# Linear models

Jeff Leek

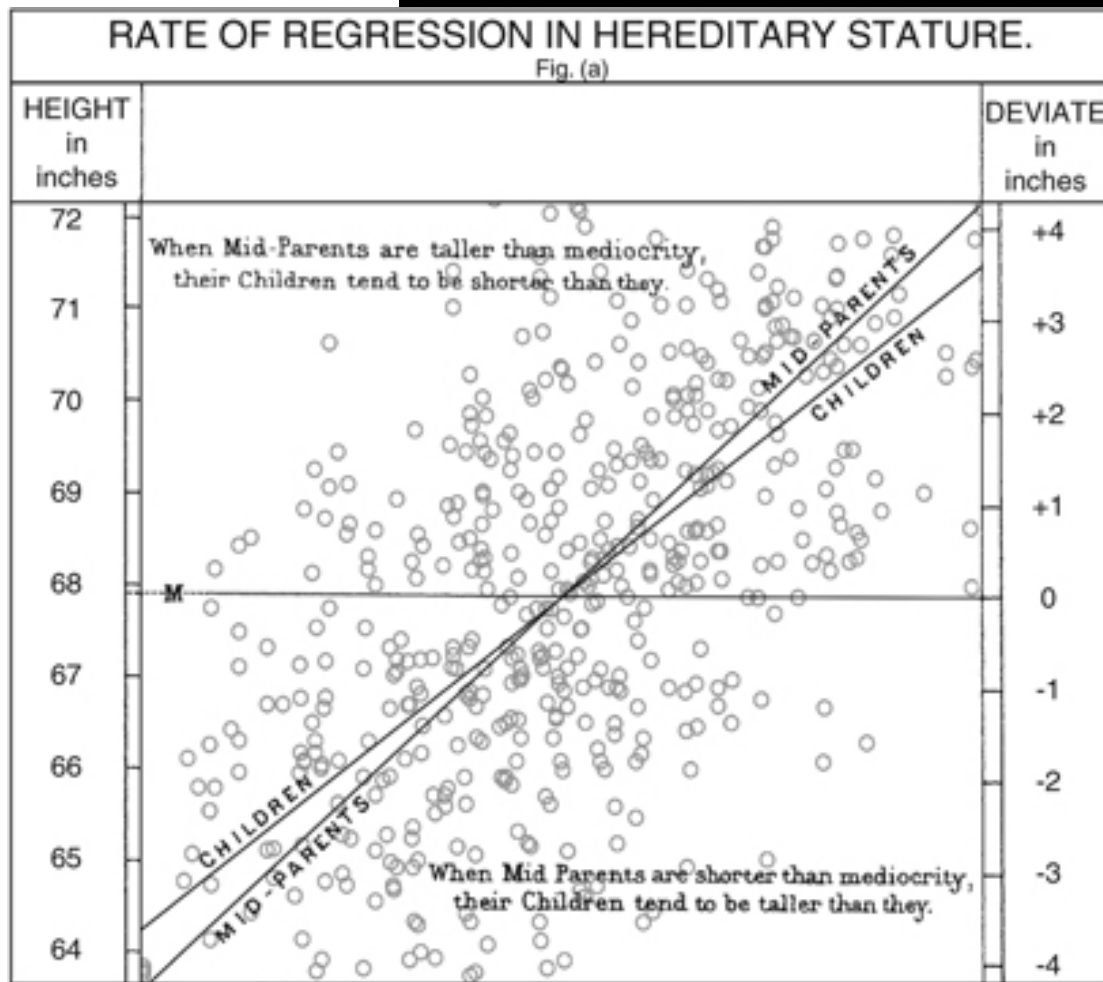@jtleek

www.jtleek.com

# Basic idea

- Fit the "best line" relating two variables

- In math we are minimizing the relationship $(Y - b_0 - b_1X)^2$

- You can always fit a line, the question is whether it is a good fit or not

# An old, but really useful idea!

RATE OF REGRESSION IN HEREDITARY STATURE.

Fig. (a)

HEIGHT in inches | DEVIATE in inches

When Mid-Parents are taller than mediocrity, their Children tend to be shorter than they.

When Mid Parents are shorter than mediocrity, their Children tend to be taller than they.

# Still relevant everywhere in genomics

# Article

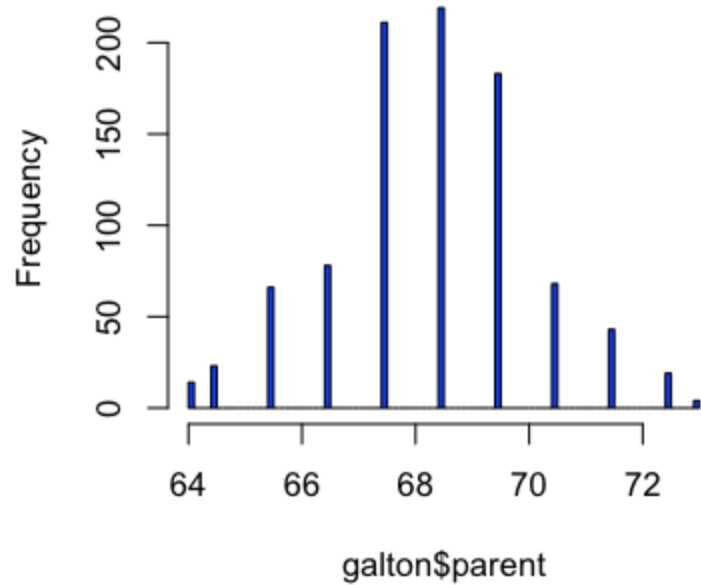# Predicting human height by Victorian and genomic methods

Yurii S Aulchenko[1,2,7], Maksim V Struchalin[1,3,7], Nadezhda M Belonogova[2,4], Tatiana I Axenovich[2], Michael N Weedon[5], Albert Hofman[1], Andre G Uitterlinden[6], Manfred Kayser[3], Ben A Oostra[1], Cornelia M van Duijn[1], A Cecile J W Janssens[1] and Pavel M Borodin[2,4]
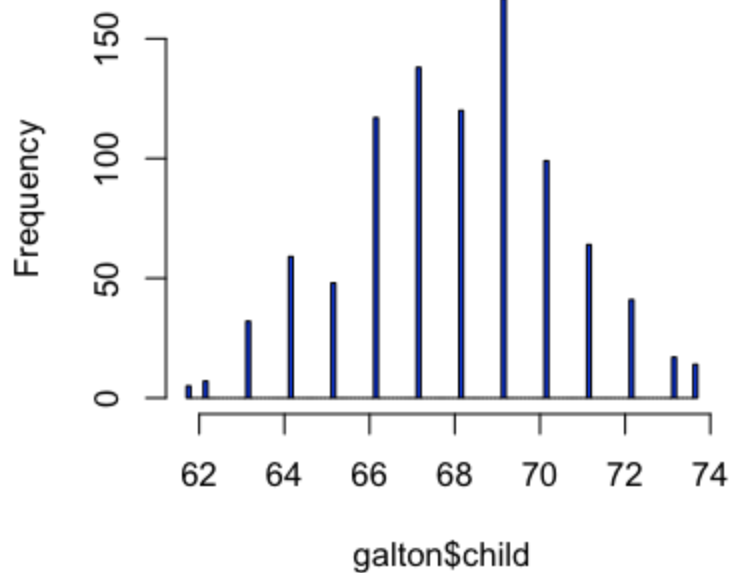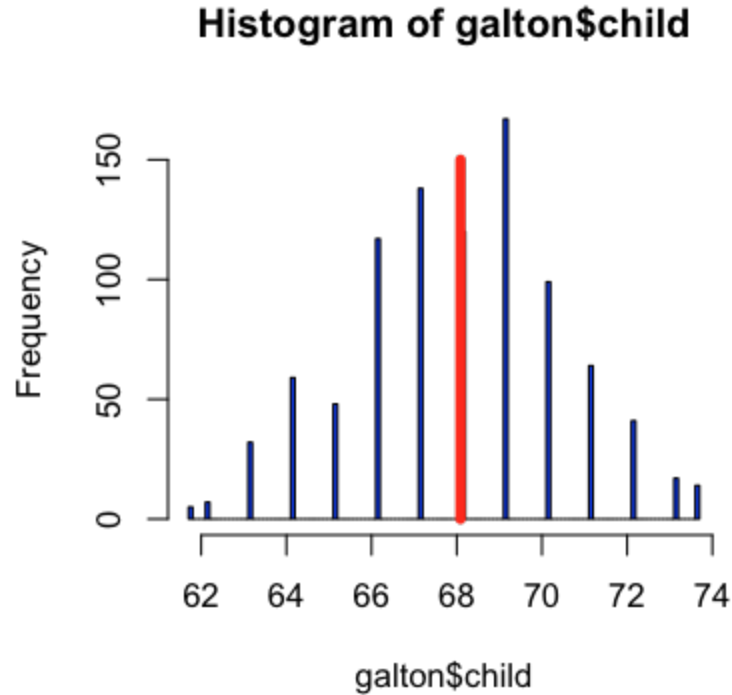
# The Galton example explained

Histogram of galton$child
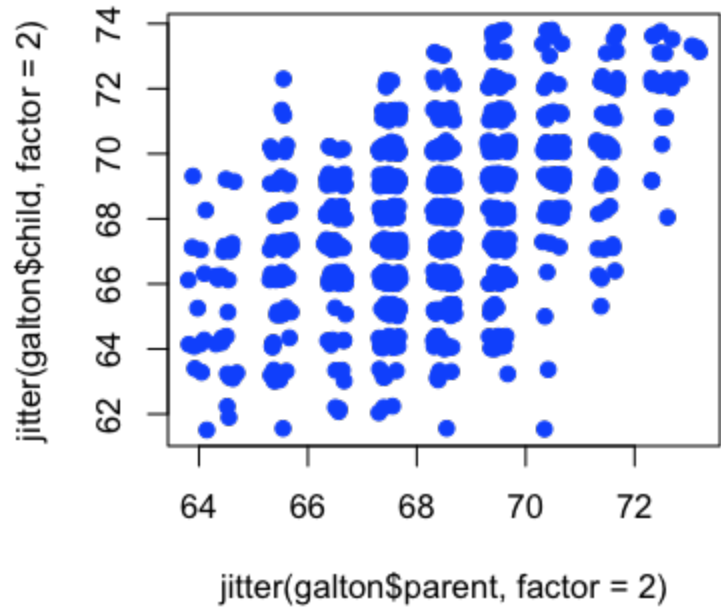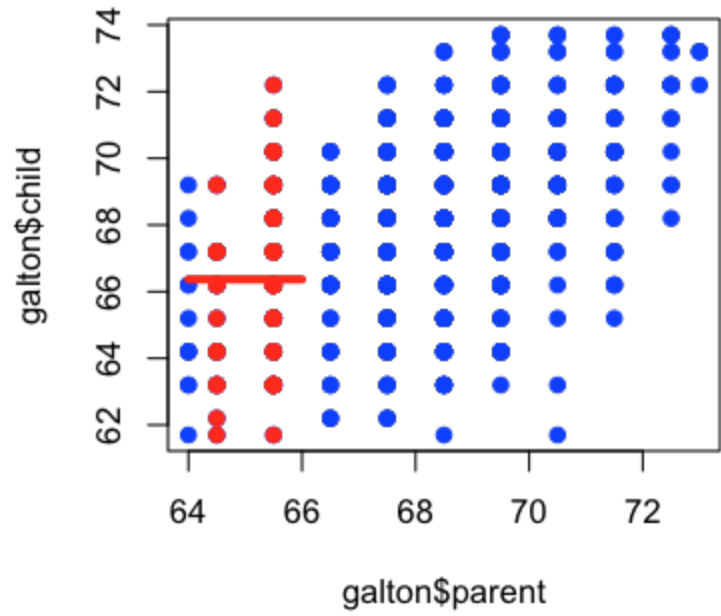
Histogram of galton$parent

Histogram of galton$child

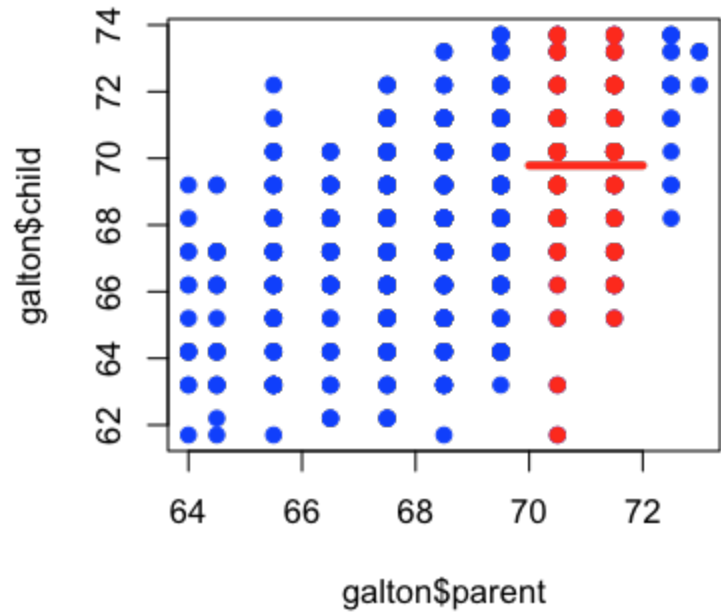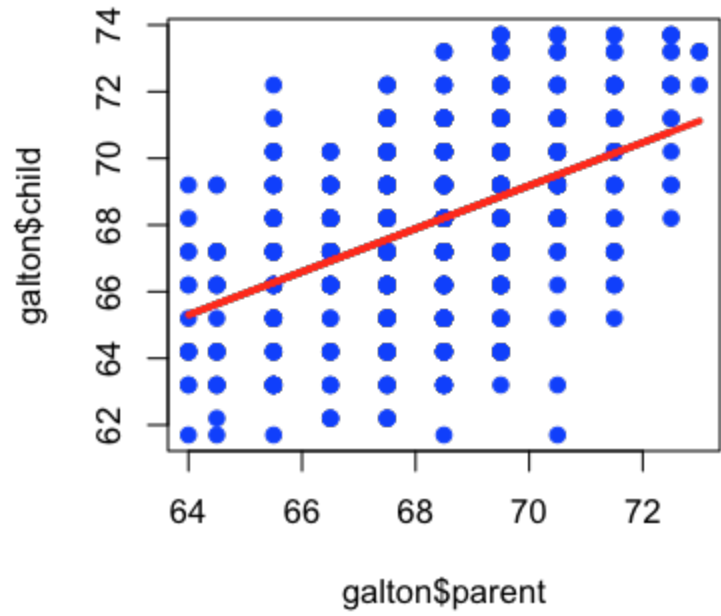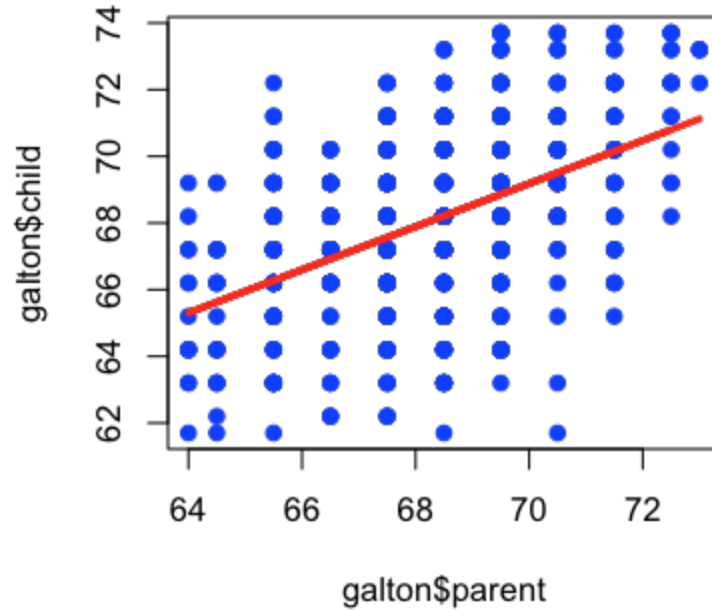**Histogram of galton$child**

Minimizes $\sum (Y - c)^2$

# Equation for a line:

$$C = b_0 + b_1 P$$

# Not all points exactly on the line

# Equation with noise
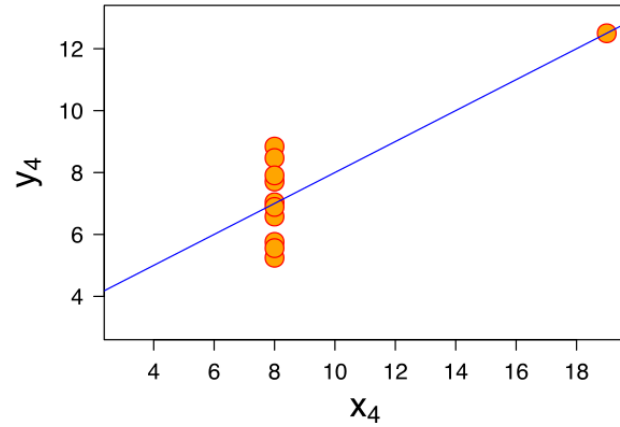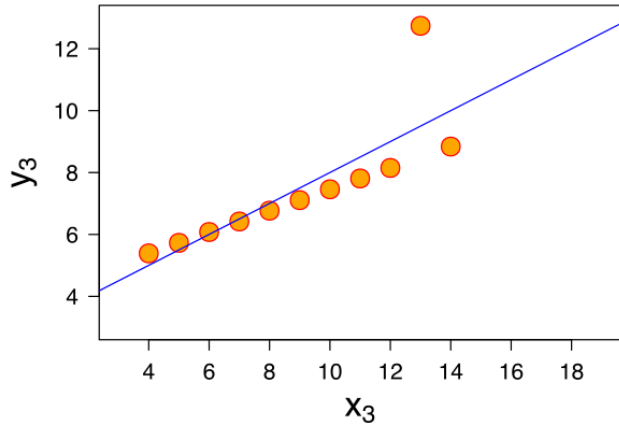
$$C = b_0 + b_1 P + \textcolor{red}{e}$$

# Fit by minimizing
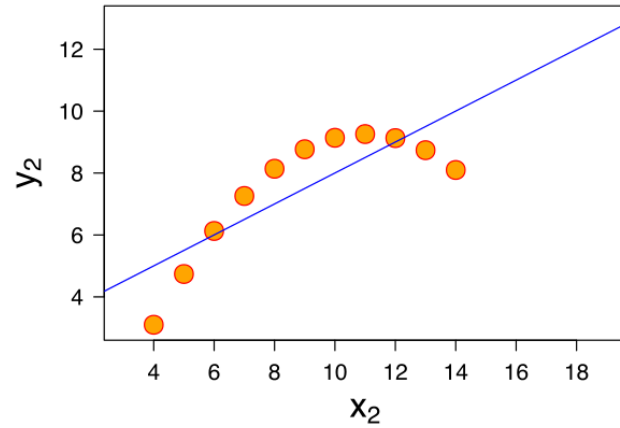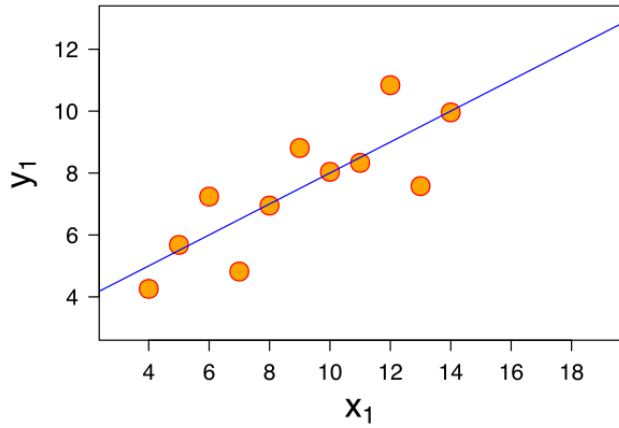
$$\sum(C - b_0 - b_1 P)^2$$

# Check residuals

You can always fit a line
 But it might not be the right thing

# Notes and further reading

- Linear models is a whole class (no joke): https://www.coursera.org/course/regmods

- Basic thing to keep in mind is does the line fit?

- Great additional notes in Chapter 2 here: http://genomicsclass.github.io/book/