

# **Introduction to data science**

## **Linear Regression analysis**

### **2023/2024**

**This lecture is based on:**

Robert West lectures, EPFL  
and the book

A. Gelman and J. Hill,  
Data Analysis Using Regression and  
Multilevel/Hierarchical models,  
Cambridge, 2007.

**Chapter 3 and 4.**

**Please read these chapters!**




Linear  
regression

# Linear regression as you know it

- **Given:**  $n$  data points  $(X_i, y_i)$ , where  $X_i$  is  $k$ -dimensional vector of predictors (a.k.a. features), and  $y_i$  is scalar outcome, of  $i$ -th data point
- **Goal:** find the optimal coefficient vector  $\beta = (\beta_1, \dots, \beta_k)$  for approximating the  $y$ 's as a linear function of the  $X$ 's:

$$\begin{aligned} y_i &= X_i \beta + \epsilon_i \\ &= \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \epsilon_i, \quad \text{for } i = 1, \dots, n \end{aligned}$$

Scalar product (a.k.a. dot product) of 2 vectors

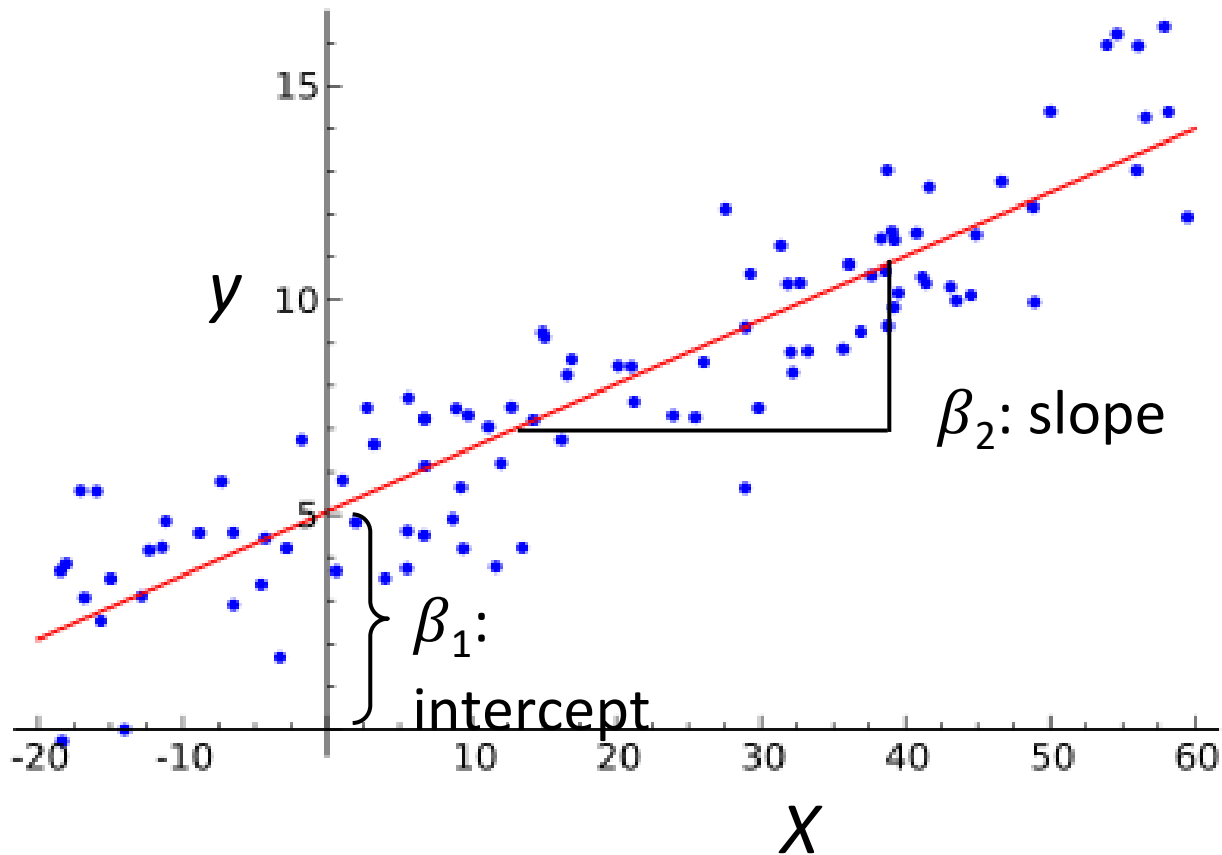


where  $\epsilon_i$  are error terms that should be as small as possible

- $X_{i1}$  usually the constant 1  $\Rightarrow \beta_1$  a constant intercept

# Example with one predictor

$$y \approx \beta_1 + \beta_2 X$$



# Linear regression as you know it

- **Given:**  $n$  data points  $(X_i, y_i)$ , where  $X_i$  is  $k$ -dimensional vector of predictors (a.k.a. features), and  $y_i$  is scalar outcome, of  $i$ -th data point

- **Goal:** find the optimal coefficient vector  $\beta = (\beta_1, \dots, \beta_k)$  for approximating the  $y$ 's as a linear function of the  $X$ 's:

$$y_i = X_i \beta + \epsilon_i$$

$$= \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \epsilon_i, \quad \text{for } i = 1, \dots, n$$

where  $\epsilon_i$  are error terms that should be as small as possible

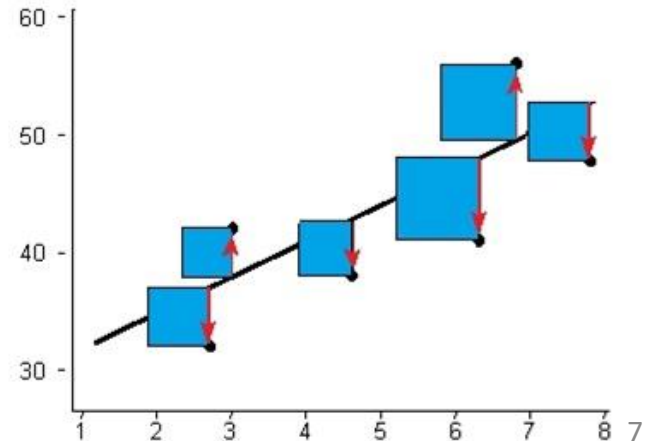
- $X_{i1}$  usually the constant 1  $\rightarrow \beta_1$  a constant intercept

# Optimality criterion: least squares

$$y_i = X_i\beta + \epsilon_i \quad \text{for } i = 1, \dots, n$$

- Intuitively, want errors  $\epsilon_i$  to be as small as possible
- Technically, want sum of squared errors as small as possible  
 $\Leftrightarrow$  find  $\hat{\beta}$  such that we minimize

$$\sum_{i=1}^n (y_i - X_i\hat{\beta})^2$$



# Use cases of regression

- **Prediction:** use fitted model to estimate outcome  $y$  for a new  $X$  not seen during model fitting (if you've seen regression before, then probably in the context of prediction)
- **Descriptive data analysis:** compare mean outcomes across subgroups of data (today!)
- **Causal modeling:** understand how outcome  $y$  changes if you manipulate predictors  $X$  (next lecture is about causality, although not primarily using regression)



# Regression as comparison of mean outcomes

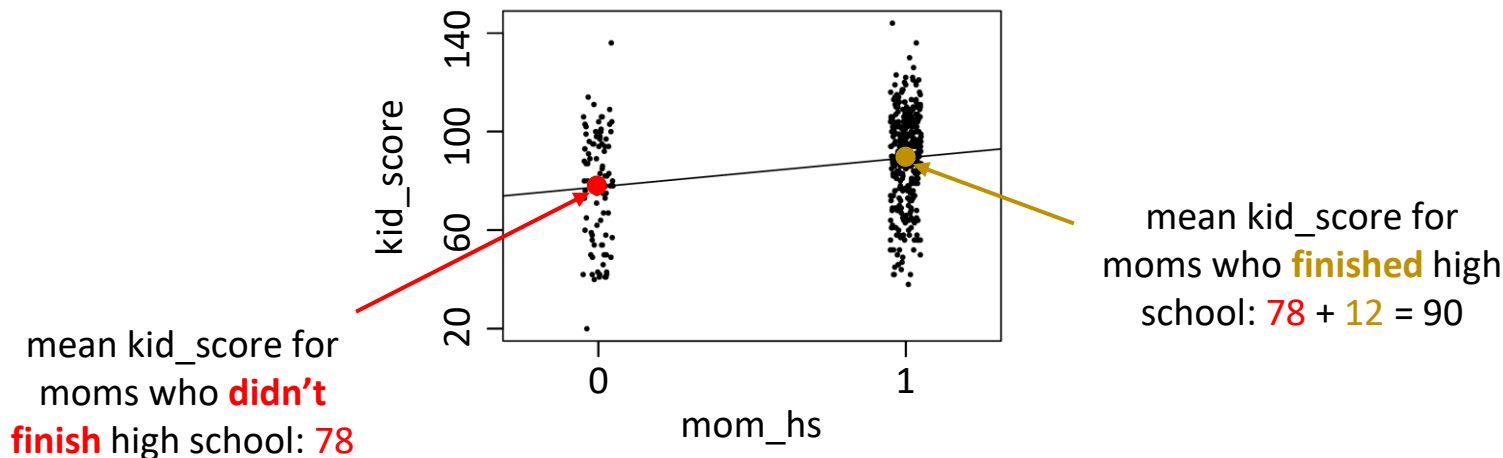
# Example with one binary predictor $X_i$

No Yes

- $X_i = \text{mom\_hs} = \text{"Did mother finish high school?"} \in \{0, 1\}$
- $y_i = \text{kid\_score} = \text{child's score on cognitive test} \in [0, 140]$

$$y_i = \beta_1 + \beta_2 X_i + \epsilon_i$$

$$\text{kid\_score} = 78 + 12 \cdot \text{mom\_hs} + \text{error}$$

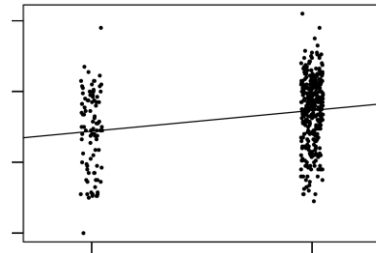


# One binary predictor $X_i$ :

## Interpretation of fitted parameters $\beta$

$$y_i = \beta_1 + \beta_2 X_i + \epsilon_i$$

- **Intercept**  $\beta_1$ : mean outcome for data points  $i$  with  $X_i = 0$
- **Slope**  $\beta_2$ : difference in outcomes between data points with  $X_i = 1$  and data points with  $X_i = 0$
- Reason: means minimize least-squares criterion

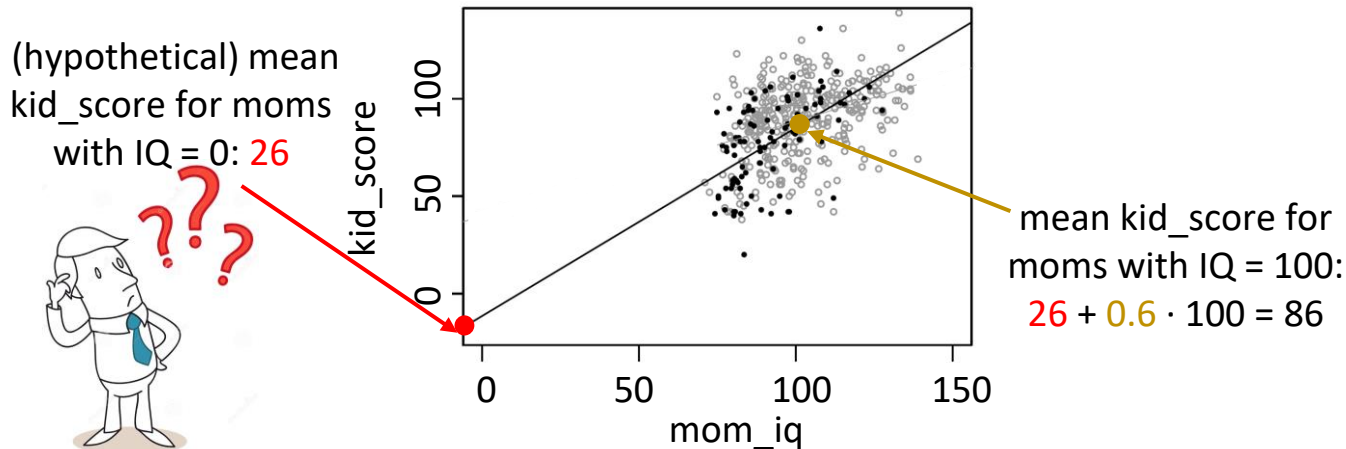


# Example with one continuous predictor $X_i$

- $X_i = \text{mom\_iq} = \text{mother's IQ score} \in [70, 140]$
- $y_i = \text{kid\_score} = \text{child's score on cognitive test} \in [0, 140]$

$$y_i = \beta_1 + \beta_2 X_i + \epsilon_i$$

$$\text{kid\_score} = 26 + 0.6 \cdot \text{mom\_iq} + \text{error}$$



# One continuous predictor $X_i$ :

## Interpretation of fitted parameters $\beta$

$$y_i = \beta_1 + \beta_2 X_i + \epsilon_i$$

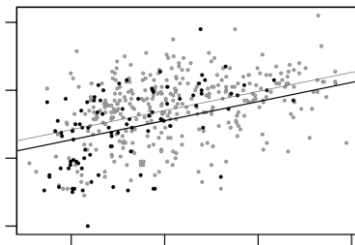
- **Intercept**  $\beta_1$ : average outcome for data points  $i$  with  $X_i = 0$
- **Slope**  $\beta_2$ : difference in outcomes between data points whose  $X_i$ 's differ by 1

# Example with multiple predictors

- ( $X_{i1} = 1 = \text{constant}$ )
- $X_{i2} = \text{mom\_hs} = \text{“Did mother finish high school?”} \in \{0, 1\}$  No Yes
- $X_{i3} = \text{mom\_iq} = \text{mother’s IQ score} \in [70, 140]$
- $y_i = \text{kid\_score} = \text{child’s score on cognitive test} \in [0, 140]$

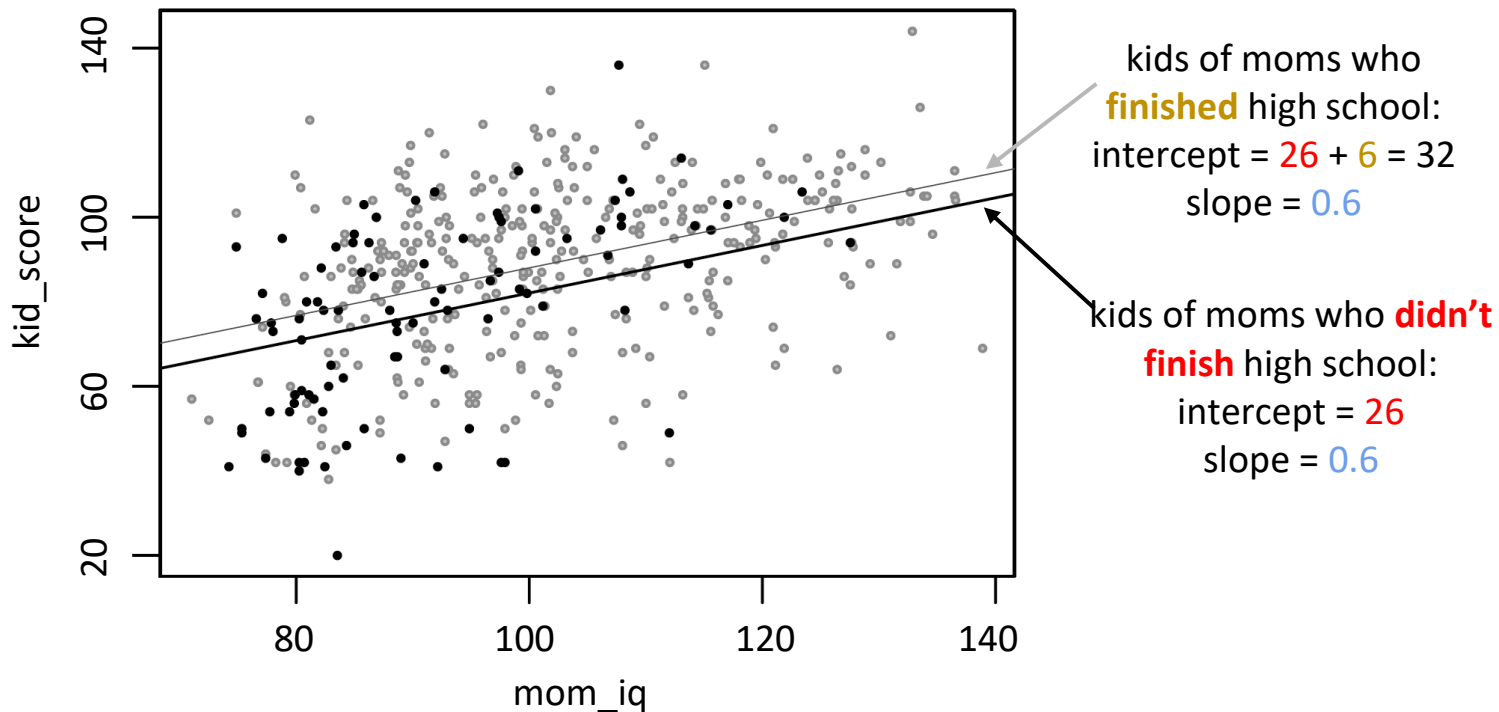
$$y_i = \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3} + \epsilon_i$$

$$\text{kid\_score} = 26 + 6 \cdot \text{mom\_hs} + 0.6 \cdot \text{mom\_iq} + \text{error}$$



# Example with multiple predictors

$$\text{kid\_score} = 26 + 6 \cdot \text{mom\_hs} + 0.6 \cdot \text{mom\_iq} + \text{error}$$



# Example with **interaction** of predictors

- $X_{i2} = \text{mom\_hs} = \text{“Did mother finish high school?”} \in \{0, 1\}$  No Yes
- $X_{i3} = \text{mom\_iq} = \text{mother’s IQ score} \in [70, 140]$
- $y_i = \text{kid\_score} = \text{child’s score on cognitive test} \in [0, 140]$

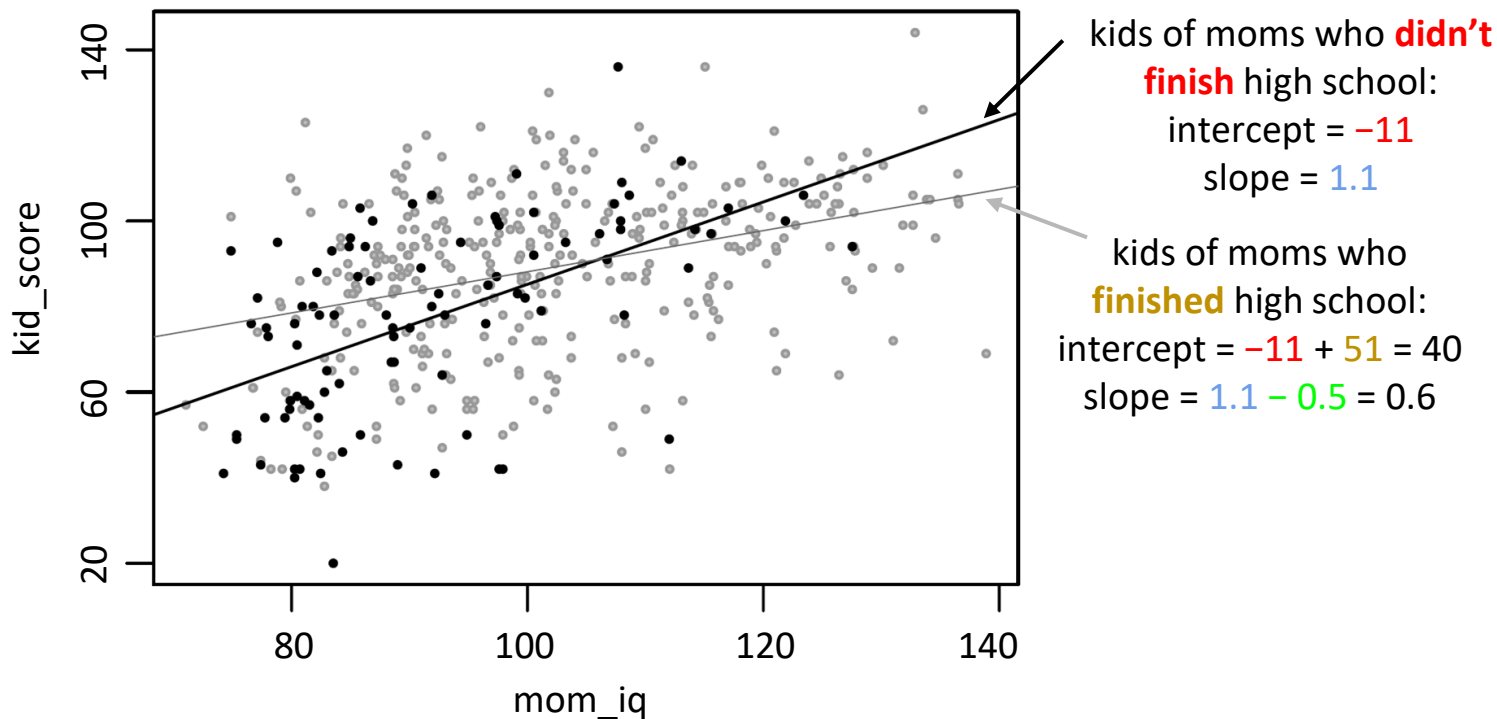
$$y_i = \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i2} X_{i3} + \epsilon_i$$

$$\text{kid\_score} = -11 + 51 \cdot \text{mom\_hs} + 1.1 \cdot \text{mom\_iq} - 0.5 \cdot \text{mom\_hs} \cdot \text{mom\_iq} + \text{error}$$

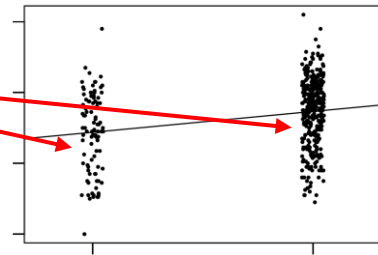


# Example with multiple predictors

$$\text{kid\_score} = -11 + 51 \cdot \text{mom\_hs} + 1.1 \cdot \text{mom\_iq} - 0.5 \cdot \text{mom\_hs} \cdot \text{mom\_iq} + \text{error}$$



So why not just compute  
the two means separately  
and then compare them?



Mom drives Mercedes    Mom doesn't drive Mercedes

Mom  
finished  
high school

avg kid\_score

90

avg kid\_score

90

Mom  
didn't finish  
high school

avg kid\_score

78

avg kid\_score

78

Mom drives Mercedes    Mom doesn't drive Mercedes

Mom  
finished  
high school

990

women

10

women

Mom  
didn't finish  
high school

10

women

990

women

	Mom drives Mercedes	Mom doesn't drive Mercedes		Mom drives Mercedes	Mom doesn't drive Mercedes
Mom finished high school	avg kid_score 90	avg kid_score 90	Mom finished high school	990 women	10 women
Mom didn't finish high school	avg kid_score 78	avg kid_score 78	Mom didn't finish high school	10 women	990 women

# CHAT ROULETTE!

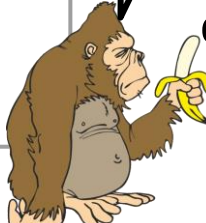
Think for 1 minute:

**What is the mean outcome for Mercedes- vs. non-Mercedes-driving moms?**

**Compare the two means! What does the comparison tell you about the two groups?**

- Then: chat with a fellow student for 3 minutes
  - Rolex Forum: Talk to neighbor (priority: left, right)
  - Zoom: You'll be randomized into small "breakout rooms"

- Mean kid\_score for Mercedes drivers:  $0.99 \cdot 90 + 0.01 \cdot 78 \approx 90$
- Mean kid\_score for non-Mercedes drivers:  $0.01 \cdot 90 + 0.99 \cdot 78 \approx 78$
- But really driving Mercedes makes no difference (for fixed high-school predictor)!
- Root of evil: **correlation** between finishing high school and driving Mercedes
- **Regression** to the rescue:  $\text{kid\_score} = 78 + 12 \cdot \text{mom\_hs} + 0 \cdot \text{mercedes} + \text{error}$

	Mercedes	No Mercedes		Mercedes	No Mercedes
Mom finished high school	mean kid_score 90	mean kid_score 90	<div>Aha!</div> 	990 women	10 women
Mom didn't finish high school	mean kid_score 78	mean kid_score 78		10 women	990 women

# Quantifying uncertainty

# Quantifying uncertainty

- Statistical software gives you more than just coefficients  $\beta$ :

Residuals:

Min	1Q	Median	3Q	Max
-52.873	-12.663	2.404	11.356	49.545

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	25.73154	5.87521	4.380	1.49e-05 ***
mom.hs	5.95012	2.21181	2.690	0.00742 **
mom.iq	0.56391	0.06057	9.309	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.14 on 431 degrees of freedom  
Multiple R-Squared: 0.2141, Adjusted R-squared: 0.2105  
F-statistic: 58.72 on 2 and 431 DF, p-value: < 2.2e-16

**p-value:** probability of estimating such an extreme coefficient if the true coefficient were zero (= null hypothesis)

# Residuals and $R^2$

- **Residual** for data point  $i$ : estimation error on data point  $i$ :

$$r_i = y_i - X_i \hat{\beta}$$

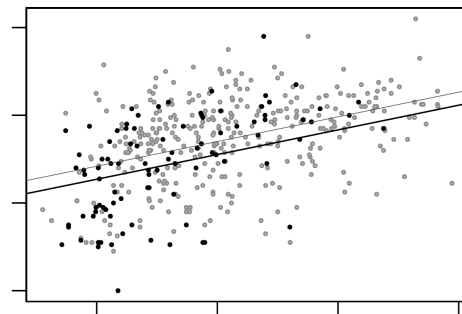
- Mean of residuals = 0  
(total overestimation = total underestimation)

- Standard deviation of residuals  
≈ average distance of predicted value from observed value  
= “unexplained variance”

- Fraction of variance explained by the model:

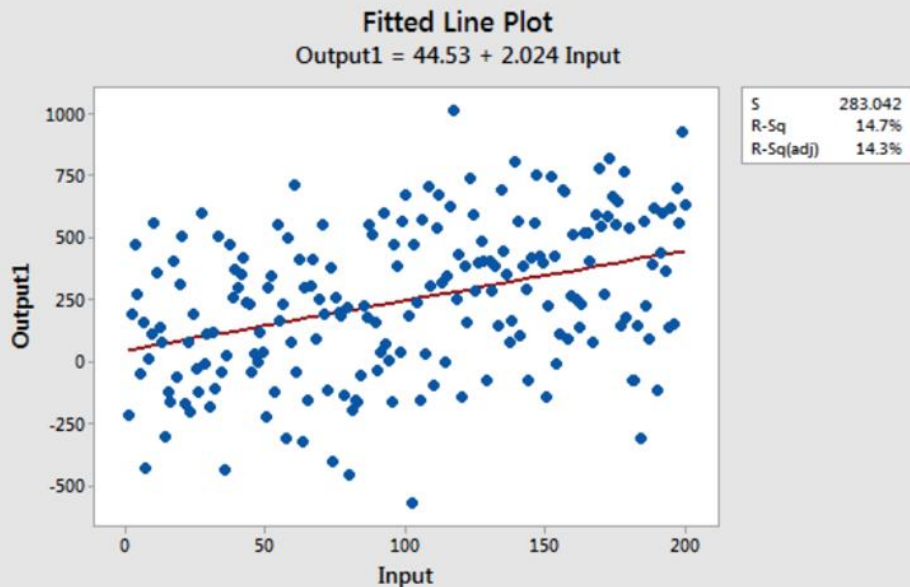
$$R^2 = 1 - \hat{\sigma}^2 / s_y^2$$

Variance of  
outcomes  $y$

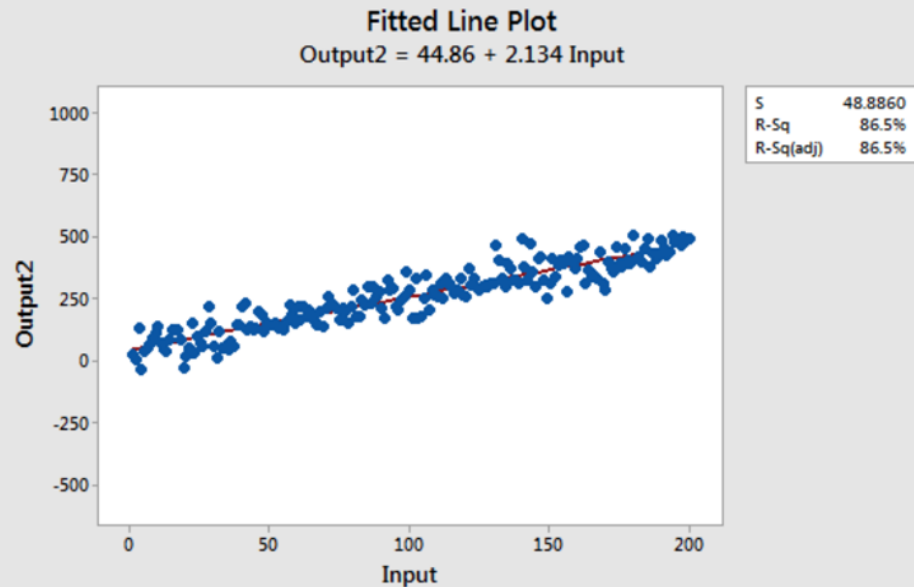


# Coefficient of determination: $R^2$

$$R^2 = 1 - \hat{\sigma}^2 / s_y^2$$



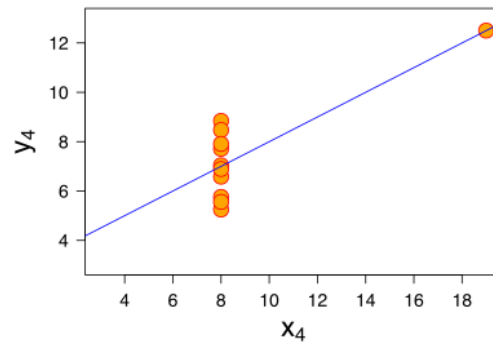
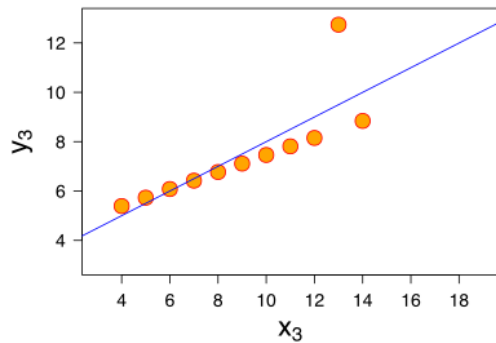
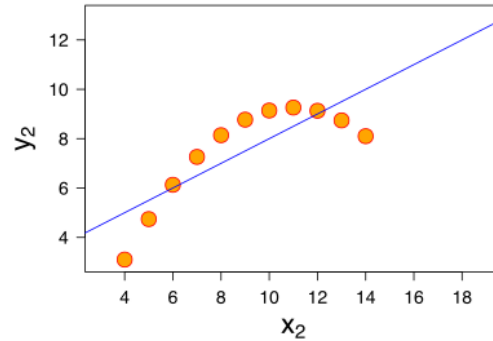
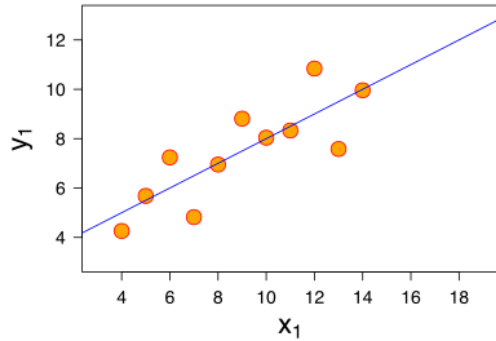
$$R^2 = 0.147$$



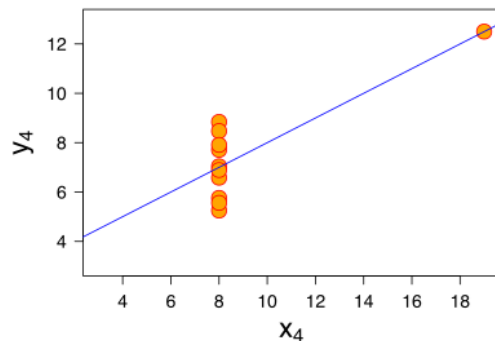
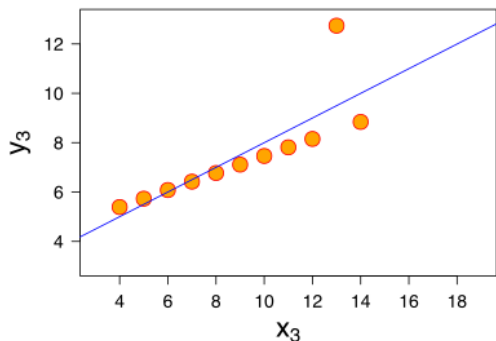
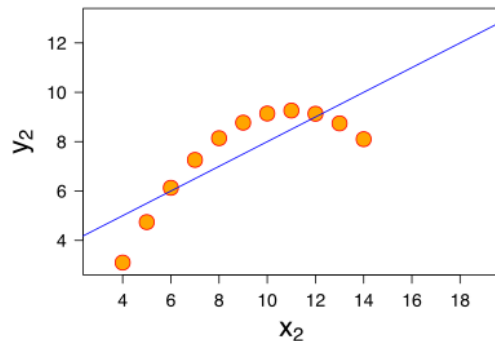
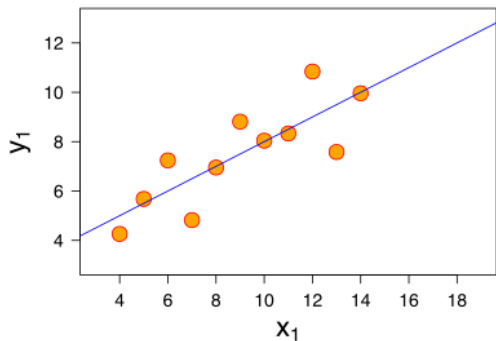
$$R^2 = 0.865$$



# Coefficient of determination: $R^2$



# Coefficient of determination: $R^2$



$R^2 = 0.67$  everywhere!

# Assumptions made in regression modeling

# Assumptions for regression modeling

## 1. Validity:

- a. Outcome measure should accurately reflect the phenomenon of interest
- b. Model should include all relevant predictors
- c. Model should generalize to cases to which it will be applied

# Assumptions for regression modeling (2)

## 2. Additivity and linearity:

$$\begin{aligned}y_i &= X_i\beta + \epsilon_i \\ &= \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + \epsilon_i, \quad \text{for } i = 1, \dots, n\end{aligned}$$

But very flexible: linear in predictors/coefficients (not in necessarily in raw inputs); predictors can be arbitrary functions of raw inputs, e.g.,

- logarithms, polynomials, reciprocals, ...
- interactions (i.e., products) of multiple inputs
- discretization of raw inputs, coded as indicator variables

# Assumptions for regression modeling (3)

- 3. Independence of errors: no interaction between data points
  - 4. Equal variance of errors
  - 5. Normality (Gaussianity) of errors
- } less important  
in practice

# Transformations of predictors and outcomes

# Transformations of predictors

- When we apply linear (technically: affine) transformations to predictors, the model stays linear
- The fitted coefficients may change, but predicted outcomes and model fit won't change
- For instance,

$$\text{earnings} = -61000 + 51 \cdot \text{height (in millimeters)} + \text{error}$$

$$\text{earnings} = -61000 + 81000000 \cdot \text{height (in miles)} + \text{error}.$$



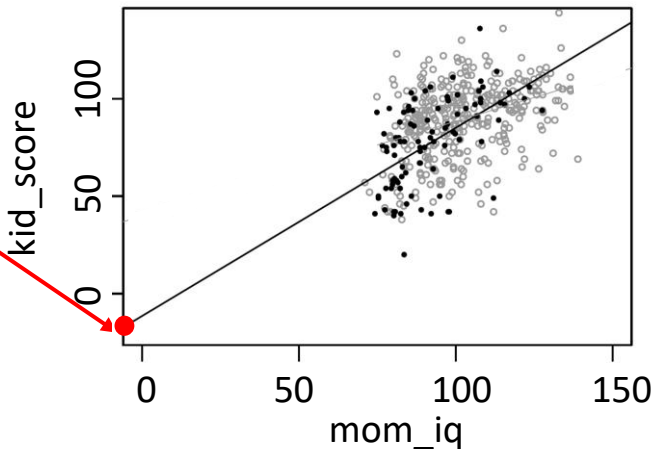
# Mean-centering of predictors

- Compute the mean value of a predictor over all data points, and subtract it from each value of that predictor:

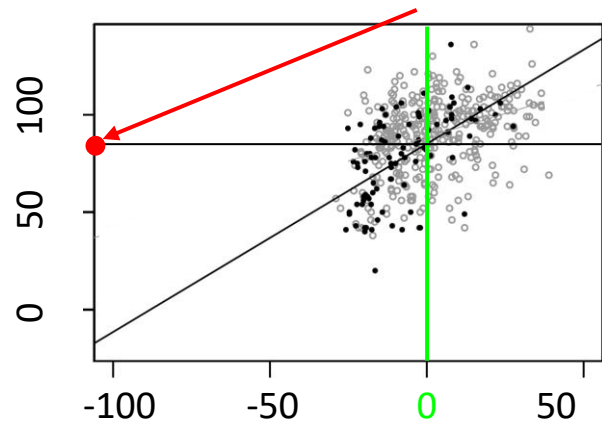
$$X_{ik} \leftarrow X_{ik} - \text{mean}(X_{1k}, \dots, X_{nk})$$

- $\Rightarrow$  the predictor  $X_{ik}$  now has mean 0

(hypothetical) mean  
kid\_score for moms  
with IQ = 0: 26



mean kid\_score for  
moms with mean IQ: 80



# After mean-centering of predictors, ...

... you have a convenient interpretation of coefficients of main predictors (main predictors == non-interaction predictors):

$\beta_k$  = mean increase in outcome  $y$  for each unit increase in  $X_{ik}$   
**when all other predictors take on their mean values**

# Standardization via z-scores

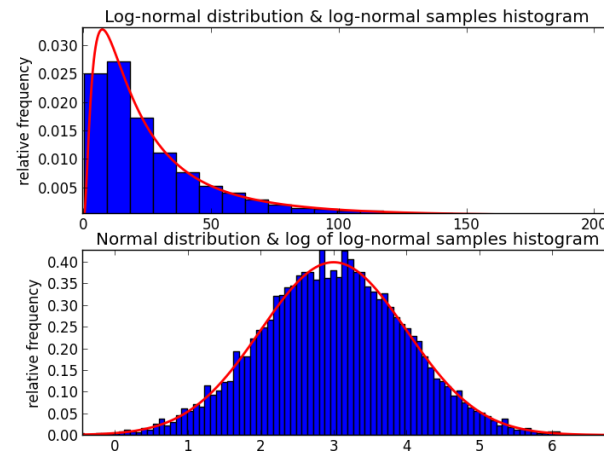
- First **mean-center** all predictors, then **divide them by their standard deviations**:

$$X_{ik} \leftarrow [X_{ik} - \text{mean}(X_{1k}, \dots, X_{nk})] / \text{sd}(X_{1k}, \dots, X_{nk})$$

- All predictors now have the same units (called “**z-scores**”): distance (in terms of standard deviations) from the mean
- This lets us compare coefficients for predictors with previously incomparable units of measurement, e.g., IQ score vs. earnings in Swiss francs vs. height in centimeters

# Logarithmic outcomes

- **Practical:** makes sense if the outcome follows a heavy-tailed distribution
- Only works for non-negative outcomes
- **Theoretical:** turns an additive model into a **multiplicative model**:



$$\log y_i = b_0 + b_1 X_{i1} + b_2 X_{i2} + \cdots + \epsilon_i$$

Exponentiating both sides yields

$$\begin{aligned} y_i &= e^{b_0 + b_1 X_{i1} + b_2 X_{i2} + \cdots + \epsilon_i} \\ &= B_0 \cdot B_1^{X_{i1}} \cdot B_2^{X_{i2}} \cdots E_i \end{aligned}$$

# Logarithmic outcomes: Interpreting coefficients

$$\begin{aligned}y_i &= e^{b_0 + b_1 X_{i1} + b_2 X_{i2} + \dots + \epsilon_i} \\ &= B_0 \cdot B_1^{X_{i1}} \cdot B_2^{X_{i2}} \dots E_i\end{aligned}$$

- An **additive** increase of 1 in predictor  $X_{.1}$  is associated with a **multiplicative** increase of  $B_1 = \exp(b_1)$  in the outcome
- If  $b_1 \approx 0$ , we can immediately interpret  $b_1$  (without needing to exponentiate it first to get  $B_1$ !) as the **relative increase** in outcomes, since  $\exp(b_1) \approx 1 + b_1$
- E.g.,  $b_1 = 0.05 \Rightarrow B_1 = \exp(b_1) \approx 1.05$   
 $\Rightarrow$  “+1 in predictor  $X_{.1}$ ” is associated with “+5% in outcome”

# How to know if your model is appropriately specified?

- Train/test split: split the data set in two parts
  - Fit the model on “training set”
  - Evaluate its accuracy on “testing set”
- If errors are not much larger on testing than on training set, your model

# Going beyond linear regression for comparing means

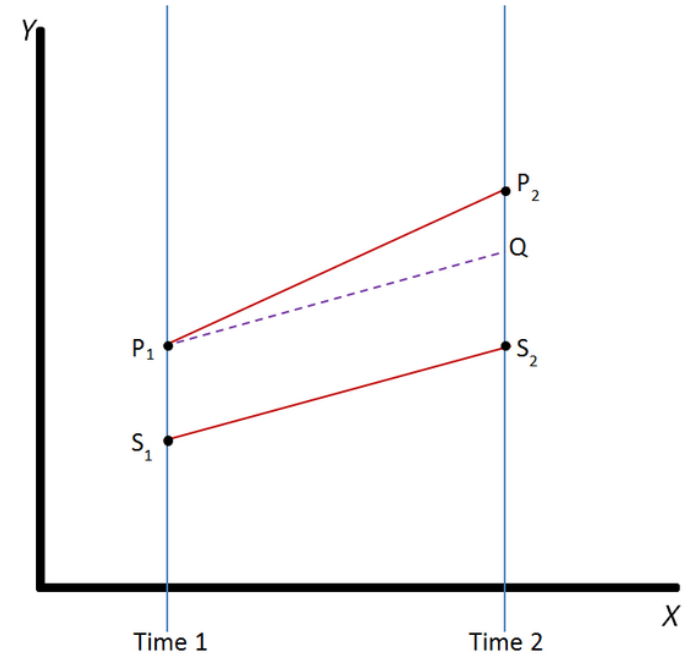
# Beyond linear regression: generalized linear models

- Logistic regression: binary outcomes
- Poisson regression: non-negative integer outcomes (e.g., counts)



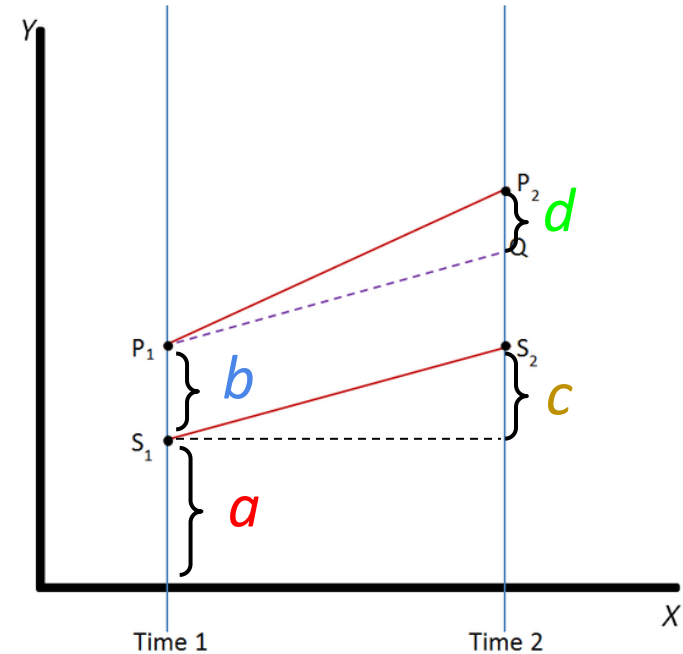
# Beyond comparing means; or, A taste of causality: “Difference in differences”

- Two groups:  $P$ ,  $S$
- At time 2, group  $P$  receives a **treatment**, group  $S$  doesn't
- Question: Did the treatment have an **effect**? If so, how large was it?
- $P$  and  $S$  don't start out the same at time 1
- There is a temporal “baseline effect”



# Beyond comparing means; or, A taste of causality: “Difference in differences” (2)

- Elegant linear model with binary predictors:  
$$y_{it} = a + b \cdot \text{treated}_i + c \cdot \text{time2}_t + d \cdot (\text{treated}_i \cdot \text{time2}_t) + \text{error}$$
- $d$  = treatment effect
- All of this with one single regression!
- You get quantification of uncertainty (significance) for free!



# Summary

- Linear regression as a tool for comparing means across subgroups of data
- How? Read group means off from fitted coefficients
- Advantages over plain comparison of means “by hand”:
  - Accounting for correlations among predictors
  - Quantification of uncertainty (significance) “for free”
  - Additive or multiplicative model: all it takes is a log
- Caveat emptor:
  - Model must be appropriately specified, else nonsense results → stay critical, run diagnostics (e.g.,  $R^2$ )