

## Regression

$y \in \mathbb{R}$  → target variable,  $x \in \mathbb{R}^p$  → features

### Def (Regression function)

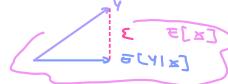
$$\mathbb{E}[Y|X]: \mathbb{R}^p \rightarrow \mathbb{R}$$

Obs: if  $(Y, X) \sim N_{p+1}$   $\Rightarrow \mathbb{E}[Y|X] = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$  (linear regression)

Data:

$$X = \begin{bmatrix} x_1 & y_1 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix} \quad x_i \in \mathbb{R}^p, y_i \in \mathbb{R} \quad i=1, \dots, n \quad \Rightarrow \hat{f}: \mathbb{R}^p \rightarrow \mathbb{R}$$

General model:  $y = \mathbb{E}[Y|X] + \varepsilon \quad \varepsilon \perp \!\!\! \perp X$

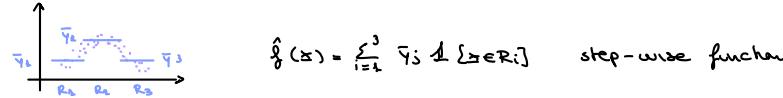


Two different approaches to deal with it:

1. Totally data driven, no model for  $f$ : "Let the data speak!" → CART (good in prediction, bad for interpretation)
2. Structured models → Linear Regression Models (good for interpretation)

### CART (Classification and regression tree)

Idea: Split the feature space  $\mathbb{R}^p$  in  $R_1, \dots, R_S$  partition of  $\mathbb{R}^p$  and predict  $y$  in  $R_i$  by means of  $\bar{y}_i = \frac{1}{|R_i|} \sum_{y_j \in R_i} y_j$



Problems and Goals: how to find  $R_1, \dots, R_S$  s.t.

1.  $\bigcup_{i=1}^S R_i = \mathbb{R}^p$
2.  $R_i \cap R_j = \emptyset \text{ for } i \neq j \quad i, j = 1, \dots, S$
3.  $R_i \cap \{x_i, i=1, \dots, n\} \neq \emptyset \quad i=1, \dots, S$

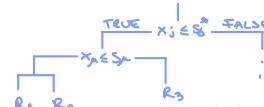
Optimal criterion: find  $R_1, \dots, R_S$  s.t.  $\sum_{i=1}^S \sum_{j \in R_i} (y_j - \bar{y}_i)^2$  is minimal (too many subsets = computationally unsolvable)

CART looks for rectangles  $R_1, \dots, R_S$



### Algorithm:

1. Look at  $x_1$  and find  $S_1^*$  (split) which maximizes:  $\sum_{i=1}^n (y_i - \bar{y})^2 - [\sum_{x_{ij} \leq S_1^*} (y_i - \bar{y}_-) + \sum_{x_{ij} > S_1^*} (y_i - \bar{y}_+)^2]$  with  $\bar{y}_- = \frac{1}{|R_1|} \sum_{x_{ij} \leq S_1^*} y_i$  and  $\bar{y}_+ = \dots$
2. repeat (1) for  $x_2, x_3, \dots, x_p \rightarrow S_2^*, S_3^*, \dots, S_p^*$
3. choose  $S_i^*$  which maximizes (\*). We now have more datasets:  $R_1 = \{x_{ij} : x_{ij} \leq S_1^*\}, R_2 = \{x_{ij} : x_{ij} > S_1^*\}$
4. Repeat (1), (2), (3) for all  $R_i$ .
5. Stop splitting on  $R_i$  if number of data points in  $R_i$  is below threshold



Whenever I stop earlier I get anyways a separation of the feature space. A goal could be to maximize purity in my rectangles  
Useful for categorical and continuous variables

To control overfitting: introduce penalization for  $S$  and minimize:  $\sum_{i=1}^S \sum_{j \in R_i} (y_j - \bar{y}_i)^2 + \alpha S = W(\alpha)$

⇒ Introduce an extra param  $\alpha$  ⇒ chosen by cross-validation

## Linear Model for Regression

Date

$$X = \begin{bmatrix} (x_{11}, \dots, x_{1p})^T & y_1 \\ (x_{21}, \dots, x_{2p})^T & y_2 \\ \vdots & \vdots \\ (x_{n1}, \dots, x_{np})^T & y_n \end{bmatrix}$$

Design matrix:

$$Z = \begin{bmatrix} 1 & z_{11} & \dots & z_{1p} \\ 1 & z_{21} & \dots & z_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{np} \end{bmatrix} \quad n \times (p+1)$$

$$\text{with: } z_1 = h_1(x_1, \dots, x_p) \dots \quad z_p = h_p(x_1, \dots, x_p)$$

and  $h_1, \dots, h_p$  are known

$$(\text{special case } r=p \quad z_1 = x_1, \dots, z_p = x_p)$$

$$\text{We will have: } E[Y|X_1, \dots, X_p] = \beta_0 + \beta_1 h_1(x_1, \dots, x_p) + \beta_2 h_2(x_2, \dots, x_p) + \dots + \beta_p h_p(x_p, \dots, x_p) = \beta_0 + \beta_1 z_1 + \dots + \beta_p z_p$$

$$\text{Model for the phenomenon: } Y = \beta_0 + \beta_1 z_1 + \dots + \beta_p z_p + \varepsilon \quad \sum \text{r.v. } E[\varepsilon] = 0, \quad \text{Var}(\varepsilon) = \sigma^2$$

$$\text{Model for the data: } \hat{Y} = (y_1, \dots, y_n)^T \quad Y = Z\beta + \varepsilon \quad \text{with: } Z \in \mathbb{R}^{n \times (p+1)} \text{ known}$$

$$\sum \text{r.v. s.t. } E[\varepsilon] = 0, \quad \text{Cov}(\varepsilon) = \sigma^2 I$$

$$\text{i.e. for } i=1, \dots, n \quad y_i = \beta_0 + \beta_1 z_{i1} + \dots + \beta_p z_{ip} + \varepsilon_i \quad \varepsilon_i \text{ s.t. } \text{Cov}(\varepsilon_i, \varepsilon_j) = \begin{cases} \sigma^2 & i=j \\ 0 & i \neq j \end{cases} \quad E[\varepsilon_i] = 0$$

## EXAMPLE (ANOVA + linear regression)

$$Y = \begin{cases} x_{11}, \dots, x_{1n_1} & \text{iid } \sim N(\mu, \sigma^2) \\ \vdots \\ x_{g_1}, \dots, x_{gn_g} & \text{iid } \sim N(\mu, \sigma^2) \end{cases}$$

$$Y = (x_{11}, \dots, x_{1n_1}, x_{21}, \dots, x_{2n_2}, \dots, x_{gn_g})^T \in \mathbb{R}^n \quad n = n_1 + n_2 + \dots + n_g$$

$$Z = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \end{bmatrix} \quad \beta = (\mu, \tau_1, \dots, \tau_g) \quad \varepsilon \sim N_n(0, \sigma^2 I)$$

$$\Rightarrow Y = Z\beta + \varepsilon \iff \begin{cases} x_{ij} = \mu + \tau_i + \varepsilon_{ij} & i=1, \dots, g, \quad j=1, \dots, n_i \\ \varepsilon_{ij} \quad \text{iid } \sim N(0, \sigma^2) \end{cases}$$

$\Rightarrow$  it's an ANOVA as LRM

notice:  $Z$  is not full rank  $\rightarrow$  problem of identification of parameters.

Geometrically speaking is not a problem, but if I want a unique solution it is!

To get it as full rank is "easy": including the constraint:  $\sum_{i=1}^g n_i \tau_i = 0$ ,  $\tau_g = \sum_{i=1}^{g-1} \frac{n_i}{n_g} \tau_i$ ,  $n_g = 0$

We have:

$$Z = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \\ -n_1/n_g & -n_2/n_g & -n_3/n_g & \dots & -n_{g-1}/n_g \\ 1 & -n_1/n_g & -n_2/n_g & \dots & -n_{g-1}/n_g \end{bmatrix}_{n \times g} \quad \beta = (\mu, \tau_1, \dots, \tau_{g-1}) \in \mathbb{R}^{g-1}$$

$$Y = Z\beta + \varepsilon$$

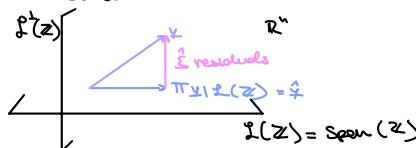
full rank makes computation easier, but not required

How to estimate  $\beta$  and  $\sigma^2$ ?

$$Y = Z\beta + \varepsilon \quad Z\beta = E[Y|Z]$$

$$Z = [c_1, \dots, c_m] \quad c_i \in \mathbb{R}^n, \quad Z\beta = \beta_0 c_1 + \dots + \beta_m c_m$$

We visualize:



$$\hat{Y} = \Pi_{\perp} Z(z) = Z\hat{\beta}, \quad Z\hat{\beta} \text{ is the solution of OLS (ordinary least square):}$$

$$\arg \min_{\beta \in \mathbb{R}^{p+1}} \|Y - Z\beta\|^2 = \arg \min_{\beta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 z_{i1} - \dots - \beta_p z_{ip})^2 = \hat{\beta}$$

## Proposition

If  $Z$  is full rank:

1.  $\hat{Y} = \Pi_{\perp} Z(z) = Z(Z^T Z)^{-1} Z^T Y = H Y$
2.  $\hat{\beta} = (Z^T Z)^{-1} Z^T Y$

Proof:

$$Z^T Z = \sum_{i=1}^{p+1} \lambda_i e_i e_i^T \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{p+1} > 0$$

$$(Z^T Z)^{-1} = \sum_{i=1}^{p+1} \frac{1}{\lambda_i} e_i e_i^T$$



For  $i = 1, \dots, r+s$   $q_i = \frac{1}{\sqrt{\lambda_i}} Z e_i$  s.t.

- $q_i \in \mathcal{L}(Z)$ ,  $i = 1, \dots, r+s$
- $q_i^T q_j = \frac{1}{\sqrt{\lambda_i \lambda_j}} e_i^T Z^T Z e_j = \sqrt{\frac{\lambda_j}{\lambda_i}} e_i^T e_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

$\Rightarrow \{q_1, q_2, \dots, q_{r+s}\}$  orthonormal basis for  $\mathcal{L}(Z)$

$$\Pi_{\mathcal{L}(Z)} = \sum_{i=1}^{r+s} \frac{q_i q_i^T}{q_i^T q_i} Y = \sum_{i=1}^{r+s} \frac{1}{\lambda_i} Z e_i e_i^T Z^T Y = Z \left( \sum_{i=1}^{r+s} \frac{1}{\lambda_i} e_i e_i^T \right) Z^T Y = Z \underbrace{(Z^T Z)^{-1}}_{B} Z^T Y$$

$\Rightarrow \hat{Y} = M Y$  and  $\hat{B} = (Z^T Z)^{-1} Z^T Y$

- Obs:
- If  $\text{rank}(Z) = k < r+s \leq n$  in the proof:  $\lambda_1 > \lambda_2 > \dots > \lambda_k > 0$ ,  $\lambda_{k+1} = \lambda_{k+2} = \dots = \lambda_{r+s} = 0$   
 $\Rightarrow$  cannot divide  $\Rightarrow$  stop summation at step  $k$   
 $\Rightarrow$  I get the generalized inverse  $\Rightarrow M = Z(Z^T Z)^{-1} Z^T$  and  $\hat{B} = (Z^T Z)^{-1} Z^T Y$
  - If  $\text{rank}(Z) = r+s = n \Rightarrow Y = \hat{Y}$   
 NOT GOOD, enormous generalization error (unless  $Y = Z\beta + \varepsilon$  and  $\text{Cov}(\varepsilon) = 0$  stupid!!)

### EXAMPLE

$\hat{Y} = H Y$   
 $\hat{\varepsilon} = (I - H) Y$   
 $H = Z(Z^T Z)^{-1} Z^T$  orth. proj. on  $\mathcal{L}(Z)$  ( $\dim = r+s$ )  
 $I - H$  orth. project. on  $\mathcal{L}^\perp(Z)$   $\dim = n - (r+s)$   
 $\|\hat{Y}\|^2 = \|\hat{Y}_\parallel\|^2 + \|\hat{\varepsilon}\|^2$   
 $SS_{\text{tot}} = SS_{\text{reg}} + SS_{\text{res}}$  1st decomposition of variance formula

$\Pi_{\mathcal{L}(Z)} = \bar{Y} \cdot 1$   
 $\Pi_{\mathcal{L}^\perp(Z)} = \frac{1 \cdot \bar{Y}}{1^T \bar{Y}} \hat{Y} = \frac{1 \cdot \bar{Y}}{1^T \bar{Y}} Y = \frac{1 \cdot \bar{Y}}{1^T \bar{Y}} Y = \bar{Y} \cdot 1$   
 $(\delta^T H^T) = (H^T \delta)^T = (H^T 1)^T = 1^T$

$\|\hat{Y} - \bar{Y} \cdot 1\|^2 = \|\hat{Y} - \bar{Y} \cdot 1\|^2 + \|\hat{\varepsilon}\|^2$   
 $\sum (y_i - \bar{Y})^2 = \sum (\hat{y}_i - \bar{Y})^2 + \sum \hat{\varepsilon}_i^2$   
 $SS_{\text{reg}} = SS_{\text{reg}} + SS_{\text{res}}$  2nd decomposition of variance formula  
 $R^2 = \frac{\sum (\hat{y}_i - \bar{Y})^2}{\sum (y_i - \bar{Y})^2} + \frac{\sum \hat{\varepsilon}_i^2}{\sum (y_i - \bar{Y})^2}$   
 $R^2$  coeff. of determination  $\rightarrow R^2 = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum (y_i - \bar{Y})^2} = 1 - \sin^2 \theta = \cos^2 \theta$

- Obs:
- $R^2 = 1 \Rightarrow \theta = 0 \Rightarrow \hat{Y} = \hat{Y}$
  - $R^2 = 0 \Rightarrow \theta = \frac{\pi}{2} \Rightarrow \bar{Y} \cdot 1 = \hat{Y}$

Rank:  
 $Z = \begin{bmatrix} z_{11} & \dots & z_{1n} \\ z_{21} & \dots & z_{2n} \\ \vdots & \ddots & \vdots \\ z_{m1} & \dots & z_{mn} \end{bmatrix}$  If  $\hat{\varepsilon} \notin \mathcal{L}(Z)$ , e.g. regression through the origin

$$\tilde{R}^2 = 1 - \frac{\|\hat{\varepsilon}\|^2}{\|\hat{Y}\|^2} = \cos^2 \eta$$

$$R^2_{\text{adj}} = 1 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2 / (n - (r+s))}{\sum_{i=1}^n (y_i - \bar{Y})^2 / (n - r)} \quad R^2 \text{ adjusted}$$

### Properties (of $\hat{B}$ and $\hat{\varepsilon}$ )

Assume  $Z$  to be full rank ( $r+s$ )

$$\hat{B} = (Z^T Z)^{-1} Z^T Y$$

$$\hat{\varepsilon} = (I - H) Y$$

$$H = Z(Z^T Z)^{-1} Z^T$$

$$1. E[\hat{B}] = B \quad (\text{unbiased})$$

$$2. \text{Cov}(\hat{B}) = \sigma^2 (Z^T Z)^{-1}$$

$$3. E[\hat{\varepsilon}] = 0 \quad (\text{recall } E[\varepsilon] = 0)$$

$$4. \text{Cov}[\hat{\varepsilon}] = \sigma^2 (I - H) \quad (\text{recall } \text{Cov}(\varepsilon) = \sigma^2 I)$$

$$5. E[\hat{\varepsilon}^T \hat{\varepsilon}] = E[\|\hat{\varepsilon}\|^2] = (n - (r+s)) \sigma^2 \Rightarrow E\left[\frac{\hat{\varepsilon}^T \hat{\varepsilon}}{n - (r+s)}\right] = \sigma^2 \quad (\text{unbiased})$$



**Proof:** 1.  $E[\hat{\beta}] = E[(Z'Z)^{-1}Z'\gamma] = (Z'Z)^{-1}Z'E[\gamma]$

$$\gamma = Z\beta + \varepsilon$$

$$E[\gamma] = Z\beta + E[\varepsilon] = Z\beta$$

$$E[\hat{\beta}] = (Z'Z)^{-1}Z'Z\beta = \beta$$

2.  $Cov(\hat{\beta}) = (Z'Z)^{-1}Z' Cov(\gamma) Z (Z'Z)^{-1} = \sigma^2 (Z'Z)^{-1}Z'Z (Z'Z)^{-1} = \sigma^2 (Z'Z)^{-1}$   
 $Cov(\gamma) = Cov(\varepsilon) = \sigma^2 I$

3.  $E[\hat{\Sigma}] = (I - H)E[\varepsilon] - (I - H)Z\beta = 0$

4.  $Cov(\hat{\Sigma}) = (I - H)Cov(\gamma)(I - H)' = \sigma^2 (I - H)(I - H)' = \sigma^2 (I - H)$

5.  $E[\hat{\Sigma}'\hat{\Sigma}] = \text{tr}[E[\hat{\Sigma}'\hat{\Sigma}]] = E[\text{tr}(\hat{\Sigma}'\hat{\Sigma})] = E[\text{tr}(\hat{\Sigma}\hat{\Sigma}')] =$

$$\hat{\Sigma}' = (I - H)\gamma = (I - H)(Z\beta + \varepsilon) = (I - H)\varepsilon$$

$$E[\hat{\Sigma}'\hat{\Sigma}] = E[\text{tr}((I - H)\varepsilon\varepsilon'(I - H))] = \text{tr}[(I - H)E[\varepsilon\varepsilon']'(I - H)] = \text{tr}(\sigma^2(I - H)) =$$

$$= \sigma^2 \text{tr}(I - H) = \sigma^2 (n - \text{tr}(H))$$

$$\text{tr}(H) = \text{tr}(Z(Z'Z)^{-1}Z') = \text{tr}(Z'Z(Z'Z)^{-1}) = \text{tr}(I_{n-r}) = r$$

$$E[\hat{\Sigma}'\hat{\Sigma}] = \sigma^2(n - (r - s))$$

this tr trick is going to be in  
the exam to be used

**Obs:**  $Cov(\hat{\beta}) = \sigma^2 (Z'Z)^{-1}$

$\rightarrow \hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_r)' \rightarrow$  are not uncorrelated unless  $Z'Z = dI$

$$\gamma = \hat{\beta}_0 + \hat{\beta}_1 z_1 + \dots + \hat{\beta}_r z_r$$

$\hookrightarrow \beta_0/z$  like that if they are correlated, then also the other ones are changed

$\rightarrow$  You can control the varab. of  $\hat{\beta}$  by controlling the design matrix

**Assumption until new notice:**  $\varepsilon \sim N_n(0, \sigma^2 I)$

$$\gamma = Z\beta + \varepsilon$$

### Proposition

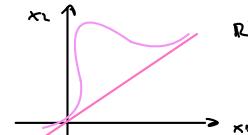
1.  $\hat{\beta}$  and  $\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-r}$  are the est. of  $\beta$  and  $\sigma^2$

2.  $\hat{\beta} \sim N_{n+r}(0, \sigma^2 (Z'Z)^{-1})$

3.  $\hat{\varepsilon} \sim N_n(0, \sigma^2 (I - H))$

4.  $\hat{\varepsilon} \perp \hat{\beta}$

5.  $\hat{\varepsilon}'\hat{\varepsilon} \sim \sigma^2 \chi_{(n-(r+s))}^2$



**Proof:** Note:  $\gamma = Z\beta + \varepsilon \sim N_n(Z\beta, \sigma^2 I)$

1. Write the L.K and start differentiating

2/3/4 
$$\begin{pmatrix} \hat{\beta} \\ \hat{\varepsilon} \end{pmatrix} = \underbrace{\begin{pmatrix} (Z'Z)^{-1} \\ I - H \end{pmatrix}}_A \gamma \sim N \left( \begin{pmatrix} \beta \\ 0 \end{pmatrix}, \begin{bmatrix} (Z'Z)^{-1} & 0 \\ 0 & I - H \end{bmatrix} \right) \sim N(AZ\beta, \sigma^2 AIA')$$

$$\hat{\varepsilon} \sim N_n(0, \sigma^2 (I - H))$$

$$\frac{1}{\sigma^2} (\hat{\varepsilon}' - 0)' (I - H)^{-1} (\hat{\varepsilon}' - 0) \sim \chi_{(n-(r+s))}^2$$

**Exercise:**  $x \sim N_p(\mu, \Sigma)$  rank  $(\Sigma) = k < p \Rightarrow (x - \mu)' \Sigma^{-1} (x - \mu) \sim \chi^2(k)$



## Linear models

$$\mathbb{R}^n \ni y = Z\beta + \varepsilon$$

$Z$  design matrix  $(r+1) \times n$  (full rank)

$$\beta \in \mathbb{R}^{r+1}$$

$$\varepsilon \sim N_n(\mathbf{0}, \sigma^2 I)$$

$$\hat{\beta} = (Z'Z)^{-1}Z'y$$

$$S^L = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-(r+1)}$$

$$\hat{\varepsilon} = (I - H)y$$

$$H = Z(Z'Z)^{-1}Z'$$

$$\hat{y} = Hy$$

Prop:

1.  $\hat{\beta} \sim N_{r+1}(\beta, \sigma^2(Z'Z)^{-1})$
2.  $\hat{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2(I-H))$
3.  $\hat{\varepsilon}'\hat{\varepsilon} \sim \sigma^2 \chi^2(n-(r+1))$

Note:

$$\frac{1}{\sigma^2} (\hat{\beta} - \beta)' (Z'Z) (\hat{\beta} - \beta) \sim \chi^2(r+1) \quad \left. \right\} \text{I}$$

$$\frac{1}{\sigma^2} \hat{\varepsilon}'\hat{\varepsilon} \sim \chi^2(n-(r+1)) \quad \left. \right\} \text{II}$$

$$(Z'Z)^{-1} (Z'Z) (\hat{\beta} - \beta) / r+1 \sim F(r+1, n-(r+1)) \quad \Rightarrow \quad \frac{1}{S^L} (\hat{\beta} - \beta)' (Z'Z) (\hat{\beta} - \beta) / (r+1) \sim F(r+1, n-(r+1))$$

$$\alpha \in (0, 1)$$

$$CI_{1-\alpha}(\beta) = \{ \beta \in \mathbb{R}^{r+1} : (\hat{\beta} - \beta)' Z'Z (\hat{\beta} - \beta) \leq (r+1) S^L \} \quad F_{r+1, n-(r+1)}$$

$$\hat{\beta} \in \mathbb{R}^{r+1}$$

$$\alpha \in \mathbb{R}^{r+1}$$

$$\alpha' \hat{\beta} \sim N_1(\alpha' \beta, \sigma^2 \alpha' (Z'Z)^{-1} \alpha)$$

$$\Rightarrow \frac{\alpha' (\hat{\beta} - \beta)}{\sigma \sqrt{\alpha' (Z'Z)^{-1} \alpha}} \sim t_{n-(r+1)} \Rightarrow \frac{\alpha' (\hat{\beta} - \beta)}{S^L \sqrt{\alpha' (Z'Z)^{-1} \alpha}} \sim t_{n-(r+1)}$$

$$\alpha \in (0, 1) \quad CI_{1-\alpha}(\alpha' \beta) = \left[ \alpha' \hat{\beta} \pm S^L \sqrt{\alpha' (Z'Z)^{-1} \alpha} t_{1-\frac{\alpha}{2}}(n-(r+1)) \right] \text{ are the true intervals for } \alpha' \beta$$

In particular

$$\alpha_i = (0 \dots 0 \overset{i}{\underset{\downarrow}{1}} 0 \dots 0)$$

$$CI_{1-\alpha}(\beta_i) = \left[ \hat{\beta}_i \pm S^L \sqrt{\text{diag}(Z'Z)^{-1}} t_{1-\frac{\alpha}{2}}(n-(r+1)) \right]$$

$$\alpha \rightarrow \frac{\alpha}{r+1} \quad (\text{Bonferroni})$$

Test:

$$H_0: \beta_i = 0 \quad \text{vs} \quad H_1: \beta_i \neq 0$$

$$\text{Reject } H_0 \text{ at level } \alpha \in (0, 1) \quad \text{if} \quad T_i = \frac{|\hat{\beta}_i|}{S^L \sqrt{\text{diag}(Z'Z)^{-1}}} > t_{1-\frac{\alpha}{2}}(n-(r+1)) \quad \Rightarrow \text{p-value}$$

$$\max_{\alpha \in \mathbb{R}^{r+1}} \frac{[\alpha' (\hat{\beta} - \beta)]^2}{S^L \sqrt{\alpha' (Z'Z)^{-1} \alpha}} = \max_{\alpha \in \mathbb{R}^{r+1}} \frac{[\alpha' (\hat{\beta} - \beta)]^2}{S^L [\alpha' (Z'Z)^{-1} \alpha]} = \frac{1}{S^L} (\hat{\beta} - \beta)' (Z'Z) (\hat{\beta} - \beta) \sim F(r+1, n-(r+1))$$

$$\text{Since } CI_{1-\alpha}(\alpha' \beta) = \left[ \alpha' \hat{\beta} \pm S^L \sqrt{\alpha' (Z'Z)^{-1} \alpha} t_{1-\frac{\alpha}{2}}(n-(r+1)) \right]$$

$$\frac{1}{S^L} \hat{\varepsilon}'\hat{\varepsilon} \sim \chi^2[n-(n+1)]$$

$$\frac{(n-(r+1))}{S^L} \sim \chi^2(n-(r+1))$$

$$P\left[ \chi_{\frac{n}{2}}^2(n-(r+1)) < \frac{(n-(r+1))S^L}{S^L} < \chi_{1-\frac{\alpha}{2}}^2(n-(r+1)) \right] = 1-\alpha$$

$$\alpha \in (0, 1)$$

$$CI_{1-\alpha}(\sigma^2) = \left[ \frac{(n-(r+1))S^L}{\chi_{1-\frac{\alpha}{2}}^2(n-(r+1))}, \frac{(n-(r+1))S^L}{\chi_{\frac{n}{2}}^2(n-(r+1))} \right]$$



$C \propto (r+L)$  matrix

$H_0: C\beta = 0 \text{ vs } C\beta \neq 0 \quad (\text{Ex: } H_0: C\beta = 0 \text{ vs } C\beta \neq 0)$

$$C\beta = \begin{bmatrix} C_{11}\beta_1 + \dots + C_{1,r+L}\beta_{r+L} \\ \vdots \\ C_{p1}\beta_1 + \dots + C_{pr+L}\beta_{r+L} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$C\hat{\beta} \sim N_p(C\beta, \sigma^2 C(z'z)^{-1}C')$$

$$\text{Under } H_0: \frac{\frac{1}{\sigma^2} (C\hat{\beta})[C(z'z)^{-1}C']^{-1}(C\hat{\beta})}{p} \sim \chi^2(p) \quad \frac{\frac{1}{\sigma^2} \sum_{i=1}^p \hat{\beta}_i^2}{\frac{1}{\sigma^2} (p-(r+L))} \sim \chi^2_{p-(r+L)}$$

$$F = \frac{1}{\sigma^2} (C\hat{\beta})[C(z'z)^{-1}C']^{-1}(C\hat{\beta}) \sim F(p, n-(r+L))$$

Reject  $H_0: C\beta = 0$  at level  $\alpha \in (0, 1)$  if  $F > F_{1-\alpha}(p, n-(r+L))$

$H_0: \beta_0 = \beta_1 = \dots = \beta_{r+L-p} = 0 \quad \text{vs} \quad H_1: \exists \beta_i \neq 0, i = r+L-p+1, \dots, r$

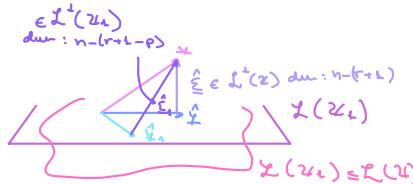
$$C = \left[ \begin{array}{cccccc} 0 & \dots & 0 & 1 & \dots & 0 \\ \vdots & & \vdots & 0 & \dots & 0 \\ \vdots & & \vdots & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 1 \end{array} \right]_p = [0 \mid I_p] \quad p \times (r+L)$$

$H_0: \beta_0 = \beta_{r+L} = \dots = \beta_{r+L-p} = 0$

is testing: full model  $y = z'\beta + \varepsilon$   
vs reduced model  $y = z_0 + \beta_1 z_1 + \varepsilon$

Reject  $H_0$  if  $S S_{\text{res}}(z_1) - S S_{\text{res}}(z_0)$  is large

$$\text{but } \frac{S S_{\text{res}}(z_1) - S S_{\text{res}}(z_0)}{\frac{s^2}{n-(r+L)}} \sim F(p, n-(r+L))$$



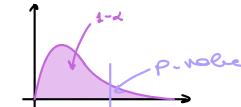
Very particular case:

$H_0: \beta_0 = \beta_1 = \dots = \beta_r = 0 \quad \text{vs} \quad H_1: \exists \beta_i \neq 0$

$$\Rightarrow Z_i = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad S S_{\text{res}}(z_i) = \sum (y_i - \bar{y})^2$$

$$\frac{S S_{\text{res}}(z_i) - S S_{\text{res}}(z_0)}{S^2 \cdot r} = \frac{\sum (y_i - \bar{y})^2 - \sum \hat{\varepsilon}_i^2}{S^2 \cdot r} = \frac{\sum (\hat{y}_i - \bar{y})^2 / r}{\sum (\hat{\varepsilon}_i^2) / (n-(r+L))} \sim F(r, n-(r+L))$$

Reject  $H_0$  (i.e. there is something in the model) if  $F > F_{1-\alpha}(r, n-(r+L))$



## PREDICTION

Model for the phenomenon  $y_0 = z_0' \beta + \varepsilon_0$

$\hookrightarrow z_0 + (z_1, z_2, \dots, z_m)$  new case, not in the training set

predict:  $E[y_0 | z_0] = z_0' \beta$

Unbiased predictor of  $z_0' \beta$ :  $\hat{z}_0' \hat{\beta} = \hat{z}_0' (z'z)^{-1} z' y$

Indeed:  $E[\hat{z}_0' \hat{\beta}] = \hat{z}_0' \beta$

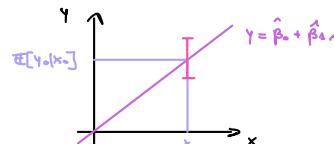
## Gauss-Markov theorem

$\hat{z}_0' \hat{\beta}$  is BLUE i.e. Best Linear Unbiased Estimator of  $z_0' \beta$

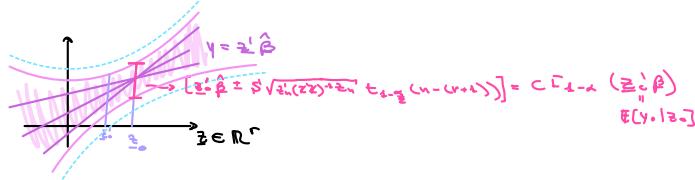
## Point prediction

$$\hat{z}_0' \hat{\beta}$$

$$CI_{1-\alpha}(\hat{z}_0' \hat{\beta}) = [\hat{z}_0' \hat{\beta} \pm S \sqrt{\hat{z}_0' (z'z)^{-1} \hat{z}_0} t_{1-\alpha/2}(n-(r+L))]$$



## Note



"The bands contain the real model with conf. 1- $\alpha$ "

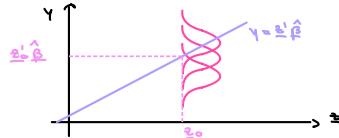
↳ wrong! each interval 95% not together  $\rightarrow$  sum 95%

$[z_0 \hat{\beta} \pm S(\sqrt{z_0'(Z'Z)^{-1} z_0}) \sqrt{F_{\alpha/2}(n-2+1, n-(2+1))}]$  all possible simultaneously

In this case not Bonferroni because  $\rightarrow +\infty$ , so F-test

Sometimes we want to estimate for  $E[y_0|z_0] \rightarrow E[\hat{y}_0|z_0]$

$$y_0 = z_0' \hat{\beta} + \varepsilon_0 \\ E[\hat{y}_0|z_0] = z_0' \hat{\beta} \\ \varepsilon_0 \perp \varepsilon_1, \varepsilon_2, \dots$$



not true  $\rightarrow$  estimate has distribution

Question: I interval s.t.  $P[y_0 \in I | z_0] = 1-\alpha$

We know  $y_0 \sim N(z_0' \hat{\beta}, \sigma^2)$ ,  $\perp \varepsilon_0$ ,  $z_0' \hat{\beta} \sim N(z_0' \hat{\beta}, \sigma^2 (Z'Z)^{-1} z_0)$

so  $y_0 - z_0' \hat{\beta} \sim N(0, \sigma^2 (I + z_0' (Z'Z)^{-1} z_0))$

$$\text{so } \frac{y_0 - z_0' \hat{\beta}}{\sqrt{\sigma^2 (I + z_0' (Z'Z)^{-1} z_0)}} \sim t(n-(2+1)) \implies I = [z_0' \hat{\beta} \pm S \sqrt{1 + z_0' (Z'Z)^{-1} z_0}] t_{1-\frac{\alpha}{2}}(n-(2+1))$$

## Lin model

$$y = z\beta + \varepsilon \\ E[\varepsilon] = 0, \text{ Cov}(\varepsilon) = \sigma^2 I \quad \text{i.e.:} \quad \begin{cases} \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 & \text{if } i \neq j \\ \text{Var}(\varepsilon_i) = \sigma^2 \end{cases}$$

$$\text{? Cov}(\varepsilon) = \sigma^2 \sum_i \quad (\varepsilon_i \text{ non cov.})$$

Unknown      Known

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} (y - z\beta)' \Sigma^{-1} (y - z\beta) \quad \text{GLS: Generalized Least Squares}$$

$$(y - z\beta)' \Sigma^{-1/2} \Sigma^{-1/2} (y - z\beta) = [\Sigma^{-1/2} y - \Sigma^{-1/2} z\beta]' [\Sigma^{-1/2} y - \Sigma^{-1/2} z\beta] = \|y - z\beta\|^2$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|y - z\beta\|^2 \implies \hat{\beta} = (z'z)^{-1} z'y = (z' \Sigma^{-1/2} \Sigma^{-1/2} z)^{-1} z' \Sigma^{-1/2} \Sigma^{-1/2} y = (z' \Sigma^{-1} z)^{-1} z' \Sigma^{-1} y$$

Note:  $y = z\beta + \varepsilon \quad \text{cov}(\varepsilon) = \sigma^2 \Sigma$   
 $\underbrace{\Sigma^{-1/2} y}_{\tilde{y}} = \underbrace{\Sigma^{-1/2} z\beta}_{\tilde{z}} + \underbrace{\Sigma^{-1/2} \varepsilon}_{\tilde{\varepsilon}} \quad \rightarrow \quad \tilde{y} = \tilde{z}\beta + \tilde{\varepsilon}$   
 $\text{cov}(\tilde{\varepsilon}) = \sigma^2 \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = \sigma^2 I$

- Note: If  $\Sigma$  is unknown:
- model  $\Sigma$  parametrically;
  - iteratively estimate  $\Sigma$ , from the residuals, initializing  $\Sigma_0 = I$
  - Transform  $y$  or  $\Sigma$

## EXAMPLES

1.  $y_i$  is the mean of  $n_i$  obs indep. with same var  $\sigma^2$

$$\text{Var}(y_i) = \frac{\sigma^2}{n_i}$$

$\Sigma$  known  $\Sigma = \begin{bmatrix} n_1 & & & \\ & n_2 & & \\ & & \ddots & \\ & & & n_n \end{bmatrix} \rightarrow \Sigma^{-1} = \begin{bmatrix} 1/n_1 & & & \\ & 1/n_2 & & \\ & & \ddots & \\ & & & 1/n_n \end{bmatrix}$  wLS

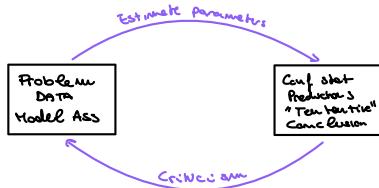


2.  $y_i$  sum of  $n_i$  obs indep. with same var  $\sigma^2$

$$\text{Var}(y_i) = n_i \sigma^2$$

$$\sum_{\text{known}} \Sigma = \begin{bmatrix} n_1 & & & \\ & n_2 & & \\ & & \ddots & \\ & & & n_n \end{bmatrix} \rightarrow \Sigma^{-1} = \begin{bmatrix} 1/n_1 & & & \\ & 1/n_2 & & \\ & & \ddots & \\ & & & 1/n_n \end{bmatrix} \text{ WLS}$$

Box



### Diagnostic for LM

- Residuals analysis, outliers, heteroscedasticity, normality, autocorrelation, ..
- Influential cases
- collinearity
- others ..

Model

$$y = z\beta + \varepsilon$$

$$\varepsilon : E[\varepsilon] = 0$$

$$\text{Cov}(\varepsilon) = \sigma^2 I$$

$$\varepsilon \sim N_n(0, \sigma^2 I)$$

distr on  $\mathbb{R}^n$

### Fitted model

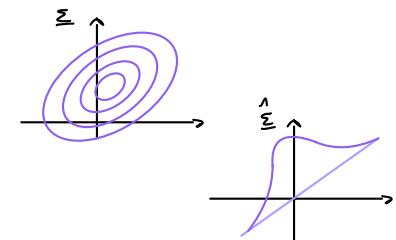
$$\hat{y} = z\hat{\beta} \implies y = z\hat{\beta} + \hat{\varepsilon}$$

$$\hat{\varepsilon} : E[\hat{\varepsilon}] = 0$$

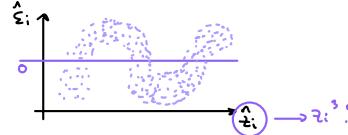
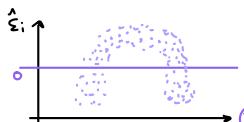
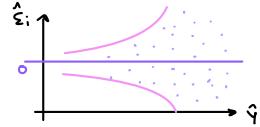
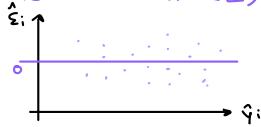
$$\text{Cov}(\hat{\varepsilon}) = \sigma^2 (I - H)$$

$$\hat{\varepsilon} \sim N_n(0, \sigma^2 (I - H))$$

distr on  $L^\perp(z)$



### Residual An ( $\hat{\varepsilon}$ )



$$\text{But: } \text{Cov}(\hat{\varepsilon}) = \sigma^2 (I - H)$$

$$\text{Var}(\hat{\varepsilon}_i) = \sigma^2 (1 - h_{ii})$$

$$h_{ii} \text{ diag } H$$

$$H = Z(Z'Z)^{-1}Z'$$

$$\hat{\varepsilon}_i \longrightarrow \frac{\hat{\varepsilon}_i}{s\sqrt{1-h_{ii}}}$$

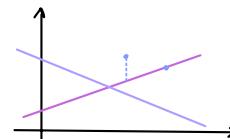
Standardized residuals

If assume  $\varepsilon \sim N$   $\rightarrow$  check Gauss of  $\hat{\varepsilon}$ : QQ plots, testing, ..

### h-ii leverage

1.  $0 \leq h_{ii} \leq 1$  because  $H = H^T$  and  $H \cdot H = H$

$$h_{ii} \uparrow \perp \implies \text{Var}(\hat{\varepsilon}_i) \downarrow 0 \quad \left. \begin{array}{l} \text{Var}(\hat{\varepsilon}_i) = 0 \\ E[\hat{\varepsilon}_i] = 0 \end{array} \right\} \implies \hat{\varepsilon}_i = 0$$



$$\text{But: } \text{Cov}(\hat{\varepsilon}) = \sigma^2 (I - H)$$

$$\text{Var}(\hat{\varepsilon}_i) = \sigma^2 (1 - h_{ii})$$

$$h_{ii} \text{ diag } H: H = Z(Z'Z)^{-1}Z'$$

$$\hat{\varepsilon}_i \longrightarrow \frac{\hat{\varepsilon}_i}{s(1-h_{ii})} \quad \text{standardized residuals}$$

### Influential case

$$X = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix} \implies Z = \begin{bmatrix} 1 & x_1 & \dots & x_n \\ 1 & x_2 & \dots & x_n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & x_n \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

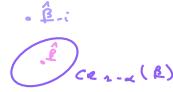
$$Z_i = \begin{bmatrix} 1 & z_{i1} & \dots & z_{ir} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix}$$

Hold out unit i  
 $y_i - \sum_{j \neq i} \hat{\beta}_{-i} z_{-ij} + \varepsilon_i \rightarrow \hat{\beta}_i$   
 if  $\hat{\beta}_i$  and  $\hat{\beta}_{-i}$  are very different  $\Rightarrow$  unit i is influential

$$D_i = \frac{(\hat{\beta}_i - \hat{\beta}_{-i})(z'_i Z)(\hat{\beta}_i - \hat{\beta}_{-i})}{S^2(r+2)} \quad \text{Cook's distance for case i}$$

$D_i$  large  $\Rightarrow$  influential cases

- Compare with  $\alpha$  quantile of  $\chi^2(r+2, n-(r+2))$
- $D_i > 1$



$$\text{In fact } D_i = \left( \frac{\hat{\varepsilon}_i}{S\sqrt{1-h_{ii}}} \right)^2 \cdot \frac{h_{ii}}{1-h_{ii}} \cdot \frac{1}{r+2}$$

$\uparrow \text{std res.}$        $\downarrow \text{with } h_{ii}$

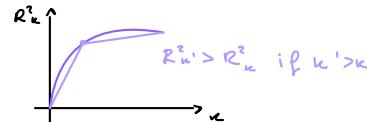
### Model selection

$$y = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r + \varepsilon$$

Can I remove some of them?

2<sup>r</sup> possible models with r regressors

- For  $k = 0, \dots, r$  fit all  $\binom{n}{k}$  with k regressors
- Choose the best one (e.g. the one with max  $R^2$ )



### Collinearity and variable selection

$$y = Z\beta + \varepsilon \quad Z \text{ design matrix, } n \times (r+2)$$

$$\text{ols. } \hat{\beta} = (Z'Z)^{-1} Z' y \quad \text{if } Z \text{ is full rank}$$

$$\text{Var}(\hat{\beta}) = \sigma^2 (Z'Z)^{-1}$$

If 2 or more regressors are close to be linearly dependent  $\Rightarrow \text{Cor}(\hat{\beta})$  will explode

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_{i=1}^r (z_{ij} - \bar{z}_j)^2} \cdot \frac{1}{1-R_j^2} = \text{diag}_j (\sigma^2 (Z'Z)^{-1})$$

$R_j^2$  is coeff of det. when  $z_j$  is regressed on  $z_1, z_2, \dots, \cancel{z_j}, \dots, z_r$

If  $\sum_{i=1}^r (z_{ij} - \bar{z}_j)^2 \uparrow \Rightarrow \text{Var}(\hat{\beta}_j) \downarrow$

If  $R_j^2 \uparrow \uparrow \Rightarrow \text{Var}(\hat{\beta}_j) \uparrow$  Collinearity

Red alarm if:  $\frac{1}{1-R_j^2} = VIF = \text{Variance Inflation Factor} > 5$

Fitted model:  $y_0 = \bar{z}_0 \hat{\beta} \quad \bar{z}_0 = (1 \bar{z}_{01} \dots \bar{z}_{0r})$

$$\bar{z}_0 = (1 \bar{z}_1 \bar{z}_2 \dots \bar{z}_r) = \frac{Z' \bar{z}}{n \cdot 1}$$

$$y_0 = \frac{1' Z}{n} (Z' Z)^{-1} Z' y = \frac{1' \bar{z}}{n \cdot 1} = \frac{1' \bar{y}}{n \cdot 1} = \bar{y} = \bar{Y}$$

$$(z'^H) = (z z)^{-1} = z'$$

That is:  $\bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{z}_1 + \hat{\beta}_2 \bar{z}_2 + \dots + \hat{\beta}_r \bar{z}_r$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{z}_1 - \hat{\beta}_2 \bar{z}_2 - \dots - \hat{\beta}_r \bar{z}_r$$

$$y_0 - \bar{Y} = \hat{\beta}_1 (z_{01} - \bar{z}_1) + \dots + \hat{\beta}_r (z_{0r} - \bar{z}_r)$$

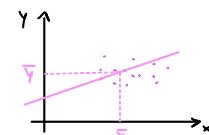
as fitted models pass through the barycenter  $(\bar{z}_1, \dots, \bar{z}_r, \bar{Y})$

$$\text{Centering: } y \rightarrow y^* = y - \bar{Y} \cdot 1 = \begin{bmatrix} y_1 - \bar{Y} \\ \vdots \\ y_n - \bar{Y} \end{bmatrix}$$

$$Z \rightarrow Z^* = \begin{bmatrix} z_{11} - \bar{z}_1 & z_{12} - \bar{z}_2 & \dots & z_{1r} - \bar{z}_r \\ \vdots & \vdots & & \vdots \\ z_{n1} - \bar{z}_1 & z_{n2} - \bar{z}_2 & \dots & z_{nr} - \bar{z}_r \end{bmatrix} \quad n \times r$$

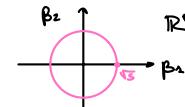
$$\text{ols: } \begin{cases} \hat{\beta}^* = \underset{\beta \in \mathbb{R}^r}{\text{argmin}} \|y^* - Z^* \beta\|^2 \\ \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{z}_1 - \dots - \hat{\beta}_r \bar{z}_r \\ \hat{\beta}_i = \hat{\beta}_i^* \quad i = 1, \dots, r \end{cases}$$

Caution: From now on  $y$  and  $Z$  ( $n \times r$ ) are centered



## Ridge regression

$$\begin{cases} \underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - z\beta\|^2 \\ \|\beta\|_2^2 \leq s \end{cases}$$

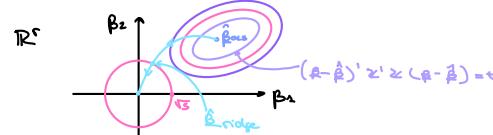


$$\|y - z\beta\|^2 = \|y - \hat{y} - z(\beta - \hat{\beta}_{OLS})\|^2 = \|z\hat{\beta}_{OLS}\|^2 = \left\| \sum_{i=1}^n z_i (\beta - \hat{\beta}_{OLS}) \right\|^2 = \|\hat{z}\|^2 + \|z(\beta - \hat{\beta}_{OLS})\|^2$$

$$\hat{y} = \hat{y} - z(z'z)^{-1}z'y$$

$$\begin{cases} \underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|z(\beta - \hat{\beta}_{OLS})\|^2 \\ \|\beta\|_2^2 \leq s \end{cases}$$

$$\boxed{\begin{cases} \underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} (\beta - \hat{\beta}_{OLS})' z' z (\beta - \hat{\beta}_{OLS}) \\ \|\beta\|_2^2 \leq s \end{cases}}$$



Consider the Lagrangian:  $\underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - z\beta\|^2 + \lambda \|\beta\|_2^2 \implies \hat{\beta}_{\text{ridge}} = (z'z + \lambda I)^{-1} z'y$

- Obs:
- $\hat{\beta}_{\text{ridge}}$  is biased
  - For any regression problem  $\Rightarrow \exists \lambda^*$  such that  $E[\|\hat{\beta}_{\text{ridge}}\|^2] \leq E[\|\hat{\beta}_{OLS}\|^2]$
  - Find the  $\lambda^*$  by cross validation

If the problem is collinear, a different sol'n could have been  $\Rightarrow$  orthog. regressor

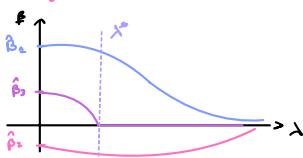
$$Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \xrightarrow{\text{PCA}} Z^* = \begin{bmatrix} p_{11} & \dots & p_{1r} \\ S_{11} & \dots & S_{1r} \\ \vdots & & \vdots \\ S_{n1} & \dots & S_{nr} \end{bmatrix}$$

$$Z^* = \begin{bmatrix} p_{11} & p_{1r} \\ S_{11} & \dots & S_{1r} \\ \vdots & & \vdots \\ S_{n1} & \dots & S_{nr} \end{bmatrix} \quad y = Z^* \beta + \varepsilon \quad \beta \in \mathbb{R}^r$$

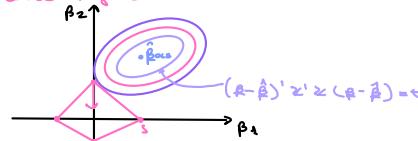
$$\begin{aligned} \hat{\beta}_{\text{PCA}} &= \underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - Z^* \beta\|^2 \\ &\begin{cases} \text{PC}_1 = e_{11} z_1 + \dots + e_{rr} z_r \\ \text{PC}_k = e_{kk} z_k + \dots + e_{rr} z_r \end{cases} \\ \hat{y} &= z_1 (\hat{\beta}_1 e_{11} + \hat{\beta}_2 e_{21} + \dots + \hat{\beta}_r e_{r1}) + \\ &+ z_2 (\hat{\beta}_1 e_{12} + \hat{\beta}_2 e_{22} + \dots + \hat{\beta}_r e_{r2}) + \dots \\ &+ z_r (\hat{\beta}_1 e_{1r} + \hat{\beta}_2 e_{2r} + \dots + \hat{\beta}_r e_{rr}) \\ &= z_1 \gamma_1 + z_2 \gamma_2 + \dots + z_r \gamma_r \end{aligned}$$

Obs: PCA-reg., ridge reg., don't have spars sol'n in terms of  $z_1, \dots, z_r$

## Ridge regression:



## Lasso regression:

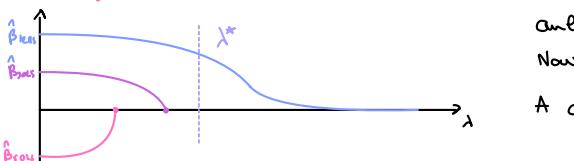


$$\begin{cases} \underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} (\beta - \hat{\beta}_{OLS})' z' z (\beta - \hat{\beta}_{OLS}) \\ \|\beta\|_1 \leq s \quad (\exists |\beta_i| < s) \end{cases}$$

Lagrangian:  $\underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - z\beta\|^2 + \lambda \sum |\beta_i|$

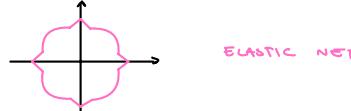
Problem: Lagrangian  $\underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - z\beta\|^2 + \lambda \sum |\beta_i|$  not identically solvable

## Lasso regression:



only one to select  
Now plenty possibilities  $\rightarrow$  boundary even more spiky  $\rightarrow \star \rightarrow \star \rightarrow \dots \rightarrow \star$   
A compromise:  $\underset{\beta \in \mathbb{R}^r}{\operatorname{argmin}} \|y - z\beta\|^2 + \lambda + \|\beta\|_1 + \lambda_2 \|\beta\|_2$

Two variables:



## ENSEMBLE METHODS

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad y_i \in \mathbb{R}^p$$

$y_i \in \mathbb{R}$   
labels (regressor)

$M_X : \mathbb{R}^p \longrightarrow \mathbb{R}^{\text{labels}}$ , given a new vector of features  $x \in \mathbb{R}^p$  we predict  $M_X(x)$

Learning  $M_X$  means dealing with bias-variance trade-off

How to reduce the variance of the prediction without losing too much in terms of bias?

We must separate concepts of measurement  $y$  and true measure  $\mu$ .

$$y = \mu + \varepsilon, \quad \mathbb{E}[\varepsilon] = 0, \quad \text{Var}(\varepsilon) = \sigma^2$$

Take one measurement:  $y_1 \leftarrow y_1 = \mu + \varepsilon_1$  with  $\mathbb{E}[y_1] = \mu$ ,  $\text{Var}(y_1) = \sigma^2$

To reduce variability we can take more meas., we take  $B$  indep. meas.

$y_1, y_2, \dots, y_B$  (realization of)

$y_1, y_2, \dots, y_B$

$$y_i = \mu + \varepsilon_i \quad i=1, \dots, B, \quad \varepsilon_i \text{ iid: } \mathbb{E}[\varepsilon_i] = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2$$

Take the average:  $\bar{y} = \frac{1}{B} \sum_{i=1}^B y_i$  realization of  $\bar{Y} = \frac{1}{B} \sum_{i=1}^B Y_i$   
 with  $\mathbb{E}[\bar{y}] = \mu$  but now  $\text{Var}(\bar{y}) = \frac{1}{B} \sigma^2$  (reduced variability)

This was the basic idea to apply to our argument

$X$   $n$  units  $\mapsto M_X$  and I generate  $X_1 \mapsto M_{X_1}$ ,  $X_2 \mapsto M_{X_2}$ , ...,  $X_B \mapsto M_{X_B}$   
 and we aggregate them  $M = \frac{1}{B} \sum_{i=1}^B M_{X_i}$

Why don't I use a big data  $X$   $B \times n$  units to learning? Not practical

### Bootstrap:

$X$  numbers training set, want to generate  $B$  date sets, pseudocode:

For  $i = 1, \dots, n$

sample the unit (row)  $u$  of  $X$

$$\rightarrow (x_i^*, y_i^*)$$

sample with replacement

$\rightarrow X^*$  obtained by resampling with replacement from  $X$ .

Repeat it for all  $B$  sets, from each of them we can extract a model and aggregate them

$\rightarrow$  BAGGING = BOOTSTRAP + AGGREGATION

check some prob:

$$\text{Unit } u \in X: \quad \mathbb{P}[u \notin X^*] = (1 - 1/n)^n \xrightarrow{n \rightarrow \infty} e^{-1} \approx 1/3$$

So  $2/3$  of sample in  $X$  can appear in  $X^*$   $\Rightarrow$  not possible to say independent

$$\text{Var}(M) > \frac{\sigma^2}{B} \quad \text{if} \quad \sigma^2 = \text{Var}(M_{X_i}) \quad \text{but in any case} \quad \sigma^2 > \text{Var}(M)$$

