

Lecture 1

Statistics 103 fall 2017.

Instructor:

Ethan Anderes

TAs:

Michael Bissell

Maxime Pouokam

What this course is about

- 1) a study of random variables and the dependence btwn two random variables.
(probability theory)
- 2) inferring truth based on incomplet/random data.
(statistics)

We interact with both every day.

Although some problems associated with 1) or 2) can be solved intuitively, many can be non intuitive and require formal tools to solve.

①

e.g. if I bet \$1 on the flip of a coin, I "expect" to win about 5 times in 10 plays.

②

- Easy

e.g. I go to the doctor and get a strep test... its positive.

doc says:

If you have strep the test works 98% of the time.
If you don't have strep the test works 92% of the time.

This is good to know but what I actually want to know is this:
what is the chance that I have strep?

- For this one intuition fails on we need to learn formal or mathematical tools to help us solve it.

Syllabus

(3)

- Summation notation $\sum_{i=1}^n$
- Random variables
 - Characterizing them with PMF or Box models
 - $E(X)$, $\text{var}(X)$, $\text{sd}(X)$
- Dependence btwn random variables
 - Characterizing them with the joint PMF or Box models
 - Conditional PMFs
 - $E(X|Y=y)$, $P(X \leq z | X+Y=3)$
 $\text{var}(X|Y=y)$
 - independence
 - covariance & linear properties of $E(\cdot)$ & $\text{var}(\cdot)$
- Continuous random variables
 - Characterizing them with PDFs
 - Gaussian random variables
 - Z-scores
 - CLT
- Correlation and Bivariate regression

- MSE, var & bias of an estimator.

(4)

- Estimation & Testing

- Comparing two populations
- regression
- confidence intervals.

-
- If you want to see more details I have old course notes up online:

github.com/EthanAnderes/STA-103-lecture-notes

- No book required. I will post these notes online.

- No graded HWKs.

Each week I will post practice problems for the Discussion sections and for you to practice at home.

- Your grade will be based on 3 midterm exams & the final exam.

Midterm 1: Oct. 16 }
Midterm 2: Nov. 6 }
Midterm 3: Nov. 29 } the lowest
midterm score dropped.

Final: Dec 14, 10:30am - 12:30am

e.g. Hypothetical scores

(5)

$$\left. \begin{array}{l} \text{Midterm 1: } \frac{10}{12} \\ \text{Midterm 2: } \cancel{\frac{9}{12}} \\ \text{Midterm 3: } \frac{11}{12} \end{array} \right\} \text{Top 2 gives } 21 \text{ out of 24}$$

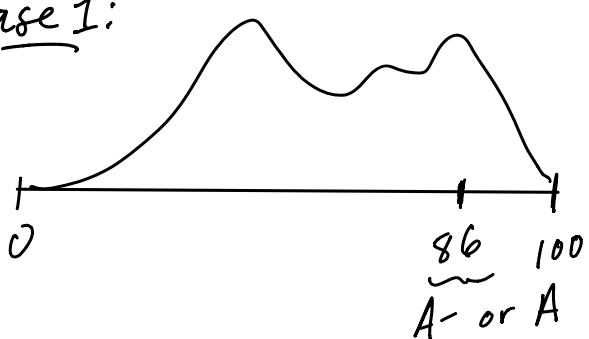
$$\text{Final: } \frac{20}{24}$$

Total grade out of 100:

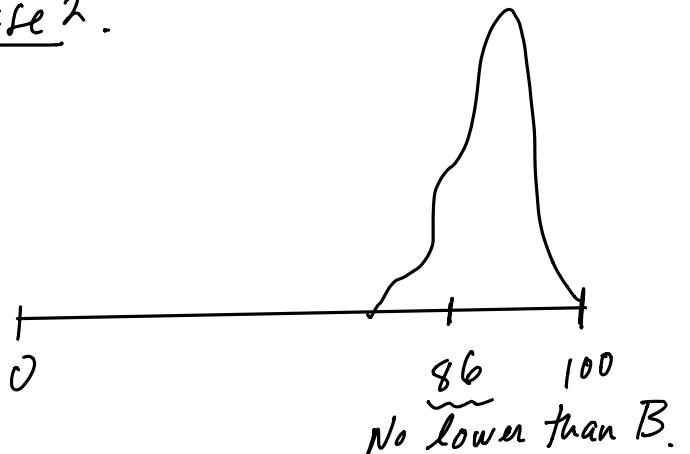
$$\underbrace{65\left(\frac{21}{24}\right)}_{\text{Midterms counts 65\%}} + \underbrace{35\left(\frac{20}{24}\right)}_{\text{Final counts 35\%}} = 86$$

To determine your letter grade
I curve based on the
distribution of "Total grade"

Case 1:



Case 2:



Note: You must take the final to pass this class!!!

(6)

Summation notation

(7)

Later in the quarter we will encounter things like this:

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

$$\hat{x} = \bar{Y} - \hat{\beta} \bar{X}$$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - (\hat{x} + \hat{\beta} X_i))^2$$

summation notation

Upper case to indicate they are random numbers

We will start the class by going over this notation.

X, Y, Z, etc will generally denote R.V.s

(8)

a, b, c, x, y ... will denote fixed non-random numbers.

X_1, X_2, \dots, X_n denote n R.V.s

x_1, x_2, \dots, x_n denotes n nm-random numbers.

$\sum_{i=1}^n x_i$ is shorthand notation for $x_1 + x_2 + \dots + x_n$

So $\sum_{i=1}^n x_i$ is just a mathy way to write $x_1 + \dots + x_n$.

In words

$\sum_{i=1}^n x_i$ = "the sum of the x_i 's as i ranges from 1 to n"

Note: $\sum_{k=1}^n x_k = \sum_{i=1}^n x_i$ This number tells you when to stop summing

This number tells you when to start.

example

Let $x_1 = 1$

$$x_2 = 5$$

$$x_3 = 7$$

$$x_4 = 2$$

Then

$$\sum_{i=1}^4 x_i = 1 + 5 + 7 + 2$$

$$\sum_{k=1}^3 x_k = x_1 + x_2 + x_3 = 1 + 5 + 7$$

$$\sum_{k=2}^3 x_k = x_2 + x_3 = 5 + 7$$

example

with the same x_1, x_2, x_3, x_4 in previous example.

$$\pi \sum_{i=1}^3 \sin(x_i) = \pi (\sin(1) + \sin(5) + \sin(7)) \\ = \pi (\sin(1) + \sin(5) + \sin(7))$$

example

$$\frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} (x_1 + \dots + x_n)$$

$= \bar{x}$ ← the regular average of x_1, x_2, \dots, x_n .

(9)

Linear properties of Σ

(10)

We all know:

$$5(7+9) = 5 \cdot 7 + 5 \cdot 9$$

$$& (7+2) + (9+12) = 7+2+9+12$$

This is abstractly expressed with Σ as follows

$$a \left(\sum_{i=1}^n x_i \right) = \sum_{i=1}^n (ax_i)$$

$$\sum_{i=1}^n (x_i + y_i) = \left(\sum_{i=1}^n x_i \right) + \left(\sum_{i=1}^n y_i \right)$$

Putting these together

$$\sum_{i=1}^n (ax_i + by_i) = a \left(\sum_{i=1}^n x_i \right) + b \left(\sum_{i=1}^n y_i \right)$$

One application of these rules is when changing units.

example

Let x_1, x_2, \dots, x_{30} denote the daily low temp (in Farenheit) in November at the south pole.

Suppose I told you $\bar{x} = -10$

Do you have enough information to find the average temp in Celsius?

For each x_i define y_i as follows

$$y_i = \frac{5}{9} (x_i - 32)$$

↑ ↑
the temp on day i in C° the temp on day i in F°

We know $\bar{x} = -10$.

We want $\bar{y} = ?$

There are plenty of tricky problems with \sum_i

example

i.e. discussion.
In lab they will show how to derive:

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - (\bar{x})^2$$

$$\bar{y} = \frac{1}{30} \left(\sum_{i=1}^{30} y_i \right)$$

$$= \frac{1}{30} \left(\sum_{i=1}^{30} \underbrace{\frac{5}{9} (x_i - 32)}_{\text{}} \right)$$

$$= \frac{1}{30} \frac{5}{9} \left(\sum_{i=1}^{30} (x_i - 32) \right)$$

$$= \frac{1}{30} \frac{5}{9} \left(\left(\sum_{i=1}^{30} x_i \right) + \underbrace{\sum_{i=1}^{30} (-32)}_{\substack{\text{"sum } (-32) \\ \text{30 times}}} \right)$$

$$= \frac{1}{30} \frac{5}{9} \left(\left(\sum_{i=1}^{30} x_i \right) - 30 \cdot 32 \right)$$

$$= \frac{5}{9} \left(\underbrace{\left(\frac{1}{30} \sum_{i=1}^{30} x_i \right)}_{\text{}} - \frac{1}{30} \cdot 30 \cdot 32 \right)$$

$$= \bar{x} = -10$$

$$= \frac{5}{9} (-10 - 32)$$

Box Models

(13)

- Box models are essentially a way to analyze randomness without math.
- The hard part is setting up the Box model... once that is done the answer is usually easy.

e.g. I go to the doctor and get a strep test... it's positive.

doc says:

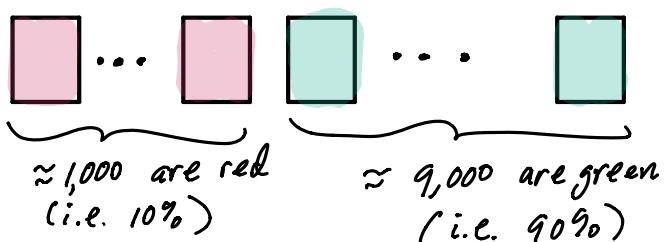
- If you have strep the test works **98%** of the time.
- If you don't have strep the test works **92%** of the time.
- At any given time about **10%** of the population has strep.

What is the chance that I have strep?

Consider a population of 10,000.

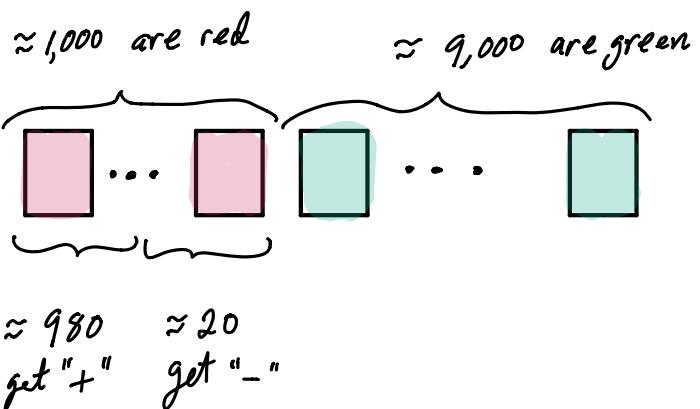
For each person make a ticket & color it red if they have strep (green otherwise).

10,000 tickets

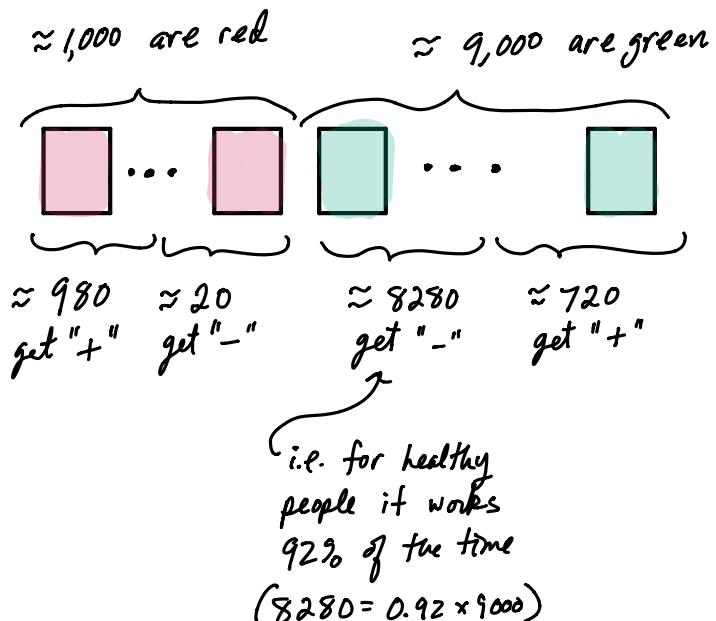


Now imagine giving each person the strep test. (14)

Those who have strep the test will work 98% of the time.

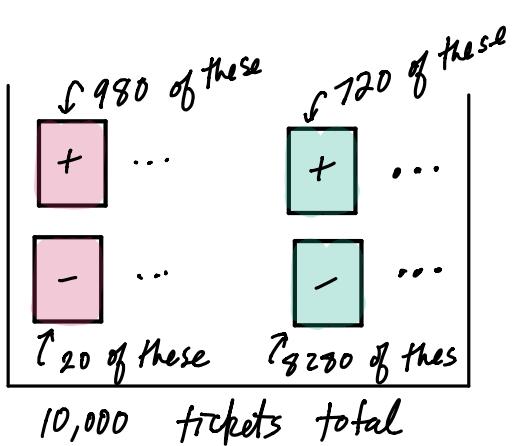


Those who don't have strep the test will work 92% of the time

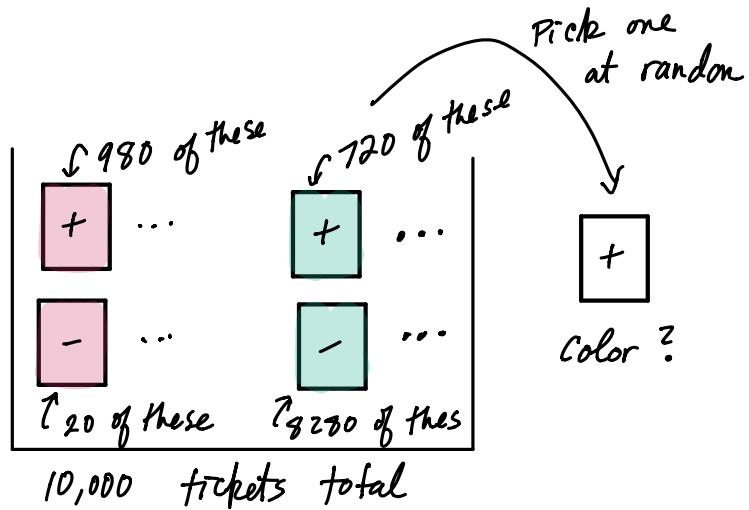


Write "+" or "-" on the ticket corresponding to the outcome of the strep test... then put them in a box.

(15)



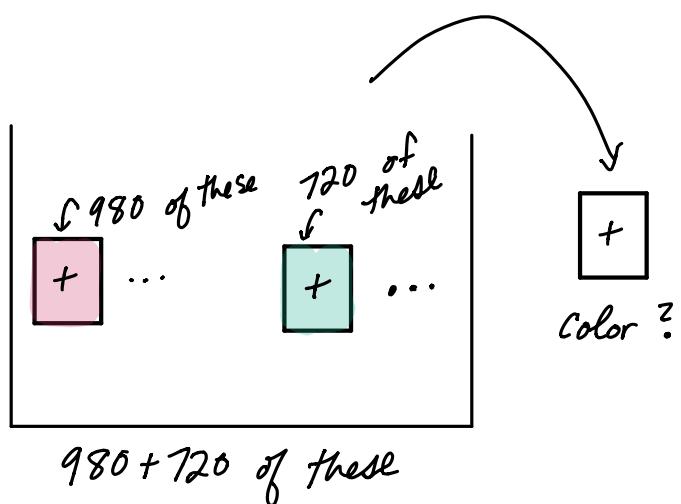
For the box model imagine someone reaches in and picks a ticket at random. They see a "+" or "-" but not the color.



This random represents your state of information: you got a "+" on the strep test but don't know if you have it or not.

(16)

Notice: since you got a "+", and each ticket was equally likely to pick so it behaves like a draw from



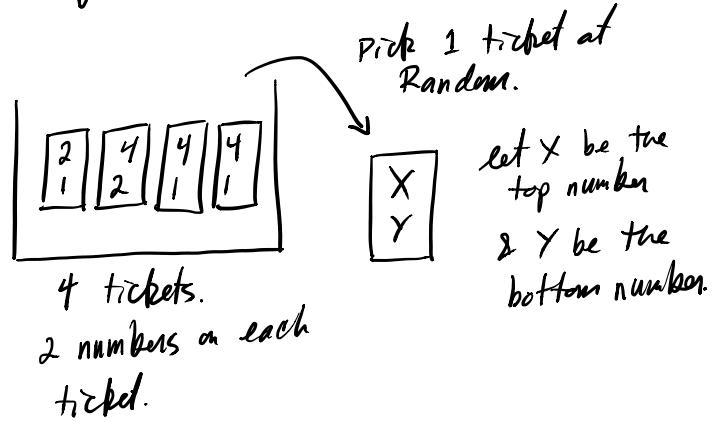
So the chance your ticket is red (i.e. you have strep) is

$$\frac{980}{980+720} \approx 0.576.$$

We will push this technique much further and use them to analyze dependence b/wn random variables.

e.g.

(17)



let X be the top number
& Y be the bottom number.

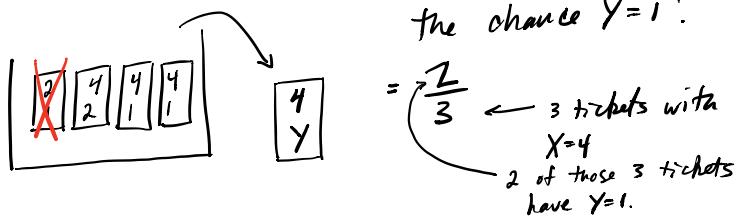
X & Y are R.V.s

Probability Calculations are very easy with Box Models.

$$P(X=4) = \text{"the probability that } X \text{ is 4"} \\ = \frac{3}{4} \quad \begin{matrix} \leftarrow 3 \text{ tickets with } X=4 \\ \leftarrow 4 \text{ tickets total.} \end{matrix}$$

$$P(X=4 \text{ and } Y=1) = \frac{2}{4}$$

$P(Y=1 \mid X=4) = \text{"if someone told you that } X=4 \text{ but not what } Y \text{ is what is the chance } Y=1".$



= $\frac{2}{3} \quad \begin{matrix} \leftarrow 3 \text{ tickets with } X=4 \\ \leftarrow 2 \text{ of those 3 tickets have } Y=1. \end{matrix}$