# Chapter 1:

None

# Chapter 2:

**2.1**  $\underset{k}{\operatorname{argmax}}\ \hat{y} = \underset{k}{\operatorname{argmin}}\ |t_k - \hat{y}|_2$

$$|t_k - \hat{y}|^2 = (1 - \hat{y})^2 + \sum_{j \neq k} \hat{y}_j^2$$

$$= 1 - 2\hat{y} + |\hat{y}|^2$$

1 by assumption

$$\Rightarrow \underset{k}{\operatorname{argmin}}\ 1 - 2\hat{y} = \operatorname{argmax} \hat{y}$$

**2.2**  $b_{1\cdots 10} \sim N\left(\begin{smallmatrix}1\\0\end{smallmatrix}, \mathbb{I}\right)$   blue

$o_{1\cdots 10} \sim N\left(\begin{smallmatrix}0\\1\end{smallmatrix}, \mathbb{I}\right)$   orange

$$N(b_i, \mathbb{I}/5) \qquad N(o_i, \mathbb{I}/5)$$

$$\Rightarrow p(x \mid blue) = \sum_{i=1}^{10} \frac{1}{10} N(b_i, \mathbb{I}/5)$$

$$\Rightarrow p(blue \mid x) = \frac{\left(\frac{1}{10} \sum_{i=1}^{10} N(b_i, \mathbb{I}/5)\right) \cdot \frac{1}{2}}{\frac{1}{10} \sum_{j=1}^{10} N(b_j, \mathbb{I}/5) + \frac{1}{10} \sum_{i=1}^{10} N(o_i, \mathbb{I}/5)}$$

$$p(blue \mid x) = p(orange \mid x)$$

$$\Rightarrow \sum N\left(b_i, \frac{\mathbb{I}}{5}\right) = \sum N\left(o_i, \frac{\mathbb{I}}{5}\right)$$

$$\sum \exp\left((x-b_i)^2 \cdot \frac{s}{2}\right) = \sum \exp\left((x-0_i)^2 \cdot \frac{s}{2}\right)$$

**2.3** $\quad Vol(B_r^p) = \nu_p r^p$

$$Prob\left(x \sim B_1^p \notin B_r^p\right) = 1 - \frac{Vol\, B_r^p}{Vol\, B_1^p} = 1 - r^p$$

$$P(r,N) := Prob\left[(x_1 \cdots x_N \text{ iid} \sim B_1^p) \in B_r^p\right] = (1-r^p)^N$$

Median is when $P_N = \frac{1}{2}$

$$\Rightarrow (1-r^p)^N = \frac{1}{2} \qquad \Rightarrow r = \left(1 - \frac{1}{2^{1/N}}\right)^{1/p}$$

by def of pdf

$$\int_{-\infty}^{med} = \int_{med}^{\infty}$$

**2.4** $\quad$ WLOG $\quad a_i = \frac{1}{\sqrt{N}} \mathbb{1}$

$$\Rightarrow z = \frac{1}{\sqrt{N}} \sum x_i \sim N(0,1) \qquad \text{if} \quad x_i \in N(0, \mathbb{1})$$

$\Rightarrow$ for $p=10$ avg distance$^2$ is

$$Mean\left(\Gamma\left(\frac{\nu}{2}, \frac{1}{2}\right)_{\nu=p}\right) = p \Rightarrow RMS \text{ dist} \sim \sqrt{p}$$

While along any axis its $1$

**2.5** $\quad$ Test point $x_0$, we know $Y = X^T\beta + \varepsilon$

$$\hat{y}_0 = x_0^T \hat{\beta}$$

$$\hat{y}_0 = x_0^T \beta + \sum_{i=1}^{N} l_i(x_0)\varepsilon_i$$

$$\left[x_0(X^TX)^{-1}X^T\right]_i$$

a) $EPE(x_0) = \mathbb{E}_{y_0|x_0} \mathbb{E}_{\tau} (y_0 - \hat{y}_0)$

$$= Var(y_0|x_0) + \left(x_0^T\beta - \mathbb{E}_{\tau}\hat{y}_0\right)^2 + \mathbb{E}_{\tau}\left(\hat{y}_0 - \mathbb{E}_{\tau}\hat{y}_0\right)^2$$

$$\underbrace{\qquad}_{Bayes} \qquad \underbrace{\qquad}_{Bias} \qquad \underbrace{\qquad}_{Var}$$

$$= \sigma^2 \quad + \quad 0 \quad + \mathbb{E}_{\tau} x_0^T(X^TX)^{-1}x_0\,\sigma^2$$

b)

$$N \to \infty \;\Rightarrow\; X^TX \to N\,Cov\,X \qquad \textcolor{red}{(empirical\ cov\ concentrates)}$$

$$\Rightarrow \mathbb{E}_{x_0} EPE(x_0) = \sigma^2\left(1 + \frac{1}{N}\mathbb{E}_{x_0}\left[x_0^T\,Cov(X)^{-1}x_0\right]\right)$$

$$\mathbb{E}_{x_0} x_0 x_0^T = Cov\,x_0$$

$$= \sigma^2\left(1 + \frac{1}{N}Tr\left[Cov(X)^{-1}Cov(x_0)\right]\right)$$

$$\textcolor{red}{assuming} \xrightarrow{\hspace{1cm}} = \sigma^2\left(1 + \frac{1}{N}Tr\left[\mathbb{1}_p\right]\right)$$
$$\textcolor{red}{x_0\ is}$$
$$\textcolor{red}{in\text{-}distribution} \qquad = \sigma^2\left(1 + \frac{p}{N}\right)$$

2.6 $\quad RSS(\theta) = \sum_{i=1}^{N}\sum_{\ell=1}^{N_i}\left(y_{i\ell} - f_\theta(x_i)\right)^2$

$y_{i\ell}, \ \ell \in 1,\dots,N_i \qquad = \sum_{i,\ell}\left(y_{i\ell} - \bar{y}_i + \bar{y}_i - f_\theta(x_i)\right)^2$

$$\textcolor{red}{i\neq e}$$

$$= \sum_{i,\ell}\left(y_{i\ell} - \bar{y}_i\right)^2 + \sum_{i} N_i\left(\bar{y}_i - f_\theta(x_i)\right)^2$$

2.7 $\quad$ a) $\quad$ L.R. : $\quad x_0(X^TX)^{-1}X^T$

$\quad$ kNN : $\ell_i(x_0;X) = \begin{cases} \frac{1}{k} & \text{if } x_i \in N_k(x_0) \\ 0 \end{cases}$

$$\textcolor{red}{Bias}$$

b) $\quad \mathbb{E}_{y|x}\left|f(x_0) - \hat{f}(x_0)\right|^2 = \left(f(x_0) - \mathbb{E}_{y|x}\hat{f}(x_0)\right)^2$

$$+ \mathbb{E}_{y|x}\left(\hat{f}(x_0) - \mathbb{E}_{y|x}\hat{f}(x_0)\right)^2$$

$$\underbrace{\qquad}_{Var}$$

c) $\underset{Y,X}{E} \left| \bar{f}(x_0) - \hat{f}(x_0) \right|^2 = \left| \bar{f}(x_0) - \underset{Y,X}{E} \hat{f}(x_0) \right|^2$ <span style="color:salmon">Bias</span>

$$+ \underset{Y,X}{E} \left| \hat{f}(x_0) - \underset{Y,X}{E} \hat{f}(x_0) \right|^2$$

<span style="color:salmon">Var</span>

d) <span style="color:salmon">bias:</span>

$$\bar{f}(x_0) - \underset{Y|X}{E} \hat{f}(x_0) = f(x_0) - \sum_i \ell_i(x_0; X) \, f(x_i)$$

<span style="color:salmon">Var:</span>

$$\underset{Y|X}{E} \left( \hat{f}(x_0) - \underset{Y|X}{E} \hat{f}(x_0) \right)^2 = \underset{\varepsilon_i}{E} \left( \sum_i \ell_i(x_0; X) \, \varepsilon_i \right)^2$$

$$= \sigma^2 \sum_i \ell_i(x_0; X)^2$$

$$S := \begin{pmatrix} \ell_1(x_0; X) \\ \vdots \\ \ell_n(x_0; X) \end{pmatrix}$$

$\Rightarrow$ Bias $= f(x_0) - S^T f$ $\Rightarrow$ Bias$^2 = \bar{f}(x_0)^2 - 2f(x_0) S^T f + f^T S S^T f$
Var $= \sigma^2 S^T S$

$$\Rightarrow \text{Bias}^2 = f(x_0)(f(x_0) - 2 S^T f) + \frac{f^T \text{Var } f}{\sigma^2}$$

for c)

<span style="color:salmon">bias:</span>

$$f(x_0) - \int \pi \, dx_i \, h(x_i) \sum_i \ell_i(x_0; X) \, f(x_i)$$

<span style="color:salmon">$\mathcal{S}(f; x_0)$</span>

<span style="color:salmon">Var:</span>

$$\sum_i \left[ \ell_i(x_0; X) \, y_i - \int dx_i \, h(x_i) \, \ell_i(x_0; x_0) \, f(x_i) \right]$$

Do they just want

$$\text{Bias}^2 \hat{f}(x_0) + \text{Var } \hat{f}(x_0) = \underset{X}{E} \left[ \text{Bias}\left(\hat{f}(x_0) | X\right)^2 + \text{Var}\left(\hat{f}(x_0) | X\right) \right]$$

## 2.8

```
Test error rate of Linear Regression is 4.12%
Train error rate of Linear Regression is 0.58%
k-NN Model: k is 1, train/test error rates are 0.00% and 2.47%
k-NN Model: k is 2, train/test error rates are 0.58% and 2.47%
k-NN Model: k is 3, train/test error rates are 0.50% and 3.02%
k-NN Model: k is 4, train/test error rates are 0.43% and 2.75%
k-NN Model: k is 5, train/test error rates are 0.58% and 3.02%
k-NN Model: k is 6, train/test error rates are 0.50% and 3.02%
k-NN Model: k is 7, train/test error rates are 0.65% and 3.30%
k-NN Model: k is 8, train/test error rates are 0.58% and 3.30%
k-NN Model: k is 9, train/test error rates are 0.94% and 3.57%
k-NN Model: k is 10, train/test error rates are 0.79% and 3.57%
k-NN Model: k is 11, train/test error rates are 0.86% and 3.57%
k-NN Model: k is 12, train/test error rates are 0.72% and 3.57%
k-NN Model: k is 13, train/test error rates are 0.86% and 3.85%
k-NN Model: k is 14, train/test error rates are 0.86% and 3.85%
k-NN Model: k is 15, train/test error rates are 0.94% and 3.85%
```

## 2.9 $\quad \mathbb{E}\, R_{tr}(\hat{\beta}) \le \mathbb{E}\, R_{te}(\hat{\beta})$

$$\mathbb{E}\left[\frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{\beta}\cdot x_i)^2\right] \le \mathbb{E}\left[\frac{1}{N}\sum_{i=1}^{N}(y_i - \beta x_i)^2\right]$$

$$= \mathbb{E}\, R_{tr}(\beta)$$

$$\underset{x,y,\tilde{x},\tilde{y}}{\mathbb{E}}\, R_{te}(\hat{\beta}) = \underset{x,y}{\mathbb{E}}\,\underset{\tilde{x},\tilde{y}}{\mathbb{E}}\, \frac{1}{M}\sum_{i=1}^{M}(\hat{\tilde{y}}_i - \hat{\beta}\tilde{x}_i)^2$$

$$\ge \underset{\tilde{x},\tilde{y}}{\mathbb{E}}\, \frac{1}{M}\sum_{i=1}^{M}(\tilde{y} - \beta^*\tilde{x}_i)^2 \qquad \color{red}{\beta^* \text{ is optimal for test set MSE}}$$

$$= \underset{x,y}{\mathbb{E}}\, \frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{\beta}x_i)^2$$

$$= \underset{x,y}{\mathbb{E}}\, R_{tr}(\hat{\beta})$$

## Chapter 3

1. We will show that the F-statistic for adding/dropping a single term is the square of its z-score

$$\hat{\sigma}^2 = \frac{RSS_1}{N - p_1 - 1} \qquad (N - p - 1)\hat{\sigma}^2 \sim \chi^2_{N-p-1}\cdot\sigma^2$$

$$\hat{\beta} \sim N(\beta,\, (X^TX)^{-1}\sigma^2) \qquad v_j = \left[(X^TX)^{-1}\right]_{jj}$$

$$z_j = \frac{\hat{\beta}_j}{\sigma\sqrt{v_j}} \sim N(0,1)$$

Show: F statistic is given by $z^2$

$$F := \frac{RSS_0 - RSS_1}{\frac{RSS_1}{N-p-1}}$$

<span style="color:pink">bigger model</span>

$$\min_{\beta} \ (y - X\beta)^2 + \lambda(\beta^T e_j - 0)$$

$$\Rightarrow \beta^0 = (X^TX)^{-1}(X^Ty - \lambda e_j)$$

$$\beta^0 \cdot e_j = 0 \Rightarrow \lambda = \frac{e_j^T(X^TX)^{-1}X^Ty}{e_j^T(X^TX)^{-1}e_j}$$

$$\Rightarrow \beta^0 = \beta - \frac{e_j^T(X^TX)^{-1}X^Ty}{e_j^T(X^TX)^{-1}e_j}(X^TX)^{-1}e_j$$

$$\Rightarrow RSS_0 = (y - X\beta^0)^2 = (y - X\beta' - \lambda X(X^TX)^{-1}e_j)^2$$

$$= RSS_1 - 2\lambda \underbrace{(y - X\beta^0)^T}_{0} X(X^TX)^{-1}e_j$$

$$+ \lambda^2 e_j (X^TX)^{-1}e_j$$

$$(y - X\beta^0)^T X = yX - yX(X^TX)^{-1}X^TX = 0$$

$$\Rightarrow RSS_0 - RSS_1 = \frac{(e_j(X^TX)^{-1}X^T)^2}{e_j(X^TX)^{-1}e_j} \quad \color{pink}{\S(\hat{\beta}_j)}$$
$$\color{pink}{\S v_j}$$

$$\Rightarrow \frac{RSS_0 - RSS_1}{\frac{RSS_1}{N-p_1-1}} = \frac{\hat{\beta}_j^2}{v_j} \frac{1}{\hat{\sigma}_j}$$

$$= (z_j)^2 \qquad \color{pink}{\ddot\smile}$$

2. Method 1 gives: $Var(y_0) = x_0^T \, Var\beta \, x_0$

$$= x_0^T (X^TX)^{-1} x_0 \, \sigma^2$$

$$\Rightarrow y_0 = \hat{y}_0 \pm \sigma \sqrt{x_0 (X^TX)^{-1} x_0}$$

Method 2 gives: $\beta \sim N(\hat{\beta}, (X^TX)^{-1}\sigma^2)$

$$C_\beta = \{\beta \mid (\beta - \hat{\beta})\frac{X^TX}{\sigma^2}(\beta - \hat{\beta}) \leq \chi^2_{q, 0.05}\}$$

$$X^T X = U^T U \qquad U \text{ is upper } \triangle$$

$$\tilde{v} := U(\beta - \hat{\beta}) \quad \text{lies in ball of}$$
$$\text{radius } r = \sigma \sqrt{x^2_{4,0.05}}$$

$$\beta = \hat{\beta} + \sigma \sqrt{x^2_{0.05}} \; U \cdot f \qquad\qquad f \in S^3$$

$\underbrace{\phantom{\beta = \hat{\beta} + \sigma \sqrt{x^2_{0.05}} U \cdot f}}_{\text{\textcolor{pink}{elliptical}}}$

$\Rightarrow$ Method 2 gives tighter bounds

b.c. it bounds all $\beta_i$ at the same time

3.3 a) Let $\theta = c^T y$ another estimate of $a^T \beta$

WLOG $\quad c = a(X^T X)^{-1} X^T + d \quad$ ← arbitrary, possibly $x, y$ dep

$$\mathbb{E}_y[c^T y] = \mathbb{E}_\varepsilon \left( a(X^T X)^{-1} X^T + d \right)(X\beta + \varepsilon)$$
$$= a^T \beta + dX\beta$$

unbiased $\Leftrightarrow$ $dX\beta = 0 \Rightarrow dX = 0$ $\forall$ vecs

$$\text{Var}(c^T y) = \sigma^2 \left( a(X^T X)^{-1} X^T + d \right)\left( a(X^T X)^{-1} X^T + d \right)^T$$
$$= \sigma^2 \underbrace{dd^T}_{\textcolor{pink}{\geq 0}} + \text{Var } a^T \hat{\beta}$$

b) As before but now $d \to D$

$$\Rightarrow DX = 0$$
$$\Rightarrow \text{Var } C^T y = \sigma^2 DD^T + \text{Var}(\hat{\beta})$$
$$\Rightarrow \text{Var } C^T y \succeq \text{Var } \hat{\beta}$$
$$\Rightarrow \tilde{V} \succeq \hat{V}$$

3.4 $\quad X = QR \quad \Rightarrow \quad X^T X = R^T R$

$$(R^T R)^{-1} R Q^T y$$
$$= R^{-1} Q^T y$$

$Q$ is $N \times (p+1)$

$Q^T Q = \mathbb{I}_{p+1}$

$R$ is $(p+1) \times (p+1)$

calculate $q_k^T y$ to fill column vec of $Q^T y \in \mathbb{R}^{p+1}$

solve $R^t Q^T y$ by backsub

$$\hat{y} = Q Q^T y$$

## 3.5 In 3.41

$$\underset{\beta}{\text{argmin}} \sum_{j=1}^{N} (y_i - \beta_0 - X \cdot \beta)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

take $x \to x - \bar{x}_j$

$$\implies \sum_i \left( y_i - \beta_0 - \sum_j (x_{ij} - \bar{x}_j) \beta_j \right)^2 + \lambda |\beta|^2$$
$$- \sum_j \bar{x}_j \beta_j$$

$$\implies \beta_0^c = \beta_0 - \sum_j \bar{x}_j \beta_j$$
$$\beta_j^c = \beta_j$$

$$\frac{\partial \mathcal{L}}{\partial \beta_0} = \sum_i y_i - N \beta_0^c - \sum_{ij} (x_{ij} - \bar{x}_j) \beta_j = 0$$

<span style="color:red">0 as a sum</span>

$$\implies \beta_0^c = \bar{y}$$

$$\tilde{y} = y_i - \beta_0^c$$
$$\tilde{x}_{ij} = x_{ij} - \bar{x}_j \qquad \implies \underset{\beta}{\min} (\tilde{X}\beta - \tilde{y})^2 + \lambda |\beta|^2$$

$$\implies \hat{\beta}_c = (X^T X + \lambda \mathbb{1})^{-1} \tilde{X}^T y$$

## 3.6 $P(\theta | X) \propto \exp\left[ -\frac{1}{\sigma^2} \sum_{j=1}^{N} (x_i \cdot \beta - y_i)^2 - \frac{1}{\tau} \mathbb{1} \beta^2 \right]$

<span style="color:red">extensive in N</span>

$$\implies \lambda = \frac{\sigma^2}{\tau} \qquad \implies \underset{\beta}{\min} (X\beta - y)^2 + \lambda \beta^2$$

$$\implies (X^T X + \lambda \mathbb{1})^{-1} X^T y$$

**3.7** $P(\beta|y) = \dfrac{P(\beta)\,P(y|\beta)}{P(y)}$

$$\Rightarrow -\log P(\beta|y) = \frac{1}{c^2}|\beta|^2 + \frac{1}{\sigma^2}\sum_{i=1}^{N}\left(y_i - \beta_0 - \sum_j x_{ij}\beta_j\right)^2 + \log Z$$

$$\lambda = \frac{\sigma^2}{c^2}$$

**3.8** let $\vec{x}_i = \begin{pmatrix} x_{1i} \\ \vdots \\ x_{Ni} \end{pmatrix} \in \mathbb{R}^N$ and $\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ be the columns of $X$

$$\vec{q}_j = \begin{pmatrix} q_{1j} \\ \vdots \\ q_{Nj} \end{pmatrix} \in \mathbb{R}^N \quad \text{the columns of } Q$$

$$X = Q R = \square \text{⬚}$$

$$\mathbf{1} = r_{00}\cdot \vec{q}_0$$

$$\Rightarrow r_{00} = \sqrt{N} \qquad \vec{q}_0 = \frac{1}{\sqrt{N}} \qquad \text{⭐}$$

$$\Rightarrow \bar{q}_j = \sum_j q_{ij}/N = \frac{1}{\sqrt{N}} q_0^T q_j = 0$$

$$\Rightarrow \bar{x}_j = r_{0j}\, q_0 = \frac{r_{0j}}{\sqrt{N}}$$

$$\Rightarrow \tilde{x} = \vec{x}_j - \bar{x}_j \mathbf{1}$$

$$= \sum_{k\neq 0} q_k\, r_{kj}$$

take $Q_2 = (q_1 \cdots q_p)$

$$\tilde{X} = (\tilde{x}_1 \cdots \tilde{x}_p) = U\Sigma V^T \in \mathbb{R}^{N\times p}$$

$Q_2$ spans the p-dim subspace spanned by the data
as does $\text{col}(\tilde{X})$

$$Q_2 R_2 = \tilde{X} = UDV^T$$

$$R_2 = Q_2^T U DV^T$$

$$Q_2 = U \Rightarrow \qquad R_2 = \overset{\text{diag}}{D}V^T$$

$$\Rightarrow V \text{ is diag w elems } \pm 1 \quad \Rightarrow \text{ so is } R$$
$$\Rightarrow \text{take all } +1 \text{ WLOG}$$

If $\tilde{X}$ has ortho columns

$\Rightarrow \tilde{X} = QR$ has $R$ diag w strictly pos entries

$\Rightarrow$ this is SVD w/ $Q = U$

So $QR = SVD$ when $X$ has orthogonal columns

## 3.9

Claim: $\underset{k}{\text{argmax}} \quad q_k^T r$

$X = Q_1 R_1$

add a predictor $x_k$

$\text{Proj}_{X_1} x_k = \sum_j (x_k^T q_j) q_j$

$r_k = x_k - \text{Proj}_{X_1} x_k$

$q_k = \dfrac{r_k}{|r_k|}$

$\hat{y}_1 \rightarrow \hat{y}_2 = \hat{y}_1 + (q_k^T y) q_k$

$\qquad = \hat{y}_1 + (q_k^T r) q_k$

$\Rightarrow RSS_2 = RSS_1 - (q_k^T r)^2$

$\Rightarrow$ pick max $k$ for $(q_k^T r)^2$

## 3.10  ⭐  Z score is $\dfrac{\beta_k}{\hat{\sigma}\sqrt{v}} = \dfrac{R_p^{-1} Q^T y}{\hat{\sigma} \, R_{pp}^{-1}} = \dfrac{q_p^T y}{\sigma} \Leftarrow$ pick smallest z-score

For least predictor added!

$\text{Var } \beta_j = \dfrac{\sigma^2}{|z_p|^2} =$

$R_{p\cdot}^{-1} = \begin{pmatrix} 0 \\ \vdots \\ R_{pp}^{-1} \end{pmatrix}$

## 3.11  $\text{Tr}\,(Y - XB)^T \, \Sigma^{-1} \, (Y - XB)$

$\Sigma = \mathbb{1}\sigma^2 \Rightarrow B = (X^T X)^{-1} X^T Y$

else $\quad S := \sqrt{\Sigma} \qquad Y \to YS, \quad B \to BS$

$$BS = (X^TX)^{-1} X^T YS \implies B = (X^TX)^{-1} X^T Y$$

No closed formula for general $j$-dependent $\Sigma$
but still solvable by quadratic programming

3.12 $\quad \tilde{X} = \begin{pmatrix} X \\ \sqrt{\lambda}\, \mathbb{1}_{p\times p} \end{pmatrix} \implies \tilde{X}^T \tilde{X} = \left( X^T \ \sqrt{\lambda}\, \mathbb{1}_{p\times p} \right) \begin{pmatrix} X \\ \sqrt{\lambda}\, \mathbb{1}_{p\times p} \end{pmatrix}$

$\qquad Y = \begin{pmatrix} y \\ 0_p \end{pmatrix} \qquad\qquad = X^TX + \lambda \mathbb{1}$

$\qquad \hat{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$

$\qquad = (\tilde{X}^T \tilde{X})^{-1} X^T y$

$\qquad = (X^TX + \lambda \mathbb{1})^{-1} X^T y$

3.13 $\quad z_m = X v_m \implies$ regress $y$ on $z_1 \cdots z_m$

$\qquad\qquad z_1 \cdots z_m$ orthogonal

$\qquad\qquad \implies y^{pcr}_m = \bar{y}\mathbb{1} + \sum_{m=1}^{M} \hat{\theta}_m z_m$

$\qquad\qquad\qquad \hat{\theta}_m = \dfrac{\langle y, z_m \rangle}{\langle z_m, z_m \rangle}$

$\implies \beta^{pcr}(M) = \sum \hat{\theta}_m \vec{v}_m$

$\quad X = UDV^T \qquad\qquad z_m = X\vec{v}_m = d_m \vec{u}_m$

$\quad \hat{\theta}_m = \dfrac{\langle u_m, y \rangle}{d_m} \implies \beta^{pcr} = V \cdot D^{-1} \cdot U^T y = (X^TX)^{-1} X^T y$

3.14 a) $z_m = \sum_j \hat{\varphi}_{mj}\, x_j^{(m-1)} \qquad \hat{\varphi}_{mj} = \langle y, x_j^{m-1} \rangle$

b) $\theta_m = \dfrac{\langle z_m, y \rangle}{\langle z_m, z_m \rangle}$

c) $\hat{y}^m = \hat{y}^{m-1} + z_m \theta_m$

d) $x_j^m = x_j^{m-1} - \dfrac{\langle z_m, x_j^{m-1} \rangle}{\langle z_m, z_m \rangle}$

$$\langle z_1, z_1 \rangle = \sum_{i,j}^{p} \hat{\theta}_{1i} \hat{\theta}_{1j} \langle x_i^0, x_j^0 \rangle$$

$$= \sum_{i=1}^{p} (\hat{\theta}_{ii})^2$$

$$\langle z_1, y \rangle = \sum_{j}^{p} \hat{\theta}_{1j} \langle y, x_j^0 \rangle \qquad \Rightarrow \quad \theta_1 = 1$$

$$= \sum_{i=1}^{p} (\hat{\theta}_{ii})^2 \qquad \Rightarrow \quad \hat{y}^1 = \hat{y}^0 + \sum_{j=1}^{p} \theta_{1j} x_j^0$$

$$x_j' = x_j^0 - \frac{\langle z_1, x_j^0 \rangle}{\langle z_1, z_1 \rangle} z_1$$

$$= x_j^0 - \frac{\theta_{1j}}{\sum \theta_{1i}^2} \sum \theta_{1j} x_i^0$$

$$\theta_{2j} := \langle y, x_j' \rangle = \langle x_0^0, y \rangle - \frac{\theta_{1j}}{\sum \theta_{1i}^2} \sum_{ij} \theta_{ij} \langle x_i^0, y \rangle$$

$$= \langle x_j^0, y \rangle - \theta_{1j} = 0$$

---

**3.15** $\quad Corr^2(y, X\alpha) \, Var(X\alpha) = \dfrac{Cov^2(y, X\alpha)}{Var(y)}$

mth PLS $\hat{\theta}_m$ yields

$$\Rightarrow \max_{\alpha} (y^T X\alpha)^2 \qquad s.t. \quad |\alpha| = 1$$
$$\alpha^T S \hat{\theta}_\ell = 0$$
$$S = X^T X$$

$\alpha^T X^T y y^T X \alpha + \lambda (\alpha^T \alpha - 1)$

$$\Rightarrow \quad (X^T y y^T X)\alpha = \lambda \alpha$$
$$\Rightarrow a = X^T y$$

$\ell = 1 \Rightarrow \alpha_1 = \dfrac{X^T y}{|X^T y|_2} \quad \alpha \, \hat{\theta}_1$

<span style="color:salmon">Not done</span>

$\ell = 2 \Rightarrow \quad \alpha_2 \propto X^T y - \dfrac{y^T X S X^T y}{y^T X S^2 X^T y} S X^T y$

$$\Rightarrow \alpha_2^T S \alpha_1 = 0 \checkmark$$

Continuum regression:

$$\max_\alpha \ (y^T X \alpha)^2 (\alpha X^T X \alpha)^{\frac{r}{1-r}-1}$$

$$\text{s.t.} \ |\alpha| = 1 \qquad \alpha^T S \hat{\alpha}_\ell \quad \ell = 1, \dots, m-1$$

3.16  $X^T X = \mathbb{1} \Rightarrow \hat{\beta}_j = \tilde{x}_j^T y$

$\Rightarrow \hat{\beta}_j \to \dfrac{\hat{\beta}_j}{1+\lambda}$  for ridge

$\hat{\beta}_j \to \hat{\beta}_j \cdot \mathbb{1}(\text{rank } \beta_j \le M)$  for Subsets

$\beta_j$ $\qquad \beta^* = \underset{\beta}{\text{argmin}} \ \frac{1}{2}(\hat{\beta}-\beta)^2 + \lambda|\hat{\beta}|$

$$F = \mathcal{L}'(\beta) = \begin{array}{ll} (\beta-\hat{\beta}) - \lambda & \beta < 0 \\ (\beta-\hat{\beta}) + \lambda & \beta > 0 \end{array}$$

$\Rightarrow \beta = \text{sgn}\,\hat{\beta}\,(\hat{\beta}-|\lambda|)_+$

3-17  Colab

3.18

3.19  $X = UDV \Rightarrow |\beta^{ridge}|^2 = \displaystyle\sum_{j=1}^{p} \frac{d_j^2 (U^T y)_j}{(d_j^2+\lambda)^2}$

For LASSO use dual form

320

Motivation   For CCA:

If   $Y_k = f(X) + \varepsilon_k$   } pool these
$\quad Y_\ell = f(X) + \varepsilon_\ell$

Goal:  Successively  maximize   $\text{Corr}^2(Yu_m, Xv_m)$

ie   $Yu_1$  most correlated  to  $Xv_1$

$X, Y$  centered,  look  @  $\dfrac{Y^T X}{N}$

$$\text{Corr}^2(Yu_m, Xv_m) = \frac{(u_m^T Y^T X v_m)^2}{\text{Var } Xv_m \ \text{Var } Yu_m} = \frac{(u_m^T Y^T X v_m)^2}{v_m^T X^T X v_m \ u_m^T Y^T Y u_m}$$

$$\mathcal{L} = v\, Y^T X v - \frac{\lambda_1}{2}(v^T X^T X v - 1) - \frac{\lambda_2}{2}(u Y^T Y u - 1)$$

$\Rightarrow \quad Y^T X v = \lambda Y^T Y u$

$\qquad X^T Y u = \lambda_2 X^T X v$  } $\Rightarrow \lambda_1 u^T Y^T Y u = \lambda_2 v^T X^T X v$

$\qquad \Rightarrow u Y^T X = \lambda_2 v(X^T X)$ $\qquad \lambda_1 = \lambda_2$

take   $M = (Y^T Y)^{-1/2}(Y^T X)(X^T X)^{-1/2} = \text{Corr}(Y, X)$

$M \cdot \underbrace{(X^T X)^{1/2} v}_{v, *} = \lambda \underbrace{(Y^T Y)^{1/2} u}_{u_1*}$

$u(Y^T Y)^{1/2} M = \lambda(X^T X)^{1/2} v \qquad \Rightarrow \quad \boxed{\begin{array}{l} v = (X^T X)^{-1/2} v_* \\ u = (Y^T Y)^{-1/2} u_* \end{array}}$

$M = u^* D^* V^{*T}$
$\qquad\quad \uparrow \qquad \uparrow$
$\qquad\quad u_* \qquad v_*$  top sing vals  of covariance matrix

For   $k = 2, \cdots, \min(K, p)$

$$\text{max} \quad u^T Y^T X v \Rightarrow \mathcal{L} = u^T Y^T X v - \frac{\lambda_1}{2}(v^T X^T X v - 1) - \frac{\lambda_2}{2}(u Y^T Y u - 1)$$

$u^T Y^T Y u = 1$
$v^T X^T X v = 1$
$u^T_i u_j = 0$
$v^T_i v_j = 0$

$$- \sum_{j<k} \alpha_j u^T u_j - \sum_{j<k} \beta_j v^T v_j$$

$$\Rightarrow \quad Y^T X v - \lambda_1 Y^T Y u - \sum_j \alpha_j u_j$$

$$X^T Y u - \lambda_2 X^T X v - \sum_j \beta_j v_j$$

$$\Rightarrow \quad u^T Y^T X v - \lambda_1 u^T Y^T Y u = 0 \quad \Rightarrow \quad \textcolor{salmon}{\text{same eqs}}$$
$$u^T Y^T X v - \lambda_1 v^T X^T X v = 0 \quad \textcolor{salmon}{\Rightarrow \text{subsequent sing vals}}$$
$$\textcolor{salmon}{\text{of } Corr(Y,X)}$$

3.21 $$B^{rr} = \underset{\text{rank } B = m}{\text{argmin}} \; \text{Tr}\left[(Y - B^T X)\left(\frac{Y^T Y}{N}\right)^{-1}(Y - BX)^T\right]$$

$$Y \to Y^* = Y \Sigma^{-1/2} \Rightarrow \text{Tr}\left[(Y^* - B^T X \Sigma^{-1/2})(Y^* - B^T X \Sigma^{-1/2})^T\right]$$

$$\Rightarrow \text{Tr}\left[\underset{\textcolor{salmon}{\text{indep}}}{\textcolor{salmon}{\cancel{Y^* Y^{*T}}}} + \Sigma^{-1/2} B^T X^T X B \Sigma^{-1/2} - 2\Sigma^{-1/2} Y^T X B \Sigma^{-1/2}\right]$$

$$\textcolor{salmon}{\underbrace{\qquad}_{\text{complete square}}}$$

$$\underset{B}{\text{min}} \; \left\| \Sigma^{-1/2} B^T (X^T X)^{1/2} - \Sigma^{-1/2} Y^T X (X^T X)^{-1/2} \right\|$$

<span style="color:salmon">Eckhardt - Young - Mirsky</span>

$$\underset{\text{rk } \hat{D} = r}{\text{min}} \; \| \hat{D} - D \|_F = \sum_{i=1}^{r} \sigma_i u_i v_i^T$$

$$\text{in} \quad SVD: \quad D = U \Sigma V^T \qquad \textcolor{salmon}{K \times M}$$

<span style="color:salmon">let</span> $\textcolor{salmon}{\hat{D} = \Sigma^{-1/2} B^T (X^T X)^{1/2}} \quad \textcolor{salmon}{U D V^T =: \Sigma^{-1/2} Y^T X (X^T X)^{-1/2}}$

$$\textcolor{salmon}{\Rightarrow} \quad \textcolor{salmon}{\Sigma^{-1/2} \hat{B}^T (X^T X)^{1/2} = \sum_{i=1}^{m} d_i v_i v_i^T = U_m D V^T}$$

$$\textcolor{salmon}{= U U_m^T U D V}$$

$$\textcolor{salmon}{\Rightarrow \quad \hat{B}_m = (X^T X)^{-1/2}\left(\sum_{i=1}^{m} d_i v_i u_i^T\right) \Sigma^{1/2}}$$

$$\textcolor{salmon}{(X^T X)^{-1}(X^T Y)\Sigma^{-1/2} \underbrace{U_m}_{\tilde{U}_m} \underbrace{U_m^T \Sigma^{1/2}}_{\tilde{U}_m^-}}$$

$$Y \Sigma^{-1/2} \to Y^* $$
$$\Rightarrow \min_{B_m} | B^T (X^TX)^{1/2} - Y^{*T} X (X^TX)^{-1/2} |^2 = B \, \tilde{U}_m \tilde{U}_m^-$$

Full rk

$$\underbrace{U \Sigma V^T}_{K \cdot m}$$

$$B^T (X^TX)^{1/2} = U_m U_m^T \left( Y^{*T} X (X^TX)^{-1/2} \right)$$

$$\Rightarrow B = \underbrace{(X^TX)^{-1} X^T Y^*}_{\text{Full rank}} \underbrace{U_m U_m^T}_{\text{Proj}} \qquad U_m = \Sigma^{-1/2} U_m^*$$
$$U_m^- = U_m^{*T} \Sigma^{1/2}$$

**3.22**   Prev was true $\forall \Sigma$

**3.23**   $\frac{1}{N} |\langle x_j, y \rangle| = \lambda \qquad j = 1, \cdots, p$

a)   take $\hat\beta = (X^TX)^{-1} X^T y = (X^TX)^{-1} \lambda \frac{1}{N}$

$$u = \alpha X \hat\beta \quad \Rightarrow \quad u \text{ moves a fraction } \alpha \text{ of the LS fit}$$

$$\Rightarrow \quad \frac{1}{N} | \langle x_j, \, y - \alpha X (X^TX)^{-1} X^T y \rangle |$$

$$= | \lambda - \alpha \langle x^T X (X^TX)^{-1} \rangle \underbrace{\langle x_j y \rangle}_{N} |$$

$$= \lambda (1 - \alpha)$$

b)   $\text{Corr} = \dfrac{\frac{1}{N} \langle x_j, y - u_\alpha \rangle}{\sqrt{\dfrac{\langle x_j, x_j \rangle}{N}} \sqrt{\dfrac{\langle y - u_\alpha, y - u_\alpha \rangle}{N}}}$

$$u_\alpha = \alpha X (X^TX)^{-1} X^T y$$
$$y - u_\alpha = (1-\alpha) y + \alpha (y - \hat{y})$$

$(1-\alpha)^2 \langle y, y \rangle + \alpha^2 RSS$

$$\Rightarrow \langle y - u_\alpha, y - u_\alpha \rangle = y^T y + \alpha^2 y^T X (X^TX)^{-1} X^T y$$
$$- 2\alpha y^T X (X^TX)^{-1} X^T y$$

$$\langle y - \alpha \hat{y}, \, y - \alpha \hat{y} \rangle$$

Needed to

$$= y^T y + 2\alpha (\alpha - 2) \, y^T \hat{y}$$

$$= N + \quad \alpha(2-\alpha) \; RSS + \alpha(\alpha-2) \; N$$

$$= (\alpha-1)^2 N + \frac{\alpha(2-\alpha)}{2} RSS$$

$$\Rightarrow \quad Corr = \frac{\lambda(1-\alpha)}{\sqrt{(1-\alpha)^2 + \dfrac{\alpha(2-\alpha)\;RSS}{N}}}$$

$$\alpha = 0 \Rightarrow \quad Corr = \lambda$$
$$\alpha = 1 \Rightarrow \quad Corr = 0$$

<span style="color:red">$\Big\}$ always tied & decreasing</span>

**3.24** $\quad \delta_k = (X_{A_k}^T X_{A_k})^{-1} X_{A_k}^T r_k$ <span style="color:red">$\in \mathbb{R}^{p^-} \Leftarrow \# \text{ active}$</span>

<span style="color:red">$N \times p^-$</span> $\qquad$ <span style="color:#6aa3e0">residual before $k$ is added</span>

$$\beta_{A_k}(\alpha) = \beta_{A_k}(0) + \alpha \delta_k \quad \Rightarrow \quad \hat{f}_k(\alpha) = \hat{f}_{A_k}(0) + \alpha X_{A_k} \delta_k$$

<span style="color:red">$\underbrace{\qquad}_{u_k}$</span>

$$\Rightarrow \quad \text{direction is} \quad u_k = X_{A_k} \delta_k \quad \in \mathbb{R}^N$$

<span style="color:red">Claim is $u_k$ makes smallest & equal angle w/ each col in $X_{A_k}$</span>

$$\Rightarrow \quad X_{A_k}^T u_k = X_{A_k}^T r_k$$

<span style="color:red">$\underbrace{\qquad}$</span>
<span style="color:red">$\vec{x}_q \cdot r_k = \vec{x}_j \cdot r_k$ by assumption of when we added $x_k$</span>

<span style="color:#6aa3e0">also all $r_k \notin A_k$ have $\ge$ corr by assumption</span>

**3.25** $\quad \hat{f}_k(\alpha) = \hat{f}_k(0) + \alpha u_k \qquad u_k = X_{A_k} \delta_k$

$$\left| c_a(\alpha) \right| = \left| x_a^T (y - \hat{f}_k(\alpha)) \right|$$

<span style="color:red">$a \in A_k$</span>
$$= x_a^T (r_k - \alpha u_k)$$
$$= x_a^T r_k - \alpha \, x_a^T u_k \quad$$ <span style="color:red">$\Leftarrow$ by 3.24 this is the same $\forall_q \in A_k$</span>

$$= \hat{C} - \alpha A$$

For $b \notin A_k$
$$|x_b^T r_k| \leq |x_j^T r_k|$$

$\Rightarrow$ for small $\alpha$
$$|c_b(\alpha)| < |c_j(\alpha)|$$
until some $\alpha^*$

Pick $b = \text{argmax}_b |c_b(\alpha)|$. Def $k(\alpha) = \max_{b \in A_k} |c_b(\alpha)|$

$$|c(\alpha^*)| = |x_b^T r_k - \alpha \, x_b^T u_k| = |x_a^T r_k - \alpha^* x_a^T u_k| = |c_a(\alpha^*)|$$

$$\Rightarrow \frac{(x_b \mp x_a) r_k}{(x_b \mp x_a)^T u_k} = \frac{x_b r_k \mp C}{x_b u_k \mp A} = \alpha^*$$

$$\Rightarrow \alpha^* = \min_{b \notin A_k} \left\{ \frac{C + x_b^T r_k}{A + x_b^T u_k}, \frac{C - x_b^T r_k}{A - x_b^T u_k} \right\}$$

$$b = \text{argmin} \quad \gamma$$

**3.26**    From 3.9    For fwd stepwise, at each step we choose $k$ with

$$r_k = x_k - \text{Proj}(x_k), \quad q_k = \frac{r_k}{|r_k|}$$

$$\hat{y} \rightarrow \hat{y} + (q_k^T y) q_k = \hat{y} + (q_k^T r) q_k$$
$y - \hat{y}$

$$y - \hat{y} = r \quad \Rightarrow \quad r \rightarrow (I - q_k q_k^T) r$$

$$\text{RSS} \rightarrow (q_k^T r)^2 \quad \text{reduction}$$

$\Rightarrow$ $j$ for which $\text{Corr}(x_{j, u_{A_c}}, r)$ is largest in magnitude

**3.27 a)** $\mathcal{L} = \mathcal{L}(\beta) + \lambda \sum_j |\beta_j|$

$$\beta_j = \beta_j^+ - \beta_j^- \qquad \beta_j^+, \beta_j^- \geq 0$$

$$\Rightarrow L = \mathcal{L}(\beta) + \lambda \sum_j (\beta_j^+ \mp \beta_j^-) - \sum_j (\lambda_j^+ \beta_j^+ + \lambda_j^- \beta_j^-)$$

$$\nabla \mathcal{L}_j(\beta) + \lambda - \lambda^+ = 0$$
$$-\nabla \mathcal{L}_j(\beta) + \lambda - \lambda^- = 0$$

**b)** $\quad \lambda^+ + \lambda^- = 2\lambda \geq 0$

$\quad \nabla 2 = \frac{1}{2}(\lambda^- - \lambda^+)$

$\quad \Rightarrow |\nabla 2| = \frac{1}{2}|\lambda^- - \lambda^+| \leq \frac{1}{2}(\lambda^- + \lambda^+) = \lambda$

$\quad \Rightarrow \lambda = 0 \Rightarrow \nabla 2_j = 0$

**KKT** $\quad \beta_j^+ > 0, \ \lambda > 0 \Rightarrow \lambda^+ = 0 \quad \nabla 2 = -\lambda < 0 \quad \beta_j^- = 0$

$\quad \beta_j^- > 0, \ \lambda > 0 \Rightarrow \lambda^- = 0 \quad \nabla 2 = \lambda > 0 \quad \beta_j^+ = 0$

$\quad \Rightarrow \quad \nabla 2 = -\lambda \ \text{sign}(\beta_j)$

$\quad \Rightarrow \quad \lambda = x_j^T(y - X\beta) = x_j^T r$

$\quad \underbrace{\qquad\qquad\qquad}$
$\quad$ all active preds $\ (\beta_j \neq 0)$ have same corr $= \lambda$

**c)** $\quad X^T(y - X\hat{\beta}(\lambda)) = \theta(\lambda)$

$$\theta(\lambda)_j = \begin{cases} -\lambda \ \text{sgn}(\beta_j) & \beta_j \neq 0 \quad x_j \in S \\ x_j^T y & \beta_j = 0 \quad \text{ie } j \notin S \end{cases}$$

$\hat{\beta}_j(\lambda) = (X^T X)^{-1}[X^T y - \theta(\lambda)]$

$\hat{\beta}(\lambda) - \hat{\beta}(\lambda_0) = (X^T X)^{-1}(\theta(\lambda_0) - \theta(\lambda))$

$$\begin{cases} (\lambda - \lambda_0) \ \text{sgn} \ \beta_j(\lambda_0) & j \in S \\ 0 & j \notin S \end{cases}$$

$\quad \underbrace{\qquad\qquad\qquad}$
$\quad\quad\quad$ linear in $\lambda$

**328** $\quad \hat{\beta}^{new} = \underset{\beta}{\text{argmin}} \ \|y - X\beta - x_j^T \beta_j^*\|$

$\quad\quad\quad\quad \text{s.t.} \quad |\beta|_1 + |\beta_j^*| \leq t$

$\quad\quad\quad \tilde{\beta}_j = \beta_j + \beta_j^*, \quad \tilde{\beta} \ \text{has} \ \beta_j \rightarrow \tilde{\beta}_j$

$$\implies \underset{\beta}{\text{argmin}} \ \| y - X\tilde{\beta} \|$$

$$\text{s.t.} \ |\tilde{\beta}|_1 + \underbrace{(|\beta_j| + |\beta_j^*|) - (\tilde{\beta}_j|)}_{\geq 0} \leq t$$

$\implies$ *More stringent*

By symmetry, set $\beta_j = \beta_j^* = \dfrac{\beta^{orig}}{2} \implies |\tilde{\beta}|_1 = |\beta^{orig}|_1$

$\&$ objective is the same

$\implies$ sol'n

$\implies$ *Final effect reduces $\beta_j, \beta_j^*$ by $\frac{1}{2}$*

**3.29**

$$\beta = \frac{X^T y}{X^T X + \lambda} \qquad\qquad X \in \mathbb{R}^{N,1}$$

$$\tilde{X} = (X, X) \in \mathbb{R}^{N,2} \implies \beta = (\tilde{X}^T \tilde{X} + \lambda \tilde{I})^{-1} \tilde{X}^T y$$

$$\tilde{\lambda} = \lambda I_{2\times 2} \qquad\qquad \implies \beta_1 = \beta_2$$

$$\begin{pmatrix} X^T \\ X^T \end{pmatrix} (X \ X) = \begin{pmatrix} X^T X & X^T X \\ X^T X & X^T X \end{pmatrix} = \tilde{X}^T \tilde{X}$$

$$\| y - 2X\beta \|^2 + 2\lambda \|\beta\|^2 \implies \beta = \frac{2 X^T y}{4 X^T X + 2\lambda}$$

$$= \frac{X^T y}{2 X^T X + \lambda}$$

for $m$ copies

$$\frac{X^T y}{m X^T X + \lambda}$$

3.30 $\quad \hat{X} = \begin{pmatrix} X \\ \gamma \mathbf{1}_{p+1} \end{pmatrix} \in \mathbb{R}^{N+p+1, \, p+1} \quad \tilde{y} = \begin{pmatrix} Y \\ 0 \end{pmatrix} \in \mathbb{R}^{N+p+1, \, 1}$

$$\gamma = \sqrt{\lambda \alpha}$$

$\Rightarrow$ Lasso for

$$\|\tilde{y} - X\beta\|^2 + \lambda(1-\alpha)\|\beta\|_1$$

## Chapter 4

4.1 $\qquad \max \quad a^T B a$
$\qquad a^T W a = 1$

$$\mathcal{Z} = a^T B a - \lambda a^T W a$$

$$\Rightarrow \quad B a = \lambda W a \quad \Rightarrow \quad W^{-1} B a = \lambda a$$

4.2 a) likelihood of data from 2 Gaussians w $\quad \gamma_1 = \frac{N_1}{N}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \gamma_2 = \frac{N_2}{N}$

$$(x-\mu_1)^T \Sigma (x-\mu_1) - (x-\mu_2)^T \Sigma (x-\mu_2) + \log \frac{N_1}{N_2}$$

$$x^T \Sigma(\mu_2-\mu_1) - (\mu_2+\mu_1)\Sigma^{-1}(\mu_2-\mu_1) + \log \frac{N_1}{N_2}$$

$$x^T \Sigma^{-1}(\mu_2-\mu_1) = (\mu_2+\mu_1)\Sigma^{-1}(\mu_2-\mu_1) - \log \frac{N_2}{N_1}$$

b) $\qquad X = \begin{pmatrix} 1 & \tilde{x}_i \end{pmatrix} \Rightarrow X^T X = \begin{pmatrix} N & N\bar{x} \\ N\bar{x} & \sum_i x_i x_i^T \end{pmatrix}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \bar{x} \in \mathbb{R}^p$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad x_i \in \mathbb{R}^p$

$\qquad X^T X \begin{pmatrix} \beta_0 \\ \beta \end{pmatrix} = X^T Y$

$\qquad \Rightarrow N\beta_0 + N\bar{x}\cdot\beta = N\bar{y} \Rightarrow \beta_0 = \bar{y} - \bar{x}\cdot\beta$

$\qquad N\beta_0 \bar{x} + \sum_i x_i x_i^T \cdot \beta = \sum y_i x_i$

$\qquad \left(\frac{1}{N}\sum_i x_i x_i^T - \bar{x}\,\bar{x}^T\right)\cdot\beta = \frac{1}{N}\sum y_i x_i - \bar{y}\bar{x}$

Take $\mu_1 = \frac{1}{N_1} \sum_{G_1} x_i$ $\mu_2 = \frac{1}{N_2} \sum_{G_2} x_i$

$$y \in \left\{ \underset{G_1}{\frac{-N}{N_1}} , \underset{G_2}{\frac{N}{N_2}} \right\} \qquad \Rightarrow \bar{y} = 0 \qquad \bar{x} = \mu_1 + \mu_2$$

$$\sum x_i y_i = -N\mu_1 + N\mu_2$$

$\Rightarrow \quad \beta_0 = (-N + N)^0 - (N_1 \underset{2}{\mu_1} + N_2 \mu_2) \cdot \beta$

$* \Rightarrow \quad \left( \sum x_i x_i^T - N \bar{x} \bar{x}^T \right) \cdot \beta = N(\mu_2 - \mu_1)$

$\hat{\Sigma}$ is estimate of $Var(X|k)$

$$(N-2) \hat{\Sigma} = \left( \sum_{G_1} (x_i - \mu_1)(x_i - \mu_1) + \sum_{G_2} (x_j - \mu_2)(x_j - \mu_2) \right)$$

$$= \sum x_i x_i^T - N_1 \mu_1 \mu_1^T - N_2 \mu_2 \mu_2^T$$

$\Rightarrow \sum x_i x_i^T - N\bar{x}\bar{x}^T = (N-2)\hat{\Sigma} + N_1 \mu_1 \mu_1^T + N_2 \mu_2 \mu_2^T - \frac{1}{N}(N_1\mu_1 + N_2\mu_2)(N_1\mu_1 + N_2\mu_2)^T$

$= (N-2)\hat{\Sigma} + \underbrace{\frac{N_1 N_2}{N}(\mu_1 - \mu_2)(\mu_1 - \mu_2)^T}$

$N \Sigma_B$ , $\Sigma_B = \frac{N_1}{N}\frac{N_2}{N}(\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$

$*$ is as desired

c) $\hat{\Sigma}_B \beta \propto \mu_2 - \mu_1$ by its structure

$\Rightarrow \quad \hat{\Sigma} \beta \propto \mu_2 - \mu_1$

$\Rightarrow \quad \beta \propto \underbrace{\hat{\Sigma}^{-1}(\mu_2 - \mu_1)}_{\text{LDA coeff}}$

d) Replacing $y \in \{t_1, t_2\}$

$\beta_0 = \frac{1}{N}(N_1 t_1 + N_2 t_2) - \mu \cdot \beta$ $\qquad \mu = \mu_1 + \mu_2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \Delta = \mu_1 - \mu_2$

$\left[ (N-2)\hat{\Sigma} + N\hat{\Sigma}_B \right] \cdot \beta = (N_1 t_1 \mu_1 + N_2 t_2 \mu_2) - \frac{1}{N}(N_1 t_1 + N_2 t_2)(N_1 \mu_1 + N_2 \mu_2)$

$\qquad\qquad\qquad\qquad = \frac{N_1 N_2}{N}(t_1 - t_2)(\mu_1 - \mu_2)$

e) $\quad \beta_0 = -\mu \cdot \beta$

$\Rightarrow \quad \hat{f}(x) = (x - \mu) \cdot \beta$

$\quad (x - \mu) \Sigma^{-1} (\mu_2 - \mu_1) = 0 \qquad$ is decision bdy

$\quad \cong$ to LDA when $N_1 = N_2$

else, we have extra $\log \frac{N_1}{N_2}$ const in LDA

**4.3** $\quad \pi^{new} = \pi^{old}$

$\mu^{new} = B^T \mu^{old}$

$\hat{\Sigma}^{new} = B^T \hat{\Sigma}^{old} B \qquad \Rightarrow \quad (\Sigma^{new})^{-1} = B^{-1} (\Sigma^{old})^{-1} (B^T)^{-1}$

$\quad (x - \mu^{new})^T (\Sigma^{new})^{-1} (\mu_2^{new} - \mu_1^{new})$

$= (x - \mu^{old})^T B B^{-1} (\Sigma^{-1})^{old} (B^T)^{-1} B^T (\mu_2^{old} - \mu_1^{old})$

$= (x - \mu^{old}) (\Sigma^{-1})^{old} (\mu_2^{old} - \mu_1^{old})$

**4.4** $\quad \beta \in \mathbb{R}^{(p+1), (K-1)} \qquad\qquad \beta \in \binom{\beta_0}{\beta}, \quad x^T \in \binom{1}{x^T}$

$P(G = k \mid X = x) = \dfrac{\exp(\beta_k^T x)}{1 + \exp(\beta_k^T x)} \qquad\qquad k = 1 \cdots, K-1$

$\hat{P}(G = K \mid X = x) = \dfrac{1}{1 + \exp(\beta_k^T x)}$

$\ell(\beta) = \sum_{i=1}^{N} \log P(g_i \mid x_i ; \beta)$

$= \sum_{j=1}^{N} \sum_{k=1}^{K-1} \mathbb{1}(y_i = k) \beta_k^T x_{ik} - \log \left( 1 + \sum_{i} \exp \beta_k^T x_i \right)$

$\dfrac{\partial \ell}{\partial \beta} = \sum_{i} \left[ \mathbb{1}(y_i = k) - \dfrac{e^{\beta_k x_i}}{1 + \sum \exp(\beta_k x_i)} \right] \vec{x_i}$

$$\frac{\partial^2 \ell}{\partial \beta_k \, \partial \beta_\ell} = \sum_{i=1}^{N} p_k(x_i;\beta) \, p_\ell(x_i;\beta) \; x_{ik} \; x_{i\ell} \qquad k \neq \ell$$

$$= -\sum_i p_k(1-p_k) \; x_{ik} \, x_{ik} \qquad k = \ell$$

$$= -\sum_i \left[ diag(p_k(x_i;\beta)) - p_k(x_i;\beta)(1-p_\ell(x_i;\beta)) \right] \vec{x_i} \vec{x_i}^T$$

$$\beta^{new} = \beta^{old} - H^{-1} g$$

$$g = X^T(y_k - p_k) \quad \in \mathbb{R}^{p+1, \, k-1}$$

$$= \begin{pmatrix} X^T & & \\ & X^T & \\ & & \ddots \end{pmatrix} \begin{pmatrix} y_i - p_i \\ \vdots \\ y_{k-1} - p_{k-1} \end{pmatrix}$$

$$H = -X^T W X$$

$$W = \begin{pmatrix} P_1 & R_1 R_2 \cdots R_1 R_{k-1} \\ R_2 R_1 & \ddots & \vdots \\ \vdots & & \vdots \\ R_{k-1} R_1 & \cdots & P_{k-1} \end{pmatrix} \qquad [R_k]_{ij} = \left[ diag\,\overset{M \times N}{p_k(x_i;\beta)} \right]_{ij}$$

$$[P_k]_{ij} = \left[ diag\,\overset{M \times N}{p_k(x_i;\beta)(1-p_k(x_i;\beta))} \right]$$

$$\beta^{new} = \beta^{old} + (X^T W X)^{-1} X^T (y-p)$$

$$= (X^T W X)^{-1} (X^T W) \left[ X\beta^{old} + W^{-1}(y-p) \right]$$

$$\underbrace{\phantom{X\beta^{old} + W^{-1}(y-p)}}$$

z    is    our    "target"

**4.5** log likelihood:

$$l(\beta) = \sum_i \left[ \beta_i x_i y_i - \log(1 + e^{\beta_i x_i}) \right]$$

$$= \sum_i \left[ (\beta_0 + \beta_1 x_i) y_i - \log Z \right]$$

$$= \sum_i \left[ (\beta_0 + \beta_1 x_0 + \beta_1(x_i - x_0)) y_i - \log Z \right]$$

choose $\beta_0$ at this is 0

$$\Rightarrow \sum_i \left[ y_i \beta(x_i - x_0) - \log(1 + e^{\beta(x_i - x_0)}) \right]$$

$$\Rightarrow \sum_{i \in N_1} \beta(x_i - x_0) - \log(1 + e^{\beta(x_i - x_0)}) - \sum_{i \in N_2} \log(1 + e^{\beta(x_i - x_0)})$$

const $\to 0$        arb neg

a) Same story but w/ plane

$$l(\beta) \to \infty$$

b) 

$$l(\beta) = \sum_{j=1}^{N} \left[ \sum_{k=1}^{K-1} \mathbb{1}(y_j = k) \beta_k \cdot x_j - \log\left( 1 + \sum_{\ell=1}^{K-1} e^{\beta_\ell \cdot x_j} \right) \right]$$

$$= \sum_k \sum_{i \in S_k} \left[ \beta_k \cdot x_i - \log\left( 1 + \sum_{\ell=1}^{K-1} e^{\beta_\ell \cdot x_i} \right) \right]$$

$$+ \sum_{i \in S_K} \left[ - \log\left( 1 + \sum_{\ell=1}^{K-1} e^{\beta_\ell \cdot x_i} \right) \right]$$

$$\exists \beta_k \quad \text{for } k = 1, \dots K-1$$

$$\text{s.t.} \quad \beta_k \cdot x > 0 \quad \forall x \in S_k$$

$$\Rightarrow l(\beta) \to \infty$$

**4.6**  **a)** $\exists \beta$  s.t.  $\begin{array}{l}\beta^T x_i > 0 \quad \text{if } y_i = 1 \\ \beta^T x_i < 0 \quad \text{if } y_i = -1\end{array}$  $\Rightarrow y_i \beta^T x_i > 0$

$$\Rightarrow y_i \beta^T z_i > 0 \qquad z_i = \frac{x_i}{|x_i|}$$

if  $m = \min y_i \beta^T z_i$  then  $\frac{1}{m} y_i \beta^T z_i \geq 1 \quad \forall i$

set  $\beta \to \frac{\beta}{m}$

**b)** $\| \beta_{new} - \beta_{sep} \|^2 = \| \beta_{old} - \beta_{sep} + y_i z_i \|^2$

$$= \| \beta_{old} - \beta_{sep} \|^2 + \| y_i z_i \|^2 + 2 y_i (\beta_{old} - \beta_{sep})^T z_i$$

$$= \| \beta_{old} - \beta_{sep} \|^2 + 1 + \underbrace{2 y_i \beta_{old} z_i}_{\substack{<0 \\ \text{since } z_i \\ \text{was misclassified} \\ \text{before}}} - \underbrace{2 y_i \beta_{sep} z_i}_{-2}$$

$$\leq \| \beta_{old} - \beta_{sep} \|^2 - 1$$

$$\implies \leq \| \beta_{start} - \beta_{sep} \|^2 \quad \text{steps}$$

**4.7**  $D(\beta, \beta_0) = - \sum_i \underbrace{y_i (x_i^T \beta + \beta_0)}_{\substack{\text{signed dist} \\ \text{to hyperplane}}}$  $\qquad \| \beta \| = 1$

No, because no need for optimal sep
in case of class imbalance

**4.8**  $\ell(\mu, \Sigma) = -\frac{1}{2} \sum_{\ell=1}^{K} \sum_{g(i)=\ell} (x_i - \mu_\ell)^T \Sigma^{-1} (x_i - \mu_\ell) - N \log |\Sigma|$

**4.9**  See  Colab

Test error $= Err_T = \mathbb{E}[L(Y, \hat{f}(X)) \mid T]$
/generalization error
$\quad X, Y \sim D_{test}$

Expected prediction error
/ie Expected test error $\quad Err = \mathbb{E}_T Err_T = \mathbb{E}_{X,Y} L(X, \hat{f}(X))$

↑ easier

$$\overline{err} = \frac{1}{N} \sum_i L(y_i, \hat{f}(x_i))$$

↗ bad estimate of test

take $T = \{(x_1, y_1), \cdots, (x_N, y_N)\}$

$$Err_T = \mathbb{E}_{X^0, Y^0}\left[L(Y^0, \hat{f}(X^0)) \mid T\right]$$

if $X^0 \neq x_1 \cdots x_n \Rightarrow$ "extra sample"

in-sample takes $x_1 \cdots x_N$ & new responses $Y_1^0 \cdots Y_N^0$

$$Err_{in} = \frac{1}{N} \sum_i \mathbb{E}_{Y^0}\left[L(Y_i^0, \hat{f}(x_i)) \mid T\right]$$

$$op = Err_{in} - \overline{err}$$

$$\mathbb{E}_Y \; op =: \omega$$

take $L = |\cdot|^2 \Rightarrow Err_{in} - \overline{err} = \frac{1}{N} \sum_i \mathbb{E}_{Y_i^0}\left[(Y_i^0)^2\right] - \mathbb{E}_Y (Y^2)$

$$- 2 \mathbb{E}_{Y_i^0}[Y_i^0] \, \mathbb{E}_{Y_i}[\hat{y}_i]$$

$=\mathbb{E}_Y[Y]$

$$+ 2 \mathbb{E}_{Y_i} Y_i \hat{y}_i$$

$$= \frac{2}{N} \sum_i cov(y_i, \hat{y}_i)$$

**7.1** For lin reg

$$\sum_i \text{Cov}(\hat{y}_i, y_i) = \underset{Y}{E}\left[ y^T X^T (X^T X)^{-1} X y^T \right]$$

$$= Tr\left[ H \; \text{Cov}(y, y) \right]$$
$$\underbrace{\qquad}_{\sigma_\varepsilon^2 \, \delta_{ij}}$$

$$= Tr \, H \; \sigma_\varepsilon^2$$

$$= d \sigma_\varepsilon^2$$

$$\Rightarrow \frac{2}{N}\sum_i \text{Cov}(\hat{y}_i, y_i) = \frac{2d}{N}\sigma_\varepsilon^2$$

$$\Rightarrow E_y \, \overline{\text{Err}}_{in} = \underset{Y}{E} \, \overline{\text{err}} + \frac{2d}{N}\sigma_\varepsilon^2 \qquad \leftarrow AIC$$

$$AIC: \; -2\,\mathbb{E}\, \log p_\theta(Y) = -2\,\mathbb{E}\, \log lik + \frac{2d}{N}$$

**7.2** $\quad Pr(Y=1 \mid x_0) = f(x_0) \qquad \hat{G} = \mathbb{1}\left[ f(x_0) > \tfrac{1}{2} \right]$

First let $G = 1 \Rightarrow f(x_0) > \tfrac{1}{2}$

$$\text{Err}(x_0) = Pr\left( Y \neq \hat{G}(x_0) \mid X = x_0 \right)$$

$$= Pr\left( Y = 1 \mid x_0 = x_0 \right) Pr\left( \hat{G} = 0 \mid X = x_0 \right)$$

$$= f(x_0) \, Pr\left( \hat{G} = 0 \mid X = x_0 \right) + (1 - f(x_0))\left( 1 - Pr(\hat{G} = 0 \mid X = x_0) \right)$$

$$= 1 - f(x_0) + (2 f_0 - 1) \, Pr\left( \hat{G} = 0 \mid X = x_0 \right)$$

$$= \text{Err}_{Bayes} + |2 f_0 - 1| \, Pr\left( \hat{G} \neq G \mid X = x_0 \right)$$

<span style="color:salmon">general form</span> $\longrightarrow$

<span style="color:salmon">take again $f > \tfrac{1}{2}$</span>
$$Pr\left( G \neq \hat{G} \mid X = x_0 \right) = P\left( \hat{G} = 0 \mid X = x_0 \right)$$
$$= P\left( f(x_0) < \tfrac{1}{2} \right)$$

$$= Pr\left[\frac{\hat{f}_{(x_0)} - E\hat{f}(x_0)}{\sqrt{Var(x_0)}} < \frac{\frac{1}{2} - E\hat{f}(x_0)}{\sqrt{Var\hat{f}(x_0)}}\right]$$

$$= \Phi\left[\frac{|\frac{1}{2} - E\hat{f}(x_0)|}{\sqrt{Var\hat{f}(x_0)}} \cdot sgn(\frac{1}{2} - EF)\right]$$

<span style="color:pink">bby bias</span>

<span style="color:pink">true pred</span>

<span style="color:pink">if $\hat{f}$ is on the wrong side its better to increase the var</span>

**7.3** $\quad \hat{f} = Sy$

a) $\quad S_{ij} = x_i^T(X^TX + \lambda\Omega)^{-1}x_j^T$

$$\hat{f}(x_i) = x_i^T(X^TX + \lambda\Omega)^{-1}X^Ty$$

$$\Rightarrow \hat{f}^{-i}(x_i) = x_i(X_{-i}^T X_{-i} + \lambda\Omega)^{-1}X_{-i}^Ty_{-i}$$

$$= x_i(X^TX - x_ix_i^T + \lambda\Omega)^{-1}(X^Ty - x_iy_i)$$

$$\underbrace{(X^TX + \lambda\Omega)^{-1}}_{\color{pink}A} + \color{pink}\frac{A^{-1}x_ix_i^TA^{-1}}{1 - x_i^TA^{-1}x_i}$$

$$\hat{f}^{-i}(x_i) = \hat{f}(x_i) - y_iS_{ii} - \frac{x_iA^{-1}x_ix_i^TA^{-1}x^Ty}{1 - x_i^TA^{-1}x_i} \overset{\color{pink}S_{ii}T}{+} \frac{x_i^TA^{-1}x_ix_i^TA^{-1}x_iy_i}{1 - x_i^TA^{-1}x_i}\color{pink}S_{ii}$$

$$= \hat{f}(x_i) - y_iS_{ii} + \frac{S_{ii}\hat{f}(x_i)}{1 - S_{ii}} - \frac{S_{ii}^2y_i}{1 - S_{ii}}$$

$$= \frac{\hat{f}(x_i)}{1 - S_{ii}} - \frac{y_iS_{ii}}{1 - S_{ii}}$$

$$= \frac{\hat{f} - y_iS_{ii}}{1 - S_{ii}}$$

$$\Rightarrow y^i - \hat{f}^{-i}(x_i) = \frac{y_i - \hat{f}(x_i)}{1 - S_{ii}}$$

b) $\quad |y^i - \hat{f}^{-i}(x_i)| = \frac{|y_i - \hat{f}(x_i)|}{1 - S_{ii}}$

$$> |y_i - \hat{f}(x_i)|$$

$$S = (X^TX + \lambda\Omega)^{-1}$$

$$S^2 = S \Rightarrow S_{ii}^2 = \sum_{k \neq i} S_{ik}^2 + S_{ii}^2$$

$$\leq S_{ii}$$

$$\Rightarrow 0 \leq S_{ii} \leq 1$$

c) Replace $y_i$ with $\hat{f}^{-i}(x) =: y'$

$$\Rightarrow \hat{f}^{-i} = Sy'$$

$$= \sum_{j \neq i} S_{ij} y_j + S_{ii} \hat{f}^{-i}$$

$$= \hat{f}(x_i) - S_{ii} y_i + S_{ii} f^{-i}$$

$$\Rightarrow y^i - \hat{f}^{-i}(x_i) = \frac{y - f(x_i)}{1 - S_{ii}}$$

## 7.4

take $L = |\cdot|^2 \Rightarrow Err_{in} - \overline{err} = \frac{1}{N} \sum_i \underset{Y_i^0}{\mathbb{E}}[(Y_i^0)^2]^{=0} - \underset{Y}{\mathbb{E}}(Y^2)^{=0}$

$$-2 \underset{Y_i^0}{\mathbb{E}}[Y_i^0] \underset{Y_i}{\mathbb{E}}[\hat{y}]$$

$= \mathbb{E}[Y]$

$$+ 2 \underset{Y_i}{\mathbb{E}} Y_i \hat{y}_i$$

$$= \frac{2}{N} \sum_i cov(y_i, \hat{y}_i)$$

## 7.5

$$\sum_i Cov(y_i, S_{ij} y_j)$$

$$= Tr[S_{ij} cov Y_i y_j]$$

$$= Tr[S_{ij}] \sigma^2$$

## 7.6

$$\hat{f} = \frac{1}{k} \sum_{j: x_j \in N_k(x)} y_j = \frac{1}{k} \sum_j \gamma_j y_j$$

$$\Rightarrow \hat{Y} = Sy$$

$$\Rightarrow d.o.f. = Tr S = \frac{1}{k} Tr \mathbb{1} + off\ diag = \frac{N}{k}$$

$$7.7 \quad GCV(\hat{f}) = \frac{1}{N} \sum_{j=1}^{N} \left( \frac{y_i - \hat{f}(x_i)}{1 - \frac{TrS}{N}} \right)^2$$

<span style="color:pink">GCV<br>Approx to $S_{ii}$</span>

$$\approx \frac{1}{N} \sum_{j=1}^{N} |y_i - \hat{f}(x_i)|^2 \left( 1 + 2 \frac{TrS}{N} \right)$$

$$= err + \frac{2 TrS}{N} \hat{\sigma_\epsilon}^2 \quad \longleftarrow \text{using MLE}$$

<span style="color:pink">rather than<br>an unbiased one</span>

<span style="color:pink">$C_p$</span>