

STAC58 - Chapter 1

Textbook: Probability and statistics - Evans & Rosenthal,
Chapter 5-9

Mastering Statistical Evidence using relative belief

- M. Evans

website: <http://www.mathstat.yorku.ca/~mikevans/stac58/stac58.html>

Evaluations - Midterm 40%
- final 60%

January 7, 2019
11:15 AM

STAC 58 - Statistical Inference

Basics: Introduction

- Statistical inference is not so much about the methods of statistics but the "why".
- What is statistics as a subject all about?
- statistical methods are used in:
 - Finance
 - Machine learning
 - medicine
 - quantum physics
 - :
more!
- Furthermore, "statistical reasoning" is becoming more and more important!
- It is being used as a tool to reason about reality.
- Note: significant decisions are made based on statistical analyses.
- So we want the rules of statistical reasoning to be sound = logical, free of contradictions, paradoxes, etc... So we feel confident that whatever the conclusion/inference we draw makes sense.

- Current state of statistics
 - Many different points of view about what the correct statistical reasoning is.
 - This makes learning the subject hard.

- Purpose of this course (STAT58 - Statistical Inference)

- 1.) Survey the various approaches
- 2) present the outline of a logical way to develop a theory of statistical reasoning.

- Some phenomenon / context in the real world that we have questions about
- Questions like:
 - 1) what is the value of some quantity of interest?
e.g. mean half life length of a neutron
Answer: An estimate of assessment of its error
 - 2) Does a certain quantity take a particular value?
Answer: hypothesis assessment - evidence for or against and a measure of strength.
- when can statistical inference play a role?

Theory tells
you how
accurate
estimate is.

Statistical Problems

- The first thing we need to do is be very clear about what a statistical problem is.
- It is all based on "measuring" and counting.
- We have a population Ω = a finite set of objects of interests.

Eg. Ω = set of all students enrolled at UoT on Jan 7, 2019

$\#(\Omega) < \infty$

cardinality / # of items in the set

- we have a measurement(w) defined on Ω

$$X: \Omega \rightarrow \mathcal{X} \quad \forall \omega \in \Omega \quad x(\omega)$$

- for $\omega \in \Omega$ = set of students at UoT.

Define

$x_1(\omega)$ = height of ω in cm (interval)

$x_2(\omega)$ = weight of ω in kg (interval)

$x_3(\omega)$ = gender of ω (categorical)

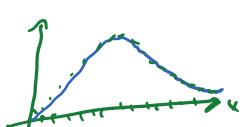
$$x = (x_1, x_2, x_3): \Omega \rightarrow R \times R \times \{M_1, F\}$$

Ω and X define relative frequency distribution over x .

$$x(\omega) = \begin{pmatrix} x_1(\omega) \\ x_2(\omega) \\ x_3(\omega) \end{pmatrix}$$

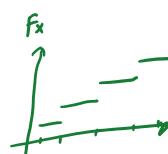
$x(\omega)$ is 3d measurements

$$K: \Omega \rightarrow R' \times R' \times \{M_1, F\}$$



simplify by introducing continuous approximation

discrete is too hard form!
approx it.



Step function

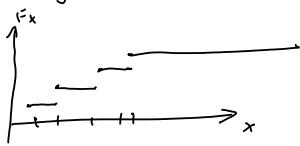
- When $x = R$ (or an interval)

$$F_X(x) = \frac{\#\{\omega: x(\omega) \leq x\}}{\#\Omega} = \text{cumulative distribution function of } X \quad (\text{CDF of } X)$$

$$= \sum_{z \in \mathcal{X}} f_Z(z)$$

$$f_Z(x) = F_X(x) - F_X(x-\epsilon) \quad \text{where } F_X(x-\epsilon) = \lim_{\epsilon \downarrow 0} F_X(\epsilon)$$

- So F_x and f_x are two equivalent ways of presenting a frequency distribution.

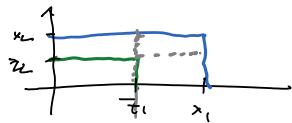


- when $X = \mathbb{R}^2$

★ $F_x(x_1, x_2) = \frac{\#\{\omega : x_1(\omega) \leq x_1, x_2(\omega) \leq x_2\}}{\#\Omega}$

$$= \sum_{\substack{z_1 \leq x_1 \\ z_2 \leq x_2}} f_x(z_1, z_2)$$

$$f_x(x_1, x_2) = \lim_{z_1 \nearrow x_1} \left[F_x(x_1, z_2) - F_x(x_1, z_1) - F_x(z_1, z_2) + F_x(z_1, z_1) \right]$$



So, $F_x \leftrightarrow f_x$

- The whole point of any statistical analysis is to learn something about F_x .
- how do we do this?
- If possible we do an estimate, namely compute $x(\omega)$ & $w(\omega)$ of the form f_x .
- Typically count (return to this in a moment)
- why do we want to know F_x ?

e.g. relationships among variables.

- Suppose (x, y) , where $x: \Omega \rightarrow X$, $y: \Omega \rightarrow Y$
and we want to know if there is a relationship between $x \& y$ on Ω .

- form the conditional relative frequency distribution.

★ $f_{y|x}(y|x) = \frac{\#\{\omega | x(\omega) = x, y(\omega) = y\}}{\#\{\omega | x(\omega) = x\}}$

$$= \frac{f(x, y)(x, y)}{f_x(x)}$$

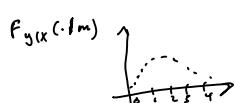
★ Definition: x and y are related variables over Ω if $f_{y|x}(\cdot|x)$ changes as x changes.

- The "form" of the relationship between x and y is given by how $f_{y|x}(\cdot|x)$ changes as x changes.

e.g. $\Omega = 1^{\text{st}}$ year students at UoT

$y = \text{GPA as of Dec 31, 2015}$

$x = \text{gender}$



- often simplifying assumptions are introduced.

- regression assumption: $f_{y|x}(\cdot|x)$ changes at most through its mean as x changes.



$$\begin{aligned} \frac{t(y|x)(\omega)}{\#\{\omega | x(\omega) = x\}} &= \frac{1}{\#\{\omega | x(\omega) = x\}} = x \quad \#\{\omega | x(\omega) = x\} = x^3 \\ E_k &= \sum_{\omega} y f_{y|x}(y|x) \end{aligned}$$

$f_{g(x)}(F)$

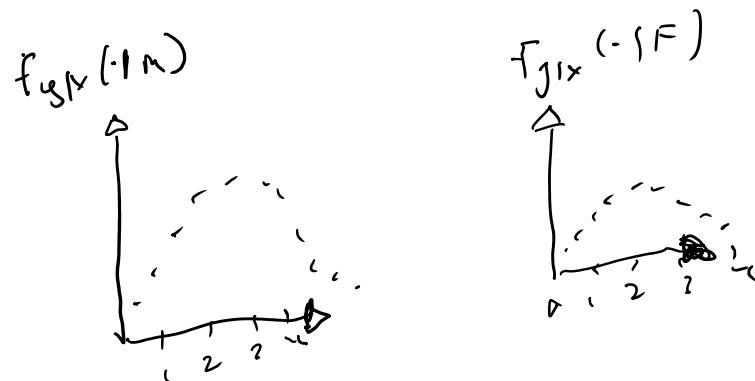
e.g. relationship among variables

- Suppose $(x, y) : \Omega \rightarrow X \times Y$ and we want to know if there is any relationship between X & y .
- Form the conditional relative frequency distribution

$$\star F_{y|x}(y|x) = \frac{\#\{w: x(w)=x, y(w)=y\}}{\#\{w: x(w)=x\}} = \frac{F_{(x,y)}(x,y)}{f_x(x)}$$

\star Def X & y are related variables over a population Ω , if $F_{y|x}(\cdot|x)$ changes as x changes.

eg. Ω = students at USTT
 $x(w)$ = gender $y(w)$ = GPA



- The form of the relationship

\star if x and y have no relationship then
 $f_{y|x}(y|x) = f_y(y)$
 $\Leftrightarrow F_{(x,y)}(x+y) = F_x(x)F_y(y)$

- the form of the relationship between X & y when it exists is given by how $F_{y|x}(\cdot|x)$ changes with x

- Often simplifying assumptions are made

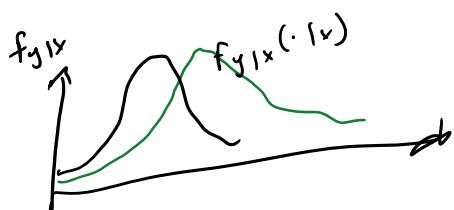
regression assumptions

- here $y \in \mathbb{R}$ and we assume

$f_{y|x}(\cdot | x)$ changes at most though

★ $E(y|x)(x) = \frac{\sum_{\omega: x(\omega)=x} y(\omega)}{\#\{\omega: x(\omega)=x\}}$ \geq average value of y in the sub population

$\sum \omega: x(\omega)=x$



- regression assumption

- same distribution but shifted.

Linear regression assumption

$$E(y|x) \in \{g_1, \dots, g_K\} \text{ where}$$

$$g_i: x \rightarrow \mathbb{R} \quad i=1, \dots, k.$$

$$\text{i.e. } E(y|x)(x) = p_1 g_1(x) + \dots + p_k g_k(x) \text{, for some } p_1, \dots, p_k \in \mathbb{R}$$

$$\text{e.g. } g_1(x) = 1, \quad g_2(x) = x, \quad g_3(x) = x^2$$

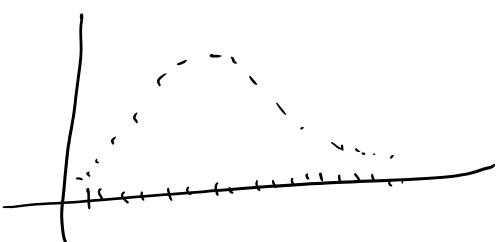
Infinity and Continuity

- Recall, all populations are finite

- As a result, all relative frequency distributions are positive on a finite # of points

e.g. $\Omega = \text{students at the university of Toronto}$

$x(\omega) = \text{height in centimeters.}$

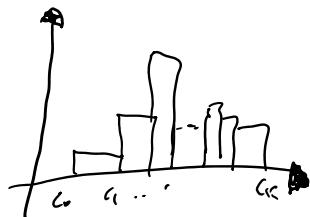


- Sometimes measurements can be thought of as being taken to an ever finer accuracy.

- histogram - divide up range of $x: \Omega \rightarrow \mathbb{R}'$ into k intervals: $c_0 < c_1 < \dots < c_k$

- Density histogram

$$\star f_x^{\text{dense}}(x) = \sum_{z \in (c_i, c_{i+1})} f_x(z) \frac{\chi_{(c_i, c_{i+1})}(z)}{(c_{i+1} - c_i)} = \frac{\# \{w: c_i < x(w) \leq c_{i+1}\}}{\#(\Omega)(c_{i+1} - c_i)}, \text{ when } x \in$$



Then for $c_i > c_j$

$$F_x(c_j) - F_x(c_i) = \int_{c_i}^{c_j} f_x^{\text{dense}}(x) dx$$

- For such measurements we can introduce the idea of continuous approximation.

- Let f be a probability density function for x so

$$(i) f(x) \geq 0$$

$$(ii) \int_{-\infty}^{\infty} f(x) dx = 1$$

- then f approximation f_x^{dense} if $\forall a < b$

$$\int_a^b f(x) dx \approx \int_a^b F_x^{\text{dense}}(x) dx$$

make precise via limit
as measurement acc.
make to finer and
finer accuracy.

- so continuous distribution (and infinite sample pass) arise as approximation

e.g. $\Omega = \text{students of U of T}$

- $X(\omega) = \text{height in cm}$

- plausibly we cause the "approximation"

$$X(\omega) \sim N(\mu, \sigma^2), \quad \forall \mu \in \mathbb{R}, \sigma^2 \in (0, \infty)$$

$$\text{actually } \mu = \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} X(\omega)$$

$$\sigma^2 = \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} (X(\omega) - \mu)^2$$

e.g. $\Omega = \text{students at U of T}$

- $Y(\omega) = \text{weights in grams}$

- $X(\omega) = \text{height in cm}$

$$- \begin{pmatrix} X(\omega) \\ Y(\omega) \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 \end{pmatrix} \right)$$

$$\mu_y = \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} Y(\omega), \quad \mu_x = \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} X(\omega)$$

$$\sigma_x^2 = \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} (X(\omega) - \mu_x)^2$$

$$r_{xy} = \frac{1}{\sigma_x \sigma_y} \frac{1}{\#(\Omega)} \sum_{\omega \in \Omega} (X(\omega) - \mu_x)(Y(\omega) - \mu_y) = \frac{\sum_{\omega \in \Omega} (X(\omega) - \mu_x)(Y(\omega) - \mu_y)}{\sqrt{\sum_{\omega \in \Omega} (X(\omega) - \mu_x)^2} \sqrt{\sum_{\omega \in \Omega} (Y(\omega) - \mu_y)^2}}$$

$$- Y(\omega) | X(\omega) = x_1 \sim N(\mu_y + \sigma_y r_{xy} \frac{x_1 - \mu_x}{\sigma_x}, \sigma_y^2 (1 - r_{xy}^2))$$

- So a relationship between X and Y exists if $r_{xy} \neq 0$

- how strong is the relationship?
- the closer ρ_{xy} is to 1 the stronger the stronger theoretically.
- note we can write

$$\begin{aligned} E(x_1 x_2) &= \mu_1 + \frac{\sigma_1}{\sigma_2} \rho (\mu_2 - \mu_1) \\ &= (\mu_1 - \frac{\sigma_1}{\sigma_2} \rho \mu_2) + \frac{\sigma_1}{\sigma_2} \rho \mu_2 \\ &= B_0 + B_1 x_2 \end{aligned}$$

where $B_0 = \mu_1 - \frac{\sigma_1}{\sigma_2} \rho \mu_2$, $B_1 = \frac{\sigma_1}{\sigma_2} \rho$

so alternatively we can write

$$y = B_0 + B_1 x + \varepsilon, \text{ where } \varepsilon \sim N(0, \sigma^2)$$

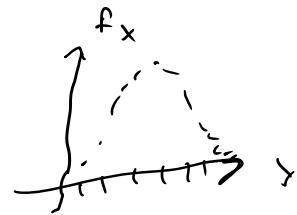
note - if we know f_x we know everything from a statistical point-of-view

note - when is a relationship between X & Y causal
 - must be able to assign any value of X to any $w \in \Omega$

- Ω = population

$x: \Omega \rightarrow X$ a measurement ($\begin{pmatrix} x \\ y \end{pmatrix}$)

for $x \in X$ $f_x(x) = \frac{\#\{w: x(w) = x\}}{\#\Omega}$



- with a census we know everything

Sampling

- Typically we can't conduct a census

- so select $n < \#\Omega$ and $\{w_1, \dots, w_n\} \subseteq \Omega$

obtaining $x_1 = x(w_1), \dots, x_n = x(w_n)$

- based on the data $x = (x_1, x_2, \dots, x_n)$ we make inference about Ω

- how do we select $\{w_1, \dots, w_n\}$?

- clearly we want to avoid selection effects as this will "bias" our results/inferences.

e.g. if $w_i \in \Omega = \text{student's pass}$ is always female then x would not be representative.

Solution

• We use a "random mechanism" to select $\{w_1, \dots, w_n\}$

e.g. - put $\#\Omega$ chips in a bowl each labelled w^i a number in $\{1, 2, \dots, \#\Omega\}$ where i corresponds to the i^{th} number of Ω

- stir up the chips and pick one without looking and say the one labelled i .

- continue without replacement to obtain $x_1 = x(w_1), \dots, x_n = x(w_n)$

- then select the population element

We model this via probability:

$$P\{x(\omega_1) = x_1\} = \frac{\#\{\omega : x(\omega) = x_1\}}{\#\Omega}$$

$$= f_x(x_1)$$

$$P\{x(\omega_2) = x_2 | x(\omega_1) = x_1\}$$

$$= \begin{cases} \frac{\#\{\omega : x(\omega) = x_2\} - 1}{\#\Omega - 1}, & x_2 = x_1 \\ \frac{\#\{\omega : x(\omega) = x_2\}}{\#\Omega - 1}, & x_2 \neq x_1 \end{cases}$$

$$= \begin{cases} \frac{f_x(x_1) - \frac{1}{\#\Omega}}{1 - \frac{1}{\#\Omega}}, & x_2 = x_1 \\ \frac{f_x(x_2)}{1 - \frac{1}{\#\Omega}}, & x_2 \neq x_1 \end{cases}$$

, given you removed the first person what's the probability of getting the second person.

we can express it in terms of original frequency function,

$$\approx f_x(x_2)$$

$$\therefore P(x(\omega_1) = x_1, x(\omega_2) = x_2) = P(x(\omega_2) = x_2 | x(\omega_1) = x_1) P(x(\omega_1) = x_1)$$

$$\approx f_x(x_1) f_x(x_2) \Rightarrow \perp \quad , \text{when } \#\Omega \text{ is large}$$

- Provided $n \ll \#\Omega$

$$P(x(\omega_1) = x_1, \dots, x(\omega_n) = x_n) \approx f_x(x_1) \cdots f_x(x_n)$$

so x_1, \dots, x_n are approximately i.i.d \rightarrow

Note:

- if we know f_x we know everything from a statistical point of view.

Note:

- if we have generated the data via a random mechanism then we are justified in referring to the data as being objective

Statistical Models

- for inference about f_x , we have the data $x = (x_1, \dots, x_n)$
- some basic inference justified by convergence results.

e.g. $A \subseteq \mathcal{X}$

$$- P(A) = \sum_{x \in A} f_x(x) = \text{proportion of } \underbrace{\omega \in \Omega}_{\text{of individuals}} \text{ s.t. } x(\omega) \in A$$

$$\begin{aligned} - \text{estimate by } P(A) &= \frac{1}{n} \sum_{i=1}^n I_A(x(\omega_i)) \\ &= \text{proportion of sample values s.t. } x_i \in A \\ &\xrightarrow{\text{SLLN}} P(A) \text{ as } n \rightarrow \infty \end{aligned}$$

- $I_A(x(\omega)) \sim \text{Bernoulli}(P(A))$, when ω is randomly selected.

$$\begin{aligned} \text{Var}(\hat{P}(A)) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(I_A(x(\omega_i))) \\ &= \frac{1}{n} \text{Var}(I(x(\omega))) \\ &= \frac{P(A)(1-P(A))}{n} \end{aligned}$$

- actually $n \hat{P}(A) \sim \text{Binomial}(n, P(A))$

$$\begin{aligned} - \frac{\hat{P}_n(P(A) - P(A))}{\sqrt{n \hat{P}(1-\hat{P}(A))}} &\xrightarrow{\text{CLT}} N(0, 1) \\ &\text{since } P(A) \text{ with} \end{aligned}$$

$$\begin{aligned}
 & - \frac{\hat{P}(A) - P(A)}{\sqrt{\hat{P}(A)(1-\hat{P}(A))}} \xrightarrow{\text{CLT}} N(0, 1) \\
 & - \text{so } \hat{P}(A) \pm 3 \sqrt{\frac{\hat{P}(A)(1-\hat{P}(A))}{n}}
 \end{aligned}$$

contains $P(A)$ with
 virtual confidence

- but requires in "large"
- similarly for other characteristics of f_Y
e.g. $\mu_Y, \sigma^2_Y, \frac{c_X}{\text{av}}, \text{etc...}$
- want inference methods that don't require large n .
- Typically this involves making assumptions
as assumption in value making (subjective) choices.

Principle of Empirical Criteria

- Any ingredient (assumption) we use as part of a statistical analysis must be checkable against the (objective) data to ensure that it makes sense.
- The most important assumption made is the choice of a statistical model.
- we assume $f_X \in \{f_\theta : \theta \in \Theta\} = M$
where f_θ is a density \propto for each $\theta \in \Theta$
- Θ is called the parameter of the model and Θ is called the parameter space.
- Typically Θ indexes; to each value of θ there corresponds a unique density f_θ (no non-identifiability)
- we have to check the assumption

$$f_X \in \{f_\theta : \theta \in \Theta\} = M$$

how? later
- So suppose we have checked M against X and have decided M is okay
(note: in general we never say it is correct)
- then we want rules to apply to (M, X) to make inferences about true value of θ
(and thus equivalent true f_Y)

Eg. $\Omega = \text{students at UofT}$

- $X(\omega) = \text{ht of } \omega \text{ in cm}$
- assume $X(\omega) \sim N(\mu, \sigma^2)$
- $\Theta = (\mu, \sigma^2) \in \mathbb{D} = \mathbb{R} \times (0, \infty)$
- To check the model look at standard residuals
- $r_i = \frac{x_i - \bar{x}_i}{s_i} \stackrel{\text{fast}}{\sim} \text{Uniform}(-\sqrt{n}, \sqrt{n})$

- propose test statistic that look for patterns

$$\text{eg. skewness } T(r) = \frac{1}{n} \sum_{i=1}^n r_i^3 \xrightarrow{\text{CLL}} 0$$

- so should be close to zero

- why make the assumption $f_x \in M$

(1) small $n \Rightarrow$ assumption replaces data for more accurate info.

(2) enables development of a theory of inference as opposed to asymptotes.

Types of Inferences

- Often we don't need to know θ but only $\mu = \mathbb{E}(\theta) \in \bar{\Theta}$

$$\text{eg} \quad - x = (x_1, \dots, x_n) \sim N(\mu, \sigma^2)$$

$$(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$$

- $\mathbb{E}(\mu, \sigma^2) = \mu$, inference about mean

- $\mathbb{E}(\mu, \sigma^2) = \sigma$, inference about standard deviation

- $\mathbb{E}(\mu, \sigma^2) = \frac{\sigma}{\sqrt{n}}$

- $\mathbb{E}(\mu, \sigma^2) = \mu + \sigma z_{0.75}$ (3rd quartile)

- $\mathbb{E}(\mu, \sigma^2) = \Phi\left(\frac{x_2 - \mu}{\sigma}\right) - \Phi\left(\frac{x_1 - \mu}{\sigma}\right)$

= probability of interval (x_1, x_2)

- two types of inferences (this is what we want the theory of inference to give)

(1) Estimate $\hat{\psi}(x)$ together with an assessment of the accuracy of the estimates as given by a set $(c(x) \subseteq \bar{\Theta}$ with $\hat{\psi}(x) \in C(x)$ and size of $C(x)$ gives assessment.

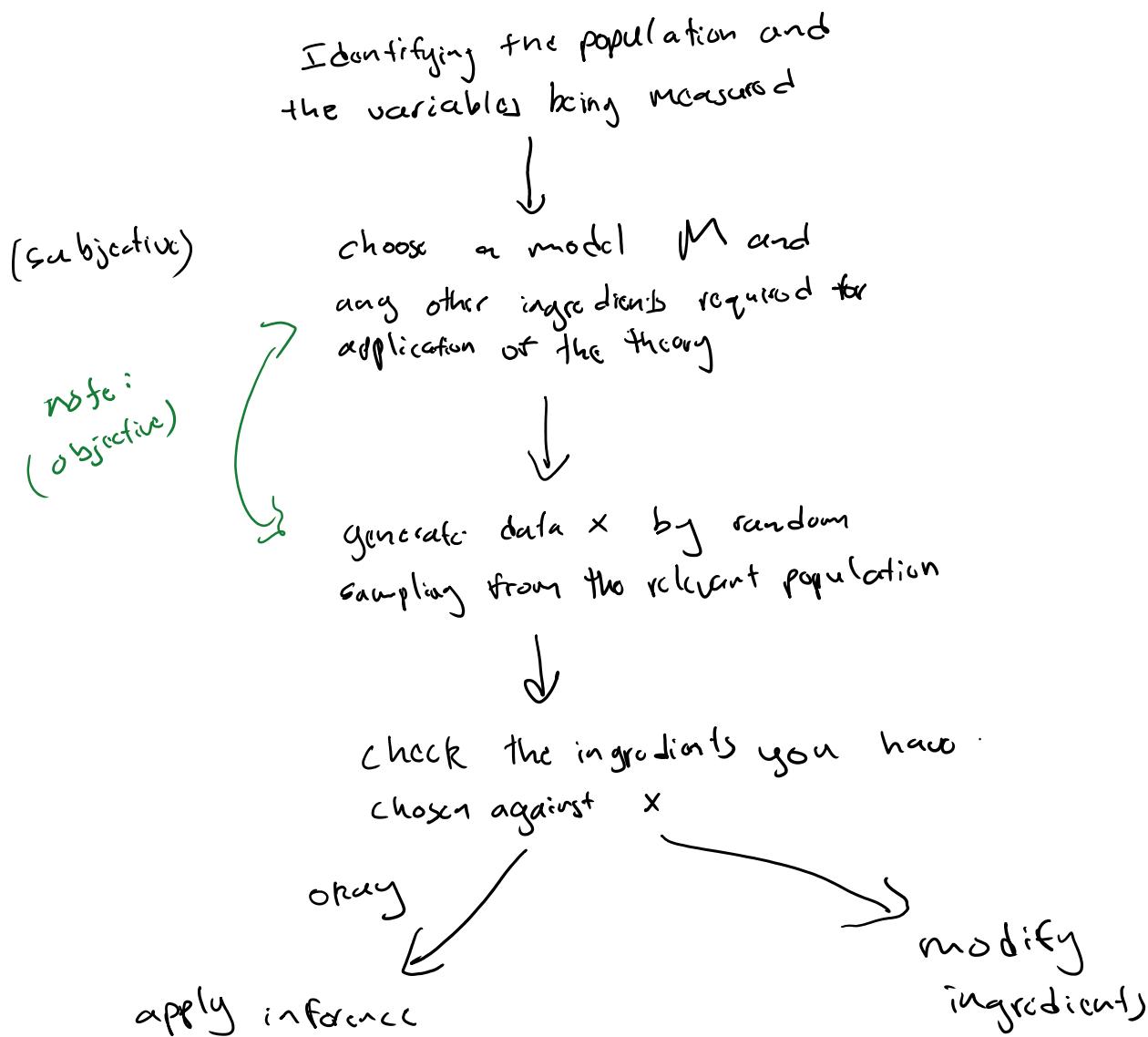
(2) assess hypothesis $H_0 = \{\psi_0\} \subseteq \bar{\Theta}$ by a quoting a measure of the evidence

- Basically means we have a hypothesis H_0 that the true value of ψ is $\psi_0 \in \bar{\Theta}$ that H_0 is true together with a measure of the strength of this evidence.

- all theories of statistical inference attempt to answer 1 or 2.

Statistical Analysis

- We can conceive of a statistical analysis as the following steps.



Chapter 2

February 2, 2019 11:38 AM

