

(1)

SESSIONS

03/05/19

ANNOUNCEMENTS:

• sessions are on saturdays - how to do things in R (in class by hand)

• double session on last week (saturday)

- ml - python ; stats - R

• attendance sheet

• project available until spring quarter schedule (optional)

• free writing counseling = uchicago grad

• flexible program → need more time etc. → best

• grading is absolute

• a done before lecture - 0 = after

• director / founders

QUESTION

WHAT IS STATISTICS?

• communicating data → descriptive statistics summarizing data

• controlling outcomes → make good decisions in the presence of uncertainty

• uncertainty → calculated risk → how things depend on each other

something we don't know → variation = differences = identify → predict → control

• if no variation we don't need statistics? → not debatable -

• Event → put a # to this risk/probability (assign $p(n)$)

- mutually exclusive = when one occurs the other can't

- exhaustive = overall possible outcomes

• ex. SEO → conversion rates

- use theoretical prob dist to assign prob in real life

• communication → tell more/less

- ex. website w/ terabytes of data

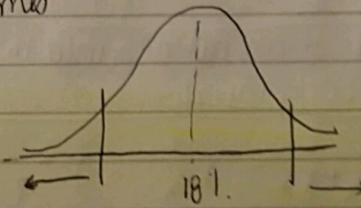
↳ metrics

• "manage var concurrently & proactively" → knowable caused eq.

↳ goes after systematic variation → statistic → identify variation in reqs

random variation → probability

↳ inherent to the process causing



ask
→

how is probability different from chance?

① VARIATION DUE TO SAMPLING

- random: samples can be representative of pop; no bias; esp. chance

↳ sampling variation

➢ stats: gives no # and confidence of accuracy

↳ reduces 1 source of var (var smaller)

- stratified: credit risk model (minority strata groups)

- probability: match variation w/ probable effect

↳ sampling variance more than population variance (?)

↳ confidence interval

$$Y = f(X) + \epsilon$$

TYPES OF STATISTICS:

- Descriptive statistics
 - Business intelligence
- Inferential statistics → focus of program
- Probability

◦ manage variance concurrently & proactively

Variance = systematic + random

of STATISTICS: SAMPLE → POPULATION

of PROBABILITY: POPULATION → SAMPLE

- assume a distribution then explain a probability
- or observe the sample then guess the distribution

VARIABLES:

◦ independent = measurements we use to explain y dependent var

◦ modeling frameworks

- supervised → have an ~~independent~~ dependent var

- unsupervised → no dependent var (y)

} Machine learning!

discrete of ◦ Nominal - labels only; no order; no degree → $p(x)$

◦ ordinal - w/ order → don't need to be evenly spaced

⊕ % - bound scale → not technically ordinal?

→ garbage in, garbage out

continuous of ◦ Interval - don't have 0 (?) - linearly related - (Age) more analytically friendly

◦ Ratio - real 0 points → Kelvin = Ratio

→ relating to the power of 2

→ probability of intervals

DISTRIBUTION OF VARIABLES:

◦ Normal distribution (continuous) ⊕ → Read about this later

- μ, σ

↳ Arithmetic mean

← σ = how far from the mean → "2 σ "

← σ^2 = variance → no "are" interpretation like σ

- symmetry

↳ kurtosis → fat tail

↳ skew

median = inconvenient
mathematically

→ doesn't rep pop if skewed

mean = biased to outliers

mode = more ordinal

CENTRAL TENDENCY MEASURES:

- stem & leaf

- mean

- median

- mode

* this graph is good

- ④ Healthins ✓
- ① STABREV
- ② REWSYSF } }
- ③ Proj For Analyt }

DISPERSION MEASURES

• σ^2 → "statistician's bread & Butter"

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

→ square it to get a meaningful metric
"sums of square"

of L1 norm - absolute val
of L2 norm - squared

④ Residual sum of squares

④ Average sum of squares (ARSE)

• SKewness

- long tail

- happiness is relative

- kurtosis = extreme events sensitivity

platy kurtiz = Chicago

eukurtiz = Manila

- modality

- implies diff segments in population → separate to get unimodal

- separation by age

- nominal
- pie chart: labels must be mutually exclusive, small sample size, people who like other stuff
 - frequency tables
 - M_w = weighted average
 - histogram

LINEAR TRANSFORMATION OF VARIABLE

$$y = a + bx$$

• correlation

• covariance

• trimmed mean = w/ data deleted

03/06/19 session 1 notes ^{concept} Review

Ch 7: NORMAL DISTRIBUTIONS → continuous distributions

1) INTRODUCTION TO NORMAL DISTRIBUTIONS

- "bell curve", "gaussian curve"
- μ, σ

$$\text{density of normal dist.} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

e = base of ln or $\log(e)$

Hypothesis testing:

$$H_0: \mu = 100$$

$$H_1: \mu \neq 100$$

- Features of Normal Distributions:

- 1) symmetric around their μ
- 2) mean = median = mode
- 3) Area under normal curve = 1.0
- 4) Denser in center & less dense in tails
- 5) Defined by 2 parameters: μ & σ
- 6) 68% of its area is within 1σ of μ
- 7) ~95% of its area is within 2σ of μ

Decision	H_0	H_1
H_0	confidence	Type II β
H_1	Type I α	power

2) HISTORY OF THE NORMAL DISTRIBUTION

- Abraham de Moivre - "as the # of events increased, the shape of the binomial dist approached very smooth curve"
- many natural phenomena are naturally distributed
- Laplace = central limit theorem

3) AREAS UNDER NORMAL DISTRIBUTIONS

⊕ online normal & inverse normal calculator → wtf.

⊕ varieties demonstration

4) STANDARD NORMAL DISTRIBUTION

- normal distribution with a $\mu=0$ and $\sigma=1$ is standard normal dist.
- z-table

$$z = \frac{(x-\mu)}{\sigma}$$

z = ^{z-score}
normal dist

x = value on original dist

μ = mean of orig dist

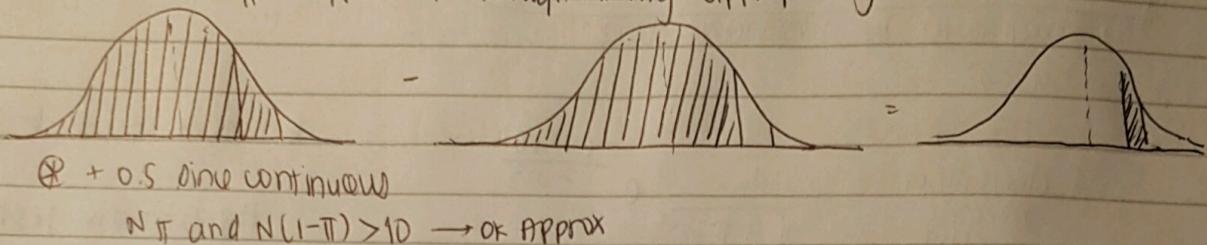
σ = σ of orig dist

- "standardizing the distribution" → transforming dist to ^{std} normal

(2)

5) NORMAL APPROXIMATION TO THE BINOMIAL

- normal dist can be used to approximate the binomial distribution (continuous dist)
- solution: round off & approximate through binning appropriately



4) STATISTICAL LITERACY DEMO

- evaluating tail risk (risk analysis)
↳ investment that moves $> 3\sigma$ from μ

↳ Larsen & Marx - 5th Ed:

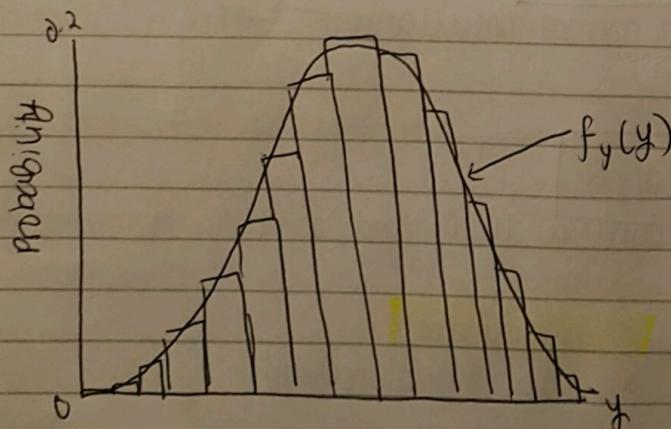
4.3 The Normal Distribution

- deMoivre:

$$f_z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad -\infty < z < \infty$$

can be used to estimate

$$P\left[a \leq \frac{X - n(\frac{1}{2})}{\sqrt{n(\frac{1}{2})(\frac{1}{2})}} \leq b\right] \quad \text{where } X = \text{binomial random var w/ large } n, p = 1/2$$



THEOREM:
4.3.1

Let X be a binomial random var defined on n independent trials for w/o p = p(success). For any #s a and b :

$$\lim_{n \rightarrow \infty} P\left[a \leq \frac{X - np}{\sqrt{np(1-p)}} \leq b\right] = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-z^2/2} dz$$

The formula for the normal probability distribution w/ μ and σ .

- Area under the curve = prob of Being within a given region (3)

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

Gaussian curve

- finding areas under the standard normal curve:

- 1) Using a normal table
2) algorithm

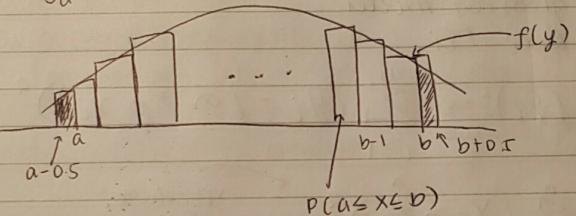
$$F_Z(z) = P(Z \leq z) \text{ associated w/ } z$$

& compute Z score
& find Z score on normal table

- the continuity correction

→ extending $a - 0.5$ to $b + 0.5$

$$\int_a^b f(y) dy \approx P(a \leq X \leq b)$$



• rule of thumb: DeMoivre-Laplace limit should only be used if magnitude of n and p are such that $n > 9 \frac{P}{1-P}$ and $n > 9 \frac{1-P}{P}$

Sum of squares

- residual sum of squares - to help you decide if a statistical model is a good fit for your data

• measures overall diff b/w your data & values predicted by estimation model

• "Residual" - distance from data point to regression line

$$\text{Total SS} = \text{Explained SS} + \text{Residual SS}$$

- Total SS = "how much variation is there in the dependent variable?"

$$\text{Total SS} = \sum (y_i - \bar{y})^2$$

$$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2 \rightarrow \text{regression.}$$

(4)

- Explained SS = "how much of the variation in the dep var does our model explain"?

$$\text{Explained SS} = \sum (\hat{y}_i - \bar{y})^2$$

- Residual sum of squares = "how much of the dependent var variation did our model NOT explain"?

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - f(x_i))^2 \\ &= \sum_{i=1}^n (\epsilon_i)^2 \end{aligned}$$

Linear simple regression model:

$$y_i = \alpha + \beta x_i + \epsilon_i$$

$$\text{RSS} = \sum_{i=1}^n (\epsilon_i)^2 = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

(means squared)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

"prediction"

↓

"square error"

- estimator that measures the square of the errors
- "risk function"
- measure of quality of an estimator (min-negative)
- closer to 0 = better

- vector of n predictions generated from a sample of n dep vars
- \hat{Y} = vector of observed values of var being predicted