

# STAT 151A Lecture 11

Henry Liev

22 September 2025

### Remark 0.1 (Quiz 1 Review)

#### 1) Goals of modeling where data comes from

- Describe associations
- Prediction  $\rightarrow$  know  $X$ , want to guess  $Y$
- Causal inference  $\rightarrow$  guess change in  $Y$  from possible change in  $X$  (not from change in some other variable)

Where does your data come from? Was it

- a random experiment?  $\rightarrow X$  assigned randomly to subjects
- a random sample from some population?  $\rightarrow$  subjects chosen randomly from population of interest

#### 2) Exploring data

- Plot the data
  - $\rightarrow$  Univariate plots (histograms, density plots)
  - $\rightarrow$  Bivariate scatter plots  $\rightarrow$  for all pairwise relationships (`pairs()`)

Pay attention to:

- $\rightarrow$  Skew (univariate)
- $\rightarrow$  Linearity (bivariate)
- $\rightarrow$  Outliers
- $\rightarrow$  Spread
- $\rightarrow$  Associations/correlations/patterns
- $\rightarrow$  Scale/measurements
- Are there transformations that would help me
  - $\rightarrow$  See the patterns
  - $\rightarrow$  Describe the data better using means, SDs, and linear models

When to transform?

- $\rightarrow$  Clean curved (Simple monotone) relationship, want to use a linear model
- $\rightarrow$  Proportion data  $\rightarrow$  maybe (logit)
- $\rightarrow$  Highly skewed data  $\rightarrow$  transform to get a clearer visualization, “protect the mean/SD”

Downsides: Interpretability, transformed variables are harder to explain/think about

$\rightarrow$  back to goals of analysis

“Bulging” rule

#### 3) Least squares regression: Observe $\vec{Y} \in \mathbb{R}^n$ , $\mathbf{X} \in \mathbb{R}^{n \times p+1}$

Solve  $\min \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2$ , ie project  $\mathbf{Y}$  onto  $\text{Col}(\mathbf{X})$

Solution  $(\mathbf{X}^\top \mathbf{X})\hat{\beta} = \mathbf{X}^\top \mathbf{Y} \rightarrow \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ ,  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{H}\mathbf{Y}$

$\vec{e} = \mathbf{Y} - \mathbf{X}\hat{\beta} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$

$\vec{e} \perp \vec{1}, \vec{e} \perp \vec{x}, \hat{y} \perp \vec{e}$

$$\|\vec{y}\|^2 = \|\hat{y}\|^2 + \|\vec{e}\|^2$$

$$\|\vec{y} - \vec{y}\vec{1}\|^2 = \|\hat{y} - \vec{y}\vec{1}\|^2 + \|\vec{e}\|^2 \rightarrow \text{TotSS} = \text{Regss} + \text{ErrSS}, R^2 = \frac{\text{RegSS}}{\text{TotSS}}$$

**Remark 0.2 (Interpreting element  $\hat{\beta}_j$  of  $\hat{\beta}$ )**

- 1) If  $\vec{1}, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_p$  are all orthogonal  
 $\hat{\beta}_j = \frac{\vec{y} \cdot \vec{x}_j}{\vec{x}_j \cdot \vec{x}_j}$  (Same as projecting  $\vec{y}$  onto  $\vec{x}_j$  alone)
- 2) If not, more complex formula, removing role of other  $\vec{x}$ 's:  
 $\hat{\beta}_j = \frac{\vec{y} \cdot \vec{e}^{(j)}}{\vec{e}^{(j)} \cdot \vec{e}^{(j)}}$ , where  $\vec{e}^{(j)}$  is the vector of residuals from regressing  $\vec{x}_j$  on other  $\mathbf{X}$ -columns  
Reminder: Regression line/plane always goes through the point  $(\bar{y}, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)$   
If you take  $\hat{\beta} \cdot (1, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_p) = \bar{y}$   
Probability: Statement about original data being sampled from some population  
Gauss-Markov probability model  $\vec{Y}, \vec{\varepsilon}$  are random variables,  $\beta$  is a vector of parameters,  
 $\mathbb{E}(\vec{\varepsilon}) = 0, \Sigma_{\vec{\varepsilon}\vec{\varepsilon}} = \sigma^2 \mathbb{I}_n$

$$\vec{Y} = \mathbf{X}\beta + \vec{\varepsilon}$$

Typically we condition on  $\mathbf{X}$

$\mathbb{E}(\hat{\beta}|\mathbf{X}) = \beta, \text{Var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$  Insights:

Adjusted  $R^2 = 1 - \frac{\text{ErrSS}/[n-(p+1)]}{\text{TotSS}/(n-1)} = 1 - \frac{n-1}{n-(p+1)}(1 - R^2)$

Formula for variance of  $\hat{\beta}_j$

In simple regression ( $\vec{x}_j$  is the only variable):  $\frac{\sigma_{\text{simple}}^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} = \frac{\sigma_{\text{simple}}^2}{\text{Var}(\beta_j)}$

In multiple regression ( $\vec{1}, \vec{x}_1, \dots, \vec{x}_p$ ):  $\frac{1}{1-R_j^2} \cdot \frac{\sigma_{\text{multiple}}^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$ , where  $R_j^2$  is from the regression of  $\vec{x}_j$  on all other  $\mathbf{X}$ -columns

Collinearity  $\rightarrow$  when  $\vec{x}_j$  is highly related to other  $\mathbf{X}$ -columns

What should we do about this?

Prediction  $\rightarrow$  how much does including a collinear variable really help your predict  $\mathbf{Y}$  better?

vs Causal inference  $\rightarrow$  knowing about collinearity is really important, but dropping variables is usually bad