

Data Analysis

1. Comp. Programming
2. Statistics
3. Applications

R package Slidyfy -

<http://ramnathv.github.com/slidyfy>

Slides - [github.com / jtleeh / dataanalysis](https://github.com/jtleeh/dataanalysis)

↑
all classes are available there
as well as slides

to compile them:

```
library(slidyfy)  
slidyfy("index.Rmd")
```

Week 1

Data - values of ~~quant~~ qualitative or
quantitative variables, belonging to
a set of items

set of items - subjects
variables - measurements

Probe

Raw data

hard to use
complex format

Processed

ready for analysis

each variable ^{forms a column}
each observation ^{forms a row}
each file stores data about one ^{test}

kind of observation

Big data - data that is too huge to manipulate on a single computer

How to represent data?

H for height, W for weight } informative

Randomness

- represents incompletely ~~var~~ measured variable
- @ random mechanism

Distributions

X - random value

P - probability

$$0 \leq x \leq 1$$

$$p_1 + \dots + p_n = 1$$

↓
continuous

↓
discrete

Parameters:

$$N(\mu, \sigma) \quad \text{Poisson}(\lambda)$$

$X \sim N(\mu, \sigma)$ means X has
the $N(\mu, \sigma)$ distribution

The most important parameters

$E[X]$ - expected value

$\text{Var}[X]$ - how 'spread out' a distribution is
measured in $(\text{unit of } X)^2$

$$\text{SD}[X] = \sqrt{\text{Var}[X]}$$

how spread out, but in the same
unit.

Conditioning

to indicate that smth is fixed

$X|\mu$ X rand variable with fixed μ

$Y|X=2$ when X is fixed at 2

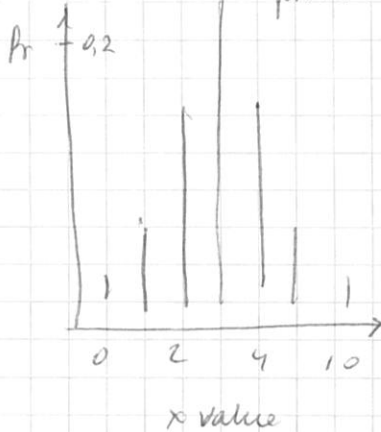
Distributions

- Binomial

$$X \sim \text{Bin}(10, 0.5)$$

number of coins

probability of Heads



dist. that describes a sum of coin flips

$$E[X] = n \cdot p$$

$$\text{Var}[X] = n \cdot p \cdot (1-p)$$

if you flip 10 coins and count when they come up Heads

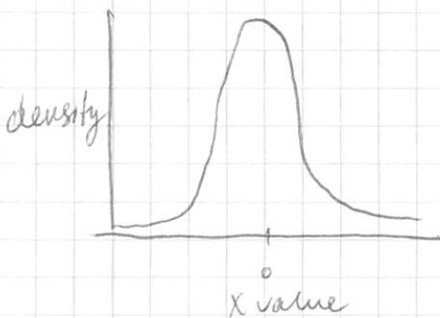
- Normal

$$N(\mu, \sigma)$$

mean

$$X \sim N(0, 1)$$

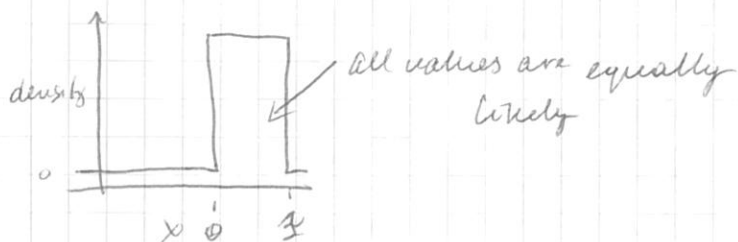
std



density - explains smth about probability for ranges of values of X .

- Uniform $U(a, b)$

$$X \sim U(0, 1)$$



Condensing

Law of total variance

Law of total expectation

\Rightarrow W.k.e

Representing data in R

my Data Frame [firstName == "jiff,"]

consider reading on style guides in R

Simulation Basics

Distributions

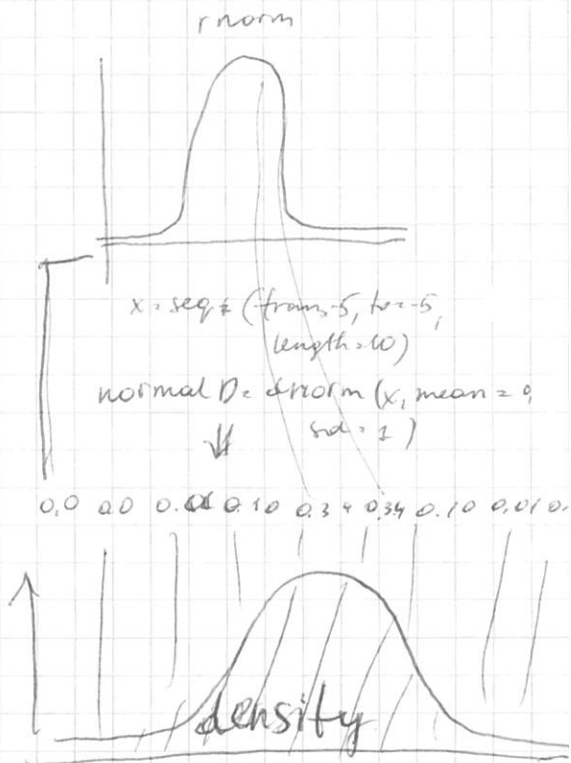
rbeta
rbinom
rgamma
rnorm
runif

Densities

dbeta
dbinom
dgamma
dnorm
dunif

Sampling

sample =



sample draws a random sample

args(sample)

x size ← number of samples
replace = F, prob = NULL

h = rnorm(10, mean = 188, sd = 3)

s = sample(heights, size = 10, replace = T)

↓
the same value
is used several
times

sample with probability

s = c(0.4, 0.3, 0.2, 0.1, 0... 0)

↓ ↓ ↓ ↓ ↓ ↓
1 2 3 4 5 ... 10 et of vec X

sample(heights, size = 10, replace = T, prob = s)

set.seed(12345)

↑
for reproducing

Types of Data Analysis

• Descriptive

goal: to describe a set of data

commonly applied to census data
(for generalizing?) (not for explaining!)

• Explanatory

goal: find relationships

ideas for following studies

• Inferential analysis

goal: use relatively small sample of data
to say smth about a bigger population

most common goal of statistics models
depends on sample you got.

• Predictive

goal: use data on some object to predict
values for another object

census - перепись, chop perepis

• Casual

goal: to find out what happens to one variable when you make another variable change

explores an average effect
randomized, to see if X causes Y

• Mechanistic

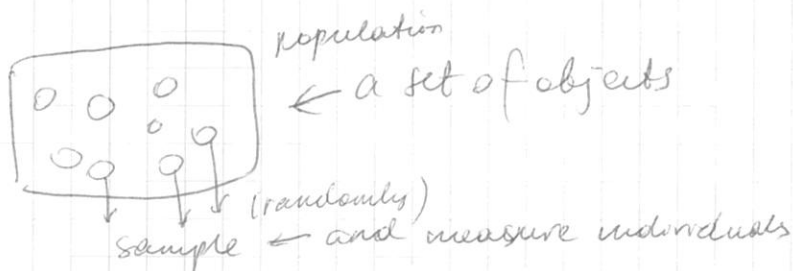
goal: understand the exact changes in variables that lead to changes in other variables for individual objects

(in one object we change $X \rightarrow$
 Y in the same object changes
in predictable way)

for finding formula, etc
by empirical analysis

Sources of Data Sets

census
observational study
randomized trials
etc



Sample (T:8, size=4, replace=F)

censal - all
observation study (sample with or replacements)
convenience (sample with probabilities) \neq
randomized no replace
→ get two groups and apply different

prediction: get a subset, do analysis,
get ~~rest~~ a subset of remaining
and check if it works

Study over time cross-sectional
sub set only at one day
longitudinal
the same individuals over time

retrospective
sampling at the end
(relationship between outcome and
exposure)