

KERNEL METHODS - PRACTICAL SESSION № 1

AMMI 2020

18/05/2020

Linear Algebra Recap

Exercise 1

Let $A \in \mathbb{R}^{m \times n}$ be a real matrix. For each function f , specify its codomain and compute ∇f

(a) $f : \begin{cases} \mathbb{R} \rightarrow ? \mathbb{R}^{\text{sym}} \\ x \rightarrow xA \end{cases}$

(b) $f : \begin{cases} \mathbb{R}^n \rightarrow ? \mathbb{R}^m \\ x \rightarrow Ax \end{cases}$

(c) $f : \begin{cases} \mathbb{R}^{m \times n} \rightarrow ? \mathbb{R} \\ X \rightarrow \text{Tr}(A^T X) \end{cases}$

Exercise 2

Let $X \in \mathbb{R}^{n \times p}$, $\lambda > 0$.

Prove that $M = X^T X + \lambda I_p$ is invertible.

Hint: Prove that the eigenvalues of M are larger than λ

Linear Regression

Using notations from the slides:

- $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ the vector of outcomes
- $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times p}$ the matrix (n rows=samples, p columns=features)
- $\beta \in \mathbb{R}^p$ the linear model's parameters

Exercise 3

- State the loss function for ordinary least squares (OLS)
- Formulate OLS as a minimization problem
- Compute the solution $\hat{\beta}^{OLS}$
- What happens if $X^T X$ is singular (not invertible)?

Exercise 4 - Bias and Variance of OLS

slide 16

Exercise 5 - Optimality of OLS

slide 17

Ridge Regression

Exercise 6 - Bias and Variance

slide 33

Exercise 7 - Performance

slide 34

Ridge Logistic Regression

Exercise 8

slide 46

EXERCISE 01

① $A \in \mathbb{R}^{m \times m}$

a) $f: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^{m \times m} \\ x \mapsto Ax \end{cases}$

$$f(x) = xA ; \nabla f \in \mathbb{R}^{m \times m}$$

$$\nabla f_{ij} = f'_{ij}(x) ; f'_{ij}(x) = f(x)_{ij} = A_{ij}x \Rightarrow \nabla f_{ij} = A_{ij} \Rightarrow \boxed{\nabla f = A}$$

b) $f: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^{m \times 1} \\ x \mapsto \underbrace{Ax}_{m \times n} \end{cases}$

m dimensions, m variables

$$\nabla f_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\frac{\partial f_i}{\partial x_j} = A_{ij}$$

$$\Rightarrow \boxed{\nabla f = A}$$

$$f_i(x) = A_i^T x \\ = \sum_k A_{ik} x_k$$

$$\nabla f \in \mathbb{R}^{m \times m}$$

c) $f: \begin{cases} \mathbb{R}^{m \times m} \rightarrow \mathbb{R} \\ X \mapsto \underbrace{\text{Tr}(A^T X)}_{\substack{\mathbb{R}^{m \times m} \times \mathbb{R}^{m \times m}} \in \mathbb{R}} \end{cases}$

$\text{Tr}(M) = \sum$ of diagonal elements of M .

$$\nabla f_{ij} = \frac{\partial f}{\partial x_{ij}} \quad \text{with } f(x) = \sum A_{ij} x_{ij}$$

$$\frac{\partial f}{\partial x_{ij}} = A_{ij} \Rightarrow \boxed{\nabla f = A}$$

If we have $f: \mathbb{R}^m \rightarrow \mathbb{R}$

$$f(x) = f(x_0) + (x - x_0)^T \nabla f(x_0) + O(\|x - x_0\|)$$

$$\approx f(x_0) + f'(x_0)(x - x_0) + O(x - x_0)$$

$$\nabla^2 f_{ij} = \frac{\partial^2 f}{\partial x_{ij}}$$

EXERCISE 2

Let $X \in \mathbb{R}^{m \times p}$, $\lambda > 0$. Prove that $M = X^T \bar{X} + \lambda I_p$ is invertible.

Let $M' = X^T X \in \mathbb{R}^{p \times p}$, we will show that M' has non-negative eigenvalues.
 Let μ be an eigenvalue of M' , and y an eigenvector.

$$M'y = \mu y \Rightarrow y^T M'y = y^T \mu y = \mu y^T y = \mu \|y\|^2$$

$$M' = X^T X \Rightarrow y^T M'y = y^T X^T X y = (Xy)^T (Xy) = \|Xy\|^2$$

$$\|Xy\|^2 = \mu \|y\|^2 \Rightarrow \mu = \frac{\|Xy\|^2}{\|y\|^2} \geq 0$$

\Rightarrow The eigenvalues of M are equal to $\lambda +$ the eigenvalues of M'

$$\mu \in \text{Eig}(M') \Leftrightarrow \mu + \lambda \in \text{Eig}(M) \quad (1)$$

Therefore the eigenvalues of M are $\geq \lambda$

So M is invertible \Rightarrow (M is invertible \Leftrightarrow its eigenvalues are $\neq 0$)

Proof of (1): M' is symmetric ($M'^T = M'$) so there exist $U \in \mathbb{R}^{p \times p}$ orthogonal ($U^T U = I$) such that $M' = U^T \Delta U$; $\Delta = \text{Diag}(\lambda_1, \dots, \lambda_m)$
 the λ_i 's are eigenvalues of M' .

$$M = M' + \lambda I = U^T \Delta U + \lambda U^T U = U^T (\Delta + \lambda) U$$

$$= U^T \begin{pmatrix} \lambda_1 + \lambda & 0 & \cdots & 0 \\ 0 & \ddots & & \\ \vdots & & \ddots & \\ 0 & & & \lambda_m + \lambda \end{pmatrix} U \Rightarrow \text{The eigenvalues of } M \text{ are } \underline{\lambda_i + \lambda}$$

EXERCISE 3

a) State the loss function for ordinary least squares (OLS)

$$y = \beta^T x + \epsilon$$

$$= \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

$$\epsilon = y - \beta^T x$$

$$\text{The criterion: } \text{MSE}(\beta) = \frac{1}{m} \sum (y_i - \beta^T x_i)^2$$

$$= \frac{1}{m} (y - X\beta)^T (y - X\beta) = \frac{1}{m} \|y - X\beta\|^2$$

b) Formulate OLS as a minimization Problem.

Ordinary Least Square

$$\min_{\beta} \text{MSE}(\beta), \text{MSE}: \mathbb{R}^p \rightarrow \mathbb{R}$$

c) Compute the solution $\hat{\beta}^{\text{OLS}}$

$$\nabla \text{MSE}(\hat{\beta}) = 0$$

$$\nabla \text{MSE} = \frac{1}{m} \left[(y - X\beta)^T (y - X\beta) \right] = \frac{2}{m} (y - X\beta)(-X) = \boxed{\frac{2}{m} X^T (X\beta - y)}$$

$$\nabla \text{MSE} = \frac{\partial \text{MSE}}{\partial \beta}$$

OLS Convex quadratic reaches optimum when $\nabla \text{MSE}(\beta^*) = 0$

$$X^T (X\beta^* - y) = 0$$

$$X^T X \beta^* = X^T y$$

* If $X^T X$ non singular $\beta^* = (X^T X)^{-1} X^T y$

$$A^T B = B^T A$$

$$\nabla \text{MSE} = \frac{2}{m} X^T (X\beta - y)$$

d) What happens if $X^T X$ is singular (not invertible)

When $X^T X$ is singular

\Rightarrow Pseudo inverse \rightarrow Find Solution when there is

\Rightarrow Regularize with small λ ; $(X^T X)^{-1}$ becomes $(X^T X + \lambda I_m)^{-1}$ no solution

Ridge regression no closed form / unstable \Rightarrow which is inv. (Ex. 2)

EXERCISE 4 - BIAS AND VARIANCE OF OLS (SLIDE 16)

Show that $\hat{\beta}^{\text{OLS}} = (X^T X)^{-1} X^T y$

Satisfies $\left\{ \begin{array}{l} E(\hat{\beta}^{\text{OLS}}) = \beta^* \\ \text{Var}(\hat{\beta}^{\text{OLS}}) = \dots \end{array} \right.$

$$\text{Var}(\hat{\beta}^{\text{OLS}}) = E(\hat{\beta}^{\text{OLS}} - \beta^*)(\hat{\beta}^{\text{OLS}} - \beta^*)^T = \sigma^2 (X^T X)^{-1}$$

For $A \in \mathbb{R}^{P \times P}$, $\text{Var}(A) \in \mathbb{R}^{P \times P}$ — is the covariance matrix.

$$\text{Var}(A) = E((A - E(A))(A - E(A))^T)$$

1-dimension $\text{Var}(A) = E(A - E(A))$

$$\text{Var}(A)_{ij} = \text{Cov}(A_i, A_j) = E(A_i A_j) - E(A_i)E(A_j)$$

In this exercise: $E(\varepsilon) = 0$; $\text{Var}(\varepsilon) = \sigma^2 I$

$$\Rightarrow \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \text{ if } i \neq j \text{ and } \text{Var}(\varepsilon_i) = \sigma^2$$

Let define $P = (X^T X)^{-1} X^T$; and assume: $y = X\beta^* + \varepsilon$

$$\begin{aligned}\hat{\beta}^{\text{OLS}} &= Py = (X^T X)^{-1} X^T (X\beta^* + \varepsilon) = (X^T X)^{-1} X^T X \beta^* + P\varepsilon \\ &= \beta^* + P\varepsilon\end{aligned}$$

E is linear

$$E[\hat{\beta}^{\text{OLS}}] = E(\beta^*) + E(P\varepsilon) = \beta^* + P E(\varepsilon) \stackrel{/\!/\!O}{=} \beta^* \Rightarrow \hat{\beta}^{\text{OLS}} \text{ is unbiased}$$

$$\begin{aligned}\Rightarrow \text{Var}(\hat{\beta}^{\text{OLS}}) &= \text{Var}(\hat{\beta}^{\text{OLS}} - \beta^*) = \text{Var}(\beta^* + P\varepsilon - \beta^*) = \text{Var}(P\varepsilon) \\ &= P \underbrace{\text{Var}(\varepsilon)}_{\sigma^2 I} P^T\end{aligned}$$

$$\begin{aligned}\text{Var}(\hat{\beta}^{\text{OLS}}) &= \sigma^2 P P^T = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \xrightarrow{P^T = X(X^T X)^{-1}} \\ &\quad \boxed{\sigma^2 (X^T X)^{-1}} \quad \left((M^{-1})^T = (M^T)^{-1} \right)\end{aligned}$$

EXERCISE 5

Gauss Markov theorem:

Least squares estimator is Best Linear Unbiased Estimator.

Assumption

$$y = X\beta^* + \varepsilon \quad \text{Among linear unbiased estimators}$$

$E(\varepsilon) = 0$ $\text{Var}(\varepsilon) = \sigma^2 I$ $\hat{\beta}^{\text{OLS}}$ has the lowest variance

if $\tilde{\beta} = C y$ unbiased ($E\tilde{\beta} = \beta^*$) then $\text{Var}(\hat{\beta}^{\text{OLS}}) \leq \text{Var}(\tilde{\beta})$

(Δ This condition does not mean $\text{Var}(\hat{\beta}^{\text{OLS}})_{ij} \leq \text{Var}(\tilde{\beta})_{ij} \forall i, j$)

i.e.:

$M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}})$ is a positive semi-definite matrix

$$\text{i.e. } \forall X_0, X_0^T M X_0 \geq 0$$

PROOF

Let $\tilde{B} = Cy$ be an unbiased estimator, $C \in \mathbb{R}^{P \times n}$

$$\text{Define: } D := C - (X^T X)^{-1} X^T \quad (*)$$

We are going to show that $\text{Var}(\hat{B}^{\text{OLS}}) \leq \text{Var}(\tilde{B})$

→ STEP 1: show that $Dx = 0$

\tilde{B} is unbiased: $E(\tilde{B}) = B^*$ (1)

$$E(\tilde{B}) = E(Cy) = C E(y) = C E(XB^* + \epsilon) = CXB^* + C E(\epsilon) = CXB^* \quad (2)$$

Combining (1) and (2) $B^* = CXB^*$

$$\Rightarrow (CX - I)B^* = 0 \quad \text{Recall } D = C - (X^T X)^{-1} X^T$$

$$= Dx B^* = 0 \quad Dx = CX - (X^T X)^{-1} X^T X = CX - I$$

This is true for all B^* $\Rightarrow \underline{Dx = 0} \quad \square$

→ step 2: show that $\text{Var}(\tilde{B}) = \text{Var}(\hat{B}^{\text{OLS}}) + \sigma^2 DD^T$

We reuse the notation $P = (X^T X)^{-1} X^T$

Since $Dx = 0$, $D P^T = P D^T = 0$

From (*) $C = D + P$

$$\Rightarrow \tilde{B} = Cy = (D + P)(XB^* + \epsilon) = \underline{DXB^*} + \underline{D\epsilon} + \underline{PB^*} + \underline{PE}$$

$$= B^* + (D + P)\epsilon$$

$$\text{Var}(\tilde{B}) = \text{Var}(B^* + (D + P)\epsilon) = (D + P) \underbrace{\text{Var}(\epsilon)}_{\sigma^2 I} (D + P)^T \quad (\text{Var}(B^*) = 0)$$

$$= \sigma^2 (D D^T + \underbrace{D P^T}_{=0} + P D^T + P P^T) \stackrel{\sigma^2 I}{=} \sigma^2 D D^T + \underbrace{\sigma^2 (X^T X)^{-1}}_{= \text{Var}(\hat{B}^{\text{OLS}})}$$

$$\Rightarrow \text{Var}(\tilde{B}) = \text{Var}(\hat{B}^{\text{OLS}}) + \sigma^2 D D^T$$

\square

Note: $y \in \mathbb{R}^P$ random vector, $C \in \mathbb{R}^{P \times P}$ fixed. Then $E(Cy) = C E(y)$
and $\text{Var}(Cy) = C \text{Var}(y) C^T$

→ Step 3 (final) $M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}})$ is positive semi-definite (P.s.d.)

• $M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}}) = \sigma^2 \Delta \Delta^T$ (from step 2)

• Let $x_0 \in \mathbb{R}^m$: $x_0^T M x_0 = \sigma^2 x_0^T \Delta \Delta^T x_0 = \sigma^2 (\Delta^T x_0)^T (\Delta^T x_0) = \sigma^2 \|\Delta^T x_0\|^2$

This is fine $\forall x_0$: M is positive semi-definite (P.s.d.) $\Rightarrow 0$

Finally: $\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta}^{\text{OLS}})$

$\tilde{\beta}$ has a "large" variance than $\hat{\beta}^{\text{OLS}}$

Why does this matter?

Suppose we observe measurements $y = \mu + \varepsilon$ what we want to estimate
noise

m measurements $y_i = \mu + \varepsilon_i$; $\mathbb{E}(\varepsilon_i) = 0$; $\text{Var}(\varepsilon) = \sigma^2 I$

We want an estimator $\hat{\mu}$ of μ .

Which is the best estimator?

① $\hat{\mu}_1 = \frac{1}{m} \sum_{i=1}^m y_i$; ② $\hat{\mu}_2 = y_i$; ③ $\hat{\mu}_3 = \frac{1}{m} \sum y_i + \frac{1}{m}$

→ Bias: $\mathbb{E}(\hat{\mu}_1) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}(y_i) = \frac{1}{m} \sum_{i=1}^m \mu = \mu = \mathbb{E}(\hat{\mu}_2)$

$\mathbb{E}(\hat{\mu}_3) = \mu + \frac{1}{m} \neq \mu$

$\Rightarrow \hat{\mu}_1$ and $\hat{\mu}_2$ are unbiased but $\hat{\mu}_3$ is biased

→ Variance: $\text{Var}(\hat{\mu}_1) = \text{Var}\left(\frac{1}{m} \sum_{i=1}^m y_i\right) = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(\mu) + \text{Var}(\varepsilon)$
 $= \frac{1}{m^2} (m \text{Var}(\mu) + m \sigma^2 I) = \frac{\sigma^2}{m} = \text{Var}(\hat{\mu}_3)$

$\text{Var}(\hat{\mu}_2) = \sigma^2$

→ Error: $\text{Er}(\hat{\mu}) = \text{bias}^2(\hat{\mu}) + \text{Var}(\hat{\mu})$

$\text{Er}(\hat{\mu}_1) = \text{bias}^2(\hat{\mu}_1) + \text{Var}(\hat{\mu}_1) = \underbrace{\mathbb{E}(\mathbb{E}(\hat{\mu}_1))}_{0} + \frac{\sigma^2}{m} = \frac{\sigma^2}{m}$

$\text{Er}(\hat{\mu}_2) = 0 + \sigma^2 = \sigma^2$

$\text{Er}(\hat{\mu}_3) = \underbrace{\mathbb{E}(\mu)\mathbb{E}(\mu)}_{0} + \underbrace{\mathbb{E}(\mu)\mathbb{E}\left(\frac{1}{m}\right)}_{0} + \underbrace{\mathbb{E}\left(\frac{1}{m}\right)\mathbb{E}(\mu)}_{0} + \underbrace{\mathbb{E}\left(\frac{1}{m}\right)\mathbb{E}\left(\frac{1}{m}\right)}_{\frac{1}{m^2}} + \frac{\sigma^2}{m} = \frac{1}{m^2} + \frac{\sigma^2}{m}$

$\Rightarrow \hat{\mu}_1$ and $\hat{\mu}_2$ are both unbiased, but $\hat{\mu}_1$ is much better than $\hat{\mu}_2$

$\Rightarrow \hat{\mu}_3$ is biased, and $\hat{\mu}_2$ is not, but $\hat{\mu}_3$ is much better than $\hat{\mu}_2$

\hookrightarrow Variance matters just as much as bias.

EXERCISE 6 (Slides 33) Ridge Regression

$$\left\{ \begin{array}{l} \text{bias}(\hat{\beta}_{\text{Ridge}}) = ? = \frac{\lambda}{1+\lambda} \beta^* \\ \text{Var}(\hat{\beta}_{\text{Ridge}}) = ? = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}_{\text{OLS}}) \end{array} \right.$$

$\lambda \geq 0$

loss

$$L(\beta) = \text{MSE}(\beta) + \lambda \|\beta\|^2 = \frac{1}{m} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|^2$$

gradient

$$\nabla J(\beta) = \frac{2}{m} \mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) + 2\lambda \beta$$

$$\nabla J(\hat{\beta}_{\lambda}) = 0 \iff \frac{2}{m} \mathbf{X}^T \mathbf{X} \beta + 2\lambda \beta = \frac{2}{m} \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \cancel{\frac{2}{m} \beta (\mathbf{X}^T \mathbf{X} + \lambda m \mathbf{I})} = \cancel{\frac{2}{m} \mathbf{X}^T \mathbf{y}} \Rightarrow \boxed{\hat{\beta}_{\lambda} = (\mathbf{X}^T \mathbf{X} + \lambda m \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}}$$

N.B.: When $\lambda > 0$, $(\mathbf{X}^T \mathbf{X} + \lambda m \mathbf{I})^{-1}$ always exist (Exercise 2)

Assumption: $\mathbf{X}^T \mathbf{X} = m \mathbf{I}$ ($\times \times$) (orthogonal vectors \mathbf{x}_i)

$$\begin{aligned} \Rightarrow \hat{\beta}_{\lambda}^{\text{Ridge}} &= \frac{1}{m(1+\lambda)} \mathbf{X}^T \mathbf{y} \quad \Rightarrow \mathbf{y} = \mathbf{X}\beta^* + \mathbf{\epsilon} \\ &= \frac{1}{m(1+\lambda)} \left(\cancel{\frac{1}{m} \mathbf{X}^T \mathbf{X} \beta^*} + \mathbf{X}^T \mathbf{\epsilon} \right) = \frac{1}{(1+\lambda)} \beta^* + \frac{1}{m(1+\lambda)} \mathbf{X}^T \mathbf{\epsilon} \end{aligned}$$

$$\Rightarrow \hat{\beta}_{\lambda}^{\text{Ridge}} = \frac{1}{(1+\lambda)} \beta^* + \frac{1}{m(1+\lambda)} \mathbf{X}^T \mathbf{\epsilon}$$

Taking expectations:

$$\mathbb{E}(\hat{\beta}_{\lambda}^{\text{Ridge}}) = \frac{1}{1+\lambda} \mathbb{E}(\beta^*) + \frac{1}{m(1+\lambda)} \mathbf{X}^T \mathbb{E}(\mathbf{\epsilon}) = \frac{1}{1+\lambda} \beta^*$$

$$\Rightarrow \text{Hence } \text{bias}(\hat{\beta}_{\lambda}^{\text{Ridge}}) = \mathbb{E}(\hat{\beta}_{\lambda}^{\text{Ridge}}) - \beta^* = -\frac{\lambda}{1+\lambda} \beta^*$$

Large λ = large bias
Small λ = small bias

\Rightarrow When $\lambda = 0$ (no regularization), bias = 0, $\hat{\beta}_{\lambda}^{\text{Ridge}} = \hat{\beta}_{\lambda}^{\text{OLS}}$

When $\lambda = +\infty$, bias = $-\beta^*$ because $\hat{\beta}_{\lambda}^{\text{Ridge}} = 0$

Assumption

$$X^T X = mI$$

$$X_i^T X_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

$$\hat{\beta}_\lambda^{\text{Ridge}} = \frac{1}{m(1+\lambda)} X^T Y$$

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{1}{m^2(1-\lambda)^2} \text{Var}(X^T Y) = \frac{1}{m^2(1-\lambda)^2} X^T \underbrace{\text{Var}(\epsilon)}_{\sigma^2 I} X$$

$$= \frac{\sigma^2}{m^2(1-\lambda)^2} \underbrace{X^T X}_{mI} = \frac{\sigma^2}{m(1+\lambda)^2}$$

Large λ = small variance
Small λ = large variance

Uncheck out assumptions

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{\sigma^2}{m(1+\lambda)^2} ; \text{Var}(\hat{\beta}^{\text{OLS}}) = (X^T X)^{-1} \sigma^2 = \frac{\sigma^2}{m}$$

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}^{\text{OLS}})$$

$\lambda \rightarrow 0 : \hat{\beta}_\lambda^{\text{Ridge}} \xrightarrow{\hat{\beta}^{\text{OLS}}} \hat{\beta}^{\text{OLS}}$: no bias, large variance

$\lambda \rightarrow +\infty : \hat{\beta}_\lambda^{\text{Ridge}} \xrightarrow{0} 0$: Large bias, no variance

EXERCISE 7. Expectation of total error

$$f(\lambda) = \underset{S, X_0}{\mathbb{E}} [\text{bias}^2(X_0^T \hat{\beta}_\lambda^{\text{Ridge}}) + \text{Var}(X_0^T \hat{\beta}_\lambda^{\text{Ridge}})]$$

$$\left| \begin{array}{l} \mathbb{E} X_0 = 0 \\ \mathbb{E} X_0 X_0^T = I \end{array} \right.$$

→ Let us compute the two terms:

$$\begin{aligned} - \text{bias}^2(X_0^T \hat{\beta}_\lambda^{\text{Ridge}}) &= (X_0^T \text{bias}(\hat{\beta}_\lambda^{\text{Ridge}}))^2 \\ &= \frac{\lambda^2}{(1+\lambda)^2} (X_0^T \beta^*)^2 \\ &= \frac{\lambda^2}{(1+\lambda)^2} \beta^T X_0 X_0^T \beta^* \end{aligned}$$

$$\text{bias}(\hat{\beta}) = \beta^*$$

$$\begin{aligned} (X_0^T \beta^*)^2 &= \sum X_{0,i} X_{0,j} \beta_i \beta_j^* \\ &= \beta^T X_0 X_0^T \beta^* \end{aligned}$$

$$- \text{Var}(X_0^T \hat{\beta}_\lambda^{\text{Ridge}}) = X_0^T \text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) X_0 = \frac{\sigma^2}{m(1+\lambda)} X_0^T X_0$$

→ Taking expectations, with $\mathbb{E}(X_0 X_0^T) = I$, we get

$$f(\lambda) = \frac{\lambda^2}{(1+\lambda)^2} \|\beta\|^2 + \frac{\sigma^2 p}{m(1+\lambda)^2}$$

$$\mathbb{E}(X_0^T X_0) = \sum_{i=1}^p \mathbb{E}(X_{0,i}^2) = p$$

Finding the optimum λ^*

Define $a = \|\beta^*\|^2 = f(+\infty)$; $b = \frac{\sigma^2 p}{m} = f(0)$

$$f(\lambda) = \frac{1}{(1+\lambda)^2} (a\lambda^2 + b) ; \quad f'(\lambda) = \frac{2}{(1+\lambda)^3} (a\lambda - b)$$

$$f'(\lambda) = 0 \Leftrightarrow \lambda^* = \frac{b}{a} = \frac{\sigma^2 p}{m \|\beta^*\|^2}$$

$$\text{Keeping our notations: } f(\lambda^*) = \frac{1}{(1+\frac{b}{a})^2} \left(a \cdot \left(\frac{b}{a} \right)^2 + b \right) = \frac{ab}{a+b}$$

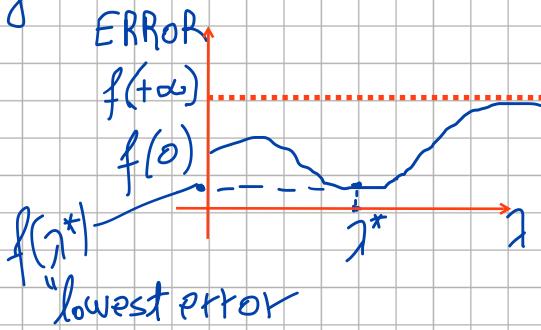
$$f(\lambda^*) = \frac{ab}{a+b} = \frac{f(0)f(+\infty)}{f(0)+f(+\infty)}$$

$$\begin{aligned} a \geq 0 \Rightarrow \frac{ab}{a+b} &\leq \frac{a(a+b)}{a+b} = a \\ b \geq 0 \Rightarrow \frac{ab}{a+b} &\leq \frac{(a+b)b}{a+b} = b \end{aligned} \quad \left. \begin{array}{l} \frac{ab}{a+b} \leq \min(a, b) \end{array} \right\}$$

$f(X)$ = error function \exists optimal λ^* that gives the best fit
cannot compute it directly. Find with cross-validation

$$\Rightarrow f(\lambda^*) \leq \min(f(0), f(+\infty))$$

Using λ^* is better than not using regularization.



Conclusion (curve)
Some regularization is better than
and better than the estimator \hat{y} .

EXERCISE 8

RIDGE LOGISTIC REGRESSION (SLIDE 46)

$$\beta^{\text{new}} = \beta^{\text{old}} - [\nabla_{\beta}^2 J(\beta^{\text{old}})]^{-1} \nabla_{\beta} J(\beta^{\text{old}})$$

Lemma: β^{new} is also the solution to $\min_{\beta} J(\beta) = \frac{1}{m} \sum_{i=1}^m \ln(1 + e^{-y_i \beta^T x_i}) + \lambda \|\beta\|^2$

where $J_q(\beta) := J(\beta^{\text{old}}) + \nabla J(\beta^{\text{old}})(\beta - \beta^{\text{old}})$
 $+ \frac{1}{2} (\beta - \beta^{\text{old}})^T \nabla_{\beta}^2 J(\beta^{\text{old}})(\beta - \beta^{\text{old}})$

J_q is the quadratic approximation to J in β^{old}

(analogy in 1 dimension: $J_q(x) = J'(x_0)(x - x_0)$

$$+ \frac{1}{2} J''(x_0)(x - x_0)^2$$

Proof:

$$\begin{aligned} \nabla J_q(\beta) &= \nabla J(\beta^{\text{old}}) + \nabla_{\beta}^2 J(\beta^{\text{old}})(\beta - \beta^{\text{old}}) \\ &= 0 \iff \beta = \beta^{\text{old}} - [\nabla_{\beta}^2 J(\beta^{\text{old}})]^{-1} \nabla_{\beta} J(\beta^{\text{old}}) \\ &= \beta^{\text{new}} \end{aligned}$$

writing all the terms of J_q that depend on β :

$$\Rightarrow \beta^T \nabla J(\beta^{\text{old}}) = \frac{1}{m} \beta^T X^T P y + 2 \lambda \beta^T \beta$$

$$\Rightarrow \frac{1}{2} \beta^T \nabla_{\beta}^2 J(\beta^{\text{old}}) \beta = \frac{1}{2m} \beta^T X^T W X \beta + \lambda \beta^T \beta$$

$$\Rightarrow -\beta^T \nabla_{\beta}^2 J(\beta^{\text{old}}) \beta^{\text{old}} = -\frac{1}{m} \beta^T X^T W X \beta^{\text{old}} - 2 \lambda \beta^T \beta^{\text{old}}$$

Putting it all together

$$\begin{aligned} J_q(\beta) &= -\frac{1}{m} \beta^T X^T W X \beta^{\text{old}} + \frac{1}{m} \beta^T X^T P y + \frac{1}{2m} \beta^T X^T W X \beta + \lambda \beta^T \beta + C \\ &= -\frac{1}{m} \beta^T X^T W \underbrace{(X \beta^{\text{old}} - W^{-1} P y)}_{\beta} + \frac{1}{2m} \beta^T X^T W X \beta + \lambda \beta^T \beta + C \\ &= \frac{1}{2m} (X \beta - \beta)^T W + \lambda \|\beta\|^2 + C \end{aligned}$$