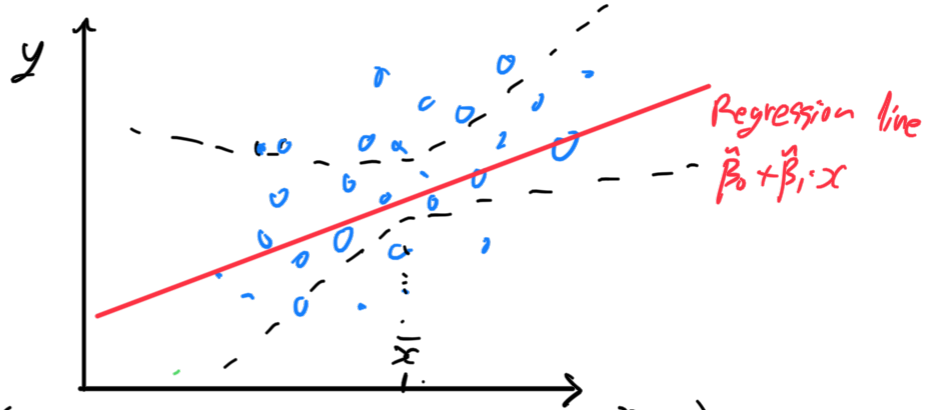


Lecture 9

17/5/2021

Recap: Simple Regression $y_i = \beta_0 + \beta_1 \cdot x_i + \epsilon_i$
 $i=1, \dots, n$

- Confidence Bands for $\beta_0 + \beta_1 \cdot x$



$$\beta_0 + \beta_1 x \underset{\text{w.p. } 1-\alpha}{\in} \left(\hat{\beta}_0 + \hat{\beta}_1 x \pm \underset{\text{w.p. } 1-\alpha}{t_{n-2}^{1-\alpha/2}} \cdot s \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}}} \right)$$

In general:

$$z_0^T \beta \underset{\text{w.p. } 1-\alpha}{\in} \left(z_0^T \hat{\beta} \pm t_{n-p}^{1-\alpha/2} \cdot s \sqrt{z_0^T (Z^T Z)^{-1} z_0} \right)$$

$z_0 \in \mathbb{R}^p$

$$s^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{\|\epsilon\|^2}{n-p}$$

- Prediction Bands (in HW)

We want to predict y_{n+1} given x_{n+1}
 and $(y_i, x_i)_{i=1}^n$.

$$y_{n+1} \underset{\text{w.p. } 1-\alpha}{\in} \left(\hat{\beta}_0 + \hat{\beta}_1 x_{n+1} \pm t_{n-2}^{1-\alpha/2} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{S_{xx}}} \right)$$

we need a line if we predict

we set a bound if we evaluate this for every $x_{n+1} \in \mathbb{R}$

- we can take avg. of m measurements at the same x_{n+1} . In this case:

$$\bar{y}_{n+1} \in \left(\hat{\beta}_0 + \hat{\beta}_1 x_{n+1} \pm t_{n-2}^{1-\alpha/2} \cdot s \sqrt{\frac{1}{m} + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{S_{xx}}} \right)$$

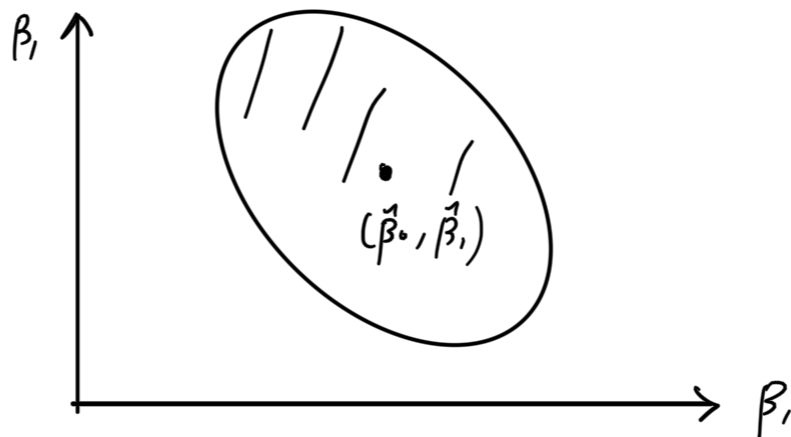
w.p. $1-\alpha$

Simultaneous Bands

- contain $(\beta_0 + \beta_1 x)_{x \in \mathbb{R}}$ with prob. $1-\alpha$
- In p dim, contain $(z^T \beta)_{z \in \mathbb{R}^p}$
- From the distribution of $\hat{\beta}$:

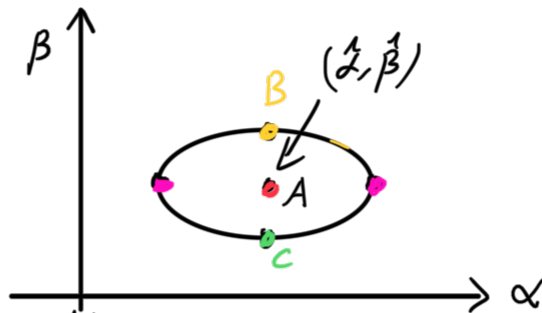
$$\Pr\left((\hat{\beta} - \beta)^T (Z^T Z)^{-1} (\hat{\beta} - \beta) \leq s^2 \cdot p \cdot F_{p, n-p}^{1-\alpha} \right) = 1-\alpha$$

This defines an ellipsoid in \mathbb{R}^p

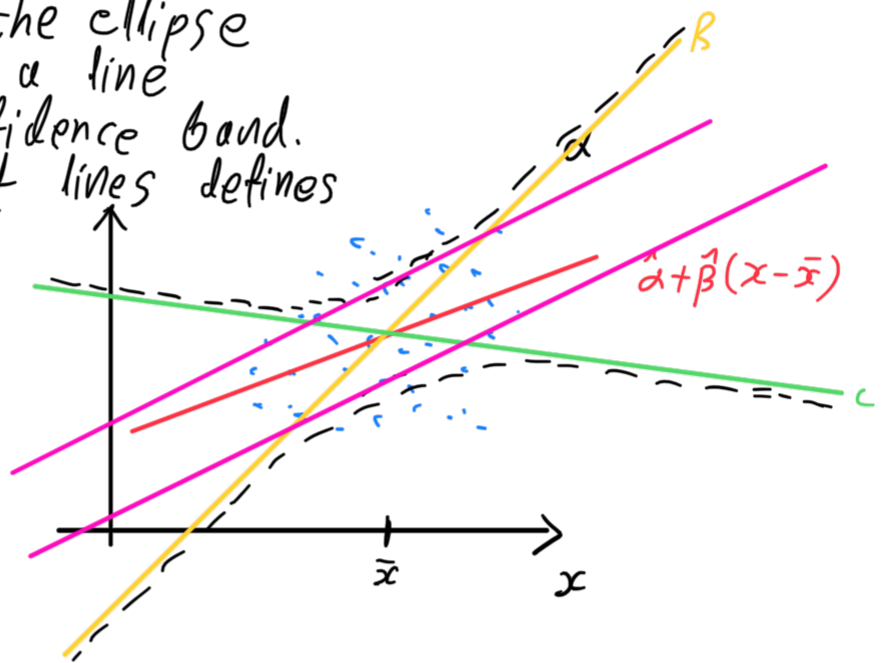


- For $p=2$, β_0 and β_1 lay in an ellipse with prob. $1-\alpha$

- If we use $\hat{\alpha} + \hat{\beta}(x - \bar{x})$ for the regression line, then $Z^T Z$ becomes diagonal and the ellipse's axes are aligned with the x & y axes:



each point in the ellipse corresponds to a line within the confidence band. The collection of lines defines the band



Working Hotelling Bands:

Confidence: $\hat{\beta}_0 + \hat{\beta}_1 x \pm \sqrt{2 F_{2, n-2}^{1-\alpha}} \cdot s \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}$

Prediction: $\hat{\beta}_0 + \hat{\beta}_1 x \pm \sqrt{2 F_{2, n-2}^{1-\alpha}} \cdot s \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}} + 1}$

INTERVAL

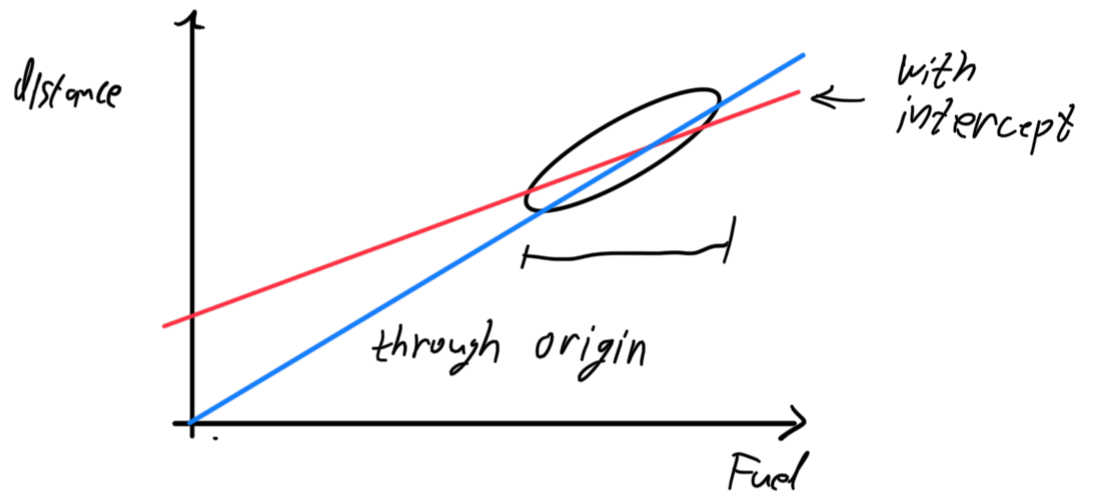
$$\hat{\beta}_0 + \hat{\beta}_1 x \pm \sqrt{2 F_{2, n-2}^{1-\alpha}} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$$

Avg. of m

Predictions: $\hat{\beta}_0 + \hat{\beta}_1 x \pm \sqrt{2 F_{2, n-2}^{1-\alpha}} \cdot s \sqrt{\frac{1}{m} + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$

(these are wider than earlier-mentioned bands; interpreted when x runs over all values)

- Regression Through the Origin



- This is dangerous if there is no data near the origin
- Also, no interpretation of R^2 in this case.

Multiple Regression:

Suppose: $Z = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \end{bmatrix} \begin{matrix} x_i \in \mathbb{R}^d \\ \in \mathbb{R}^p \\ \dots \end{matrix}$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$r = d+1$

(fraction of explained variance)

- Note: when using R^2 you must have an intercept term
- R^2 is not additive $0 \leq R^2 \leq 1$
- No rule of thumb for how large R^2 must be
- Increasing the number of predictors never decreases R^2

Some Considerations:

- A "true" β_j :
 - Depends on what variables are included.
 - A "true" β_j exists for each and every subset of the covariates used.
- Naive Face-Value Interpretation:

... " ... $\frac{p}{n}$...

- With $y_i = \sum_{j=1}^k x_{ij} \beta_j + \epsilon_i$, we may about β_j as $\frac{\partial y}{\partial x_j}$ or $\frac{\partial E[y]}{\partial x_j}$.

- This is unreliable. For example if $x_2 = x_1^2$, then changing one predictor alone is not possible.

- Wrong Signs:

- may be due to correlation within the data.

Correlation, Association, Causality

- Regression with observational data captures correlation/association, not causality.

- It could be $x \rightarrow y$, $y \rightarrow x$, or $z \rightarrow (x, y)$

Interplay Between Variables

- Competing Variables:

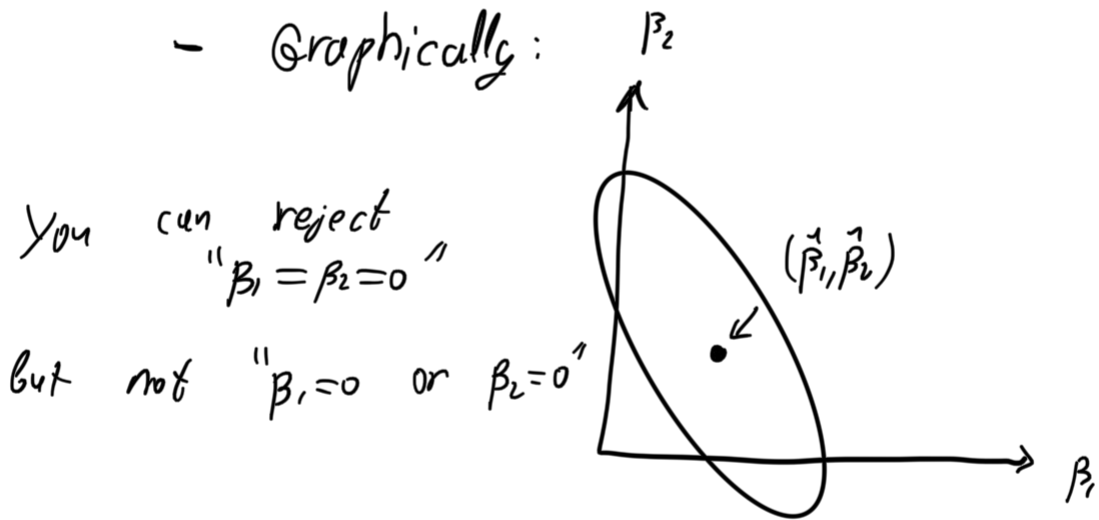
- β_1 is significant if x_2 is not in the model

- β_2 is significant if x_1 is not in the model.

- In the house prices example:

Overall Qual & Overall Cond have this property,

- Graphically:



- Occurs when there is positive correlation between x_1 & x_2

collaborating variables:

- $\hat{\beta}_1$ is sig. if x_2 is in the model
- $\hat{\beta}_2$ is sig. if x_1 is in the model

- This is much more rare than competition

- Example:

x_1 = air pressure in loc. A
 x_2 = air pressure in loc. B
 y = wind intensity $A \rightarrow B$
 (proportional to $x_1 - x_2$)

Simpson's Paradox

	Two hospitals		
	A	B	
Mild cases	S	D	
Critical cases	S	D	

Hospital	A	50	0
	B	80	10

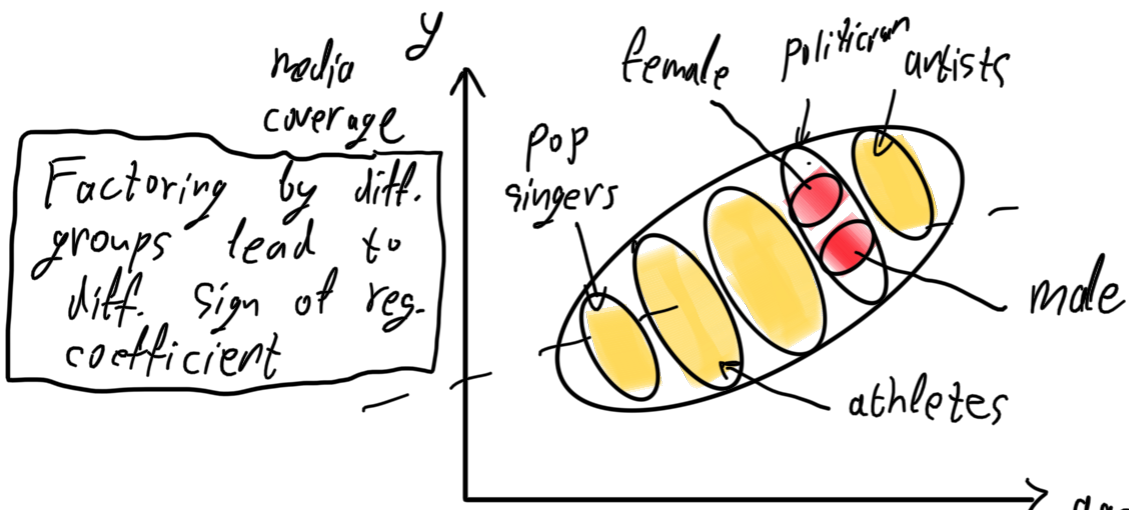
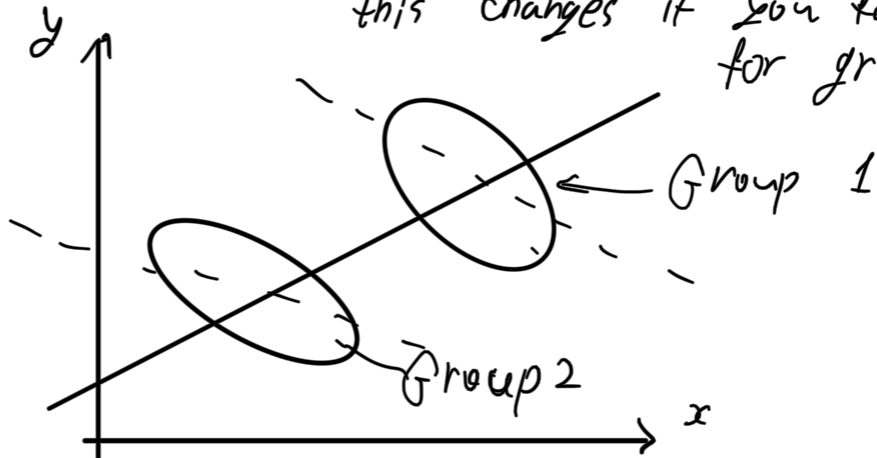
A	25	25
B	0	10

Total:

	S	D
A	75	25
B	80	20

- B looks better overall, although it loses in each category individually.

- With $y = \beta_0 + \beta_1 \cdot \text{Age}$, $\beta_1 > 0$, but this changes if you factor for groups



Partial Correlation

x y e

- We want to look for the connection between height and spelling ability of kids from ages 7-12. We get positive correlation, but kids height & spelling ability change as they grow.
- We regress height on age, and get residuals. These residuals are "age-adjusted height"
- We do the same for spelling-level, getting "age-adjusted spelling."

The correlation between the two sets of residuals is the partial correlation of height and spelling-level, adjusted for age.

Def. Partial Correlation P_{ijk}

Partial correlation of x_i x_j , adjusting for x_k is the correlation of resid. for x_i on x_k and resid. for x_j on x_k

For Gaussian data:

$$P_{ijk} = \frac{P_{ij} - P_{ik}P_{jk}}{\sqrt{(1-P_{ik}^2)(1-P_{jk}^2)}}$$

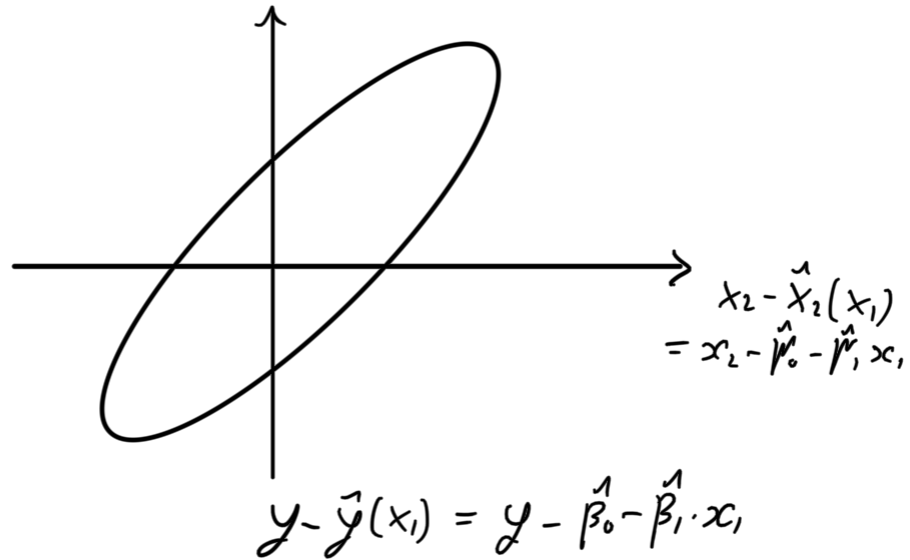
(with or without hats)

$$\hat{P}_{ij} = \frac{x_i^T x_j}{\sqrt{\dots}}$$

$$P_{ij} = \frac{\text{Cov}(x_i, x_j)}{\sqrt{\dots}}$$

$$\|x_i\| \quad \|x_j\|$$

$$\text{Var}(x_i) \quad \text{Var}(x_j)$$



- we can find partial corr. of (x_i, y) , adjusting for x_2, \dots, x_k : we regress those predictors out of both x_i and y and find the partial correlation.