

# Statistical Inference

Real Chaps  
in  $\mathbb{E} + \mathbb{R}$

①

## ① Basics

### ①. Introduction

- not so much about the methods of statistics but the why

- what is statistics as a subject all about?

- statistical methods are used

- finance.

- machine learning

- medicine

- quantum physics.

⋮

- "statistical reasoning" is becoming more important

- it is being used as a tool to reason about reality

- significant decisions are made based on statistical analysis

- so we want the rules of statistical reasoning to be sound = logical, free of contradictions, paradoxes, etc. so we feel confident that whatever conclusions/inferences

that we draw make sense

- current state of statistics

- many different points of view about what is correct statistical reasoning

- makes learning the subject harder.

- purpose of the course

(1) survey the various approaches

(2) present the outline of a logical way to develop a theory of statistical reasoning

- start from the start

- some phenomenon/context in the real world that we have questions about

- questions like (1) what is the value of some quantity of interest?  
eg. mean half life length of a neutron.

(2) does a certain quantity take a particular value?

- when can statistical inference play a role?

## ② 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  $\Omega$  = a finite set of objects of interest

Ex  $\Omega$  = set of all students enrolled at UFT on Jan. 2, 2015

- $\#(\Omega) < \infty$
  - we have a measurement(s) defined on  $\Omega$
- $X: \Omega \rightarrow \mathcal{X}$

Ex (cont'd) - for  $\omega \in \Omega$  = set of student at UFT

- define  $X_1(\omega)$  = ht of  $\omega$  in cm (interval)
- $X_2(\omega)$  = wt of  $\omega$  in kg ("")
- $X_3(\omega)$  = gender of  $\omega$  (categorical)

-  $X = (X_1, X_2, X_3): \Omega \rightarrow \mathbb{R} \times \mathbb{R} \times \{M, F\}$

- Def and  $X$  define a relative frequency function over  $\Omega$

$$f_x(x) = \frac{\#\{\omega : X(\omega) = x\}}{\#\Omega}$$

= proportion of individuals in  $\Omega$  whose  $X$  measurement is  $x \in \mathcal{X}$

- note (i)  $0 \leq f_x(x) \leq 1$

$$(ii) \sum_{x \in \mathcal{X}} f_x(x) = 1$$

and only finitely many  $x \in \mathcal{X}$  have  $f_x(x) > 0$ .

- when  $\mathcal{X} = \mathbb{R}$  (or an interval)

$$F_x(x) = \frac{\#\{\omega : X(\omega) \leq x\}}{\#\Omega} = \text{cumulative dist. fn of } X$$

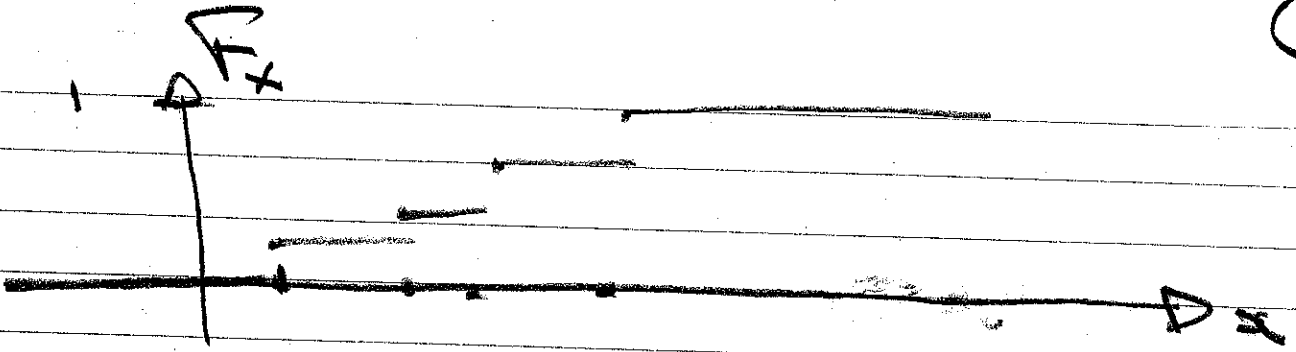
$$= \sum_{z \leq x} f_x(z)$$

$$f_x(x) = F_x(x) - F_x(x-0)$$

$$\text{where } F_x(x-0) = \lim_{z \uparrow x} F_x(z)$$

so  $F_x$  and  $f_x$  are two equivalent ways of presenting a frequency distribution.

(5)



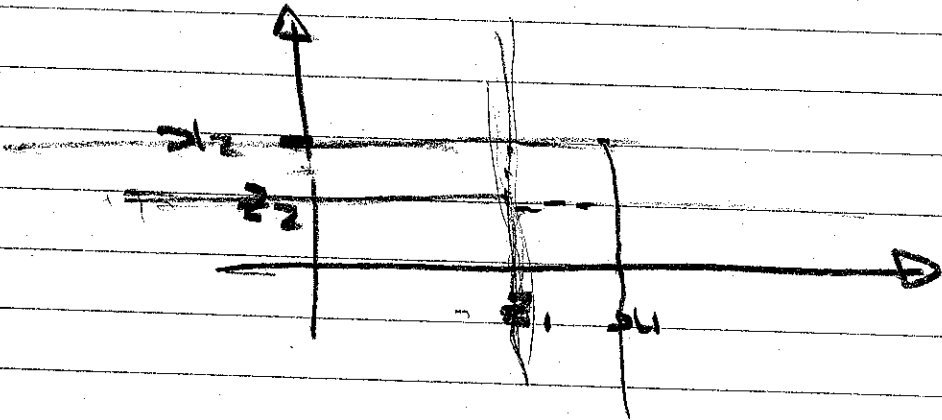
- when  $X = \mathbb{R}^2$

$$F_x(x_1, x_2) = \# \{ \omega : X_1(\omega) \leq x_1, X_2(\omega) \leq x_2 \}$$

$$\#(0)$$

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

$$F_x(x_1, x_2) = \lim_{\substack{p_1 \rightarrow x_1 \\ p_2 \rightarrow x_2}} (F_x(x_1, x_2) - F_x(p_1, x_2) - F_x(x_1, p_2) + F_x(p_1, p_2))$$



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 could do a census, namely compute  $f_X(\omega) \forall \omega \in \Omega$  at the form

- typically can't (return to this in a moment)

- why do we want to know  $F_X$ ?

relationships among variables

- suppose  $(X, Y)$  where  $X: \Omega \rightarrow \mathcal{X}, Y: \Omega \rightarrow \mathcal{Y}$  and we want to know if there is an relationship between  $X$  and  $Y$ .

- from the empirical relative frequency distribution

$$P_{Y|X}(y|x) = \frac{\#\{\omega \mid X(\omega)=x, Y(\omega)=y\}}{\#\{\omega \mid X(\omega)=x\}}$$

$$= \frac{P_{(X,Y)}(x,y)}{P_X(x)}$$

$$= \frac{P_{(X,Y)}(x,y)}{P_X(x)}$$

$$= \frac{P_{(X,Y)}(x,y)}{P_X(x)}$$

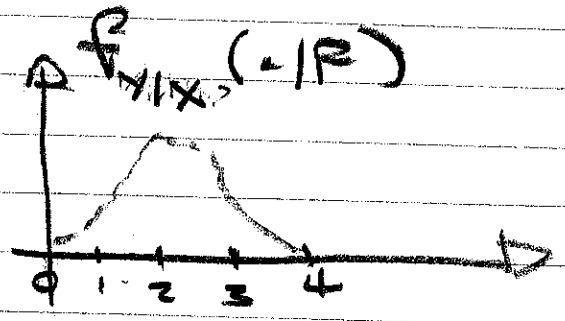
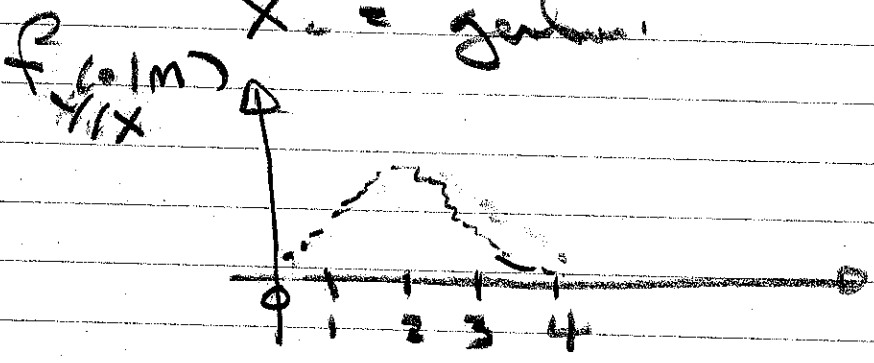
what roles  $Y$  a response and  $X$  a predictor?

Def  $X$  and  $Y$  are related variables over  $\Omega$  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.

Ex  $\Omega =$  1st yr students at UofT  
 $Y =$  GPA as of Dec 31, 2015

$X =$  gender



- often simplifying assumptions are introduced

- regression assumption  $f_{Y|X}(\cdot|x)$  changes at most through its mean as  $x$  changes,  $E(Y|X=x)$   

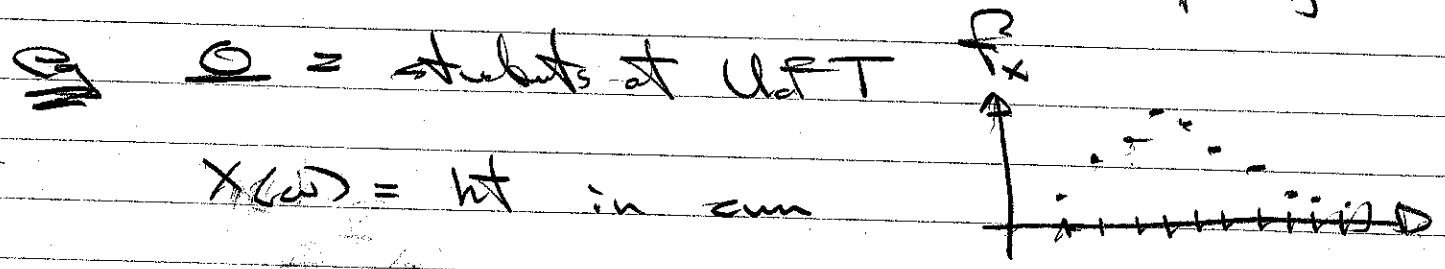
$$= \int y f_{Y|X}(y|x) dy$$

$$E_x = \int y f_{Y|X}(y|x) dy$$

- linear regression assumption  $E(Y|X) = \alpha + \beta X$

### ③ Infinity and Continuity

- all populations are finite
- so all relative frequency distributions are positive on a finite number of points



- sometimes a measurement can be thought of as being taken to an over finer accuracy

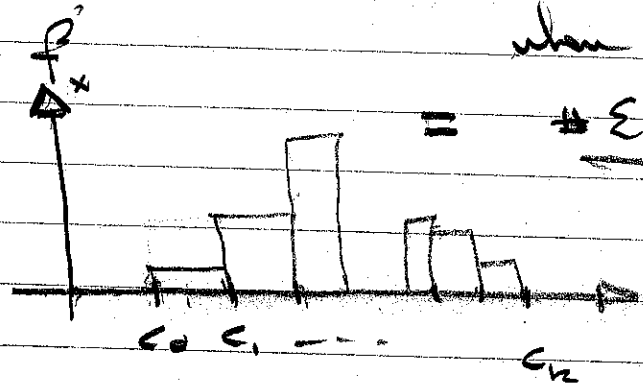
- histogram - divide up range into  $k$  intervals  
 $c_0 < c_1 < \dots < c_k$

### density histogram

$$f_x^{\text{den}}(a) = \sum_{\omega \in (c_j, c_{j+1})} f_x(\omega) / (c_{j+1} - c_j)$$

when  $a \in (c_j, c_{j+1}]$

$$= \frac{\#\{\omega : c_j < X(\omega) \leq c_{j+1}\}}{\#\{ \omega \in (c_j, c_{j+1}) \}}$$



- then for  $c_j > c_i$



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

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

- let  $f$  be a prob. density f<sub>n</sub> so  
 (i)  $f(x) \geq 0$  (ii)  $\int_{-\infty}^{\infty} f(x) dx = 1$

- then  $f$  approximates  $f_x^{disc}$  if  $a < b$

$$\int_a^b f(x) dx \approx \int_a^b f_x^{disc}(x) dx$$

make precise via limit as measurements are made by finer and finer accuracy

- so continuous distributions (and infinite sample spaces) arise as approximations.

### III Q = students at UoFT

-  $X(\omega)$  = ht in cm.

- plausibly we can use the "approximation"

$$X(\omega) \sim N(\mu, \sigma^2)$$

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

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

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

Q = students at Uaf T

Y(w) = wt in kgms

X(w) = ht in cms

(X(w), Y(w)) ~ N((mu\_1, mu\_2), (sigma\_1^2, sigma\_1 sigma\_2 rho, sigma\_1 sigma\_2 rho, sigma\_2^2))

mu\_x = 1/n sum\_{w in Q} X(w), mu\_y = 1/n sum\_{w in Q} Y(w)

sigma\_x^2 = 1/n sum\_{w in Q} (X(w) - mu\_x)^2

rho = 1/(sigma\_x sigma\_y) \* 1/n sum\_{w in Q} (X(w) - mu\_x)(Y(w) - mu\_y)

= sum\_{w in Q} (X(w) - mu\_x)(Y(w) - mu\_y)

observed value n

sqrt(sum\_{w in Q} (X(w) - mu\_x)^2) sqrt(sum\_{w in Q} (Y(w) - mu\_y)^2)

Y(w) | X(w) ~ N(mu\_y + sigma\_y rho (X(w) - mu\_x) / sigma\_x, sigma\_y^2 (1 - rho^2))

so if a relationship between X and Y exists rho != 0

- how strong is the relationship?

- the closer  $\rho^2$  is to 1 the stronger the relationship.

- note we can write

$$\begin{aligned}
E(X_1 | X_2 = x_2) &= \mu_1 + \frac{\sigma_1}{\sigma_2} \rho (x_2 - \mu_2) \\
&= (\mu_1 - \frac{\sigma_1}{\sigma_2} \rho \mu_2) + \frac{\sigma_1}{\sigma_2} \rho x_2 \\
&= \beta_0 + \beta_1 x_2
\end{aligned}$$

where  $\beta_0 = \mu_1 - \frac{\sigma_1}{\sigma_2} \rho \mu_2$ ,  $\beta_1 = \frac{\sigma_1}{\sigma_2} \rho$

so alternatively we can write

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad \text{where } \varepsilon \sim N(0, \sigma^2)$$

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

- note - when is a relationship between  $X$  and  $Y$  causal

- must be able to assign any value of  $X$  to any  $\omega \in \Omega$ .

### ③ Sampling

- typically we can't count a census
- so select  $n < \#(\Omega)$  and  $\{\omega_1, \dots, \omega_n\} \subseteq \Omega$  obtaining  $x_1 = X(\omega_1), \dots, x_n = X(\omega_n)$
- based on the data  $x = (x_1, \dots, x_n)$  we make inferences about  $f_X$
- how do we select  $\{\omega_1, \dots, \omega_n\}$ ?
- clearly we want to avoid selection effects as this will "bias" our results

eg - if  $\omega_i \in \Omega = \text{UFT students}$  was always female then  $x$  would not be representative.

- solution - use a "random mechanism" to select  $\{\omega_1, \dots, \omega_n\}$

eg - put  $\#(\Omega)$  chips in a bowl each labelled with a number in  $\{1, 2, \dots, \#(\Omega)\}$  where  $i$  corresponds to  $i$ -th member of  $\Omega$

- stir up the chips and pick one without looking say the one labelled  $i$
- continue without replacement

Sampling with and without replacement

- the joint population count;

- we model this via probability

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

$$= f_X(x_1)$$

$$P(X(\omega_2) = x_2 \mid X(\omega_1) = x_1)$$

$$= \frac{\#\{\omega : X(\omega) = x_2\} - 1}{\#\{\Omega\} - 1} \quad x_2 = x_1$$

$$= \frac{\#\{\omega : X(\omega) = x_2\}}{\#\{\Omega\} - 1} \quad x_2 \neq x_1$$

$$= \frac{f_X(x_2) - 1/\#\{\Omega\}}{1 - 1/\#\{\Omega\}} \quad x_2 = x_1$$

$$= \frac{f_X(x_2)}{1 - 1/\#\{\Omega\}} \quad x_2 \neq x_1$$

$$\approx f_X(x_2)$$

$$\therefore P(X(\omega_1) = x_1, X(\omega_2) = x_2) = P(X(\omega_1) = x_1) P(X(\omega_2) = x_2 \mid X(\omega_1) = x_1)$$

$$\approx F_{X_1}(a_1) F_{X_2}(a_2)$$

- provided  $n \ll \infty$

$$P(X_1 = a_1, \dots, X_n = a_n)$$

$$= F_{X_1}(a_1) \dots F_{X_n}(a_n)$$

$\Rightarrow$   $X_1, \dots, X_n$  are <sup>approximately</sup> i.i.d.  $F_X$

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

# ④ Statistical Models

15.

- for inference about  $f_X$  we have the data  $x = (x_1, \dots, x_n)$
- some basic inferences justified by emergence results

eg  $A \subseteq X$

-  $P(A) = \sum_{x \in A} f_X(x) =$  proportion of  $\omega \in \Omega$  s.t.  $X(\omega) \in A$

- estimate by  $\hat{P}(A) = \frac{1}{n} \sum_{i=1}^n I_A(x_i)$

$=$  proportion of sample values s.t.  $x_i \in A$

SLLN

$\rightarrow \hat{P}(A)$  as  $n \rightarrow \infty$

-  $I_A(X(\omega)) \sim \text{Bernoulli}(P(A))$   
when  $\omega$  randomly selected

-  $\text{Var}(\hat{P}(A)) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(I_A(X(\omega_i)))$   
 $= \frac{1}{n} \text{Var}(I(X(\omega))) = \underline{P(A)(1-P(A))}$

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

-  $\frac{\sqrt{n}(\hat{P}(A) - P(A))}{\sqrt{P(A)(1-P(A))}} \xrightarrow{\text{CLT}} N(0,1)$

- so  $\hat{P}(A) \pm 3 \sqrt{\frac{P(A)(1-P(A))}{n}}$  contains  $P(A)$   
with virtual certainty

- but requires  $n \rightarrow \text{large}$
- similarly for other characteristics of  $f_x$   
e.g.  $\mu_x$ ,  $\sigma_x^2$ ,  $\sigma_x/\mu_x$ , etc.
- want inference methods that don't require large  $n$ .
- typically this involves making assumptions
- assumptions involve making (subjective) choices

### Principle of Empirical Criticism

Any ingredient (assumption) we use as part of a statistical analysis must be checked against the (objective) data to ensure that it makes sense.

- the most important assumption made is the choice of a statistical model
- we assume  $f_x \in \{f_\theta : \theta \in \Theta\} = \mathcal{M}$   
where  $f_\theta$  is a density on  $\mathcal{X}$  for each  $\theta \in \Theta$
- $\theta$  is called the parameter of the model  
and  $\Theta$  is called the parameter space
- typically  $\theta$  indexes; to each value of



⊙ there corresponds a unique density  $F_0$   
(no nonidentifiability)

- we have to check the assumption

$$F_x \in \{F_0 : \theta \in \Theta\} = \mathcal{M}$$

- how? later

- so suppose we have checked  $\mathcal{M}$   
against  $x$  and have decided  $\mathcal{M}$  is  
okay (note: in general we can never  
say it is correct)

- then we want rules to apply to  
 $(\mathcal{M}, x)$  to make inferences about  
true value of  $\theta$  (and thus equivalently  
the true  $F_x$ )

eg  $\Theta = \text{subset of } U \subseteq T$

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

- assume  $X(\omega) \sim N(\mu, \sigma^2)$

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

- to check the model look at  
standardized residuals

-  $r = \frac{\sum_{i=1}^n x_i - \bar{x}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$  <sup>fast</sup>  $\sim$  Uniform  $(\frac{1}{2}, \frac{1}{2})$   $\wedge$   $N(0, 1)$  (14)

- Propose test statistics that look for patterns.

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

so should be close to 0

- why make this assumption  $R \times \epsilon M$

(1) small  $n$  so assumption replaces data for more accurate inference.

(2) enables the derivation of a theory of inference as opposed to asymptotics

## ⑤ Types of Inferences

- often we don't need to know  $\theta$  but only  $T = T(\theta) \in \mathbb{R}$

⑥ -  $x = (x_1, \dots, x_n) \sim N(\mu, \sigma^2)$

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

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

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

$$\mathbb{R}(\mu, \sigma^2) = \sigma/\mu$$

$$\mathbb{R}(\mu, \sigma^2) = \mu + \sigma z_{\alpha/2} \quad (\text{conf. interval})$$

$$\begin{aligned} \mathbb{P}(\mu, \sigma^2) &= \Phi\left(\frac{x_2 - \mu}{\sigma}\right) - \Phi\left(\frac{x_1 - \mu}{\sigma}\right) \\ &= \text{prob out of } (x_1, x_2) \end{aligned}$$

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

(1) estimate  $\hat{T}(x)$  together with an assessment of the accuracy of the estimate as given by a set  $C(x) \subset \mathbb{R}$  with  $\hat{T}(x) \in C(x)$  and size of  $C(x)$  gives assessment

(2) assess hypothesis  $H_0: \mu \leq \mu_0 \in \mathbb{R}$   
 (i.e. we have a hypothesis that the true  
 value of  $\mu$  is  $\mu_0 \in \mathbb{R}$ ) by quoting  
 a measure of the evidence that  $H_0$   
 is true together with a measure  
 of the strength of this evidence.

- all theories of statistical inference  
 attempt to answer (1) and (2)

# Statistical Analyses

- we can conceive of a statistical analysis as following these steps

identify the population and the variables being measured



(subjective)

choose a model  $M$  and any other ingredients required for application of the theory

note

(objective)



generate data  $x$  by random sampling from the relevant population.



check the ingredients you have chosen against  $x$

okay

apply theory of inference

modify ingredients