

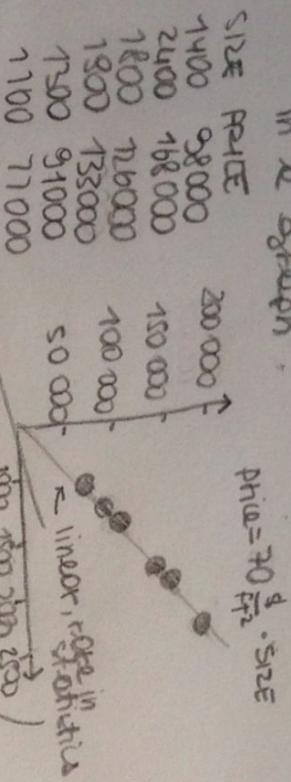
## Intro to Statistics

DATA → STATISTICS → DECISIONS

### Scatter Plot

- Each data item is a dot

in a graph.



A linear data set  
is really easy to  
predict

Even without a fixed dollar price per square foot, the  
relationship can be linear.

$$\text{SIZE PRICE}$$

1700	52000
2100	65000
1900	59000
1200	42000
1600	50000
2200	68000

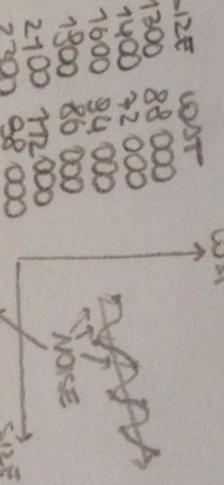
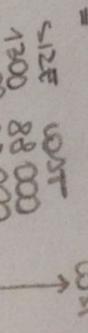
$$\text{Price} = 30 \frac{g}{\text{ft}^2} \cdot \text{size} + 20000$$

Some prices may be outliers (more water...)  
and not fit in the relationship

Scatter plots don't detect when there's what's  
called "noise". That is, the data deviation from  
the expectation, in some random, noisy way

Sort of dilutes the line of noise in data by  
pooling data points into a single summation bar

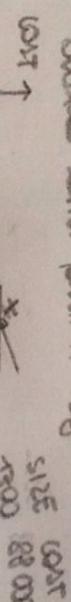
## Bar Charts



There might be factors that  
would affect the ~~fixed~~ <sup>average</sup> cost.

If those factors aren't included,  
to a statistician that's  
called "random noise".

In a bar chart, we take out raw  
data and pull it together.



It shows the count for the selection  
that fall into each category.

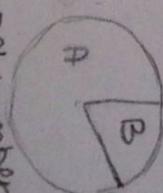
## Bar Charts

In statistics you use pie charts to  
visualize relative data!

PARTY A: 74000 votes

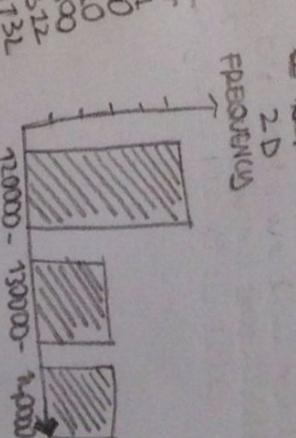
PARTY B: 181000 votes

$$\frac{181000}{255000} = 0.2 \quad \frac{74000}{255000} = 0.08$$



It is important to the total number  
of voters but it shows that A  
got many more voters than B.

Histograms - special case of  
bar chart



... so a bar chart helps you to  
pool together groups of data into  
a single bar, and understand  
global trends

Used with a lot of data points.  
A scatter plot with a lot of  
data points looks tells you very  
little

## Dot Charts

A dot chart is a much finer  
representation of the data.  
There's a much better way to  
really understand the dependence  
of cost to data.

While the bar doesn't give you the  
linear relationship, it really  
gives you a general idea of how  
and in which cases, the cost  
increases

From plotting import \*

data = [...]

histplot(data)

southplot(x, y)

Admissions Cox Study

- From UC Berkeley to New  
whether admissions were  
gender biased

Simpson's Paradox

Data is made up to illustrate  
the effect!

MALE APPLIED ADMITTED RATE ACCEPTED RATE

MAJOR A 4000 50% 700 80%

MAJOR B 100 70% 900 180

From that looks like female student

are being favored

BOTH 1000 460 46% 100 260 26%

When you look at both majors  
together males have a higher  
admissions rate than females

statistic is ambiguous. In choosing  
how to graph your data you can  
majorly impact what people believe  
to be the case

Statistics is deep and often manipulated

=

Probability

... is just the opposite of causation

DATA  $\xrightarrow{\text{probabilities}}$  CAUSES

$\leftarrow$  probability

Probability gives us a language to  
describe the relationship between data  
and underlying causes

## Flipping coins

$$P(\text{HEADS}) = 0.5$$

$$P(\text{TAILS}) = 0.5$$

MV

P(HEADS) + P(TAILS) = 1

$P(H,H) = 0.25$

TRUTH TABLE  $\leftarrow$  Down but only  
possible outcome

FUP-1 FUP-2

0.25

0.25

0.25

0.25

H T H T

T H T H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

H H H H

T T T T

## Conditional Probability

We're going to study the most interesting way  
where the outcome of the first does impact  
the outcome of the second

$$P(\text{Cancer}) = 0.1$$

$$P(\text{Cancer}) = 0.9$$

Blood test for cancer

$$P(\text{POSITIVE}|\text{CANCER}) = 0.1 \leftarrow \text{Test notes}$$

$$P(\text{NEGATIVE}|\text{CANCER}) = 0.9 \leftarrow \text{Minimize}$$

misses

Then the outcome of the test depends on

whether the patient has cancer or not

That is called a conditional probability

I know what the probability of the

test on the left given that we

knew the result on the right

is actually the one

independent

independent

disjoint depend

on the first outcome

is actually the one

positive

negative

cancer

test

P(Cancer)

P(Cancer)

P(Positive|Cancer)

P(Negative|Cancer)

P(Cancer)

P(Cancer)

P(Positive|Cancer)

P(Negative|Cancer)

P(Cancer)

P(Cancer)

P(Positive|Cancer)

P(Negative|Cancer)

P(Cancer)

P(Cancer)

$$P(H,H) = P(H) * P(H)$$

$$P(H,T) = P(H) * P(T)$$

$$P(T,H) = P(T) * P(H)$$

$$P(T,T) = P(T) * P(T)$$

(CANCER) \* (POSITIVE|CANCER)

(CANCER) \* (NEGATIVE|CANCER)



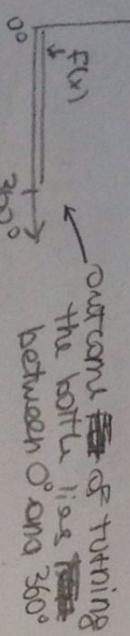
Continuous probability distributions

In continuous distribution, every outcome has probability 0!

You're dropping a package out of an airplane and it takes between 3 and 3.2 min to reach the ground. Between the prob. of density is uniform

$$P(A < x \leq b) = \frac{b-a}{360}$$

Probability that a spinning bottle stops at a given angle.



We want to assign 0 to anything outside that range

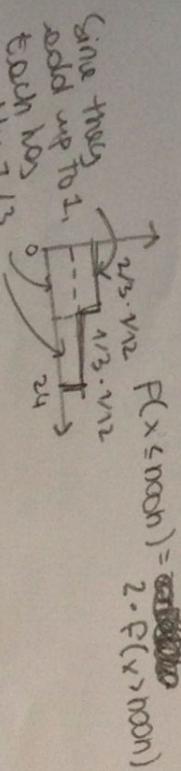
There's a function that makes every function in the range equally likely.

- Each outcome has equal value
- The area under the function/integrates to 1

$$f(x) = \frac{1}{360}$$

=

Time of day when people are born  
Twice as likely to be born before noon



(since they add up to 1)

exactly 1/13

Correlation and causation!

SICK	IN HOSPITAL	DIED
AT HOME		
8000	40	4
	40	40

Chances of dying in hospital are no longer than at home

↑ Correlation

It doesn't mean by knowing you're in a hospital increases your chance of dying

Being in a hospital increases your probability of dying by a factor of 40 ← causation

The correlation does not imply causation

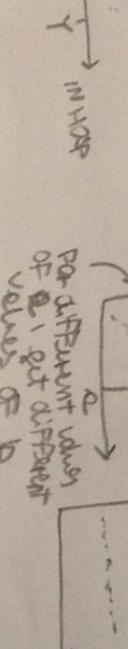
SICK	AT HOME	DIED
HEALTHY	0	4
20	40	20

Since you're sick you're more likely to die

that means the data is bivariate

correlation

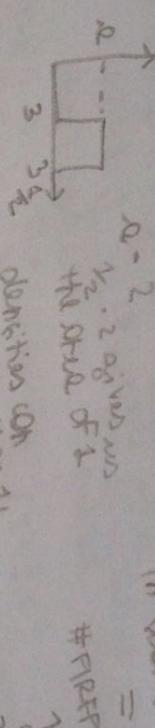
not causation!



The example admitted the sickness variable that in part also cause you to die, and who attended hospital or not.

Once you knew you were sick being in a hospital negatively correlated with you dying hospital or not.

In statistics we call it a confounding variable!

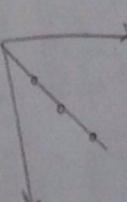


The # of Firefighters is correlated with size of fire

large correlation, Firefighters cause fire

size ↑ → #firefighters  
But the graph is divided into two directions

correlation: size causes #firefighters  
we don't know to draw conclusion



ESTIMATORS:

Flips  $\rightarrow$  1 head 0 tail

100101  $P(\text{HEADS}) = 0.5$

0111  $P(\text{HEADS}) = 0.8 \rightarrow$

Empirical Frequency

DATA  $x_1, x_2, \dots, x_n \rightarrow \frac{?}{?} x_i$  shows

$x_1, x_2, \dots, x_n \rightarrow \frac{?}{?} x_i$  between 0 and 1

Maximum Likelihood Estimator

using MLE

1663265462  $\downarrow$

$P(x_1) = 0.1 P(x_2) = 0.2 P(x_3) = 0.1$

$P(x_4) = 0.1 P(x_5) = 0.2 P(x_6) = 0.4$

Maximum Likelihood Estimator (MLE)

Estimation Problem: DATA  $\sim P$

$\leftarrow P \sim P(\text{DATA})$

multiple independent events!

101  $P = \frac{1}{2}$   $P(\text{DATA}) = 0.125$   $(0.5 \times 0.5 \times 0.5)$

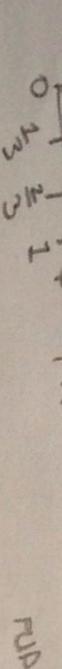
$P = \frac{1}{3}$   $P(\text{DATA}) = 0.074$   $(0.33 \times 0.66 \times 0.33)$

$P = 4$   $P(\text{DATA}) = 0$

$P(\text{DATA}) \uparrow$   $\rightarrow$  points at 2/3

that point is the one that maximizes the likelihood of the data.

Therefore it's called MLE!



Mode

DATA  $\sim P$

MLE always has the same value when we have raw data.

Add FOR DATA! DATA  $\sim P$

0.667  $\leftarrow P(\text{HEADS}) + P(\text{TAILS})$

0.4  $\uparrow$  everything towards 0.5

0.5  $\downarrow$  The last  $\rightarrow$  has now data 0.5

0.75  $\leftarrow$  therefore is further away from 0.5

... with plenty of data, the Laplace estimator gives about the same result as the maximum estimator, but when data is scarce, it works much better than the maximum one.

the Laplace est.

Roll a die, you get... Laplace estimator

1232 123456  $\rightarrow$  one data point for each outcome

MLE LAPLACE

0.25 0.2

0.5 0.3

$\frac{1}{N} \sum (1 - \bar{x}_i)$

$\frac{1}{N} K$   $\sum$  datapoints

House prices

mean median mode

Mode the day most frequently represented at the party

Mode

They are

mean median mode

median mode

mean median mode

Mode is useful if data distribution of data is multimodal

Mode

they are

mean median mode

median mode

mean median mode

Mode the day most frequently

represented at the party

Mode

they are

mean median mode

median mode

mean median mode

Mode the day most frequently

represented at the party

Mode

they are

mean median mode

median mode

mean median mode

## Variance

$\frac{1}{N} \sum (x_i - \mu)^2$  this is  
the mean

Variance is a measure of how far the data is spread

standard deviation!

The variance is the quadratic deviation from the mean

If you don't want something quadratic  
 $\text{STANDARD DEVIATION} = \sqrt{\text{VARIANCE}}$

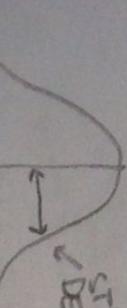
$$\text{STD DEVIATION} = \text{DEVIANCE}$$

$$\text{FAMILY ... MEAN} = 18$$

$$\text{STD DEVIATION} = 12.18$$

$$\text{FRIENDS ... MEAN} = 12$$

$$\text{STD DEVIATION} = 0.8$$



While here you expect the deviation to be just below 1 year

$$\text{STD DEVIATION} = 12.18$$

How much you would expect the age of an individual family member to deviate from the mean

$$\text{MEAN } \mu = \frac{1}{N} \sum x_i \quad \text{STD DEVIATION } \sigma = \sqrt{\frac{1}{N} \sum (x_i - \mu)^2}$$

$\sigma^2$  without  $\mu = \frac{1}{N} \sum x_i^2 - \frac{1}{N^2} (\sum x_i)^2$

Using ratio

Fixed amount 1000

$N' = N + 1000$  ← mean shift  
 $\sigma' = 12.0$  ← std deviation  
(for variance would be 1.2²)

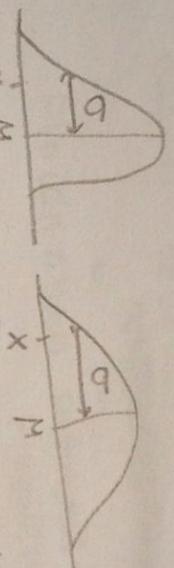
and otherwise don't change

## Standard score

For any deviation you can score how far in or out a point  $x$  is

$\frac{x - \mu}{\sigma}$   
# quartiles  
- median = 2nd quartile  
- upper quartile = 3rd quartile

Everything else you just put somethin in between



relative total mean and variance these points correspond That's called the standard score  
 $Z = \frac{x - \mu}{\sigma}$

$$Z = \frac{x - \mu}{\sigma}$$

The ~~other~~ # of elements is given by  $4^n + 3$   
↑ lowest quartile  
# quartiles = median = upper quartile

Ages	19	20	21	22	23	24	25
frequency	2	1	7	3	2	1	1

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

With these numbers we can calculate

how to calculate the median

but there is no formula

for the quartiles

## Binomial Distribution

2 COIN FLIPS

H H # HEADS = # TAILS

H T ] ↗ 5 COIN FLIPS  
# HEAD = # TAILS

T H ↗ 0

6 COIN FLIPS

HHTT ↗ # HEADS =  
HHTH ↗ # TAILS

HTTH ↗

HTHH ↗

THHT ↗

THHH ↗

TTHT ↗

TTHH ↗

TTTH ↗

TTTT ↗

10 coins 4 heads  
 $\frac{10 \cdot 9 \cdot 8 \cdot 7}{4 \cdot 3 \cdot 2 \cdot 1} = 210$

10 coins 5 heads  
 $\rightarrow \frac{10 \cdot 9 \cdot 8 \cdot 7 \cdot 6}{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = 252$

$k! \rightarrow \leq 4 \cdot 3 \cdot 2 \cdot 1$

$n! = n(n-1)(n-2)(n-3) \cdots 1$

$\frac{n!}{k!(n-k)!} \leftarrow \# \text{ ways}$

= 

FUP COINS TIMES P(HHEAD)=0.5

OCCURS →  $\frac{5}{25} = 0.2$  ORANGE

WHICH MEANS  $\frac{1}{2}$  HEADS

$\frac{5!}{2! \cdot 3!} = 10$  OUTCOMES  
 $P(\# \text{ HEADS}=3) \approx 0.3125$

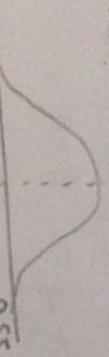
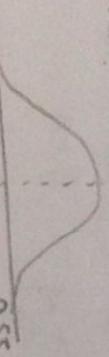


FUP COIN 2 TIMES P(HHEAD=0.8)  
P(H HEAD, 1-TAIL)  
= 0.096

FUP COIN 1000 TIMES → # HEADS ≈ 0.5  
REPEAT EXPERIMENT (↑)  
1000 TIMES

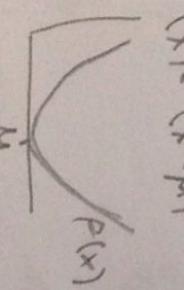
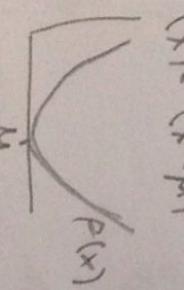
Histograms of the means, NOT

Flip coin n times  $P(\# \text{ HEADS}=k)$   
 $P(\text{HEADS})=P$

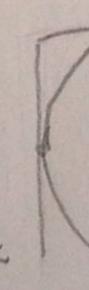
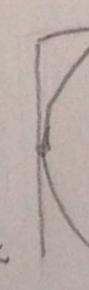
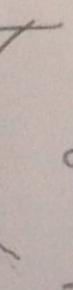
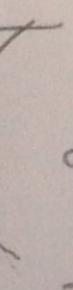


Normal distribution  
For any outcome x  
 $-F(x) = (x - \mu)^2$

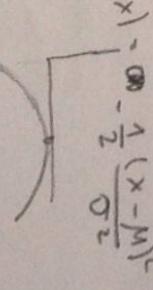
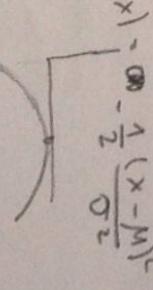
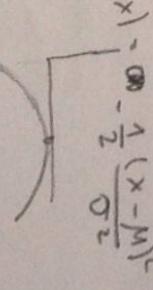
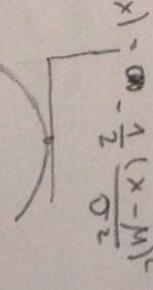
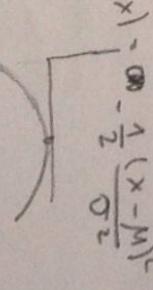
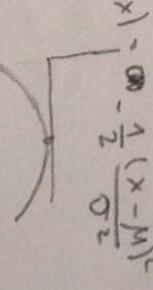
$-F(x) = \frac{(x - \mu)^2}{\sigma^2}$  notes it  
is wider if bigger  
it's smaller then closer



$-F(x) = e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$



Normal distribution  
or function



0.00

$$F(x) = \frac{1}{2} \left[ -\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right]$$

this describes  
the limit of making  
infinite rolls  
can flip



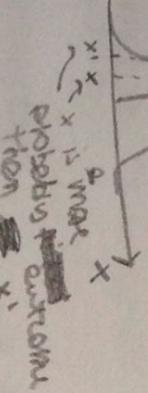
No matter what you  
do, when  $n$  is really  
large, you get a  
normal distribution

that also however  
doesn't add up to 1.  
Bell curve like this  
but to  $\sqrt{2\pi\sigma^2}$

We want all odds to add up to 1 so we  
wanted a coin flip and its complement  
to add up to 1. We normalize it  
so much

Normal distribution  $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$

$$P(x)$$



- For a single outcome use  $P$
- For a few outcomes use  $\frac{N!}{n_1! n_2! \dots n_k!} p^{n_1} (1-p)^{n_2} \dots$

Binomial distribution

- For infinitely many outcomes use  $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$

You're a golf player and you hit the  
ball really hard, and it lands  
with some uncertainty (that's random)  
where the distance is  $\mu = 100 \text{ m}$   
and your variance is  $\sigma^2 = 30 \text{ m}^2$ .  
And then you do the something again  
and then you do the something again

$\mu = 200 \text{ m}$  (combined stroke, expected distance)  
 $\sigma^2 = 30 \text{ m}^2$  (combined variance)

So, we express things in terms of std deviation

$$\sigma = 14.14 \text{ m}$$

## Manipulating normal distributions

different varieties in probability are  
normal distributed

$$\text{given mean of } 10000 \quad \mu = 10000 \quad \sigma = 10000$$

$$\mu' = 10000$$

$$\exp \left\{ -\frac{1}{2} \frac{(x-60000)^2}{10000^2} \right\}$$

$$\exp \left\{ -\frac{1}{2} \frac{(x_1-10000-60000)}{10000^2} \right\} \quad \mu \text{ becomes } 10000$$

$$\text{Double story, } \mu = 10000 \quad \sigma = 10000 \quad \sigma \text{ is not appended}$$

$$\mu = 10000 \quad \sigma = 10000$$

$$\mu' = 140000 \quad \sigma = 20000$$

$$\exp \left\{ -\frac{1}{2} \frac{(x-70000)^2}{10000^2} \right\} = \exp \left\{ -\frac{1}{2} \frac{\left(\frac{1}{2}x-70000\right)^2}{10000^2} \right\} *$$

$$= \exp \left\{ -\frac{1}{2} \frac{\left(\frac{1}{2}x_1 - \frac{1}{2}140000\right)^2}{10000^2} \right\} = \exp \left\{ -\frac{1}{2} \frac{\left(x_1 - 140000\right)^2}{10000^2} \right\}$$

$$= \exp \left\{ -\frac{1}{2} \frac{(x_1 - 140000)^2}{20000^2} \right\}$$

$$A = N(\mu, \sigma^2) \quad B = N(\mu, \sigma^2)$$

$$A+B$$

$$\mu' = \mu + b$$

$$\sigma'^2 = a^2 \sigma^2$$

$$A-B \quad \mu' = 0 \quad \sigma'^2 = 25^2$$

# Most better than average

Most drivers believe they drive better than the average driver

It's ~~possible~~ possible that even 99% of all people are smarter than the average person

4555555555555555  
avg = 4.95 95% exceed the average

Abortion's might and proof

Most are good women, but some extreme women skew up the actual

Scatterplot of today's weight vs last year's weight we get an approximation linear fit with some outliers

## Confidence intervals

... we can divine on what basis on how many items to vote  
Depend on margin of error  
60% Party A  $\pm 2\%$   
40% Party B  $\pm 3\%$   
margin of error

people  
you want  
to know  
the outcome  
can vote  
not  
right

Election day

SAMPLE  
selection of randomly  
drawn people from that  
pool that are representative  
of the pool at large

TRUE VALUE OF } P  
PERCEIVING A VOTE }  
IN A STATISTICS we  
can't know that

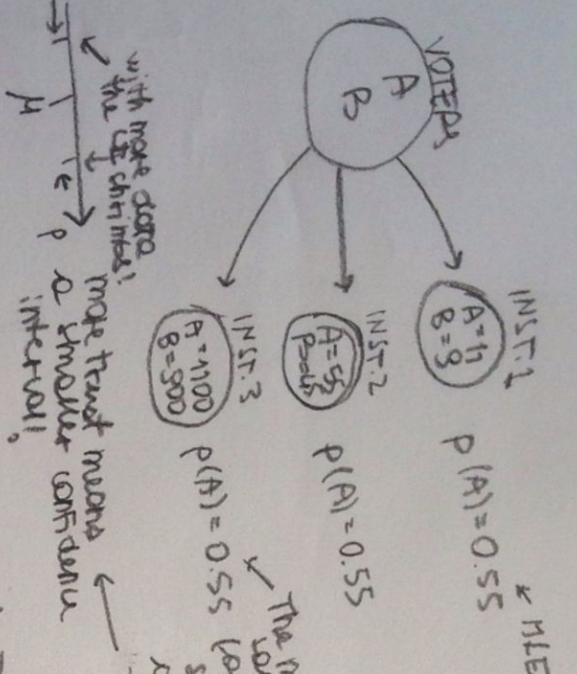
$\mu$

$\frac{1}{n} \sum x_i$

We do a sample  
 $x_1, \dots, x_n$   
95%

$\mu$

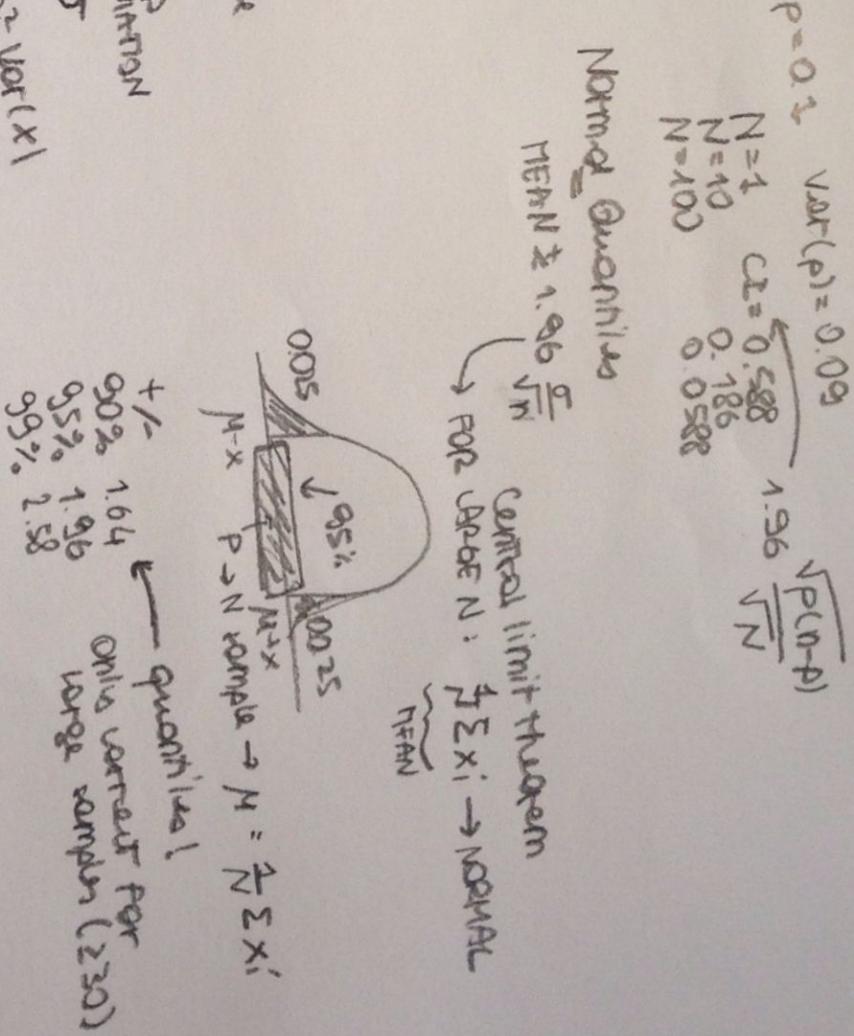
$C_i >$



$\sigma = 0.5$	$\mu = 0.5$	$\sigma^2 = 0.25$	$\text{MEAN}(\sum X_i)$	$\sqrt{\text{Var}(\frac{1}{N} \sum X_i)}$	$\text{STP. REV}$	$C_I$	$\text{STP. REV}$
$N=1$	0.5	0.25	0.5	0.25	0.5	0.98	0.196
$N=2$	1	0.25	0.5	0.125	0.35	0.69	0.147
$N=5$	2	0.25	0.5	0.05	0.26	0.31	0.078
$N=10$	5	0.25	0.5	0.025	0.16	0.21	0.038

the speed of the moon goes down as N increases: 0.05

$$\begin{aligned} P &= 0 \\ \sigma^1 &= 0 \\ \sigma^2 &= 0 \end{aligned}$$



# Hypothesis testing

SAMPLE → INFORMATION → STATISTICS → DECISION

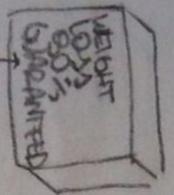
=

YES / NO

LOST WEIGHT?

YES n = 73.3%

NO 4 binomial!



Null hypothesis  
(or H<sub>0</sub>)

P = 0.9

Alternate hypothesis  
(or H<sub>1</sub>)

P < 0.9

Prediction doesn't  
work as was on  
conventional

Unless proven wrong we believe null hypothesis  
is actually correct. Or do we have

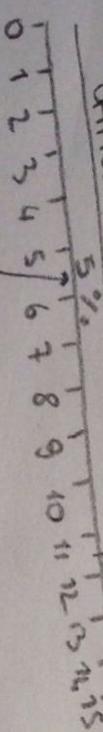
sufficient evidence to be very likely  
the null hypothesis is wrong? The alternative

hypothesis is correct. Then we reject  
the null hypothesis.

Assume H<sub>0</sub> is correct

YES n = H<sub>0</sub>: P = 0.9  
NO 9 H<sub>1</sub>: P < 0.9

Critical Region



confidence  
level

$$= \text{Corresponding probability for every outcome } \frac{N!}{k!(N-k)!} p^k (1-p)^{(N-k)}$$

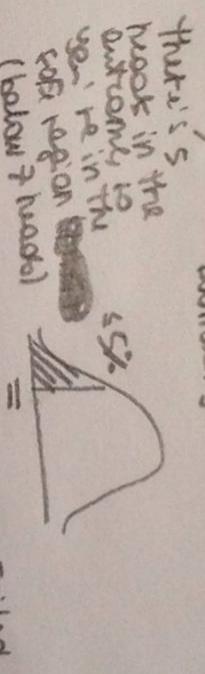
	YES n	NO 4
0	0	12.012
1	1	13.027
2	2	14.034
3	3	15.020
4	4	16.004

Outcomes up to 10 (15%) is in  
the critical region

NUL HYPOTHESIS: H<sub>0</sub>: P = 0.3  
ALT. HYPOTHESIS: H<sub>1</sub>: P > 0.3  
THHTTTHTTTHH → 0.45

02.03.12.25 12.13.06.02 0 0 0  
0 1 2 3 4 5 6, 7 8 9 10 11

≤ 2.5%  
boundaries



two-tailed  
test!

H<sub>0</sub>: P = 0.5 < 2.5%  
H<sub>1</sub>: P ≠ 0.5

N = 14 # HEADS = 3

0 0 0.005 0.02 0.06 0.12 0.18 0.24 0.18 0.12 0.06 0.022 0.05 0 0  
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14  
≤ 2.5%  
critical  
region

## Hypothesis test - part 2

HYPOTHESIS TEST  
PART 2

previous

that part  
here

WIKIPEDIA CLAIM AT  
NBA player is tall

tall measurements yourself  
199 200 204 202 203 204 205 206

$$\text{confidence interval } \left\{ \frac{1}{N} \pm \alpha \sqrt{\frac{\sigma^2}{N}} \right\} \rightarrow 202.5 \pm 1.916$$

mean

100.58

$\alpha = 2.635$

we should reject  
but depending  
on which town we  
pick different  
players, and we  
didn't pick the  
town independently

Dance club: average age 26

$$\begin{array}{ll} 4: 21 & N=30 \\ 6: 24 & \mu = 26.97 \\ 7: 26 & \sigma = 18.57 \\ 11: 29 & \alpha = 1.96 \\ 2: 40 & \text{FOR 95%} \end{array}$$

Given data  $x_1, \dots, x_n$

But taking the sample  
of one night is not  
nearly independent.  
You should go at  
random nights and  
take a person each night

$$\begin{aligned} \mu &= \frac{1}{N} \sum x_i \\ \sigma^2 &= \frac{1}{N} \sum (x_i - \mu)^2 \end{aligned}$$

$$\text{get } t(N-1, P) = \alpha$$

$$\text{size CI} = \alpha \sqrt{\frac{\sigma^2}{N}}$$

## Regression

noise or random error!

$$\begin{aligned} X &= 2 \quad 5 \quad 3 \\ Y &= 6 \quad 7 \quad 2 \quad 3 \\ b &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ &\rightarrow \frac{16 + 0 + 1 + 9}{16 + 0 + 1 + 9} \rightarrow b \text{ is } 1! \end{aligned}$$

$$a = \bar{y} - b\bar{x}$$

$$\rightarrow 3 - (1 \cdot 2) = 1$$

$$\text{noise or random error!}$$

The line is commonly described by a functional relationship between  $x$  and  $y$  in the form  $y = bx + a$ .

Given data, get  $(a, b)$  that best fit data

"we want to find a line  
that our data is the  
result of this linear  
relationship"

"we want to minimize the  
difference between  
the data and the  
line in the  $y$  direction!"

"we want to minimize  
the sum of the distances  
from the data points  
to the line"

"we want to minimize  
the sum of the distances  
from the data points  
to the line"

"we want to minimize  
the sum of the distances  
from the data points  
to the line"

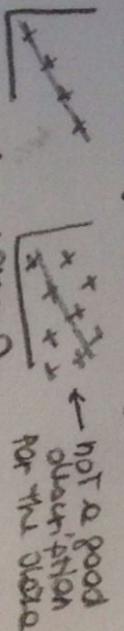
"we want to minimize  
the sum of the distances  
from the data points  
to the line"

$$b = \frac{\sum [(x_i - \bar{x})(y_i - \bar{y})]}{\sum (x_i - \bar{x})^2}$$

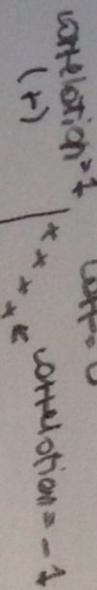
$$\bar{x} = \frac{\sum x_i}{N}$$

## Correlation!

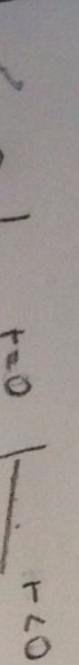
... is a measure that lies in the range  $-1 \leq r \leq 1$  that tells us how far the data is described by a line



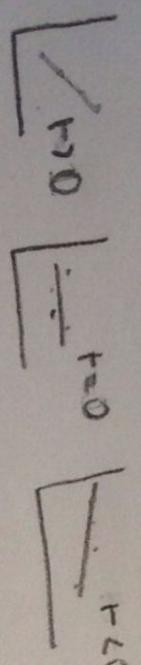
correlation  $> 0$



correlation  $< 0$



correlation  $= 0$



correlation  $\neq$  not a good description for the data

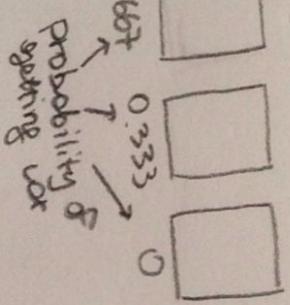
Given the linear regression  $y = \hat{b}_0 + \hat{b}_1 x$   
+ will be positive  
because of that

- between 1 and -1
- tells how how related two variables are

- 1 and  $-1$  stand for perfect linear data

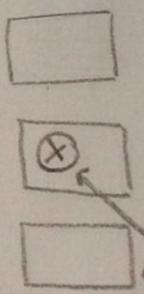
$r = \frac{\sum [(x_i - \bar{x})(y_i - \bar{y})]}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}}$

	true location	Monty shows...	prob.
1	1	1	$\frac{1}{3}$
2	2	1	$\frac{1}{3}$
3	3	2	$\frac{1}{3}$
		3	$\frac{1}{3}$



Probability of getting car

case study - regression  
prob. won't show the one I wanted  
 $\bar{x} = 81.2 \quad \bar{y} = 79.7 \quad \text{std dev}(x) = 10.6 \quad \text{std dev}(y) = 9.34$   
 $r = 0.76$   
 $b = 0.67$   $\leftarrow$  slope  
 $a = 25.2$   $\leftarrow$  intercept  
 $y = 25.2 + 0.67x$

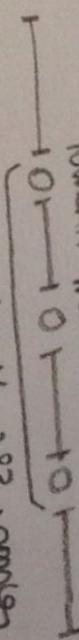


Let's make a deal!  
Answer

## Monty Hall Problem

### Weight loss study

163 complete outliers using quantiles  
remove lower.. median upper quantile



picked!  $\approx 83$  samples

$$n = 83 \quad \sum x_i = 6618.47 \quad \sum x_i^2 = 528679$$

$$\mu = 79.74$$

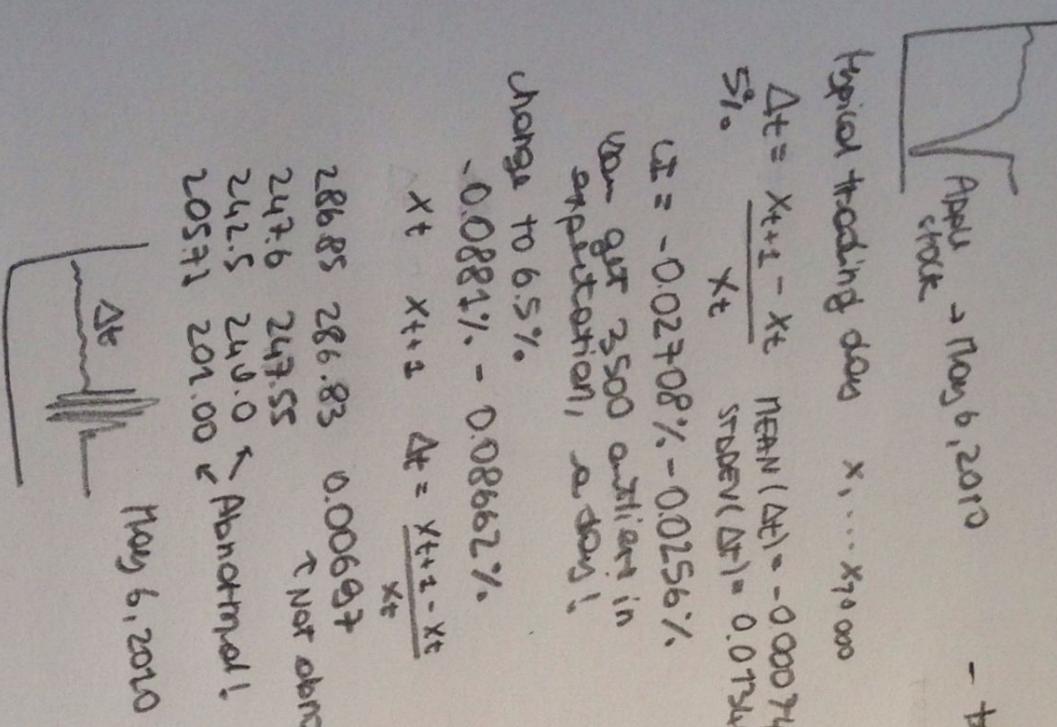
$$\pm = 0.92$$

$$Q5\% =$$

$$\frac{V\mu}{\sqrt{n}} = 1.96$$

Monty picks door 3, you  
be in either three cases  
 $P(3) = 2/3 \quad P(2) = 1/3 \quad P(1) = 0$

## Flash crash example



## Challenger example

