

Homework 1

Dai Yuehao (1800010660@pku.edu.cn)

October 12, 2020

1 No free lunch

For a loss function $l(f, h)$, we hope to have

$$\begin{aligned}\sum_f E_{ote}(\mathcal{L}_a|X, f) &= \sum_f \sum_h \sum_{x \in \mathcal{X}-X} P(x) l(f(x), h(x)) P(h|X, \mathcal{L}_a) \\ &= \sum_{x \in \mathcal{X}-X} P(x) \sum_h P(h|X, \mathcal{L}_a) \sum_f l(f(x), h(x)) \\ &= \sum_f l(f(x), h(x)) \sum_{x \in \mathcal{X}-X} P(x) \sum_h P(h|X, \mathcal{L}_a) \\ &= \sum_f l(f(x), h(x)) \sum_{x \in \mathcal{X}-X} P(x),\end{aligned}$$

thus l need to satisfy that $\sum_f l(f(x), h(x))$ is the same for any given h . An assumption is that l has the form

$$l(f, h) = g(|f - h|), \quad \sup |g| < \infty,$$

under this circumstance

$$\sum_f l(f(x), h(x)) = \sum_f [g(1)\mathbb{I}(h(x) \neq f(x)) + g(0)\mathbb{I}(f(x) = h(x))] = 2^{|\mathcal{X}|-1}[g(0) + g(1)].$$

Otherwise, we can choose

$$h_1 = \arg \max_f \sum l(f(x), h(x)), \quad h_2 = \arg \min_f \sum l(f(x), h(x)),$$

let algorithm a yield h_1 with probability 1, algorithm b yield h_2 with probability 1, then we have

$$E_{ote}(\mathcal{L}_a|X, f) > E_{ote}(\mathcal{L}_b|X, f).$$

2 ESL 3.4

Let Z_i denote the vector obtained at the i^{th} step of the Gram-Schmidt procedure, then we have

$$Z_i = X_i - \sum_{k=0}^i \frac{X_i^T Z_k}{Z_k^T Z_k} Z_k,$$

thus Z_p is the only vector that contains X_p . Now from $Z_0 = X_0$, we regress y on Z_0, Z_1, \dots, Z_{i-1} and obtain the residual r_i , then regress r_i on Z_i and obtain the coefficient $\hat{\beta}'_i$. Eventually we have $\hat{\beta}'_p = \hat{\beta}_p$,

here $\hat{\beta}_p$ is the coefficient of X_p . Now we reverse the procedure, we remove X_p from the combination of $\{Z_i\}$ and the rest of the combination of $\{Z_i\}$ must be the regression of y on $\{Z_0, \dots, Z_{p-1}\}$ (since they are orthogonal), hence the coefficient of Z_{p-1} is exact the coefficient of X_{p-1} . We continue this procedure until we get all the coefficients.

On the other hand, with the QR decomposition of X we have

$$\hat{\beta}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y = \mathbf{R}^{-1} \mathbf{Q} y,$$

thus

$$\mathbf{R} \hat{\beta} = \mathbf{Q}^T y, \quad (2.0.1)$$

equation (2.0.1) is easy to solve because \mathbf{R} is upper triangular so that we can solve the equation by back-substitution, just the same as we show above.

3 ESL 3.6

In terms of Bayes formula we have

$$p(\beta|y) = \frac{p(y|\beta)p(\beta)}{p(y)} \propto p(y|\beta)p(\beta)$$

which has normal distribution, hence the mean of the posterior distribution is the maximal of the density function. Now we take the logarithm and we need to

$$\min_{\beta} \frac{(y - \mathbf{X}\beta)^T (y - \mathbf{X}\beta)}{2\sigma^2} + \frac{\beta^T \beta}{2\tau},$$

by taking the derivative of β and let it be zero we have

$$-\frac{\mathbf{X}^T \mathbf{X} \hat{\beta}}{\sigma^2} - \frac{\hat{\beta}}{\tau} + \frac{\mathbf{X}^T y}{\sigma^2} = 0,$$

hence we obtain

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T y, \quad \lambda = \frac{\sigma^2}{\tau}.$$

Thereby we conclude the proof, the relationship between the regularization parameter λ in the ridge formula, and the variances τ and σ^2 is the equation $\lambda = \sigma^2/\tau$.

4 ESL 3.8

We consider

$$\mathbf{X} = \mathbf{Q}\mathbf{R}, \quad \tilde{\mathbf{X}} = \mathbf{U}\Sigma\mathbf{V}^T$$

where \mathbf{R} is upper triangular, hence we have

$$\mathbf{Q}_1 = \left(\frac{1}{\sqrt{N}}, \dots, \frac{1}{\sqrt{N}} \right)^T,$$

hence \mathbf{Q}_2 forms an orthogonal complement of $\mathbf{1}$ in the column space of \mathbf{X} . On the other hand we have

$$\mathbf{1}^T \tilde{X}_i = \mathbf{1}^T \left(X_{i+1} - \frac{1}{N} \mathbf{1}^T X_{i+1} \mathbf{1} \right) = 0,$$

since $\text{rank}(\tilde{\mathbf{X}}) = p$, we know that $\tilde{\mathbf{X}}$ forms an orthogonal complement of $\mathbf{1}$ in the column space of \mathbf{X} . In terms of the nature of SVD decomposition we know that the column space of \mathbf{U} is the same as $\tilde{\mathbf{X}}$'s, thereby we conclude the proof.

If $\mathbf{Q}_2 = \mathbf{U}$, we know that \mathbf{Q}_2 is the Gram-Schmidt orthogonal of $\tilde{\mathbf{X}}$ with regularization, hence we have $\tilde{\mathbf{X}} = \mathbf{Q}_2 \mathbf{R}_2$ and $\mathbf{R}_2 = \Sigma \mathbf{V}^T$, thus \mathbf{R}_2 is both upper triangular and orthogonal, which forces \mathbf{R}_2 to be diagonal. Thereby we claim that \mathbf{Q}_2 and \mathbf{U} are the same when $\tilde{\mathbf{X}}$ is orthogonal.

5 ESL 3.9

At this time we know that r is orthogonal with the column space spanned by \mathbf{X}_1 hence the variable that will reduce the RSS the most is the one whose vertical component of the column space spanned by \mathbf{X}_1 , denoted by u , is the most parallel with r , hence we need to maximize

$$\frac{|u_j^T r|}{\|u_j\|_2}, \quad q < j \leq p.$$

Now suppose we have $\mathbf{X}_1 = \mathbf{Q}\mathbf{R}$, let z_i denote the i^{th} column vector of \mathbf{Q} and we denote

$$u_j = x_j - \sum_{i=1}^q (z_i^T x_j) z_i, \quad q < j \leq p,$$

then

$$v_j = \frac{u_j}{\|u_j\|},$$

then the reduce of RSS is $\|(y^T v_j) v_j\|^2 = (y^T v_j)^2$, hence we can solve the equivalent problem

$$j^* = \arg \max_j (y^T v_j)^2,$$

then x_{j^*} is the variable we want.

6 ESL 3.11

We derive the results from minimizing

$$\text{RSS}(\mathbf{B}) = \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B})^T \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})] = \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B}) \Sigma^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})^T].$$

This is a convex optimization problem, in terms of the one order condition and since Σ^{-1} is positive definite we can write $\Sigma^{-1} = \mathbf{S}^T \mathbf{S}$, now we take the derivative of $\mathbf{B}\mathbf{S}^T$ and let

$$(\mathbf{X}^T \mathbf{X}) \hat{\mathbf{B}} \mathbf{S}^T - \mathbf{X}^T \mathbf{S}^T \mathbf{S} \mathbf{Y} \mathbf{S}^T = 0$$

yielding that

$$\hat{\mathbf{B}} \mathbf{S}^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mathbf{S}^T,$$

since \mathbf{S}^T is non-singular we have

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

If the covariance matrices are different, we need to minimize

$$\text{RSS}(\mathbf{B}) = \sum_{i=1}^N (y_i - x_i \mathbf{B})^T \Sigma_i^{-1} (y_i - x_i \mathbf{B}),$$

however this time we can no longer write a closed-form solution by taking the derivative of $\mathbf{X}\mathbf{S}^T$ as before since Σ_i^{-1} varies. In spite of this, we can still take the derivative of \mathbf{B} and obtain

$$\sum_{i=1}^N (x_i^T \mathbf{S}_i^T \mathbf{S}_i x_i) \mathbf{B} - \sum_{i=1}^N \mathbf{X}^T \mathbf{S}_i^T \mathbf{S}_i y_i = 0$$

and use some numerical methods to solve this equation.

7 ESL 3.12

Then we have

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_p \end{pmatrix}, \quad \tilde{y} = \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix}$$

thus the least square solution is

$$\hat{\beta} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \tilde{y} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T y.$$

8 ESL 3.29

For a single variable X and response y , the result is

$$a = (X^T X + \lambda)^{-1} X^T y = \frac{\sum_{i=1}^N x_i y_i}{\lambda + \sum_{i=1}^N x_i^2}.$$

Now we include a copy of X and

$$a_{(2)} = (\mathbf{X}_{(2)}^T \mathbf{X}_{(2)} + \lambda \mathbf{I}_2)^{-1} \mathbf{X}_{(2)}^T y, \quad \mathbf{X}_{(2)} = \begin{pmatrix} X & X \end{pmatrix}$$

With some calculation we have

$$a_{(2)} = \begin{pmatrix} \frac{\sum_{i=1}^N x_i y_i}{\lambda + 2 \sum_{i=1}^N x_i^2} \\ \frac{\sum_{i=1}^N x_i y_i}{\lambda + 2 \sum_{i=1}^N x_i^2} \end{pmatrix}$$

hence both coefficients are identical.

In general if there is m copies of the variable X , we also have

$$a_{(m)} = (\mathbf{X}_{(m)}^T \mathbf{X}_{(m)} + \lambda \mathbf{I}_2)^{-1} \mathbf{X}_{(m)}^T y,$$

and it is very easy to solve the equation and obtain that each element of $a_{(m)}$ is

$$\frac{\sum_{i=1}^N x_i y_i}{\lambda + m \sum_{i=1}^N x_i^2}.$$