

REGRESSION

Training set :

$$\mathbb{X} = \left[\begin{array}{cccc} x_1 & x_2 & \dots & x_p & y \\ x_{11} & x_{12} & \dots & x_{1p} & y_1 \\ x_{21} & x_{22} & \dots & x_{2p} & y_2 \\ \vdots & & & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} & y_n \end{array} \right] \quad \begin{array}{l} x_i \in \mathbb{R}^p \\ y_i \in \mathbb{R} \end{array}$$

X
covariates
Y
output variable

General goal: "explain" Y in terms of X
(= "explain" the distribution/variability of Y in terms of X)

explain/predict: Some models are used for prediction and don't look for the explainability, others go on the opposite way

More specifically regression is concerned with estimating:

$$\mathbb{E}[Y | X = x] = f(x)$$

unknown function:

REGRESSION FUNCTION

Use \mathbb{X} to estimate f by means of \hat{f} .

Two basic approaches:

1. Totally data driven (non-parametric) :

good for prediction if we take good care of the bias-variance trade-off (overfitting)

] it's good if we don't know much about the problem (about f)

2. Parametric approach (model based) :

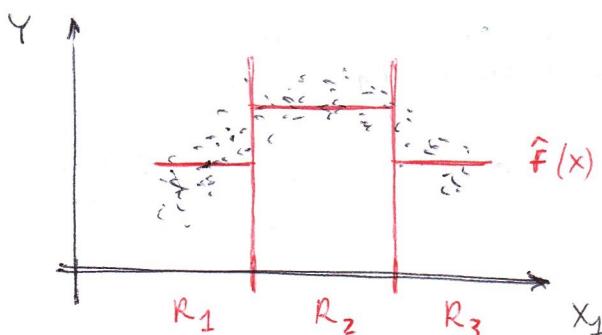
good for interpretation and prediction, quantify the uncertainty, embedding of prior knowledge

CART

CLASSIFICATION AND REGRESSION TREES (totally data driven (1))

Basic idea: \hat{f} is piecewise constant

→ partition \mathbb{R}^p (features space) in $\{R_1, \dots, R_j\}$ j finite (unknown) and assume \hat{f} to be constant over the elements of the partition



Cart takes specific values for these constants, which are the means of the y we observed in the partition.

$$\text{Let } \bar{y}_i = \frac{1}{\#R_i} \sum_{x_j \in R_i} y_i \quad i=1, \dots, j$$

$$\Rightarrow \hat{f}(x) = \sum_{i=1}^j y_i \mathbb{1}_{\{x \in R_i\}}$$

Note: $\{R_1, \dots, R_j\}$ is a finite partition of \mathbb{R}^P such that:

$$R_i \cap \{x_1, \dots, x_n\} \neq \emptyset \quad i=1, \dots, j$$

How do we find $\{R_1, \dots, R_j\}$ and j ?

$$\text{Goal: minimize}_{\{R_1, \dots, R_j\}, j} \sum_{i=1}^j \sum_{x \in R_i} (y_i - \bar{y}_i)^2$$

without overfitting the data

CART \rightarrow greedy algorithm for "solving" the problem

Iterative process: (CART)

1. Consider the cut off $s_1 \in \mathbb{R}$ such that

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{total variability}} - \left[\sum_{\{x_{j_1} \leq s_1\}} (y_i - \bar{y}_1)^2 + \sum_{\{x_{j_1} > s_1\}} (y_i - \bar{y}_2)^2 \right] \quad (*)$$

is maximized, where:

$$\bar{y}_1 = \frac{1}{\#\{j : x_{j_1} \leq s_1\}} \sum_{\{x_{j_1} \leq s_1\}} y_i$$

$$\bar{y}_2 = \frac{1}{\#\{j : x_{j_1} > s_1\}} \sum_{\{x_{j_1} > s_1\}} y_i$$

= the mean of the y 's for all the units for which the first variable (x_{j_1}) is less than s_1

What we are trying to do is:
we have the variability of the y 's,
we try to find the split with respect
to the first variable that (sort of)
makes the larger difference in
total variability *

Repeat for $x_2, x_3, \dots, x_p \Rightarrow$ it generates the splits: s_1, \dots, s_p

Choose the cutoff s_j^* which maximizes $(*)$.

Hence \mathbb{R}^P is split:

$$R_1 = \{x \in \mathbb{R}^P : x_j \leq s_j^*\}$$

$$R_2 = \{x \in \mathbb{R}^P : x_j > s_j^*\}$$

For each split we
have a value of $(*)$;
we choose the one
that maximizes $(*)$.

* We're trying to find
the split so that the
sum of the variability
within the two groups
is less than the big
variability that we
had before the split.

2. Now we iterate the same procedure on the elements of the partition. (In the first case we iterate on R_1 or R_2)

Stop partitioning an element R_j if R_j contains less than 5 units ($5/10/\dots$, just a given number of units).

The problem is overfitting.

We don't want partitions with empty regions (we wouldn't be able to compute the mean in those regions)

This reminds of K-means:
we're trying to split the x 's (to find groups among the x 's) minimizing not the variability of the x 's within the group but minimizing the variability of the y 's (the thing we want to predict)

this is the real problem
(if we create a set for each point then the objective function will be zero, this solves the optimality problem but it's obviously overfitting the data (we have to remember bias-variance trade off))

optimality problem but
it's obviously overfitting
the data (we have to
remember bias-variance
trade off))

Solution for overfitting : **PRUNING** the tree (penalize large trees) ← which are the ones that overfit data
 (we grow a large tree and then we prune it)

Choose $\alpha > 0$.

Grow a large tree following the previous procedure.

Prune bottom-up the branches of the tree with the goal of minimizing :

$$\sum_{i=1}^m \sum_{x_j \in R_i} (y_j - \bar{y}_i)^2 + \alpha m = W(m, \alpha)$$

* which mathematically means to merge the rectangles

Choose α by cross-validation.

Note: CART usually doesn't beat linear models (we can always generate a linear model that reproduces exactly the CART (it's a matter of taking the right transformations)). However it's a good thing to use CART to check what type of non-linearity the model should have.

→ this penalizes the number of branches that we got in the tree (= number of elements of the partition)

LINEAR MODELS FOR REGRESSION

$$\underline{X} = \begin{bmatrix} X_1 & X_2 & \cdots & X_p & Y \\ X_{11} & X_{12} & \cdots & X_{1p} & Y_1 \\ X_{21} & X_{22} & \cdots & X_{2p} & Y_2 \\ \vdots & & & & \\ X_{n1} & X_{n2} & \cdots & X_{np} & Y_n \end{bmatrix}$$

Design matrix :

$$\underline{Z} = \begin{bmatrix} 1 & z_1 & \cdots & z_r \\ 1 & z_{11} & \cdots & z_{1r} \\ \vdots & z_{21} & \cdots & z_{2r} \\ \vdots & & \ddots & \\ 1 & z_{n1} & \cdots & z_{nr} \end{bmatrix} \in \mathbb{R}^{n \times r}$$

RANDOM FORESTS

(Build an ensemble of trees and average them.)

We build many trees (many CARTs) by building many different training sets (by **BOOTSTRAP**: sampling from the training set with replacing). We build an ensemble of trees more or less uncorrelated (not totally correlated because we have sampled from training set and from the variables). Each of these tree will produce a prediction. We take the average of all the predictions as final prediction.

Note: We lose some interpretability of the model with this method.

z_1, \dots, z_r are known functions (so transformations) of X_1, \dots, X_p

Linear model :

$$E[Y | z_1, \dots, z_r] = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r$$

where the unknowns are β_0, \dots, β_r .

The model for Y :

given $\underline{z} = [1, z_1, \dots, z_r]^T$:
$$Y = f(\underline{z}) + \varepsilon$$

such that : $E[\varepsilon] = 0$ and $\varepsilon \perp\!\!\!\perp \underline{z}$. → ε is the part that is left out (that is not connected to informations expressed by \underline{z} (so by \underline{X}))
 with $f(\underline{z}) = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r$

Consider the "data":

$$\underline{Y} = [Y_1, \dots, Y_n]^T$$

$$\underline{Z} = \begin{bmatrix} 1 & z_{11} & z_{12} & \cdots & z_{1r} \\ 1 & z_{21} & z_{22} & \cdots & z_{2r} \\ \vdots & & & & \\ 1 & z_{n1} & z_{n2} & \cdots & z_{nr} \end{bmatrix}$$

⇒

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

where $\beta = [\beta_0, \beta_1, \dots, \beta_r]^T$,

$$\varepsilon \in \mathbb{R}^n \text{ s.t. } \left\{ \begin{array}{l} E[\varepsilon] = 0 \\ \text{Cov}(\varepsilon) = \sigma^2 I \end{array} \right.$$

! we have to say something about how the units are connected and we're saying that they're UNCORRELATED (\neq independent)

i.e.

$$Y_i = \beta_0 + \beta_1 Z_{i1} + \beta_2 Z_{i2} + \dots + \beta_r Z_{ir} + \varepsilon_i$$

with $\varepsilon_1, \dots, \varepsilon_n$:

- uncorrelated
- same variance
- mean 0
- independent of the Z 's

Linear models are very flexible.

Example: ANOVA one-way : ANOVA is a particular case of linear models

$$\begin{aligned} X_{11}, \dots, X_{1n_1} &\stackrel{iid}{\sim} N(\mu_1, \sigma^2) \\ X_{21}, \dots, X_{2n_2} &\stackrel{iid}{\sim} N(\mu_2, \sigma^2) \\ &\vdots \\ X_{g1}, \dots, X_{gn_g} &\stackrel{iid}{\sim} N(\mu_g, \sigma^2) \end{aligned} \quad \left. \right\} \text{II}$$

$$\text{let } \underline{Y} = [X_{11}, X_{12}, \dots, X_{1n_1}, X_{21}, \dots, X_{2n_2}, \dots, X_{g1}, \dots, X_{gn_g}]^\top$$

$$\underline{Y} \in \mathbb{R}^{n=n_1+n_2+\dots+n_g}$$

$$\begin{aligned} \underline{Z} &= \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & \vdots \\ \vdots & \vdots & 1 & \dots & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{bmatrix} \quad \begin{array}{c} \uparrow n_1 \\ \uparrow n_2 \\ \dots \\ \uparrow n_g \end{array} \\ (\text{design matrix}) \quad (n \times (g+1)) & \quad \text{"dummy variables"} \end{aligned}$$

$$\underline{\beta} = [\mu, \tau_1, \dots, \tau_g]^\top$$

$$\underline{Y} = \underline{Z} \underline{\beta} + \underline{\varepsilon} \quad \underline{\varepsilon} \sim N(0, \sigma^2 I)$$

i.e.

$$X_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad i=1, \dots, g \quad j=1, \dots, n_i$$

$$\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$(*) \text{ constraint: } \sum_{i=1}^g n_i \tau_i = 0 \implies \tau_g = - \sum \frac{n_i}{n_g} \tau_i$$

$$\begin{aligned} \Rightarrow \underline{Z} &= \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & \vdots \\ \vdots & \vdots & 1 & \dots & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ 0 & 1 & \vdots & \dots & \vdots \\ \vdots & \vdots & 1 & \dots & 0 \\ 0 & 0 & 1 & \dots & \vdots \\ -n_1/n_g & -n_2/n_g & -n_3/n_g & \dots & 1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -n_2/n_g & -n_3/n_g & \dots & -n_{g-1}/n_g \end{bmatrix} \quad \begin{array}{c} \uparrow n_1 \\ \uparrow n_2 \\ \dots \\ \uparrow n_{g-1} \\ \uparrow n_g \end{array} \\ (\text{$n \times g$}) & \quad \text{full rank} \end{aligned}$$

not full rank
(it corresponds to the overparametrization of the ANOVA : that's why we add a constraint (*).)

Fitting the linear model

(i.e. estimating $\beta_0, \beta_1, \dots, \beta_r, \sigma^2$)

OLS : Ordinary Least Square

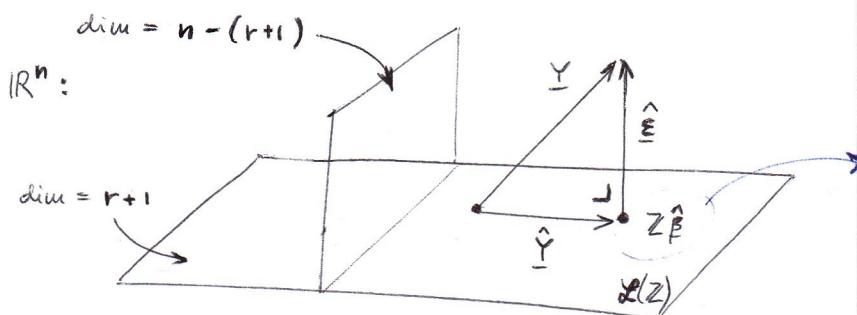
$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{r+1}} \|Y - Z\beta\|^2$$

If Z is full rank $\Rightarrow \hat{\beta} = (Z^T Z)^{-1} Z^T Y$

proof. (geometrical)

$$Y = Z\beta + \varepsilon \quad \text{model for the data}$$

$Z\beta \in \mathcal{L}(Z)$ linear space spanned by the columns of Z



(it's a linear transformation of the Y performed through a function of the design matrix)

We're trying to find the one vector $\in \mathcal{L}(Z)$ that is closest to Y in terms of Euclidean distances: this is the orthogonal projection ($\hat{Y} = Z\hat{\beta}$)

$\hat{\beta}$ is such that: $Z\hat{\beta} = \Pi_{Y|\mathcal{L}(Z)}$

in order to do the projection we need an orthonormal basis of $\mathcal{L}(Z)$

Note that:

$$Z^T Z \in \mathbb{R}^{(r+1) \times (r+1)} \quad (\text{full rank} \Rightarrow \text{spectral decomp.})$$

$$Z^T Z = \sum_{i=1}^{r+1} \lambda_i e_i e_i^T \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{r+1} > 0 \quad \text{full rank } (Z)$$

$$(Z^T Z)^{-1} = \sum_{i=1}^{r+1} \frac{1}{\lambda_i} e_i e_i^T \quad (\text{if } Z \text{ is not full rank: } Z^{-1} \rightarrow Z^* \text{ (pseudo inverse)})$$

$$\text{let } q_i = \sqrt{\frac{1}{\lambda_i}} Z e_i \quad i = 1, \dots, r+1$$

- $q_i \in \mathcal{L}(Z)$

- $q_i^T q_j = \sqrt{\frac{1}{\lambda_i \lambda_j}} e_i^T Z^T Z e_j = \frac{\lambda_j}{\sqrt{\lambda_i \lambda_j}} e_i^T e_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

$\Rightarrow \{q_1, \dots, q_{r+1}\}$ is an orthonormal basis for $\mathcal{L}(Z)$

$$\begin{aligned} \Pi_{Y|\mathcal{L}(Z)} &= \sum_{i=1}^{r+1} \Pi_{Y|q_i} = \sum_{i=1}^{r+1} \frac{q_i q_i^T}{q_i^T q_i} Y \\ &= \sum_{i=1}^{r+1} \frac{1}{\lambda_i} Z e_i e_i^T Z^T Y \end{aligned}$$

$$\begin{aligned}\Pi_{\mathcal{Z}(\mathcal{Z})} \underline{Y} &= \underline{Z} \left(\sum_{i=1}^{r+2} \frac{1}{\lambda_i} \underline{e}_i \underline{e}_i^\top \right) \underline{Z}^\top \underline{Y} \\ &\downarrow \\ &= \underline{Z} (\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y} \\ &\downarrow \\ &= \hat{\underline{Y}} \quad \text{"FITTED VALUES"}\end{aligned}$$

$$\Pi_{\mathcal{Z}(\mathcal{Z})} \underline{Y} = \underline{H} \underline{Y} = \hat{\underline{Y}} \quad H = \text{"HAT MATRIX"}$$

$$\Rightarrow \hat{\beta} = (\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y} \quad (\text{since } \underline{Z} \hat{\beta} = \hat{\underline{Y}} = \underline{Z}(\dots))$$

NOTATIONS:

$$\hat{\underline{Y}} = \underline{H} \underline{Y}$$

fitted values

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

residuals

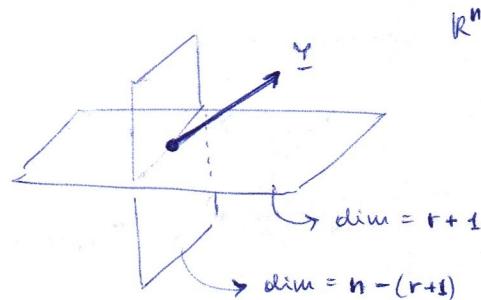
$$\Rightarrow \hat{\underline{Y}} \perp \hat{\underline{\epsilon}}$$

$$\Rightarrow \underline{Y} = \hat{\underline{Y}} + \hat{\underline{\epsilon}}$$

all we can say about \underline{Y}
in terms of the features \underline{X}

it's \perp to $\hat{\underline{Y}}$, nothing in $\hat{\underline{\epsilon}}$
can be captured by the
informations that we have from the \underline{X}

Note that:



$r+1 = n$

If we take $r+1$ too large (if we take too many variables) we're filling up \mathbb{R}^n with a linear space (generated by the columns of \underline{Z}) that will have a perfect fitting: $\underline{Y} \in \mathcal{Z}(\mathcal{Z})$, no residuals. Perfect job? ABSOLUTELY NOT. Overfitting.

Linear Models:

$$\underline{Y} = \underline{Z}\underline{\beta} + \underline{\varepsilon}$$

$\underline{Y} \in \mathbb{R}^n$ (vector of dependent variables)

$\underline{Z} \in \mathbb{R}^{n \times (r+1)}$ design matrix

$\underline{\beta} \in \mathbb{R}^{r+1}$ (unknown)

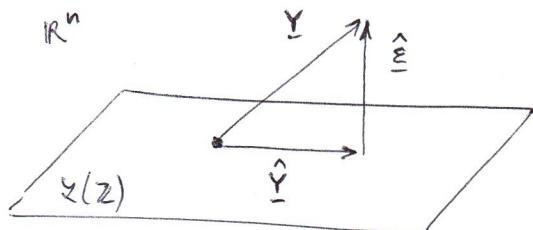
$\underline{\varepsilon}$ r.v. $\in \mathbb{R}^n$ s.t. $\begin{cases} \mathbb{E}[\underline{\varepsilon}] = \underline{0} \\ \text{(only random vector)} \end{cases}$

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

OLS:

Find the vector $\hat{\underline{Y}} \in \mathcal{L}(\underline{Z})$ "closest" to \underline{Y} .

"closest" in the sense of Euclidean distance in \mathbb{R}^n .



Why Euclidean distance?
Because $\text{Cov}(\underline{\varepsilon}) = \sigma^2 \mathbf{I}$
(otherwise Mahalanobis' distance if the $\text{Cov}(\underline{\varepsilon})$ is not diagonal)

$$\hat{\underline{Y}} = \Pi_{\mathcal{L}(\underline{Z})} \underline{Y}$$

fitted values

$$\hat{\underline{\varepsilon}} = \underline{Y} - \hat{\underline{Y}}$$

residuals

$$\hat{\underline{Y}} \perp \hat{\underline{\varepsilon}}$$

and $\underline{Y} = \hat{\underline{Y}} + \hat{\underline{\varepsilon}}$

($\neq \underline{\varepsilon}$, $\hat{\underline{\varepsilon}}$ is just what is left out after the projection. We would like $\hat{\underline{\varepsilon}}$ to be a realization of $\underline{\varepsilon}$ but we'll see that it's not possible)

Obs.

1. If rank (\underline{Z}) = $r+1 \leq n$ (full rank \underline{Z}) then:

$$\hat{\underline{Y}} = \underline{Z}(\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y} = H \underline{Y} \quad = \text{linear transfr. of } \underline{Y}$$

where H is the orthogonal projection on $\mathcal{L}(\underline{Z})$, and:

$$\hat{\underline{\beta}} = (\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y} \quad = \text{linear transfr. of } \underline{Y}$$

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

Note that:
the geometry always remains the same, it does not depend on the rank of \underline{Z}

2. If rank (\underline{Z}) = $k < r+1 \leq n$:

$$\dim(\mathcal{L}(\underline{Z})) = k < r+1$$

(the geometry is still the same!)

We have to project \underline{Y} on $\mathcal{L}(\underline{Z})$ but $\mathcal{L}(\underline{Z})$ has a dimension smaller than $r+1$: $\dim(\mathcal{L}(\underline{Z})) = k$

$$\underline{Z}^\top \underline{Z} = \sum_{i=1}^{r+1} \lambda_i e_i e_i^\top$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k = 0 = \lambda_{k+1} = \dots = \lambda_{r+1}$$

$$(\underline{Z}^\top \underline{Z})^{-1} = \sum_{i=1}^k \frac{1}{\lambda_i} e_i e_i^\top \rightarrow$$

we generate a basis for $\mathcal{L}(\underline{Z})$ which has k elements and not $r+1$

$$\Rightarrow \hat{\underline{Y}} = \underline{Z}(\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y}$$

($\exists!$ $\hat{\underline{Y}}$: projection)

$$\hat{\underline{\beta}} = (\underline{Z}^\top \underline{Z})^{-1} \underline{Z}^\top \underline{Y}$$

(this is one possible $\hat{\underline{\beta}}$)

($\exists!$ $\hat{\underline{\beta}}$: we have $r+1$ variables and k elements of the basis)

$\Rightarrow \exists \infty$ solutions : ∞ possible representations for $\hat{\underline{\beta}}$ in terms of the original variables)

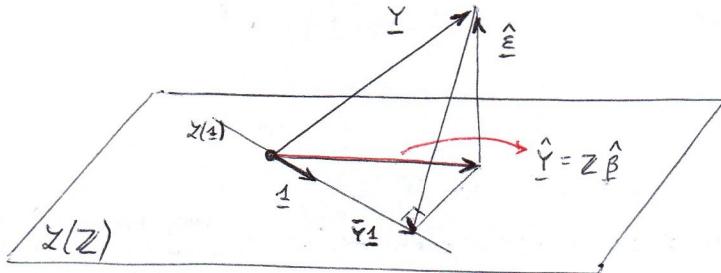
3. If $\text{rank}(Z) = r+1 = n$:

$$\Rightarrow \mathbb{X}(Z) = \mathbb{R}^n$$

$$\Rightarrow \underline{Y} = \hat{\underline{Y}} \text{ and } \hat{\underline{\varepsilon}} = \underline{0}$$

Standard way to judge how good is the fitting of the model:

COEFFICIENT OF DETERMINATION



(Assume $\text{rank}(Z) = r+1 \leq n$)

$$\|\underline{Y}\|^2 = \|\hat{\underline{Y}}\|^2 + \|\hat{\underline{\varepsilon}}\|^2$$

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \varepsilon_i^2$$

$$SS_{\text{TOT}} = SS_{\text{Reg}} + SS_{\text{Res}}$$

Note that: $\underline{1} \in \mathcal{L}(Z)$ ($\underline{1}$ is the first column of Z) and:

$$\Pi \underline{Y} | \mathcal{L}(\underline{1}) = \bar{Y} \cdot \underline{1}$$

$$\Pi \hat{\underline{Y}} | \mathcal{L}(\underline{1}) = \bar{Y} \cdot \underline{1}$$

In fact (proof.):

$$\Pi \hat{\underline{Y}} | \mathcal{L}(\underline{1}) = \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \hat{\underline{Y}} = \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} H \underline{Y}$$

$$\underline{1}^T H = (H^T \underline{1})^T = (H \underline{1})^T = \underline{1}^T \quad \text{since } \underline{1} \in \mathcal{L}(Z)$$

$$\Rightarrow \Pi \hat{\underline{Y}} | \mathcal{L}(\underline{1}) = \frac{\underline{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \hat{\underline{Y}} = \bar{Y} \cdot \underline{1}$$

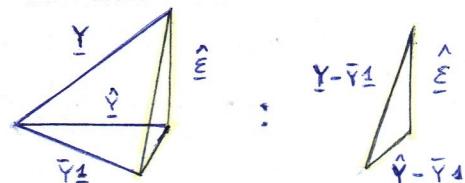
* H is an orthogonal projector
($H^T = H$)
on $\mathcal{L}(Z)$

Hence:

$$\|\underline{Y} - \bar{Y} \cdot \underline{1}\|^2 = \|\hat{\underline{Y}} - \bar{Y} \cdot \underline{1}\|^2 + \|\hat{\underline{\varepsilon}}\|^2 \rightarrow$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2$$

$$CSS_{\text{TOT}} = CSS_{\text{Reg}} + SS_{\text{Res}}$$



DECOMPOSITION OF VARIANCE
the variability of \underline{Y} is decomposed into 2 parts:

- variability accounted for by the regression
- variability that we cannot capture through the variables that we used for the model

CSS =
centered
sum of
squares

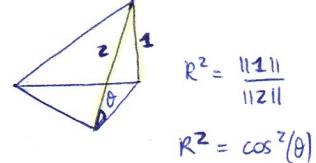
$$\Rightarrow 1 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} + \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

$$\Rightarrow R^2 = 1 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

(*) proportion of total variability explained by the regression model

$$= 1 - \sin^2(\theta)$$

since: $R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$



$$R^2 = 1 \Rightarrow \theta = 0 \Rightarrow \hat{\Sigma} = 0 \quad \text{"perfect fit"}$$

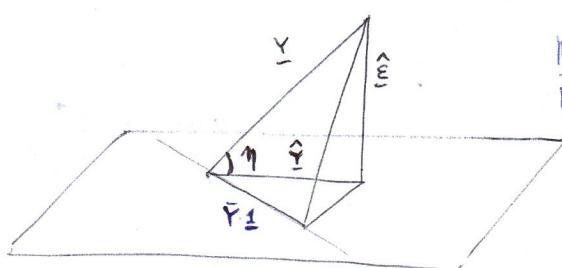
$$R^2 = 0 \Rightarrow \theta = \frac{\pi}{2} \Rightarrow \hat{Y} = \bar{Y} \cdot \underline{1} \quad Y_i = \bar{Y} \quad i = 1, \dots, n$$

Obs. If $\underline{1} \notin \Sigma(Z) \Rightarrow$ do not use R^2

(Regression through the origin: $\beta_0 = 0$)

(*) does not hold (we can compute it but it has no meaning, it may turn out negative or +∞ (for instance))

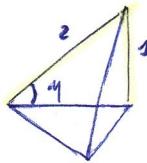
$$\text{But: } \|\underline{Y}\|^2 = \|\hat{\underline{Y}}\|^2 + \|\hat{\Sigma}\|^2$$



$$\frac{\|\hat{\Sigma}\|^2}{\|\underline{Y}\|^2} = \cos(\eta)^2 = \hat{R}^2$$

! variability around $\hat{\Sigma}$ and not variability around the mean

: measure of fit (even if it is no longer the proportion of total variability explained by the regression model): smaller the $\eta \rightarrow$ closer $\hat{\Sigma}$ to Σ



$$\text{We can compute: } \hat{R}^2 = 1 - \frac{\|\hat{\Sigma}\|^2}{\|\underline{Y}\|^2} = \frac{\|\hat{\underline{Y}}\|^2}{\|\underline{Y}\|^2}$$

Another index:

Adjusted R^2 :

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

! We want to take the ratio between two means spreading these lengths among the degrees of freedom that we have available: $\hat{\Sigma} \in \mathbb{R}^n$, $\hat{\Sigma} \in \Sigma^\perp(Z) \Rightarrow \hat{\Sigma}$ has n components but they're not all free $\dim(\Sigma^\perp(Z)) = n - (r+1)$

Properties of $\hat{\beta}$, $\hat{\Sigma}$ (Assume $\text{rank}(Z) = r+1 \leq n$)

Theorem:

$$1. E[\hat{\beta}] = \beta \quad \hat{\beta} \text{ unbiased}$$

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

$$3. E[\hat{\Sigma}] = 0$$

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

$$5. E[\hat{\Sigma}^T \hat{\Sigma}] = E\left[\sum_{i=1}^n \hat{\varepsilon}_i^2\right] = \sigma^2 (n - (r+1))$$

! $\hat{\Sigma}$ cannot be the realization of our error term Σ (the error has n free components). The same for $\underline{Y} - \bar{Y} \cdot \underline{1}$ which has n components but only $n-1$ are free ($\underline{Y} - \bar{Y} \cdot \underline{1} \in \Sigma^\perp(\underline{1})$) and $\dim(\Sigma^\perp(\underline{1})) = n-1$. (so we divide by the degrees of freedom of the vector of which we're considering the length. In this way the comparison will be fair.)

proof. (just compute)

$$1. E[\hat{\beta}] = E[(Z^T Z)^{-1} Z^T \underline{Y}] = (Z^T Z)^{-1} Z^T Z \beta = \beta$$

$$E[\underline{Y}] = E[Z \beta + \Sigma] = Z \beta + E[\Sigma] = Z \beta$$

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

$$\stackrel{\downarrow}{=} (Z^T Z)^{-1} Z^T (\sigma^2 I) \underbrace{Z}_{\text{Cov}(Y)} (Z^T Z)^{-1}$$

$$\stackrel{\downarrow}{=} \sigma^2 (Z^T Z)^{-1} \text{Cov}(Y)$$

$$3. E[\hat{\varepsilon}] = E[(I-H)Y] = E[Y] - E[HY]$$

$$\stackrel{\downarrow}{=} Z\beta - E[Z\hat{\beta}]$$

$$\stackrel{\downarrow}{=} Z\beta - Z\beta = 0$$

$$4. \text{Cov}(\hat{\varepsilon}) = \text{Cov}((I-H)Y)$$

$$\stackrel{\downarrow}{=} (I-H)(\sigma^2 I)(I-H)^T$$

$$\stackrel{\downarrow}{=} \sigma^2 (I-H)(I-H)^T \quad \text{→ } (I-H) \text{ orthogonal projector} \Rightarrow \text{idempotent and symmetric}$$

$$\stackrel{\downarrow}{=} \sigma^2 (I-H)$$

$$5. E[\hat{\varepsilon}^T \hat{\varepsilon}] = \text{Tr } E[\hat{\varepsilon}^T \hat{\varepsilon}] = E[\text{Tr}(\hat{\varepsilon}^T \hat{\varepsilon})]$$

$$\stackrel{\downarrow}{=} E[\text{Tr}(\hat{\varepsilon} \hat{\varepsilon}^T)] \quad \text{Trick : } \begin{array}{l} \hat{\varepsilon}^T \hat{\varepsilon} = \text{number} \\ \hat{\varepsilon} \hat{\varepsilon}^T = \text{matrix} \end{array}$$

$$\stackrel{\downarrow}{=} E[\text{Tr}((I-H)Y Y^T (I-H)^T)]$$

$$(I-H)Y = (I-H)(Z\beta + \varepsilon) = (I-H)\varepsilon$$

$$E[\hat{\varepsilon}^T \hat{\varepsilon}] = E[\text{Tr}((I-H)\varepsilon \varepsilon^T (I-H)^T)]$$

$$\stackrel{\downarrow}{=} \text{Tr}((I-H) \underbrace{E[\varepsilon \varepsilon^T]}_{\text{Cov}(\varepsilon) = \sigma^2 I} (I-H)^T)$$

$$\stackrel{\downarrow}{=} \sigma^2 \text{Tr}((I-H)(I-H)^T)$$

$$\stackrel{\downarrow}{=} \sigma^2 \text{Tr}(I-H)$$

$$\stackrel{\downarrow}{=} \sigma^2 (n - \text{Tr}(H))$$

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

$$\stackrel{\downarrow}{=} \text{Tr}(Z^T Z (Z^T Z)^{-1})$$

$$\stackrel{\downarrow}{=} \text{Tr}(I_{r+1}) = r+1$$

$$E[\hat{\varepsilon}^T \hat{\varepsilon}] = \sigma^2 (n - (r+1))$$

Obs.

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

If we have control of Z ($=$ of X_1, X_2, \dots) then design it in such a way that the variability of $\hat{\beta}$ is small.
(DESIGN OF EXPERIMENTS)

$$2. \mathbb{E}[\hat{\Sigma}^T \hat{\Sigma}] = \sigma^2(n-(r+1))$$

$$\mathbb{E}\left[\frac{\hat{\Sigma}^T \hat{\Sigma}}{n-(r+1)}\right] = \sigma^2 \quad \leftarrow \text{unbiased estimator for the variance}$$

$$\text{Def. } S^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n-(r+1)}$$

Corollary: S^2 is unbiased for σ^2 .

From now on : $\underline{\Sigma} \sim N_n(\underline{0}, \sigma^2 I)$

$$\underline{Y} = \underline{Z}\underline{\beta} + \underline{\Sigma} \Rightarrow \underline{Y} \sim N_n(\underline{Z}\underline{\beta}, \sigma^2 I)$$

Theorem:

If $\underline{\Sigma} \sim N_n(\underline{0}, \sigma^2 I)$:

$$1. \hat{\underline{\beta}} \text{ and } \hat{\sigma}^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n} \quad (\text{MLE estimators})$$

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

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

$$4. \hat{\underline{\Sigma}} \perp \hat{\underline{\beta}}$$

$$5. \hat{\Sigma}^T \hat{\Sigma} = \sum_{i=1}^n \hat{\Sigma}_i^2 \sim \sigma^2 \chi^2(n-(r+1))$$

since we have a model we can write the likelihood and then maximize the likelihood in terms of $\underline{\beta}$ and σ^2 (and get $\hat{\underline{\beta}}, \hat{\sigma}^2$)

The novelty is in the distributions: \underline{Y} is gaussian so $\hat{\underline{\beta}}$ and $\hat{\underline{\Sigma}}$ and σ^2 are gaussian (since they're linear functions of \underline{Y})

proof.

1. First compute derivatives of log likelihood

2-3-4. Note that:

$$\begin{bmatrix} \hat{\underline{\beta}} \\ \hat{\underline{\Sigma}} \end{bmatrix} = \begin{bmatrix} (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \\ I - \underline{Z}(\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \end{bmatrix} \underline{Y} \quad \text{and } \underline{Y} \sim N_n(\underline{0}, \sigma^2 I)$$

Verify that:

$$\text{Cov}\left(\begin{bmatrix} \hat{\underline{\beta}} \\ \hat{\underline{\Sigma}} \end{bmatrix}\right) = \sigma^2 \begin{bmatrix} (\underline{Z}^T \underline{Z})^{-1} & \underline{0} \\ \underline{0} & I-H \end{bmatrix}$$

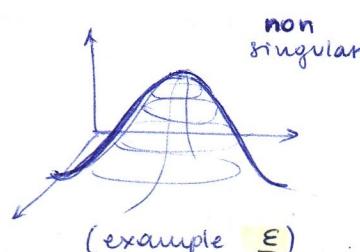
$$5. \hat{\underline{\Sigma}} \sim N_n(\underline{0}, \sigma^2(I-H)), \det(I-H) = 0$$

SINGULAR
GAUSSIAN
DISTRIBUTION

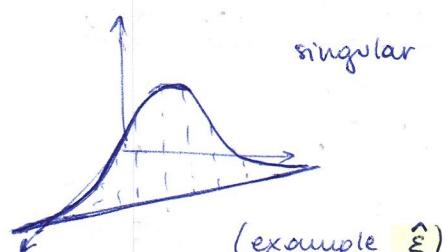
rank(I-H) = n-(r+1) : $\hat{\underline{\Sigma}}$ lies in \mathbb{R}^n but it has $n-(r+1)$ d.o.f.

Example of Singular Gaussian distribution: gaussian distribution defined only on a linear subspace of the entire space

\mathbb{R}^2 :



vs.



We would like to take: $\hat{\Sigma}^T (\mathbf{I} - \mathbf{H})^{-1} \hat{\Sigma}$ for the Mahalanobis' distance between $\hat{\Sigma}$ and its mean (which is 0)
 we cannot do that: $\det(\mathbf{I} - \mathbf{H}) = 0$
 \Rightarrow we consider: $\hat{\Sigma}^T (\mathbf{I} - \mathbf{H})^{-1} \hat{\Sigma} \sim \sigma^2 \chi^2(n - (r+1))$ (where: $\hat{\Sigma}^T (\mathbf{I} - \mathbf{H})^{-1} \hat{\Sigma} \sim \sigma^2 \chi^2(n)$)
 \Rightarrow (Mahalanobis' distance of Gaussian from its mean)
 $\hat{\Sigma}^T \hat{\Sigma} \sim \sigma^2 \chi^2(n - (r+1))$

Now: we want the inference:

$$\left. \begin{aligned} & (\hat{\beta} - \beta)^T (\sigma^2 (\mathbf{Z}^T \mathbf{Z})^{-1})^{-1} (\hat{\beta} - \beta) \sim \chi^2(r+1) \\ & \hat{\Sigma}^T \hat{\Sigma} \sim \sigma^2 \chi^2(n - (r+1)) \end{aligned} \right\} \text{Mahalanobis' distance of } \hat{\beta} \text{ from its mean}$$

$$\begin{aligned} & \frac{\frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (\mathbf{Z}^T \mathbf{Z}) (\hat{\beta} - \beta)}{r+1} \sim F(r+1, n - (r+1)) \\ & \xrightarrow[\text{since they're } \perp \perp]{} \frac{\hat{\Sigma}^T \hat{\Sigma}}{\sigma^2 (n - (r+1))} \end{aligned}$$

$$\Rightarrow \boxed{\frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (\mathbf{Z}^T \mathbf{Z}) (\hat{\beta} - \beta) \sim (r+1) F(r+1, n - (r+1))}$$

$$S^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n - (r+1)}$$

Confidence regions for β :

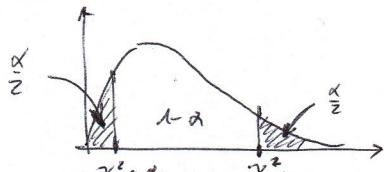
$$CR_{1-\alpha}(\beta) = \left\{ \beta \in \mathbb{R}^{r+1} : \frac{1}{\sigma^2} (\hat{\beta} - \beta)^T (\mathbf{Z}^T \mathbf{Z}) (\hat{\beta} - \beta) \leq (r+1) F_{\alpha}(r+1, n - (r+1)) \right\}$$

Confidence intervals for σ^2 :

$$\hat{\Sigma}^T \hat{\Sigma} \sim \sigma^2 \chi^2(n - (r+1)) \Rightarrow \frac{1}{\sigma^2} \hat{\Sigma}^T \hat{\Sigma} \sim \chi^2(n - (r+1))$$

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

$$\begin{aligned} CI_{1-\alpha}(\sigma^2) &= \left\{ \sigma^2 \in \mathbb{R} : \chi_{1-\frac{\alpha}{2}}(n - (r+1)) \leq \frac{(n - (r+1)) S^2}{\sigma^2} \leq \chi_{\frac{\alpha}{2}}^2(n - (r+1)) \right\} \\ &= \left\{ \sigma^2 \in \mathbb{R} : \frac{(n - (r+1)) S^2}{\chi_{\frac{\alpha}{2}}^2(n - (r+1))} \leq \sigma^2 \leq \frac{(n - (r+1)) S^2}{\chi_{1-\frac{\alpha}{2}}^2(n - (r+1))} \right\} \end{aligned}$$



Simultaneous CI for $a^T \beta$, $a \in \mathbb{R}^{r+1}$ (for linear combinations of β_i)

$$\max_{a \in \mathbb{R}^{r+1}} \frac{1}{S^2} \frac{(a^T (\hat{\beta} - \beta))^2}{a^T (\mathbf{Z}^T \mathbf{Z})^{-1} a} = \frac{1}{S^2} (\hat{\beta} - \beta)^T (\mathbf{Z}^T \mathbf{Z}) (\hat{\beta} - \beta) \sim (r+1) F(r+1, n - (r+1))$$

$$\Rightarrow \text{Sim CI}_{1-\alpha}(a^T \beta) = \left[a^T \hat{\beta} \pm \sqrt{a^T (\mathbf{Z}^T \mathbf{Z})^{-1} a} \sqrt{S^2(r+1) F_{1-\alpha}(r+1, n - (r+1))} \right]$$

$$\underline{Y} = \underline{\beta} \underline{Z} + \underline{\varepsilon}$$

$\underline{Z} \in \mathbb{R}^{n \times (r+1)}$ design matrix

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

$$\underline{\varepsilon} \sim N_n(\underline{0}, \sigma^2 \underline{I}) \Rightarrow \underline{Y} (\in \mathbb{R}^n) \sim N_n(\underline{Z}\underline{\beta}, \sigma^2 \underline{I})$$

$$\hat{\underline{\beta}} \sim N_{r+1}(\underline{\beta}, \sigma^2 (\underline{Z}^\top \underline{Z})^{-1})$$

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

$$\hat{\underline{\varepsilon}}^\top \hat{\underline{\varepsilon}} \sim \sigma^2 \chi^2(n-(r+1))$$

$$\Rightarrow \frac{1}{s^2} (\hat{\underline{\beta}} - \underline{\beta})^\top (\underline{Z}^\top \underline{Z}) (\hat{\underline{\beta}} - \underline{\beta}) \sim (r+1) F_{(r+1, n-(r+1))}$$

$$\Rightarrow CR_{1-\alpha}(\underline{\beta}), \text{ Sim CI}_{1-\alpha}(\underline{a}^\top \underline{\beta}) \quad \forall \underline{a} \in \mathbb{R}^{r+1}$$

$$\text{Sim CI}_{1-\alpha}(\underline{a}^\top \underline{\beta}) = \left[\underline{a}^\top \hat{\underline{\beta}} \pm \sqrt{\underline{a}^\top (\underline{Z}^\top \underline{Z})^{-1} \underline{a}} \sqrt{s^2(r+1) F_\alpha(r+1, n-(r+1))} \right]$$

simultaneous CI_{1-α} for any linear comb. of the β parameters

Special case:

$$\text{Sim CI}_{1-\alpha}(\beta_i) \quad i = 0, 1, \dots, r \quad (\text{it's like } \underline{a} = [0 \dots 0 \underset{i}{1} 0 \dots 0]^\top)$$

$$\text{Sim CI}_{1-\alpha}(\beta_i) = \left[\hat{\beta}_i \pm \sqrt{\text{diag}_i(\underline{Z}^\top \underline{Z})^{-1}} \sqrt{s^2(r+1) F_\alpha(r+1, n-(r+1))} \right]$$

Obs.! All the software packages for regression compute one-at-the-time confidence intervals (not simultaneous w.r.t. $\forall \underline{a}$) for β_i :

$$\left[\hat{\beta}_i \pm \sqrt{\text{diag}_i(\underline{Z}^\top \underline{Z})^{-1}} \sqrt{s^2} t_{\frac{\alpha}{2}}(n-(r+1)) \right] \leftarrow \begin{array}{l} \text{(the CI that we get} \\ \text{after fixing } \underline{a}, \text{ so} \\ \text{CI } (\underline{a}^\top \underline{\beta}) \text{ for a fixed } \underline{a} \end{array}$$

Obvious compromise: use Bonferroni CI's

d → ask for α/k (Bonferroni's correction)

k = # of β's for which we want CI's

Testing the β's (test for a set of the parameters)

$$\begin{cases} H_0: C\underline{\beta} = \underline{0} \\ H_1: C\underline{\beta} \neq \underline{0} \end{cases} \quad \left(\text{more generally: } \begin{cases} H_0: C\underline{\beta} = \underline{k}_0 \\ H_1: C\underline{\beta} \neq \underline{k}_0 \end{cases} \right)$$

where $C \in \mathbb{R}^{p \times (r+1)}$ matrix given.

We're testing p different linear combinations

$$C\underline{\beta} = \begin{bmatrix} C_{11}\beta_0 + C_{12}\beta_1 + \dots + C_{1(r+1)}\beta_r \\ C_{21}\beta_0 + C_{22}\beta_1 + \dots + C_{2(r+1)}\beta_r \\ \vdots \\ C_{p1}\beta_0 + C_{p2}\beta_1 + \dots + C_{p(r+1)}\beta_r \end{bmatrix}$$

Meaning of β_i :
 Every β_i is the derivative of our model wrt the variable:
 $\beta_i = \frac{\partial Y}{\partial Z_i}$
 = how much $E[Y]$ is increasing by increasing in 1 unit in Z_i

$\hat{C\beta}$ estimator for $C\beta$.

Under H_0 :

$$\left. \begin{array}{l} \hat{C\beta} \sim N_p(0, \sigma^2 C(Z^\top Z)^{-1} C) \\ \hat{\varepsilon}^\top \hat{\varepsilon} \sim \sigma^2 \chi^2(n-(r+1)) \end{array} \right\} \perp \quad \begin{array}{l} \text{(since the first depends on } \hat{\beta} \text{ and} \\ \text{the second depends on } \hat{\varepsilon} \text{ and the two} \\ \text{vectors are stochastically independent)} \end{array}$$

$$\frac{(C\hat{\beta})^\top (C(Z^\top Z)^{-1} C)^{-1} (C\hat{\beta})}{\sigma^2} \sim \chi^2$$

$$\Rightarrow \frac{\frac{1}{\sigma^2} (\hat{C\beta})^\top (C(Z^\top Z)^{-1} C)^{-1} (\hat{C\beta})}{\frac{\hat{\varepsilon}^\top \hat{\varepsilon}}{\sigma^2(n-(r+1))}} \sim F(p, n-(r+1))$$

$$\Rightarrow \frac{1}{\sigma^2} (\hat{C\beta})^\top (C(Z^\top Z)^{-1} C)^{-1} (\hat{C\beta}) \sim p F(p, n-(r+1))$$

Reject H_0 at level α if :

$$\frac{1}{\sigma^2} (\hat{C\beta})^\top (C(Z^\top Z)^{-1} C)^{-1} (\hat{C\beta}) > p F_\alpha(p, n-(r+1))$$

We're basically taking the Mahalanobis' distance of $\hat{C\beta}$ from its mean (0) and H_0 assumes that if this distance is too big we have evidence against H_0

Special case:

$$\left\{ \begin{array}{l} H_0: \beta_r = \beta_{r-1} = \dots = \beta_{r-(p-1)} = 0 \\ H_1: \exists \beta_j \neq 0 \quad j = r-(p-1), \dots, r-1, r \end{array} \right. : \begin{array}{l} \text{the last regressors of the model can} \\ \text{be selected out? (Taken out of the model)} \\ \text{we would get a simpler model} \end{array}$$

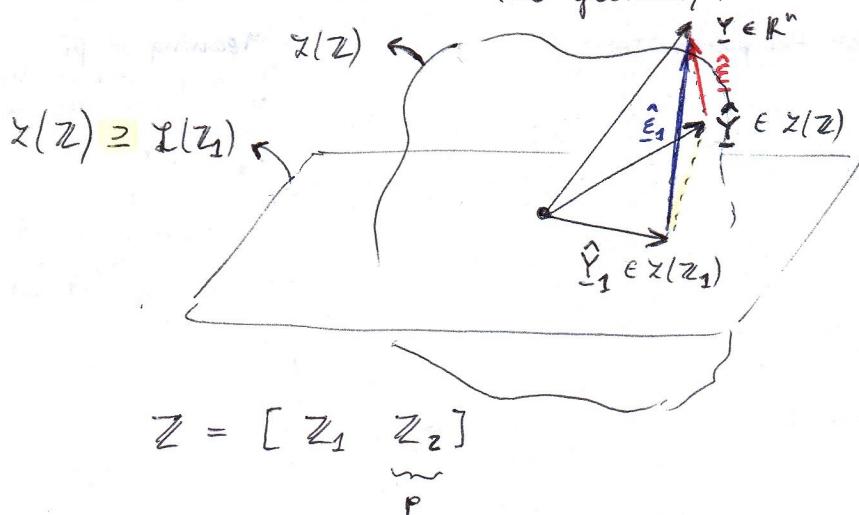
(i.e. Can we select out of the model the regressors $Z_{r-(p-1)}, \dots, Z_r$?)

Why? We would gain degrees of freedom for estimating the variability (bias-variance trade off). We would gain degrees of freedom for the space where $\hat{\varepsilon}$ lives (\Rightarrow less uncertainty about $\hat{\varepsilon}$).

\Rightarrow take C such that:

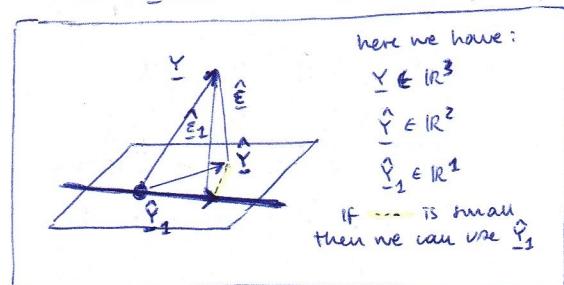
$$C = \begin{bmatrix} 0 & \dots & 0 & 1 & 1 & \dots & 1 \\ \vdots & & & & & & \\ 0 & \dots & 0 & & & & 1 \end{bmatrix} = [0 \quad I_p] \in \mathbb{R}^{p \times (r+1)}$$

Let's have a look at the geometry:



Comparing: { model $Y = Z\beta + \varepsilon$
model $Y = Z_1\beta_1 + \varepsilon_1$

When the two models are the same?
When ε is small (ε the residuals and $\hat{\varepsilon}_1$ are basically the same)



We're comparing 2 models for explaining the same Y (the two are nested: Z_1 comes from Z , obtained by setting to 0 some of the parameters)

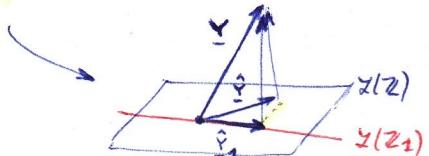
Reject $H_0: \beta_1 = \beta_2 = \dots = \beta_r = 0$ if:

$$SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z) = \hat{\Sigma}_1^T \hat{\Sigma}_1 - \hat{\Sigma}^T \hat{\Sigma}$$

In fact:

$$\frac{SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z)}{S^2 \cdot p} \sim F(p, n-(r+1))$$

$$S^2 = \frac{\hat{\Sigma}^T \hat{\Sigma}}{n-(r+1)} = \frac{SS_{\text{res}}(Z)}{n-(r+1)}$$



the two models are different (\underline{Y} and \underline{Y}_1) if S^2 is big.

Very special case:

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_r = 0 \\ H_1: \exists \beta_j \neq 0 \quad j=1, \dots, r \end{cases} \rightarrow \text{is it worth doing this regression or if we take the mean of } \underline{Y} \text{ it'll be enough? } (\underline{\beta}_0)$$

$$\rightarrow Z_1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \rightarrow \hat{\Sigma}_1^T \hat{\Sigma}_1 = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\frac{SS_{\text{res}}(Z_1) - SS_{\text{res}}(Z)}{S^2 r} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\frac{\sum_{i=1}^n \varepsilon_i^2}{n-(r+1)} \cdot r}}{\frac{\sum_{i=1}^n \varepsilon_i^2}{n-(r+1)}} \sim F(r, n-(r+1))$$

(classical F test produced by any software for regression) \rightarrow

PREDICTION

Which, btw, is not the basic goal: the basic goal of regression is to understand the link between variables (covariates) and the dependent variables (\underline{Y}) (which we want to understand how to control w.r.t. variables)

that's the test that comes out with the summary: it checks if the variability explained by the model wrt the variability that is residual is big enough to guarantee that the model has some meaning or not

$$\underline{Y} = \underline{Z} \underline{\beta} + \underline{\varepsilon}$$

model for the training set

ASSUMPTIONS:

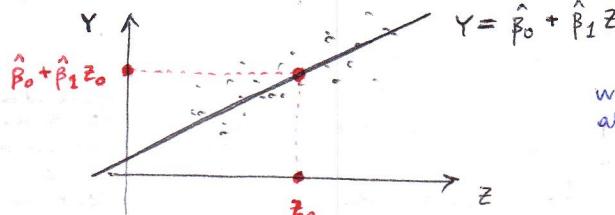
We take a new case (new statistical unit) for which we know the Z 's but we want to use the model to predict \underline{Y} or at least the mean of \underline{Y} :

Let $Z_0 = [1 \ Z_{01} \ \dots \ Z_{0r}]^T$ is a given vector capturing the values of the regressors for which we want the prediction of \underline{Y} .

$$\text{new } \underline{Y} \quad (\text{not used in the training}) \quad \underline{Y}_0 = Z_0^T \underline{\beta} + \varepsilon_0 \quad (\text{new statistical unit})$$

$$\varepsilon_0 \perp \& \text{ independent } \underline{\varepsilon}$$

$$\varepsilon_0 \sim N(0, \sigma^2)$$



what is the uncertainty about this prediction?

Recall: $\mathbb{E}[Y_0 | \mathbb{Z}_0] = \mathbb{Z}_0^T \beta$ \Rightarrow we're predicting the MEAN of Y, not the Y using the line

Natural estimator for $\mathbb{Z}_0^T \beta$ is : $\mathbb{Z}_0^T \hat{\beta}$ (Why? Gauss-Markov theorem)

Gauss-Markov theorem:

$\mathbb{Z}_0^T \hat{\beta}$ is the best estimator (minimum variance) among those:

- linear
- unbiased

$\mathbb{Z}_0^T \hat{\beta}$ is BLUE for $\mathbb{Z}_0^T \beta$.

$$\mathbb{Z}_0^T \hat{\beta} \sim N_1(\mathbb{Z}_0^T \beta, \sigma^2 \mathbb{Z}_0^T (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}_0)$$

* and $\mathbb{Z}_0^T \hat{\beta}$ is a linear transformation of $\hat{\beta}$

$$\text{because: } \hat{\beta} \sim N_{r+1}(\beta, \sigma^2 (\mathbb{Z}^T \mathbb{Z})^{-1}) *$$

$$\hat{\epsilon}^T \hat{\epsilon} \sim \sigma^2 \chi^2(n-(r+1))$$

$$|| \quad (\mathbb{Z}_0^T \hat{\beta} \perp \hat{\epsilon}^T \hat{\epsilon})$$

$$\frac{\mathbb{Z}_0^T \hat{\beta} - \mathbb{Z}_0^T \beta}{\sqrt{\sigma^2 \mathbb{Z}_0^T (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}_0}} \sim N(0, 1)$$

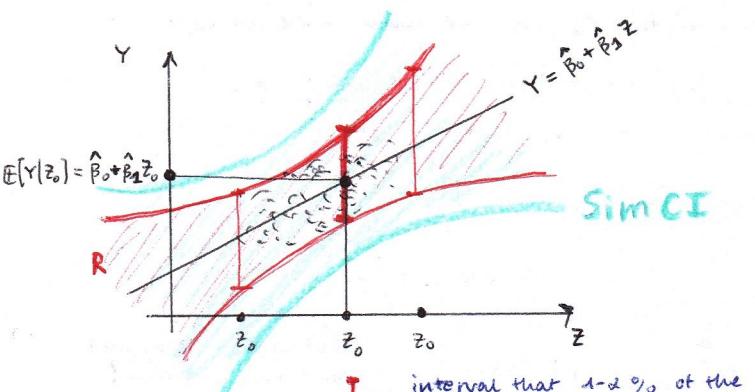
$$\sim t(n-(r+1))$$

(*) So we'll have confidence intervals but also a region that covers the linear model for all possible values of \mathbb{Z}_0

$$\frac{\mathbb{Z}_0^T \hat{\beta} - \mathbb{Z}_0^T \beta}{S \sqrt{\mathbb{Z}_0^T (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}_0}} \sim t(n-(r+1))$$

we can use this for building up confidence intervals for the MEAN of the Y corresponding to a new setting of our regressors (it's like a special case of CI for linear combinations of the β 's)

$$\begin{aligned} CI_{1-\alpha}(\mathbb{Z}_0^T \beta) &= CI_{1-\alpha}(\mathbb{E}[Y_0 | \mathbb{Z}_0]) \\ &= \left[\underbrace{\mathbb{Z}_0^T \hat{\beta}}_{\text{prediction}} \pm S \sqrt{\mathbb{Z}_0^T (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}_0} \underbrace{t_{\frac{n}{2}}}_{\text{uncertainty}} (n-(r+1)) \right] \quad \alpha \in (0, 1) \end{aligned}$$



$CI_{1-\alpha}(\mathbb{Z}_0^T \beta)$: interval that $1-\alpha$ % of the time that is used cover the right value for the mean of Y given \mathbb{Z}_0

We generate our $CI_{1-\alpha}(\mathbb{Z}_0^T \beta)$ for any possible \mathbb{Z}_0 . Note that: the closer we are to the bivariate center, the smaller are these intervals (I). That's because of " $\mathbb{Z}_0^T (\mathbb{Z}^T \mathbb{Z})^{-1} \mathbb{Z}_0$ " which makes the uncertainty larger when we go further away from the data.

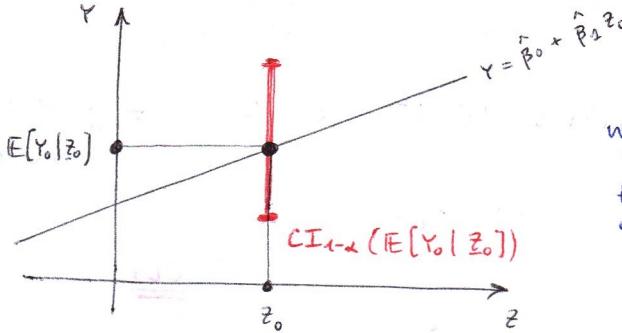
Note that the band R is the region of the CI for our predictions when \mathbb{Z}_0 varies but it doesn't mean that $1-\alpha$ % of the time the regression line will fall in R because THE IC ARE NOT SIMULTANEOUS! (they're only one-at-the-time!) Each one of them is generated by an algorithm that $1-\alpha$ % of the time covers the right value for the prediction of the mean but we're taking an α number of them here so we cannot expect that that band will cover with confidence $1-\alpha$ the the regression line.

→ SIMULTANEOUS CI (*)

$$\text{Sim CI}_{1-\alpha}(\underline{z}_0^T \hat{\beta}) = \text{Sim CI}_{1-\alpha}(\mathbb{E}[Y_0 | \underline{z}_0])$$

$$= \left[\underline{z}_0^T \hat{\beta} \pm S \sqrt{\underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0} \sqrt{(r+1) F_{\alpha}(r+1, n-(r+1))} \right]$$

This is prediction for $\mathbb{E}[Y_0 | \underline{z}_0]$, meaning:



$$Y_0 \sim N(\underline{z}_0^T \hat{\beta}, \sigma^2)$$

Can we find an interval s.t. $P(Y_0 \in I | \underline{z}_0) = 1 - \alpha$?
 (= prediction interval)

$$Y_0 \sim N_1(\underline{z}_0^T \hat{\beta}, \sigma^2)$$

$$\underline{z}_0^T \hat{\beta} \sim N_1(\underline{z}_0^T \beta, \sigma^2 \underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0)$$

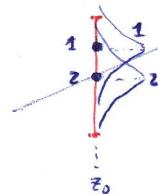
$$\rightarrow Y_0 - \underline{z}_0^T \hat{\beta} \sim N(0, \sigma^2 (1 + \underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0))$$

$$\hat{\Sigma}^T \hat{\Sigma} \sim \sigma^2 \chi^2(n-(r+1))$$

- we guarantee that, no matter how many predictions, even 1 billion of \underline{z}_0 , we're sure that the confidence level is still $1-\alpha\%$.
- we know that within the region () will be the TRUE linear model with confidence $1-\alpha\%$ of the time

What about Y_0 ? Not its mean. What do we know about Y_0 ? $Y_0 \sim N(\underline{z}_0^T \hat{\beta}, \sigma^2)$ so:

the $\text{CI}_{1-\alpha}(\mathbb{E}[Y_0 | \underline{z}_0])$ is predicting the mean of the distribution of Y_0 , we're not predicting Y_0 .



If the mean (the real mean) is in Y at z_0 then the distribution of Y_0 will be the 1, otherwise if the mean is z the distribution of Y_0 will be the 2.

→ predicting the mean is not equivalent to predicting the Y (the true observation that we would make for a statistical unit).

what this $\text{CI}_{1-\alpha}$ is capturing is the uncertainty about the mean of the distribution generating that Y , not about Y .

We know already (looking at both of pictures) that this interval will be larger than CI. We're already uncertain about where the center of this distribution is, then Y is distributed around the center ⇒ extra variability

II (since $\epsilon_0 \perp\!\!\!\perp \Sigma$)

III (since $\Sigma \perp\!\!\!\perp \hat{\beta}$)

$$\frac{Y_0 - \underline{z}_0^T \hat{\beta}}{\sqrt{\sigma^2 (1 + \underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0)}} \sim N(0, 1)$$

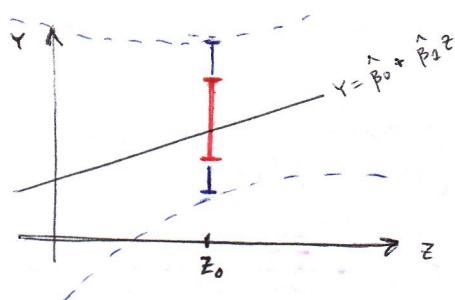
$$\frac{Y_0 - \underline{z}_0^T \hat{\beta}}{\sqrt{\frac{\hat{\Sigma}^T \hat{\Sigma}}{\sigma^2 (n-(r+1))}}} \sim t(n-(r+1))$$

$$\frac{Y_0 - \underline{z}_0^T \hat{\beta}}{S \sqrt{1 + \underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0}} \sim t(n-(r+1))$$

$$\text{PI}_{1-\alpha} = \left[\underline{z}_0^T \hat{\beta} \pm S \sqrt{1 + \underline{z}_0^T (\underline{Z}^T \underline{Z})^{-1} \underline{z}_0} t_{\frac{\alpha}{2}}(n-(r+1)) \right]$$

prediction interval of PROBABILITY $1-\alpha$

($1-\alpha$ is the probability that Y_0 belongs to $\text{PI}_{1-\alpha}$)



I: CI_{1-\alpha} ($E[Y_0 | Z_0]$)

I: PI_{1-\alpha} (Y_0)

it takes care of the fact that Y has a variability around its mean \Rightarrow the length of the prediction interval is bigger than the confidence interval (one is for Y_0 , the other for $E[Y_0 | Z_0]$)

Are these simultaneous prediction intervals? NO.
If we plot all of them (----) we will not find a prediction region for Y . If we want to find a prediction region for Y :
SIMULTANEOUS PREDICTION INTERVALS

Exercise: find the simultaneous prediction intervals (not easy)

A short excursion on **GLS** (GENERALIZED LEAST SQUARES) = what happen if we have correlated error

$$\underline{Y} = \mathbb{Z}\beta + \underline{\epsilon} \quad \left. \begin{array}{l} E[\underline{\epsilon}] = 0 \\ \text{Cov}(\underline{\epsilon}) = \sigma^2 I \end{array} \right\} \Rightarrow \text{OLS}$$

$$\hat{\beta} = \underset{\beta}{\text{argmin}} \| \underline{Y} - \mathbb{Z}\beta \|^2 = \underset{\beta}{\text{argmin}} (\underline{Y} - \mathbb{Z}\beta)^T (\underline{Y} - \mathbb{Z}\beta)$$

GLS: Assume $W \in \mathbb{R}^{n \times n}$ positive definite (that sort of control the metric that we want to use in the space of β 's instead of taking the euclidean metric)

$$\hat{\beta} = \underset{\beta}{\text{argmin}} (\underline{Y} - \mathbb{Z}\beta)^T W (\underline{Y} - \mathbb{Z}\beta)$$

$$(\underline{Y} - \mathbb{Z}\beta)^T W (\underline{Y} - \mathbb{Z}\beta) = (\tilde{W}^{1/2} \underline{Y} - \tilde{W}^{1/2} \mathbb{Z}\beta)^T (\tilde{W}^{1/2} \underline{Y} - \tilde{W}^{1/2} \mathbb{Z}\beta) \quad \begin{array}{l} \downarrow \\ \| \tilde{W}^{1/2} \underline{Y} - \tilde{W}^{1/2} \mathbb{Z}\beta \|^2 \end{array} \approx (\| \underline{Y} - \mathbb{Z}\beta \|^2)$$

$$\Rightarrow \hat{\beta} = (\mathbb{Z}^T W^{-1} \mathbb{Z})^{-1} \mathbb{Z}^T W^{-1} \underline{Y}$$

In fact: notice that:

$$W^{-1/2} \underline{Y} = W^{-1/2} \mathbb{Z}\beta + W^{-1/2} \underline{\epsilon}$$

$$\text{Cov}(W^{-1/2} \underline{\epsilon}) = \sigma^2 W^{-1}$$

Example: Assume that $\underline{\epsilon}$ is such that $\text{Cov}(\underline{\epsilon}) = \sigma^2 \Sigma$ (σ^2 unknown and Σ known).

Take $W = \Sigma$:

$$\hat{\beta} = \underset{\beta}{\text{argmin}} (\underline{Y} - \mathbb{Z}\beta)^T \Sigma^{-1} (\underline{Y} - \mathbb{Z}\beta) \quad \begin{array}{l} = \text{minimize wrt} \\ \text{the Mahalanobis' distance} \end{array}$$

$$\downarrow (\mathbb{Z}^T \Sigma^{-1} \mathbb{Z})^{-1} \mathbb{Z}^T \Sigma^{-1} \underline{Y}$$

$$\text{Note: } \text{Cov}(\Sigma^{-1/2} \underline{\epsilon}) = \sigma^2 \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = \sigma^2 I \rightarrow$$

we have transformed the variables s.t. the error term is again uncorrelated with the same variance along the components

Special cases:

$$\text{Cov}(\underline{\epsilon}) = \sigma^2 \text{diag}(w_1, \dots, w_n), \text{ i.e. } \Sigma = \begin{bmatrix} w_1 & & \\ & \ddots & \\ & & w_n \end{bmatrix} \sigma^2$$

WEIGHTED LEAST SQUARES (special case of GLS)

we believe that the statistical units do not share the same variance (they have different variances). σ^2 has a different weight along units.

Example: • suppose $Y_i \quad i=1, \dots, n$ is a mean of n_i observations -

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

$$\Rightarrow W = \text{diag} \left(\frac{1}{n_1}, \dots, \frac{1}{n_n} \right)$$

• suppose $Y_i \quad i=1, \dots, n$ is the sum of n_i observations

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

$$\Rightarrow W = \text{diag}(n_1, \dots, n_n)$$

it's not just a matter of rescaling wrt. the number of people, but also taking into account that the uncertainty that we have about these estimates is pretty different (because the variability is different)

For instance: we observe the average income but the unit is a municipality: average income for Milano, for Torino, ... We know that the number of people in those units is different so the averages have different variabilities.

For instance: in every municipality we observe a number of people having a certain disease. It's not the average for person, is the number of people having the disease (so it's a sum of Bernoulli's variables, each one may have the same variability but we're observing a different number of people (in Milano, Lecco, Torino, ...))

(# Covid in Italy cannot be directly compared with # Covid in USA (different sizes))

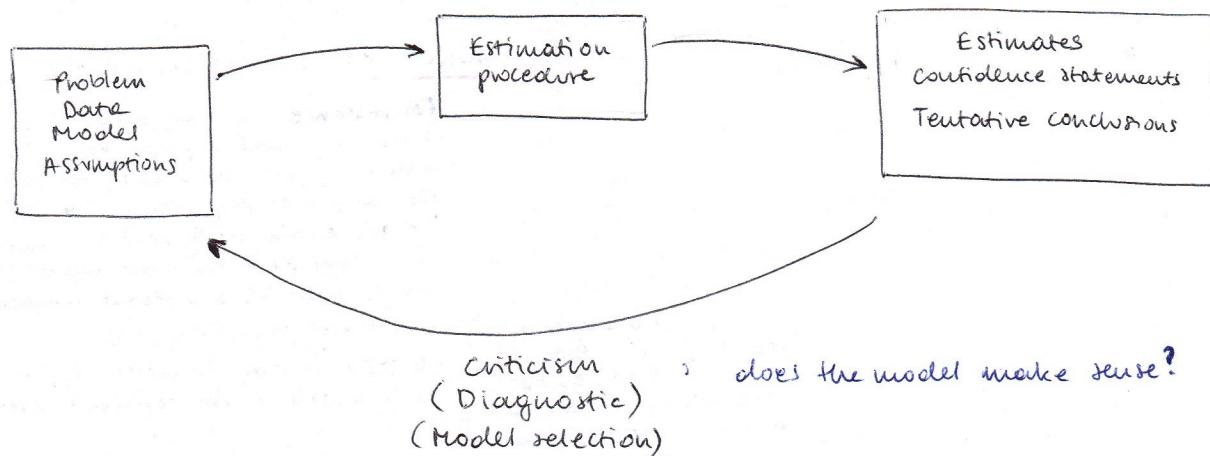
DIAGNOSTIC

: what we do after the analysis to check how good is the model

12/05

(R^2 is not enough)

Box:



Diagnostic for Linear Models

- Residual analysis: outliers, non-constant variance (heteroscedasticity), normality, autocorrelation, ...
- Influential cases or statistical units
- Collinearity among the regressors

Intro: ε vs. $\hat{\varepsilon}$

$$\underline{Y} = Z \beta + \varepsilon$$

$$\varepsilon \in \mathbb{R}^n : \begin{cases} E[\varepsilon] = 0 \\ \text{Cov}(\varepsilon) = \sigma^2 I \end{cases}$$

$$(\varepsilon \sim N_n(0, \sigma^2 I)) \text{ not necessary}$$

What do we know about $\hat{\varepsilon}$? (what changes between ε and $\hat{\varepsilon}$?)

$$\hat{\varepsilon} = (I - H) \underline{Y} = T \underline{Y} |_{\text{range}(Z)} . \quad \hat{\varepsilon} \in \text{range}(Z) \subseteq \mathbb{R}^n$$

$$\left| E[\hat{\varepsilon}] = 0 \right.$$

($\hat{\varepsilon}$ has n components but only $\dim(\text{range}(Z))$ are independent)

$$\left| \text{Cov}(\hat{\varepsilon}) = \sigma^2 (I - H) \right. \longrightarrow$$

the residuals are correlated even if the components of the error (ε) are not correlated

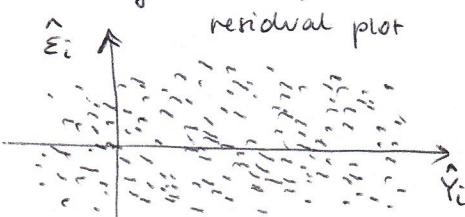
$$(\hat{\varepsilon} \sim N_n(0, \sigma^2 (I - H))) \text{ not necessary}$$

$\hat{\varepsilon}$ is gaussian on a subspace of \mathbb{R}^n (singular)

$\hat{\varepsilon}$ cannot be the realizations of ε but $\hat{\varepsilon}$ is the only thing to check the assumptions on ε (ε model)

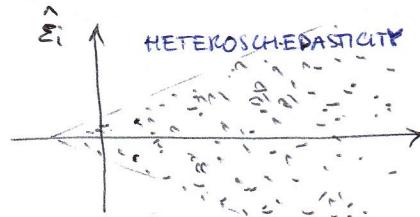
Residual analysis

Searching for large $\hat{\varepsilon}_i$ (= the fitting of the model is very poor)



1. we want to see they're not big

2. we want to check that there's no shape (if there is information (shape) we should capture it in the model)

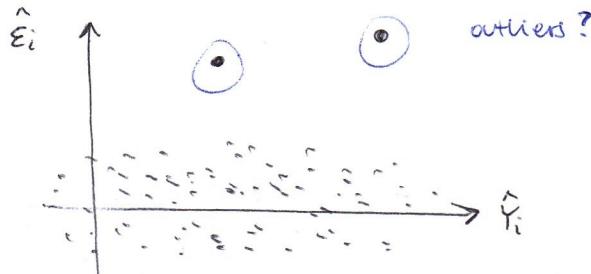


In this case there is a variability which is non constant w.r.t. the fitted values (In this case large fitted values implies large variance for $\hat{\varepsilon}_i$)

How can we fix it? (the second plot) (\Rightarrow the heteroscedasticity)

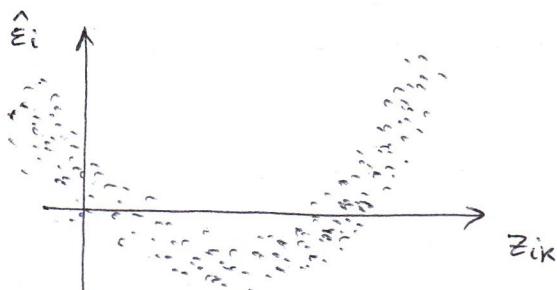
- weighted linear regression
- transform the Y's
(variance stabilizing transformations)

Another possibility for residuals:



They shouldn't be easily eliminated.
We have to understand what they represent (why they're outliers).
It might be that these outliers are a signal about existence of an other population that we didn't sample (and sometimes that's exactly what we're looking for)

We can plot $\hat{\epsilon}_i$ vs. one regressor (instead of fitted values)



$$k = 1, \dots, r$$

Again: we don't want to see patterns (the error should not be related to regressors)

\rightarrow in this case the errors are not independent with the k -regressor.
quadratic in z_{ik} ? \Rightarrow introduce z_{ik}^2



cubic transformation?

sinusoidal transformation? (sin/cos)

Note that:

$$\text{Cov}(\hat{\epsilon}) = \sigma^2(I - H) \implies \text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii}) \quad i = 1, \dots, n$$

$(h_{ii} = \text{diag}_i(H))$

$\Rightarrow \hat{\epsilon}_i$ have different variances ($i = 1, \dots$)

Repeat the residual analysis with:

$$\frac{\hat{\epsilon}_i}{\sqrt{s^2(1-h_{ii})}}$$

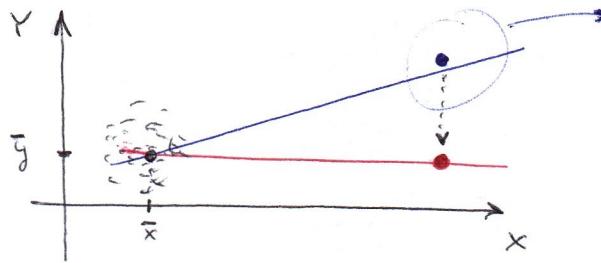
studentized residuals
(in this way the residuals will have the same variance)

- Check Gaussianity: QQ-plot of residuals or studentized residuals

- Test for autocorrelation using (for instance) Durbin-Watson test
(of course the residuals are correlated ($\text{Cov}(\hat{\epsilon}) = \sigma^2(I - H)$) but how large is this autocorrelation?)

Influential cases:

are there cases that drives all the analysis no matter how many data we have?



LEVERAGING EFFECT

if we move this point the linear regression changes drastically \Rightarrow what we read from the data (the model) is entirely determined basically by just 1 point

$$\text{Note that: } \text{Var}(\hat{\epsilon}_i) = \sigma^2(1-h_{ii}) \quad i=1,\dots,n$$

h_{ii} is called Leverage

$h_{ii} = \text{diag}_i(Z(Z^T Z)^{-1}Z^T)$: it doesn't even depend on Y , before performing the experiment we can chose Z s.t. we have small residuals

$$0 \leq h_{ii} \leq 1 \quad (\text{prove it, use: } H = H^T \text{ and } H \cdot H = H)$$

$$h_{ii} \uparrow 1 \Rightarrow \text{Var}(\hat{\epsilon}_i) \rightarrow 0 \Rightarrow \hat{\epsilon}_i \rightarrow 0$$

for a specific case that is driving the entire model

Another way to identify influential cases is to work out what happens if a specific case (a specific unit) is taken out of the dataset:

Holding out unit i to check if it is influential

$$\cancel{X} \text{ data frame} \longrightarrow Z \in \mathbb{R}^{n \times (r+1)}$$

$$\cancel{X}_{-i} \text{ data set obtained by taking out the case } i \longrightarrow Z_{-i} \in \mathbb{R}^{(n-1) \times (r+1)}$$

$$Y = Z\beta + \varepsilon \longrightarrow \hat{\beta} \in \mathbb{R}^{r+2}$$

$$Y_{-i} = Z_{-i}\beta + \varepsilon_{-i} \longrightarrow \hat{\beta}^{(i)} \in \mathbb{R}^{r+1}$$

If $\hat{\beta}^{(i)}$ and $\hat{\beta}$ are very different \Rightarrow case i is influential

Compute:

$$\frac{(\hat{\beta}^{(i)} - \hat{\beta})^T (Z^T Z) (\hat{\beta}^{(i)} - \hat{\beta})}{S^2(r+1)} := D_i$$

Cook's distance

how far $\hat{\beta}^{(i)}$ is from $\hat{\beta}$ in terms of the Mahalanobis distance

We conclude that case i is influential if D_i is large:

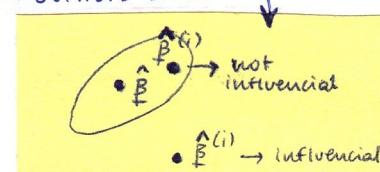
"large"?:

- compare with the quantiles of $F(r+1, n-(r+1))$: if we're very far in terms of this quantile means that we're out of any reasonable confidence region around $\hat{\beta}$
- Rule of thumb: check cases for which $D_i > 1$: the median for this distribution (for reasonable values of r and n) is always centered on 1 (or close to 1). By checking that $D_i > 1$ we're saying that the Cook's distance is larger than the median that we've expecting

Note:

$$D_i = \underbrace{\left(\frac{\hat{\epsilon}_i}{S\sqrt{1-h_{ii}}} \right)^2}_{\text{studentized residuals}} \underbrace{\frac{h_{ii}}{1-h_{ii}}}_{\text{monotone } \uparrow \text{ function of } h_{ii}} \underbrace{\frac{1}{r+1}}$$

takes care of both outliers and leveraging cases.



$\bullet \hat{\beta}^{(i)} \rightarrow \text{outlier}$

• Collinearity

One or more regressors can be expressed as a linear combination of the others, i.e. "close" to "columns of Z are not linearly independent." it's not a problem for the geometrical point of view (we can always project Y on $\mathbb{L}(Z)$), the problem is the formula: $\hat{\beta} = (\bar{Z}^T \bar{Z})^{-1} \bar{Z} Y$

$(\bar{Z}^T \bar{Z})$ is "close" to be singular

$$\Rightarrow \text{Cov}(\hat{\beta}) = \sigma^2 (\bar{Z}^T \bar{Z})^{-1} : \text{it would be a "crazy" matrix (something that jumps from } -\infty \text{ to } +\infty) \Rightarrow \text{the variability of } \hat{\beta} \text{ will be very very large} \Rightarrow \text{the uncertainty about estimating the parameters will be very large}$$

Recall that: $\hat{Y} = \bar{Y} + \mathbb{L}(Y)(Z)$; we need an orthogonal basis for $\mathbb{L}(Z)$, then we project.

Use Gram-Schmidt to build up an orthogonal basis for $\mathbb{L}(Z)$.

$$Z = \begin{bmatrix} 1 & \bar{z}_1 & \dots & \bar{z}_r \\ 1 & z_{n1} & \dots & z_{nr} \\ \vdots & & & \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix} = \begin{bmatrix} 1 & \bar{z}_1 & \dots & \bar{z}_r \\ \vdots & & & \\ 1 & \bar{z}_{r+1} & \dots & \bar{z}_r \end{bmatrix}$$

$$G-S: q_1 = 1$$

$$q_2 = \bar{z}_1 - \bar{z}_1 | q_1$$

$$q_3 = \bar{z}_2 - \bar{z}_2 | q_1, q_2$$

$$q_r = \bar{z}_{r+1} - \bar{z}_{r+1} | \bar{z}_1, \dots, \bar{z}_{r+1}$$

residuals when
we regress \bar{z}_r on
 $(1, \bar{z}_1, \dots, \bar{z}_{r-1})$] = what is left out after we
project \bar{z}_r on the linear space
generated by $1, \dots, \bar{z}_{r-1}$

$$\Rightarrow \hat{\beta}_r = \frac{q_r^T Y}{q_r^T q_r} \quad \left(\begin{array}{l} \text{how do we get } \hat{\beta}_r? \\ \text{By projecting } Y \text{ on } q_r \end{array} \right)$$

$$\text{Var}(\hat{\beta}_r) = \frac{1}{(q_r^T q_r)^2} \sigma^2 q_r^T q_r = \frac{\sigma^2}{q_r^T q_r}$$

Decomposition of variance formula applied to the regression of \bar{z}_r on $(1, \bar{z}_1, \dots, \bar{z}_{r-1})$:

$$\sum_{i=1}^n (\bar{z}_{ir} - \bar{\bar{z}}_r)^2 = \sum_{i=1}^n (\hat{\bar{z}}_{ir} - \bar{\bar{z}}_r)^2 + \underbrace{q_r^T q_r}_{\text{sum of the residuals}^2}$$

$$q_r^T q_r = \sum_{i=1}^n (\bar{z}_{ir} - \bar{\bar{z}}_r) (1 - R_r^2)$$

$\hookrightarrow R_r^2$ for the regression
of \bar{z}_r on $(1, \bar{z}_1, \dots, \bar{z}_{r-1})$

$$\text{Var}(\hat{\beta}_r) = \frac{\sigma^2}{\sum_{i=1}^n (z_{ir} - \bar{z}_r)^2} \cdot \frac{1}{1 - R_r^2}$$

$$= \frac{\sigma^2}{(n-1) S_r^2} \cdot \frac{1}{1 - R_r^2}$$

→ this works not only for the last regressor (r^{th}), it's true for any regressor (since we can change the order for the regressors)

True for $j = 1, \dots, r$

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

coefficient of determination that we get when we regress variable z_j against the other regressors

$\text{Var}(\hat{\beta}_j) \downarrow$ if $\sum_{i=1}^n (z_{ij} - \bar{z}_j)^2 \uparrow$ → this gives an indication of how we should design our matrix if we're in control of it

$\text{Var}(\hat{\beta}_j) \uparrow$ if $R_j^2 \uparrow$ → $R_j^2 \uparrow 1$ means that the variable z_j could be expressed as a linear combination of the other regressors
(that is collinearity)

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

Variance Inflation Factor $(\text{VIF}(\hat{\beta}_j))$

Red Alarm: $\text{VIF}(\hat{\beta}_j) > 5$ or > 10

Model Selection - Variable Selection

Training data: $Z = [1 \ z_1 \ \dots \ z_r]$

maybe we want to extract some variable out of the model (we keep the model and take out variables, we're not changing completely the model)

$$z_j \in \mathbb{R}^n$$

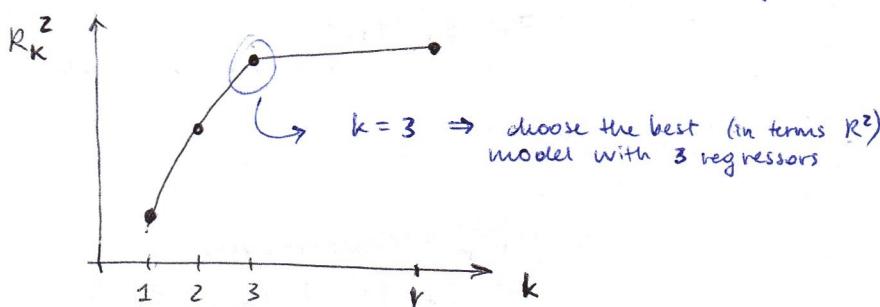
We can build 2^r models. (for every z_j we have two options: in or out the model)

1st line of attack: check them all (fails with big r)

For $k = 1, \dots, r$ fit the (k) model with k regressors and compute R^2 .

Choose the model with k regressors with highest R^2 : R_k^2

monotonic sequence. Why?
Because the best model with 1 regressor cannot beat the best model with 2 regressors (because in the model with 2 regres. we can fit one regressor to be 0 and then we go down to 1)



or plot: $R_{\text{adj}}^2, \text{AIC}, \text{BIC}, \dots$

2nd line of attack: iterative procedures (forward and backward)

Forward:

Start with the best model with **1** regressor.

Until convergence:

add 1 variable which increases the fit of the model.

Use F test to decide when to stop:

consider two models: $k, k^* \quad : \quad k^* > k$

$$\left[\frac{\frac{SS_{\text{res}}(Z_k) - SS_{\text{res}}(Z_{k^*})}{SS_{\text{res}}(Z_k)}}{\frac{n - (k^* + 1)}{n - k}} \cdot \frac{1}{k^* - k} \right]$$

stop if the p-value is large. (= it doesn't pay to move to a larger model)

(The backward selection goes the way around:

we start from the complete model and then we take out one variable, we stop taking out variables as soon as the p-values becomes to small (= the two models are different))

Linear models:

COLLINEARITY & VARIABLE SELECTION

Let's make an observation about the fitted model that we get through OLS.

$$\underline{Y} = \underline{Z} \hat{\beta} + \underline{\varepsilon}$$

model for the observed data

$$\text{OLS: } \hat{\beta} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{Y} \quad (\text{assuming } \text{rank}(\underline{Z}) = \text{full})$$

$$\text{Fitted model: } Y_0 = \underline{Z}_0^T \hat{\beta} = \text{the model that we use for} \\ (\text{for the mean of } Y) \quad E[Y_0 | \underline{Z}_0]$$

The fitted model goes through the baricenter of the observed data.

$$\text{let } \underline{Z}_0 = \frac{\underline{Z}^T \underline{1}}{\underline{1}^T \underline{1}} = [1 \ \bar{z}_1 \ \dots \ \bar{z}_r]^T = \text{vector of the means for the regressors}$$

(vector of the sample means for the columns of \underline{Z})

$$\begin{aligned} Y_0 &= \underline{Z}_0^T \hat{\beta} \\ &\stackrel{!}{=} \frac{\underline{1}^T \underline{Z}}{n} (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{Y} \\ &\stackrel{!}{=} \frac{1}{n} \underline{1}^T H \underline{Y} \\ &\stackrel{!}{=} \frac{1}{n} (\underline{H} \underline{1})^T \underline{Y} \\ &\stackrel{!}{=} \frac{\underline{1}^T \underline{Y}}{n} = \bar{Y} \end{aligned}$$

$$\underline{1} \in \mathcal{L}(\underline{Z}) \Rightarrow \underline{H} \underline{1} = \underline{1}$$

So corresponding (for the fitted model) to the mean values of the regressors we get as output always the mean value for the dependent variable (\bar{Y}).

$$\text{Hence: } \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{z}_1 + \dots + \hat{\beta}_r \bar{z}_r$$

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

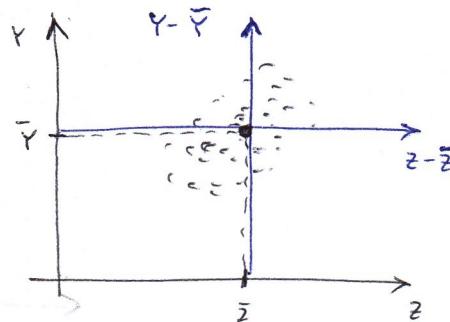
Different representation for the fitted model: (we substitute $\hat{\beta}_0$)

$$Y_0 - \bar{Y} = \hat{\beta}_1 (z_1 - \bar{z}_1) + \dots + \hat{\beta}_r (z_r - \bar{z}_r)$$

\Rightarrow the only question is to find $\hat{\beta}_1, \dots, \hat{\beta}_r$: $\hat{\beta}_0$ is automatically obtained knowing that the model will go through the origin (OLS)

\Rightarrow we can first center the variables in the baricenter and then fit the regression model:

We know that the model will pass through (\bar{z}, \bar{Y}) so we move directly the system there and we fit the model w.r.t. the new reference system



Centering:

$$\underline{Y} \longrightarrow \begin{bmatrix} Y_1 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{bmatrix} = \underline{Y}^*$$

$$\mathbb{Z} \longrightarrow \begin{bmatrix} z_{11} - \bar{z}_1 & z_{12} - \bar{z}_2 & \dots & z_{1r} - \bar{z}_r \\ z_{21} - \bar{z}_1 & \dots & z_{2r} - \bar{z}_r \\ \vdots & & \ddots & \\ z_{n1} - \bar{z}_1 & \dots & z_{nr} - \bar{z}_r \end{bmatrix} = \mathbb{Z}^* \in \mathbb{R}^{n \times r}$$

OLS becomes:

$$\left\{ \begin{array}{l} \underset{\beta \in \mathbb{R}^r}{\text{augmin}} \quad \| \mathbb{Y}^* - \mathbb{Z}^* \hat{\beta} \|^2 = \hat{\beta}^* \\ \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1^* \bar{z}_1 - \dots - \hat{\beta}_r^* \bar{z}_r \\ \hat{\beta}_i = \hat{\beta}_i^* \end{array} \right. \quad \left. \begin{array}{l} \text{Instead of solving the original problem we solve this and then we fix the } \beta_i \text{'s with } (*) \text{ in order to get the problem expressed in the original system} \\ (*) \end{array} \right.$$

From now on we assume \mathbb{Y} and \mathbb{Z} to be centered
(after getting estimates of $\hat{\beta}$, fix $\hat{\beta}_0$)

so to go back to the original system
(something that is almost always asked)

PCA REGRESSION

$$z_1, \dots, z_r$$

$$\mathbb{Z} = [z_1, \dots, z_r]$$

: the problem of collinearity is due to the fact that the regressors are not orthogonal, they can be correlated. How can we fix it? We can transform the regressors so that we have new regressors that are orthogonal. One possibility is the principal component analysis of the independent variables.

→ PCA of \mathbb{Z} just the dependent variables

→ PC_1, \dots, PC_r

Reduce dimensionality: $PC_1, \dots, PC_k, k \leq r$ (the last ones probably do not account much of the variability of the data)

$$\mathbb{Z}^* = [PC_1 \dots PC_k]$$

fit: $\mathbb{Y} = \mathbb{Z}^* \hat{\beta} + \varepsilon$ $\hat{\beta} \in \mathbb{R}^k \rightarrow$ we reduced dimensionality

fitted model: $\mathbb{Y}_0 = \mathbb{Z}_0^T \hat{\beta}$ $\mathbb{Z}_0 = [PC_1, \dots, PC_k]^T$

$$\hat{\beta} = \hat{\beta}_1 PC_1 + \dots + \hat{\beta}_k PC_k = \hat{\beta}_1 z_1 + \dots + \hat{\beta}_k z_k$$

with $z_j \in \mathbb{Z}^*$, so $z_j = PC_j$

$$PC_1 = e_{11} z_1 + e_{12} z_2 + \dots + e_{1r} z_r$$

$$PC_2 = e_{21} z_1 + e_{22} z_2 + \dots + e_{2r} z_r$$

\vdots

$$PC_k = e_{k1} z_1 + e_{k2} z_2 + \dots + e_{kr} z_r$$

substituting: $\mathbb{Y}_0 = z_1 (\hat{\beta}_1 + e_{12} \hat{\beta}_2 + \dots + e_{1r} \hat{\beta}_r) +$
 $+ z_2 (\hat{\beta}_1 + e_{22} \hat{\beta}_2 + \dots + e_{2r} \hat{\beta}_r) +$
 $+ z_r (\hat{\beta}_1 + e_{r2} \hat{\beta}_2 + \dots + e_{rr} \hat{\beta}_r)$

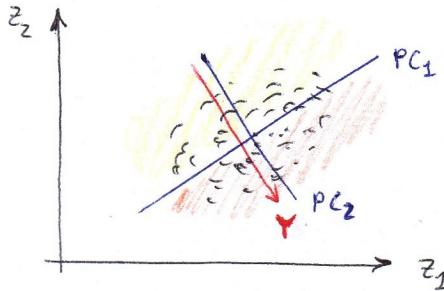
$$\Rightarrow \hat{Y}_0 := \hat{\beta}_1 z_1 + \dots + \hat{\beta}_r z_r$$

$$\hat{\beta}_j = e_{j1} \hat{\beta}_1 + \dots + e_{jk} \hat{\beta}_k$$

we'll be a fitted model with the same \hat{Y} variables (z_1, \dots, z_r) with coefficients but these coefficients are not the OLS coefficients: these are linear combinations of the k parameters so those are k coefficients that we got after PCA.

Criticism:

- solution is not sparse in terms of z_1, \dots, z_r (no variables selection) : we still got all the r variables
- we might have thrown away the information of Z correlated with Y : (doing the dimensional reduction in PCA not guarantee that the information we need to predict Y is not contained in the components that we're throwing away)



Suppose that Y increase along PC_2 (Y is higher in the yellow zone). If we consider just PC_1 (and we throw away PC_2 for dimensional reduction) we're throwing away the direction where Y is varying as a function of Z_1 and Z_2 (we're keeping a direction so that there is no increase or decrease, Y is basically constant among PC_1).

RIDGE REGRESSION

(sort of a smoother version of PCA regression)
(Hoerl & Kennard)

Collinearity \rightarrow Variability of $\hat{\beta}_{OLS}$ is large (which make predictions and interpretations all uncertain)

Ridge optimization problem:

$$\begin{cases} \text{arg min}_{\beta} \|Y - Z\beta\|^2 \\ \|\beta\|^2 \leq S \end{cases}$$

(S is a new parameter that we enter to the problem without this constraint it's OLS)

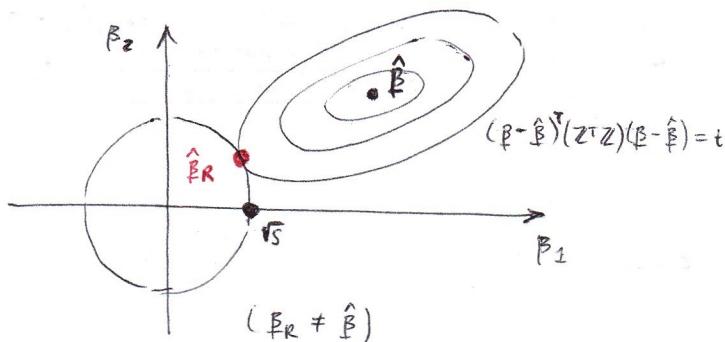
Note that:

$$\begin{aligned} \|Y - Z\beta\|^2 &= \|Y - \hat{Y} - Z(\beta - \hat{\beta})\|^2 && \text{with } \hat{Y} = H Y \\ &\stackrel{\perp}{=} \| \underbrace{\hat{Y}}_{\in Z^\perp(Z)} - \underbrace{Z(\beta - \hat{\beta})}_{\in Z(Z)} \|^2 \\ &= \|\hat{Y}\|^2 + \|Z(\beta - \hat{\beta})\|^2 \end{aligned}$$

Ridge problem:

$$\begin{cases} \text{arg min}_{\beta} \|Z(\beta - \hat{\beta})\|^2 & (\text{since } \|\hat{Y}\|^2 \perp \beta) \text{ ("it's constant w.r.t. } \beta) \\ \|\beta\|^2 \leq S \\ \hat{\beta} = (Z^T Z)^{-1} Z^T \hat{Y} & (\text{OLS solution}) \end{cases}$$

$$\Rightarrow \left\{ \begin{array}{l} \underset{\beta}{\text{arg min}} \quad (\beta - \hat{\beta})^T Z^T Z (\beta - \hat{\beta}) \\ \beta^T \beta \leq S \end{array} \right. = \text{we don't want solutions that are too big in terms of norm}$$



$\hat{\beta}_R$ is a "restricted" "shrunked" version of $\hat{\beta}_{OLS}$.

The solution is obtained via Lagrange method, which implies to solve :

$$\underset{\beta}{\text{arg min}} \quad \|Y - Z\beta\|^2 + \lambda \|\beta\|^2 \quad (\lambda = 0 \Rightarrow OLS)$$

where λ is a function of S : $\lambda \uparrow$ if $S \downarrow$

Differentiate wrt β and get:

$$\hat{\beta}_R = (Z^T Z + \lambda I)^{-1} Z^T Y$$

Obs.

→ this solves the collinearity problem : $(Z^T Z + \lambda I)$ is always invertible if $\lambda \neq 0$ and the larger the λ is the easier it is to invert

- $\hat{\beta}_R$ is biased estimator of β

- $\exists \lambda$ s.t. $E[\|\hat{\beta}_R - \beta\|^2] < E[\|\hat{\beta}_{OLS} - \beta\|^2]$

- Find the best λ by cross validation using prediction error for Y

we cross validate our predictions through many possible choices on a grid of possible λ 's and then we find the λ that minimize the prediction error for the Y

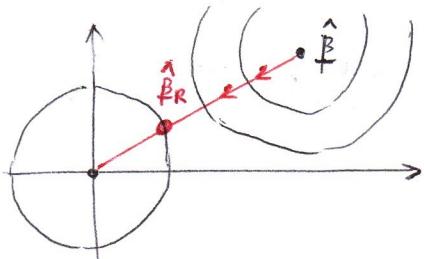
Obs.

- Standardize the regressors and the Y before fitting

Ridge Regression since $\hat{\beta}_R$ is not scale invariant

($\hat{\beta}_{OLS}$ is scale invariant)

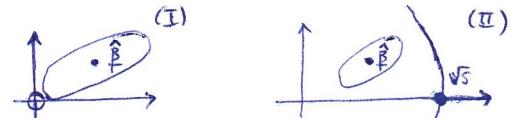
- If regressors are standardized and are orthogonal :



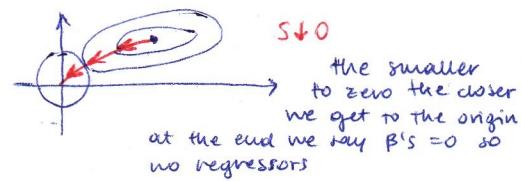
(!) it's not a curve anymore, it's a line !

$\hat{\beta}_R$ is $\hat{\beta}$ scaled.

Since S is a parameter (that control the variability of β 's) we can have:



for very large S (II) the old OLS solution ($\hat{\beta}$) is still good because satisfies the constraints. For very small S , $S \rightarrow 0$ (I) happens that the solution is being shrink toward zero:



(we penalize solutions, the larger we penalize, the closer we get to zero)

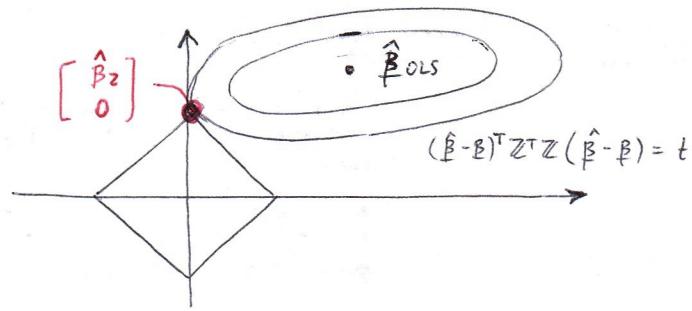
→ we're free to chose any solution but we're paying λ dollars for every kilometer that we're far from zero (so we also want to minimize the dollars that we pay at the end).

{ we're paying in terms of bias but gaining in terms of variability (bias-variance trade off)

with RIDGE:

if we measure in meters and kilograms and then we decide to switch into centimeters and grams, we cannot simply multiply our coefficients for appropriate factor in order to use the same model with a different scale

LASSO : let's change the constraint



We want a constraint with angles instead of a constraint spherical.
Why? Because we hope that the tangent point is gonna be on one of the vertices. If it's there, this means that $\beta_1 = 0$
 \Rightarrow the solution will be $[\hat{\beta}_2 \ 0]^T$, meaning that we're not only shrinking but also selecting the variables

Controlling the variability of $\hat{\beta}$ and selecting variables.

Lasso optimization problem:

$$\left\{ \begin{array}{l} \underset{\beta}{\operatorname{argmin}} \quad (\beta - \hat{\beta})^T (Z^T Z) (\beta - \hat{\beta}) \\ \|\beta\|_1 \leq s \quad \left(\sum_{i=1}^r |\beta_i| \leq s \right) \end{array} \right.$$

There is no analytical solution.

Lagrangian (Lagrange method):

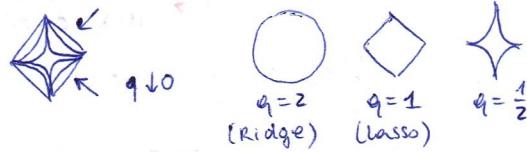
$$\boxed{\underset{\beta}{\operatorname{argmin}} \quad \|Y - Z\beta\|^2 + \lambda \|\beta\|_1}$$

Find the optimal λ through cross validation and pred. error for Y .

Variations of Lasso:

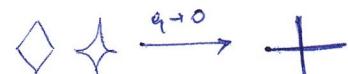
we're consider even more spiky constraints :

$$\left\{ \begin{array}{l} \underset{\beta}{\operatorname{argmin}} \quad \|Y - Z\beta\|^2 \\ \|\beta\|_q \leq s \end{array} \right.$$



$$\|\beta\|_q = \left(\sum_{i=1}^r |\beta_i|^q \right)^{1/q} \quad (q < 1)$$

For $q \downarrow 0 \Rightarrow$ model selection problem because:



it's pure selection of the variables

$$(\overset{\oplus}{P_2} \overset{+}{P_1})$$

Extension: (Elastic nets) (take both Ridge & lasso)

$$\underset{\beta}{\operatorname{argmin}} \quad \|Y - Z\beta\|^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|^2$$

General idea: (for all these problems)

fit a regression with a penalization.

This becomes super good when we already know something about the solution (in terms of model) \rightarrow we penalize solutions that are far from the ideal solutions.

\Rightarrow these methods take care of both: model selection part and the collinearity (= variability of β 's).

Ridge is not very good in selecting the variables. Lasso and variations go in the direction of model selection

Connection with GLM (Generalized linear models)

$$Y = f(\underline{z}) + \varepsilon \quad \text{basic model}$$

$$\mathbb{E}[Y|\underline{z}] = f(\underline{z})$$

linear regression: $f(\underline{z}) = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r$

In GLM: we think that (↑) holds but not for $f(\underline{z})$ but for a transformation of $f(\underline{z})$: $g(f(\underline{z}))$

$$g(\mathbb{E}[Y|\underline{z}]) = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r \quad := \text{LINK FUNCTION}$$

Linear Regression:

- $\begin{cases} \mathbb{E}[\varepsilon] = 0 \\ \text{Var}(\varepsilon) = \sigma^2 \end{cases} \Rightarrow Y \sim \begin{cases} \mathbb{E}[Y|\underline{z}] = f(\underline{z}) \\ \text{Var}(Y|\underline{z}) = \sigma^2 \end{cases}$
- $\varepsilon \sim N(0, \sigma^2) \Rightarrow Y|\underline{z} \sim N(f(\underline{z}), \sigma^2)$

GLM:

$$Y|\underline{z} \sim F \quad \left\{ \begin{array}{l} \text{Bernoulli} \\ \text{Poisson} \\ \text{Gaussian} \\ \vdots \end{array} \right.$$

((*))
the parameters β are estimated by ML (maximum likelihood) and not by least squares

the model is not only an extension of the linear model in the sense that we're taking a linear representation for a transformation of the mean (and not for the mean itself) but also we allow the error structure to be far away from the gaussian. what do we pay? least squares is no longer working. ((*))

Link functions:

- $Y \text{ Gaussian} \Rightarrow g(\mu) = \mu \quad \text{Linear Regression}$
- $Y \text{ Bernoulli}(p) \Rightarrow g(p) = \log\left(\frac{p}{1-p}\right) \quad \text{Logit Regression}$
- $Y \text{ Poisson } (\lambda) \Rightarrow g(\lambda) = \log(\lambda) \quad \text{Poisson Regression}$
- $Y \text{ Bernoulli}(p) \Rightarrow g(p) = \Phi^{-1}(p) \quad \text{Probit Regression}$
 Φ cumulative distr. of $N(0,1)$

Everything we said about Lasso, Ridge, ..., we can apply here.
(Ex. Lasso Logistic Regression, ...)