# Linear Regression Model

- **form:** $f(X) = \beta_0 + \sum_{j=1}^{P} X_j \beta_j$

$X_j$ comes from:
- input
- transformation of input ( log, square-root, square etc )
- basis expansion $X_i^2, X_i^3 \Rightarrow$ polynomial
- "dummy" coding $(0, 1, 0, 0, 0)^T$
- interaction between variables, $X_3 = X_1 \cdot X_2$

## Properties

- $RSS(\beta) = \sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2 = \sum_{j=1}^{N} (y_i - \beta_0 - \sum_{j=1}^{P} x_{ij} \beta_j)^2$

$$= (y - X\beta)'(y - X\beta)$$

$$\frac{\partial RSS}{\partial \beta} = -2X^T(y - X\beta) \qquad \frac{\partial^2 RSS}{\partial \beta \, \partial \beta^T} = 2X^T X \qquad \hat{\beta} = (X^T X)^{-1} X^T y$$

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

$$\downarrow$$

Hat Matrix

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

- $\hat{\sigma}^2 = \frac{1}{N-p-1} \sum_{j=1}^{N} (y_i - \hat{y}_i)^2, \quad (N-p-1)\hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2$

- $Y = E(Y | X_1, \cdots, X_P) + \varepsilon = \beta_0 + \sum_{j=1}^{P} X_j \beta_j + \varepsilon$

# Hypothesis Testing

- ## t-test

$$H_0 : \beta_j = 0, \quad H_1 : \beta_j \neq 0$$

Z-score: $z_j = \dfrac{\hat{\beta}_j}{\hat{\sigma}\sqrt{v_j}}$    $v_j$: $j$th diagonal element of $(X^T X)^{-1}$

$z_j \sim t_{N-p-1}$ or $N(0,1)$ (if $\sigma$ is known)   $\hat{\sigma}^2 = RSS/(n-p)$

- ## F-test

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_q = 0 \quad H_1 : \text{at least one of } \beta_1 \sim \beta_q \neq 0$$

F-stat: $F = \dfrac{(RSS_0 - RSS_1)/(p-q)}{RSS_1/(N-p-1)}$    $RSS_1$: LS fit of $p+1$ variables
                                                  $RSS_0$: LS fit of $q+1$ variables

$F \sim F_{p-q, n-p-1}$ or $\chi^2_{p_1 - p_0}$ (for large $N$)

Q: Why $F \sim F_{p-q, n-p-1}$ ? What is exercise 3.1?

- ## Confidence interval:

$$\beta_j \in \left( \hat{\beta}_j - z^{(1-\alpha)} v_j^{\frac{1}{2}} \hat{\sigma}, \ \hat{\beta}_j + z^{(1-\alpha)} v_j^{\frac{1}{2}} \hat{\sigma} \right) \quad \text{for single } \beta_j$$

$$C_\beta \in \left\{ \beta \mid (\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi^2_{p+1}{}^{(1-\alpha)} \right\} \quad \text{for entire } \beta$$

PROOF: Let $V = (X^T X)^{1/2} (\hat{\beta} - \beta)$, then

$$\begin{cases} E(V) = 0 \\ Cov(V) = (X'X)^{1/2} Cov(\hat{\beta})(X'X)^{1/2} = \sigma^2 I \end{cases}$$

since $V$ is normally-distributed, $V'V \sim \sigma^2 \chi^2_{r+1}$

$\Rightarrow (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) \sim \sigma^2 \chi^2_{r+1}$, and $(n-r-1)\hat{\sigma}^2 = \varepsilon' \varepsilon \sim \sigma^2 \chi^2_{n-r-1}$

$\Rightarrow \left[ V'V/(r+1) \right] / \hat{\sigma}^2$ is $F_{r+1, n-r-1}$ distribution

$\Rightarrow (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) \leq (r+1) \cdot \hat{\sigma}^2 F_{r+1, n-r-1}(\alpha)$

My idea: $(X'X)^{-\frac{1}{2}} (\beta - \hat{\beta}) \perp\!\!\!\perp \hat{\varepsilon}$ but why?

- ## Goodness of fit: $R^2 = 1 - \dfrac{RSS}{TSS}$, Adjusted $R^2 = 1 - \dfrac{RSS/(n-p)}{TSS/(n-1)}$

# The Gauss-Markov Theorem

- **Gauss-Markov Theorem:** $c'\hat{\beta} = c_0\hat{\beta}_0 + c_1\hat{\beta}_1 + \cdots + c_r\hat{\beta}_r$ is the smallest possible variance among all linear estimators of $c'\beta = c_0\beta_0 + \cdots + c_r\beta_r$

  **PROOF** $c'\beta = c_0\beta_0 + \cdots + c_r\beta_r$

  Estimator: $c'\hat{\beta} = c_0\hat{\beta}_0 + c_1\hat{\beta}_1 + \cdots + c_r\hat{\beta}_r = \underbrace{c'(z'z)^{-1}z'}_{\to \, a'}Y$

  $E(c'\hat{\beta}) = c'\beta$, $V(c'\hat{\beta}) = a'\sigma^2 I \, a = \sigma^2 a'a$

  Alternative Estimator: $d'Y$ with $E(d'Y) = d'Z\beta = c'\beta$, $\forall \beta \underline{\quad} \Leftrightarrow \underline{d'z = c'}$

  $V(d'Y) = d'\sigma^2 I \, d = \sigma^2 d'd = \sigma^2[(a+d-a)'(a+d-a)]$

  $\qquad = \sigma^2[a'a + (d-a)(d-a) + \underbrace{2a'(d-a)}_{}] = \sigma^2[a'a + (d-a)'(d-a)]$

  $\qquad \geq V(c'\hat{\beta})$

  $\qquad\qquad\qquad a'(d-a) = c'(z'z)^{-1}z'(z^{-1} - z(z'z)^{-1})c$

  $\qquad\qquad\qquad\qquad = c'(z'z)^{-1}c - c'(z'z)^{-1}c = 0$

  $\Rightarrow$ **$c'\hat{\beta}$ is the BLUE of $c'\beta$**

- **Mean Square Error:** $MSE(\tilde{\theta}) = E(\tilde{\theta} - \theta)^2 = Var(\tilde{\theta}) + [E(\tilde{\theta}) - \theta]^2$

  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{\text{opt.}}{\uparrow} \qquad \underset{\text{zero}}{\uparrow}$

  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{B.L.U.E.}$

# Multiple Outputs

- $Y_k = \beta_{0k} + \sum_{j=1}^{P} X_j \beta_{jk} + \varepsilon_k = f_k(X) + \varepsilon_k$

  $Y = XB + E$

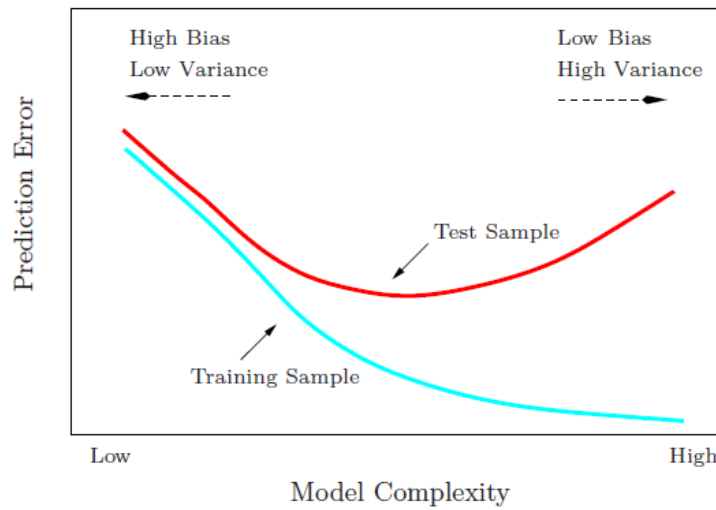- $RSS(B) = \sum_{k=1}^{K} \sum_{i=1}^{N} (y_{ik} - f_k(x_i))^2 = tr\left[(Y - XB)^\top (Y - XB)\right]$

  $\hat{B} = (X^\top X)^{-1} X^\top Y$

- if $\varepsilon_1, \cdots, \varepsilon_k$ are correlated, $Cov(\varepsilon) = \Sigma$

  $RSS(B; \Sigma) = \sum_{i=1}^{N} (y_i - f(x_i))^\top \Sigma^{-1} (y_i - f(x_i))$

# Subset Selection    discrete



High Bias
Low Variance
← - - - - - -

Low Bias
High Variance
- - - - - - →

Test Sample

Training Sample

Prediction Error

Low

High

Model Complexity

- Forward-stepwise   v.s.   Backward-stepwise
      ↓                           ↓
   start with $\beta_0$       drop the smallest Z-score
      ↓                           ↓
   add params              reduce the variance
      ↓
   improve the fit

- AIC/BIC Criterion: ??

$$AIC = -\frac{2}{N}l(\beta) + 2\frac{d}{N} , \quad BIC = -2l(\beta) + (\log N)d$$

- Shrinkage Methods <span style="color:red">continuous</span>
  - <span style="color:green">Ridge Regression</span>
    $$\hat{\beta}^{ridge} = \underset{\beta}{\text{argmin}}\left\{\sum_{i=1}^{N}(y_i-\beta_0-\sum_{j=1}^{P}x_{ij}\beta_j)^2 + \lambda\sum_{j=1}^{P}\beta_j^2\right\}$$

    or

    $$\hat{\beta}^{ridge} = \underset{\beta}{\text{argmin}}\sum_{i=1}^{N}(y_i-\beta_0-\sum_{j=1}^{P}x_{ij}\beta_j)^2$$

    $$\text{Subject to } \sum_{j=1}^{P}\beta_j^2 \le t$$

Ex. 3.5 Consider the ridge regression problem (3.41). Show that this problem is equivalent to the problem

$$\hat{\beta}^c = \underset{\beta^c}{\text{argmin}}\left\{\sum_{i=1}^{N}[y_i - \beta_0^c - \sum_{j=1}^{p}(x_{ij} - \bar{x}_j)\beta_j^c]^2 + \lambda\sum_{j=1}^{p}\beta_j^{c2}\right\}. \quad (3.85)$$

Give the correspondence between $\beta^c$ and the original $\beta$ in (3.41). Characterize the solution to this modified criterion. Show that a similar result holds for the lasso.

- $RSS(\lambda) = (y-X\beta)^T(y-X\beta) + \lambda\beta^T\beta$

  $\hat{\beta}^{ridge} = (X^TX+\lambda I)^{-1}X^Ty$

- <span style="color:purple">MLE v.s. MAP</span>

  $$\hat{\theta}_{MLE} = \text{argmax } P(X;\theta)$$
  $$= \text{argmax }\sum_{i=1}^{n}\log P(x_i;\theta)$$
  $$= \text{argmin } -\sum_{j=1}^{n}\log P(x_i;\theta)$$

  $$\hat{\theta}_{MAP} = \text{argmax } P(\theta|X)$$
  $$= \text{argmin } -\log P(\theta|X)$$
  $$= \text{argmin } -\log(X|\theta) - \log P(\theta) + \overbrace{\log P(X)}^{\color{cyan}=Const}$$
  $$= \text{argmin } -\log P(X|\theta) + \frac{1}{2}\theta'\Sigma^{-1}\theta$$

  Since $\theta \sim N(0,\Sigma)$,
  $$P(\theta) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}}e^{-\frac{1}{2}\theta'\Sigma^{-1}\theta}$$

- <span style="color:purple">log-posterior of $\beta$:</span> $\tilde{l}(\beta) = -\sum\log P(y_i;\beta) + \log P(\beta)$
  $$= -\sum_{i=1}^{n}(y_i-\beta_0-x_i^T\beta)^2/2\sigma^2 + \frac{1}{2\tau^2}\sum_{i=1}^{P}\beta_i^2$$
  $$\Leftrightarrow -\sum_{i=1}^{n}(y_i-\beta_0-x_i^T\beta) + \underset{\color{red}\lambda}{\underbrace{\frac{\sigma^2}{\tau^2}}}\sum_{i=1}^{P}\beta_i^2$$