

## Week 4

### Clustering Example

Single Value Decomposition

```
svd1 = svd(scale(samsungData,  
  [samsungData$subject != 1,  
    - c(562, 563)]))
```

```
par(mfrow = c(1, 2))
```

```
plot(svd1$u[, 1], col = numericalActivity,  
     pch = 19)
```

```
plot(svd1$u[, 2], col = ..., pch = 19) --
```

↑  
left singular vector  
explains the most variation

find maximum contributor

```
plot(svd1$v[, 2], pch = 19)
```

↑  
weight with which these  
variables contribute to  
variance

So we may calculate the variables with maximal weight

maxContrib = which.max(svd\$V[, 2])

distanceMatrix <- dist(samplingData  
[sd\$subject == 1, c(10:12, maxContrib)])

↓  
then draw it

## Regression

### Basic Least Squares

#### Goals of Statistical Modelling

- describe the distribution
- describe the relationship

how to describe the distribution

lm  
↑ linear model

$$C_i = b_0 + b_1 P_i + e_i$$

↑  
error term

lm\$(galton\$child ~ galton\$parent)  
lm\$fit

lm\$residual

↑  
left overs

$$\sum_{i=1}^n \underbrace{(C_i - \{b_0 + b_1 P_i\})^2}_{\text{smallest}}$$

how to summarize a distribution?

if you ~~would~~ ~~have~~ to choose only one number to describe some data, what would you use?



mean (for symmetric distribution)

↗  
center of distribution

also, it minimizes

$$\sum_i (C_i - \mu)^2$$

$C_i$  - observation

↓  
value that

minimizes this formula is avg of  $C$

'jitter' - add some noise to a variable

$\text{lm}(x \sim y)$

↑  
quantitative  
var

↑  
variables to include  
in your model

equation for a line

$$C_i = b_0 + b_1 P_i + e_i$$

$\uparrow$                        $\uparrow$                        $\uparrow$  (error term)  
intercept              slope              everything we didn't  
term                      measure  
(where the line  
crosses 0)

best line is:

$$\sum_i (C_i - \underbrace{\{b_0 + b_1 P_i\}}_{\text{smallest}})^2$$

lin 1 § fitted - line, found with least squares

(how many variations around this line?)

$\downarrow$

lin 1 § residuals - left overs

## Inference Basics

(brlog, brgenesis)

```
library(MuMIn); data(galton);
```

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

```
lm1 = lm(galton$child ~ galton$parent)
```

```
lines(galton$parent, lm1$fitted, col="red",  
      lwd=3)
```

lm1

Coefficients

(Intercept)

23.942

galton\$parent

0.646

$$C_i = b_0 + b_1 P_i$$

let's generate a population of families

```
newGalton = data.frame(parent = rep(NA, 1e6),  
                        child = rep(NA, 1e6))
```

- newGalton\$parent = rnorm(1e6,  
 mean = mean(galton\$parent),  
 sd = sd(galton\$parent))

with the same properties as our sample

- newGalton\$child =  $\underbrace{\text{Intercept}}_{b_0} + \underbrace{\text{slope}}_{b_1} \cdot \underbrace{\text{newGalton\$parent}}_{\text{child}}$   
+ rnorm(1e6, sd = sd(lm1\$residuals))  
error term (noise)

smooth Scatter (new Galton \$ parent,  
new Galton \$ child)

abline (lm1, col = "red", lwd = 3)

↑  
regression line from Galton's Data  
fits into newly generated data!

Let's take a sample

(from newly generated data)

sample Galton1 = newGalton [sample(1:1000,  
size = 50, replace = F), ]

50 el. from  
generated  
data

sample lm1 = lm (sample Galton1 \$ child,  
sample Galton1 \$ parent)

↑  
Linear model from ~~ge~~ sample  
and it's different from original lm1!

plot (sample Galton1 \$ parent,  
sample Galton1 \$ child, pch=9, col = "blue")

different! { lines (sample Galton1 \$ parent, sample lm1 \$ fitted,  
lwd = 3, lty = 2) ← new  
abline (lm1, col = "red", lwd = 3) ← old

if we take another sample, we'll see again  
different results

So line you get from sample isn't the same line for the population

Histogram of estimates

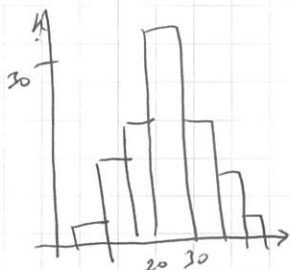
par(mfrow = c(1, 2))  
 hist(intercept coeffs, ...)  
 hist(slope coeffs, ...)

coefficients from  
 lm run many  
 times on different  
 samples

slope coeffs

```

i = list of 100
for (i in 1:100) {
  sam = sample from population
  l[i] = lm(sam)
}
  
```



How these coefficients are distributed?

From the central limit theorem

estimate of int  $\hat{b}_0 \sim N(b_0, \text{Var}(\hat{b}_0))$   
 estimate of slope  $\hat{b}_1 \sim N(b_1, \text{Var}(\hat{b}_1))$

} follow normal distribution

we don't know  
 this var. exactly, but  
 we can  
 estimate them  
 as well

$\hat{b}_0 - b_0$  estimated

estimation:

$$\hat{b}_0 \sim N(\hat{b}_0, \text{Var}(\hat{b}_0))$$

$$\hat{b}_1 \sim N(\hat{b}_1, \text{Var}(\hat{b}_1))$$

$\sqrt{\text{Var}(\hat{b}_0)}$  - "standard error" of the estimate  
 $\hat{b}_0$   
(S.E. ( $\hat{b}_0$ ))

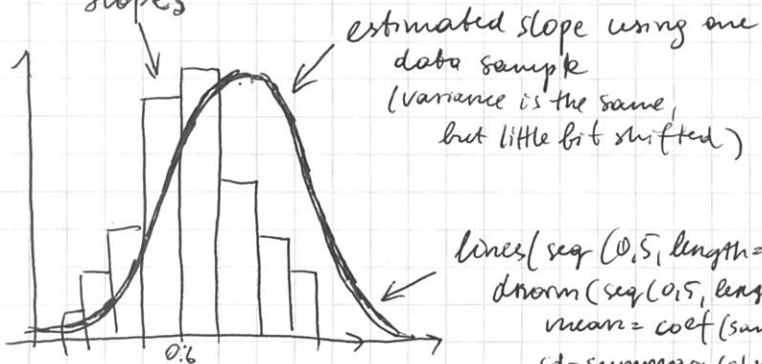
summary(sampleLm1)



gives you residuals  
and estimates for the  
coefficient

as well as S.E.s  
for  $b_0$  and  $b_1$

distribution of  
Slopes



```
lines(seq(0.5, length=100),  
dnorm(seq(0.5, length=100),  
mean=coef(summary(sampleLm1))[2,  
sd=summary(sampleLm1)$coeff[2,2],  
lwd=3,  
col="red")
```



## Standardized coefficients

$$\hat{b}_0 \approx N(b_0, \hat{\text{Var}}(\hat{b}_0))$$

dt in R

$$\frac{\hat{b}_0 - b_0}{\text{S.E.}(\hat{b}_0)} \sim t_{n-2}$$

← center  
← t-distribution (as n grows, it gets closer to  $N(0,1)$ )  
← 2 degrees of freedom  
↑  
number of samples you had

degrees of freedom tell you how much variation you have left over after estimating your parameters

we loose 2 degrees of freedom when we calculate mean and the slope

$$\hat{b}_1 \approx N(b_1, \hat{\text{Var}}(\hat{b}_1))$$
$$\frac{\hat{b}_1 - b_1}{\text{S.E.}(\hat{b}_1)}$$

## Confidence Intervals

We have an estimated  $\hat{b}_1$  and want to know how good our estimate is

⇓

Create a "level  $\alpha$  confidence interval"

set of plausible values for  $b_1$

So a confidence interval will include the real parameter  $\beta_1$  in

$\alpha\%$  of the time in repeated studies (i.e. sampling)

↓  
(so if we sample 100 times, in  $\alpha$  samples  $\beta_1$  will be in the calculated interval)

•  $\alpha$ -confidence interval

$$\rightarrow (\hat{\beta}_1 - T_{\alpha/2} \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + T_{\alpha/2} \cdot SE(\hat{\beta}_1))$$

↖  
can be calculated using `lm` command

↑  
particular quantile of  $t$ -distribution

↓  
summary(sample `lm`)\$coeff

• `confint`

~~confint~~ (sample `lm` \$ $\beta_1$ , level = 0.95)

will calculate confidence interval

## How to report the inference

A one inch increase  
in parental height  
is associated with  
a 0.77 inch increase  
in child's height  
(95%, CI: 0.42 - 1.12 inches)

confint (sampleLm,  
level = 0.95)

(Intercept) -7.8 39.5

parent 0.42 1.19

sampleLm\$coeff

(Intercept)

15.86

parent

0.7698

## P-values

most commonly reported measure of  
"statistical significance"

Ideally suppose we see nothing

how likely/unlikely that there is  
a relationship between variables?

## Approach

1. Define the hypothetical distribution of a  
data summary (statistic)

(null hypothesis)

"when nothing is going on"

2. Calculate the summary / test statistics with the data we have
3. Compare what we calculated to our hypothetical distribution and see if the value is "extreme" (p-value)

null hypothesis

$$\frac{\hat{b}_1 - b_1}{S.E.(\hat{b}_1)} \sim t_{n-2}$$

↓

$H_0$ : There is no relationship between parent and child height  
(i.e.  $b_1 = 0$ )

Under the null hypothesis the distribution is:

$$\frac{\hat{b}_1}{S.E.(\hat{b}_1)} \sim t_{n-2}$$

↑  
now we compare it with  $t_{n-2}$  and we see if there are any surprises



Some typical values

used to see if there  
are any relationships  
between 2 variables

$P < 0.05$   
statistically  
significant

← usually used

$P < 0.01$   
strongly significant

$P < 0.001$  very significant

Usually in report both confidence interval  
and p-value are mentioned

How to interpret?

`summary(lm(g$child ~ g$parent))$coeff`

	Estimate	S.E.	t	Pr(> t )
(intercept)	23.97	2.81	8.5	6.53e-17
g\$parent	0.64	0.04	15.7	(1.73e-49)

A one inch increase in parental height is  
associated with a 0.77 inch increase  
in child's height

(95% CI: 0.42 - 1.12 inches)

This difference was statistically significant  
( $P < 0.001$ )

## Regression with factor variables

rotten tomatoes score vs rating

```
plot(movies$score ~ jitter(as.numeric (
  movies$rating)), col="blue",
      xaxt="n", pch=19)
```

↑  
factor

```
meanRatings = tapply(movies$score,
  movies$rating, mean)
```

```
points(1:4, meanRatings, col="red",
      pch="-", las=5)
```

Another way to write it down:

$$\mathbb{1}(R_{ai} = "PG") = \begin{cases} 1 & \text{if } R_{ai} = "PG" \\ 0 & \text{if } R_{ai} \neq "PG" \end{cases}$$

← regression) like

$$S_i = b_0 + b_1 \mathbb{1}(R_{ai} = "PG") + b_2 \mathbb{1}(R_{ai} = "PG-13") + \\ + b_3 \mathbb{1}(R_{ai} = "R") + \dots$$

$b_0$  - avg of G

$b_0 + b_1$  = avg of PG     $b_0 + b_2$  = avg of PG-13

$b_0 + b_3$  = avg of R

`lm1 = lm(movies$score ~ as.factor(movies$rating))`  
`summary(lm1)`



and we get estimate for  
each movie type

`lm1 =` --

`anova(lm1)` ← analysis of variance table  
(check what anova is)



# Multiple Regressions

- regression with multiple covar. abcs
- still using least square / central limit theorem

$$\text{lm Both} = \text{lm}(\text{hunger} \sim \text{Numeric} + \text{hunger} \sim \text{Year} + \text{hunger} \sim \text{Sep})$$

↑ categorize by sex  
same slope

$$\text{lm Both} = \text{lm}(\text{hunger} \sim \text{Numeric} + \text{hunger} \sim \text{Year} + \text{hunger} \sim \text{Sex} + \text{hunger} \sim \text{Sex} * \text{hunger} \sim \text{Year})$$

↑  
2 slopes

# Regression in the Real World

Ideal data: galton data set  
(~~data~~ cloud - shape)

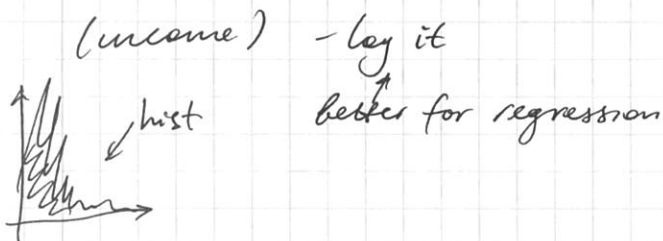
## • Confounders

a variable that is correlated with both outcome and the covariates

- confounders can change the regression line
- change the sign of the line!
- can be detected by careful exploration

↳ data visualization

## • right-skew



## • outliers

outliers - data points that do not appear to follow the pattern of the other data points

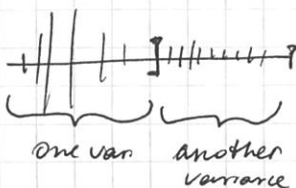
can have dramatic impact on how you visualize data and on your analysis

what can you do?

- you have to know if it's real or not
- if not - remove it
- or
  - logarithm it
  - sensitivity analysis (leave out for a while)
  - robust methods -

• A line isn't always the best summary  
(wiki - Linear Regression)

- changing variance;  
when there are 2 variances



- Box-Cox transform
- variance stabilizing transform
- weighted least squares
- etc

## • Units

Absolute vs Relative

- standardize, but keep track

but it affects

- model fits
- interpretation
- inference

## • overloading regression

too many params in regression

## • correlation and causation

be critical  
/ consider  
alternative

was  $X$  indeed the cause for  $Y$ ?

(chocolate consumption vs nobel prize winners)