} \quad \text{for } i, j \in 1, 2, \dots, N$$

Here $N_k(x_i)$ is the k-nearest neighborhood of observation x_i (x_i itself included). \mathbf{H} is a $N \times N$ matrix, and each row of it only has k elements with value $\frac{1}{k}$, and all other elements are 0. With such construction, $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$.

2.2 kNN leave-one-out

To leave-one-out, we need to reconstruct the \mathbf{H} matrix, taking the sample left out into account. Assume the i_0 th observation needs to be left out, make \mathbf{H} such that

$$H_{i,j} = \begin{cases} \frac{1}{k}, & \text{if } x_j \in N_k(x_i) \text{ and } i \neq i_0 \\ 0, & \text{otherwise} \end{cases} \quad \text{for } i, j \in 1, 2, \dots, N$$

Here $N_k(x_i)$ should be the set of x_i 's k nearest neighbors among all x 's but excluding x_{i_0} .

Denote the \mathbf{H} matrix for \mathbf{X} with the i th observation left out as \mathbf{H}_i . The square error with the i th observation left out would be

$$\begin{aligned} e_i &= (\hat{\mathbf{y}} - \mathbf{y})^T(\hat{\mathbf{y}} - \mathbf{y}) \\ &= (\mathbf{H}_i\mathbf{y} - \mathbf{y})^T(\mathbf{H}_i\mathbf{y} - \mathbf{y}) \\ &= \mathbf{y}^T(\mathbf{H}_i - \mathbf{I})^T(\mathbf{H}_i - \mathbf{I})\mathbf{y} \end{aligned}$$

Leave-one-out cross validated square error would be

$$\begin{aligned} SE &= \frac{1}{N} \sum_{i=1}^N e_i \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{y}^T(\mathbf{H}_i - \mathbf{I})^T(\mathbf{H}_i - \mathbf{I})\mathbf{y} \end{aligned}$$

2.3 SVD in regression

Plug $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ into linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, we get

$$\mathbf{Y} = \mathbf{U}\mathbf{D}\mathbf{V}^T\boldsymbol{\beta} + \mathbf{e}$$

Introduce a new vector $b = \mathbf{D}\mathbf{V}^T\beta$, the linear model becomes

$$\mathbf{Y} = \mathbf{U}b + \mathbf{e}$$

This takes the form of normal OLS linear regression model, and b can be estimated by

$$\hat{b} = (\mathbf{U}^T\mathbf{U})^{-1}\mathbf{U}^T\mathbf{y}$$

Because \mathbf{U} is an orthogonal matrix, $\mathbf{U}^T\mathbf{U} = \mathbf{I}$, \hat{b} simplifies to

$$\hat{b} = \mathbf{U}^T\mathbf{y}$$

And $\hat{\beta}$ can be estimated by

$$\hat{\beta} = (\mathbf{V}^T)^{-1}\mathbf{D}^{-1}\hat{b}$$

Again, because \mathbf{V} is also an orthogonal matrix, $(\mathbf{V}^T)^{-1} = \mathbf{V}$, thus

$$\hat{\beta} = \mathbf{V}\mathbf{D}^{-1}\hat{b}$$

Define $\mathbf{A} = \mathbf{V}\mathbf{D}^{-1}$, we get an estimator for β and

$$\hat{\beta} = \mathbf{A}\hat{b}$$

To fit the model on a dataset, first perform SVD on the design matrix \mathbf{X} , then calculate \mathbf{D}^{-1} by replacing every diagonal element with its reciprocal (because \mathbf{D} is diagonal, all other places should be 0). Then multiply \mathbf{U} and \mathbf{y} to get \hat{b} . Finally calculate $\hat{\beta}$ via $\hat{\beta} = \mathbf{A}\hat{b}$.

Reference: Mandel, John. "Use of the singular value decomposition in regression analysis." The American Statistician 36.1 (1982): 15-24.

2.4 Change rank

Take the first r elements of \hat{b} to be \hat{b}_r , take the first r columns of \mathbf{A} to be \mathbf{A}_r , then $\hat{\beta}_r = \mathbf{A}_r\hat{b}_r$.

If \mathbf{A} and \hat{b} already exist, this computation should take $O(Nrp)$ time for the matrix multiplication $\hat{\beta}_r = \mathbf{A}_r\hat{b}_r$. Meanwhile since $r \leq p \leq N$, this time complexity is $O(N^3)$.