



Introduction to regression

Regression

Brian Caffo, Jeff Leek and Roger Peng
Johns Hopkins Bloomberg School of Public Health

Questions for this class

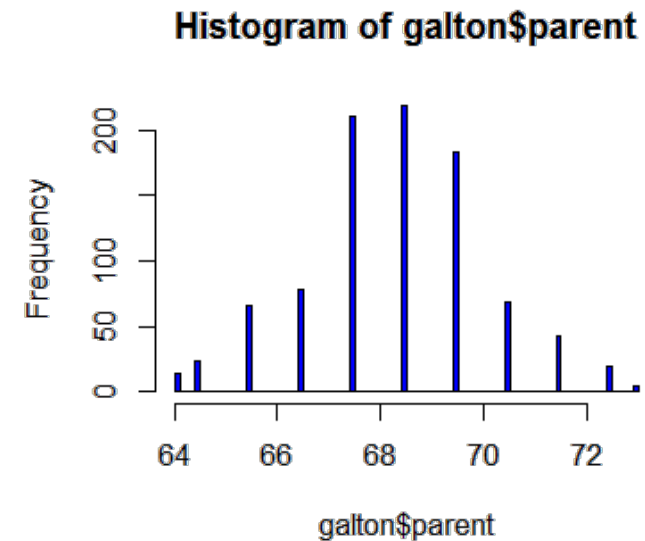
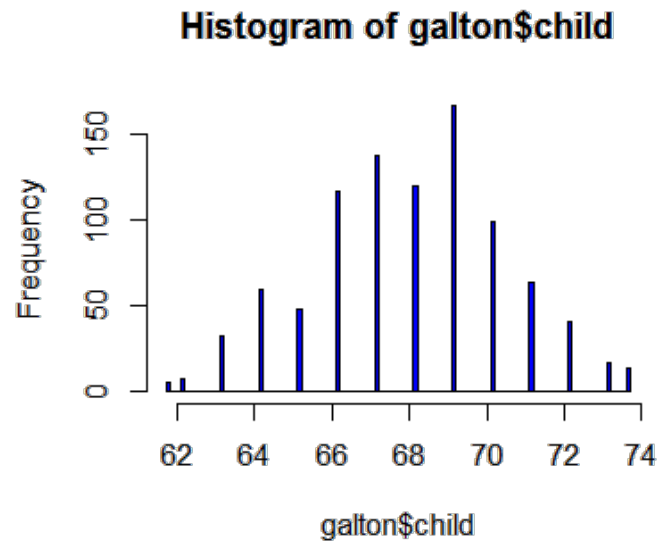
- Consider trying to answer the following kinds of questions:
 - To use the parents' heights to predict childrens' heights.
 - To try to find a parsimonious, easily described mean relationship between parent and children's heights.
 - To investigate the variation in childrens' heights that appears unrelated to parents' heights (residual variation).
 - To quantify what impact genotype information has beyond parental height in explaining child height.
 - To figure out how/whether and what assumptions are needed to generalize findings beyond the data in question.
 - Why do children of very tall parents tend to be tall, but a little shorter than their parents and why children of very short parents tend to be short, but a little taller than their parents? (This is a famous question called 'Regression to the mean'.)

Galton's Data

- Let's look at the data first, used by Francis Galton in 1885.
- Galton was a statistician who invented the term and concepts of regression and correlation, founded the journal Biometrika, and was the cousin of Charles Darwin.
- You may need to run `install.packages("UsingR")` if the `UsingR` library is not installed.
- Let's look at the marginal (parents disregarding children and children disregarding parents) distributions first.
 - Parent distribution is all heterosexual couples.
 - Correction for gender via multiplying female heights by 1.08.
 - Overplotting is an issue from discretization.

Code

```
library(UsingR); data(galton)
par(mfrow=c(1,2))
hist(galton$child,col="blue",breaks=100)
hist(galton$parent,col="blue",breaks=100)
```



Finding the middle via least squares

- Consider only the children's heights.
 - How could one describe the "middle"?
 - One definition, let Y_i be the height of child i for $i = 1, \dots, n = 928$, then define the middle as the value of μ that minimizes

$$\sum_{i=1}^n (Y_i - \mu)^2$$

Sum of squared distances
between 'middle' and each child's height

=> it's just the mean!

- This is physical center of mass of the histogram.
- You might have guessed that the answer $\mu = \bar{X}$.

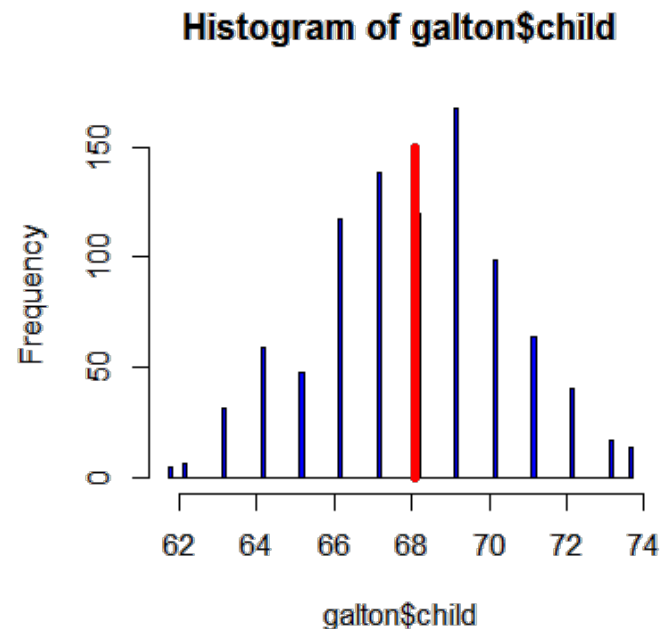
Experiment

Use R studio's manipulate to see what value of μ minimizes the sum of the squared deviations.

```
library(manipulate)
myHist <- function(mu){
  hist(galton$child,col="blue",breaks=100)
  lines(c(mu, mu), c(0, 150),col="red",lwd=5)
  mse <- mean((galton$child - mu)^2)
  text(63, 150, paste("mu = ", mu))
  text(63, 140, paste("MSE = ", round(mse, 2)))
}
manipulate(myHist(mu), mu = slider(62, 74, step = 0.5))
```

The least squares estimate is the empirical mean

```
hist(galton$child,col="blue",breaks=100)  
meanChild <- mean(galton$child)  
lines(rep(meanChild,100),seq(0,150,length=100),col="red",lwd=5)
```



The math follows as:

$$\begin{aligned}\sum_{i=1}^n (Y_i - \mu)^2 &= \sum_{i=1}^n (Y_i - \bar{Y} + \bar{Y} - \mu)^2 \\&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + 2 \sum_{i=1}^n (Y_i - \bar{Y})(\bar{Y} - \mu) + \sum_{i=1}^n (\bar{Y} - \mu)^2 \quad a^2 + 2ab + b^2 \\&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + 2(\bar{Y} - \mu) \sum_{i=1}^n (Y_i - \bar{Y}) + \sum_{i=1}^n (\bar{Y} - \mu)^2 \quad \begin{array}{l} \text{indep. von } i \\ \rightarrow \text{rausziehen} \end{array} \\&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + 2(\bar{Y} - \mu) \left(\sum_{i=1}^n Y_i - n\bar{Y} \right) + \sum_{i=1}^n (\bar{Y} - \mu)^2 \\&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + \sum_{i=1}^n (\bar{Y} - \mu)^2 \quad \begin{array}{l} \text{sum_i_bis_n}(Y_i)/n = \text{mean}(Y) \\ \text{Daher ist diese Klammer grad 0.} \end{array} \\&\geq \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad \begin{array}{l} > 0 \text{ da Summe von} \\ &\text{Quadraten} \end{array}\end{aligned}$$

Aufgrund dieser Ungleichung:

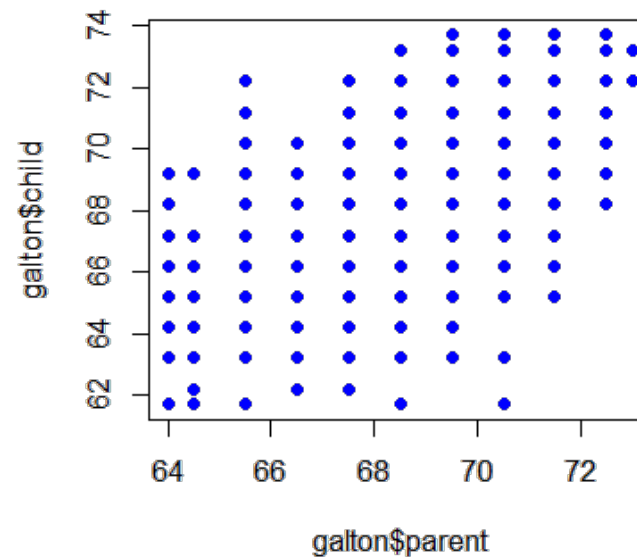
Fuer /jedes/ μ , gilt:

diese Summe der quadrierten Distanzen ist /groesser/
als wenn man statt μ $\text{mean}(Y)$ einsetzt!

Also muss $\text{mean}(Y)$ das Minimum sein. HX, huebscher Beweis ohne Analysis.

Comparing childrens' heights and their parents' heights

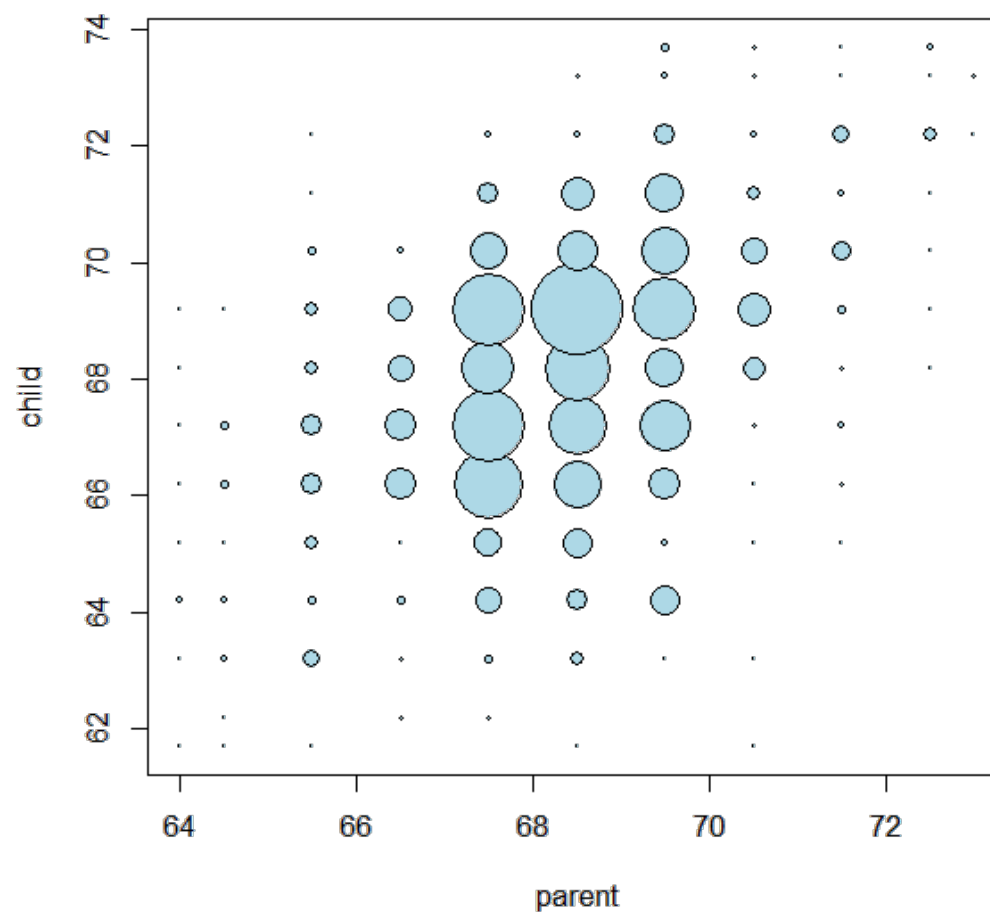
```
plot(galton$parent,galton$child,pch=19,col="blue")
```



Overplotted:

each point represents several points.
See next slide for how to represent
this better.

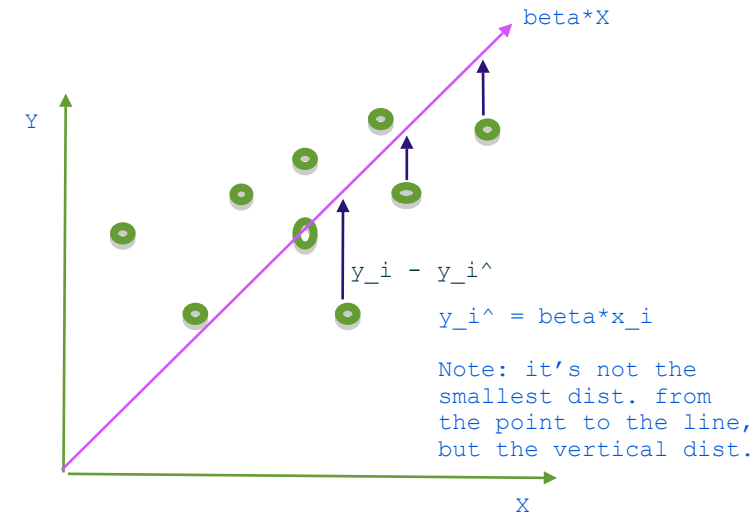
Size of point represents number of points at that (X, Y) combination (See the Rmd file for the code).



Regression through the origin

- Suppose that X_i are the parents' heights.
- Consider picking the slope β that minimizes

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



- This is exactly using the origin as a pivot point picking the line that minimizes the sum of the squared vertical distances of the points to the line
- Use R studio's manipulate function to experiment
- Subtract the means so that the origin is the mean of the parent and children's heights

das ist besser als durch (0,0) !

```

myPlot <- function(beta){
  y <- galton$child - mean(galton$child)
  x <- galton$parent - mean(galton$parent)
  freqData <- as.data.frame(table(x, y))
  names(freqData) <- c("child", "parent", "freq")
  plot(
    as.numeric(as.vector(freqData$parent)),
    as.numeric(as.vector(freqData$child)),
    pch = 21, col = "black", bg = "lightblue",
    cex = .15 * freqData$freq,
    xlab = "parent",
    ylab = "child"
  )
  abline(0, beta, lwd = 3)
  points(0, 0, cex = 2, pch = 19)
  mse <- mean( (y - beta * x)^2 )
  title(paste("beta = ", beta, "mse = ", round(mse, 3)))
}
manipulate(myPlot(beta), beta = slider(0.6, 1.2, step = 0.02))

```

The solution

In the next few lectures we'll talk about why this is the solution

In R: mit `lm`.

```
lm(I(child - mean(child))~ I(parent - mean(parent)) - 1, data = galton)
```

mean abziehen, so dass
Drehpunkt das mean ist

-1: don't use an intercept

Call:

```
lm(formula = I(child - mean(child)) ~ I(parent - mean(parent)) -  
    1, data = galton)
```

Coefficients:

```
I(parent - mean(parent))  
0.646
```

Visualizing the best fit line

Size of points are frequencies at that X, Y combination

