## 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 analysis.

- So we want the rules of statistical reasoning to be _____ = logical, free of contradictions, _____, 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 the correct statistical reasoning is.
    • This makes learning the subject hard.

- Purpose of this course (STAC5B - 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 interests?
        eg. mean half life length of a neutron

        2) Does a certain quantity take a particular value?


— When can statistical inference play a role?

## 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 $\underline{\Omega}$ = a finite set of objects of interests.

Eg. $\underline{\Omega}$ = set of all students enrolled at UofT on Jan 2, 2015

- $\#(\Omega) < \infty$

cardinality / # of items in the set

- we have a measurement(s) defined on $\Omega$

$$X : \Omega \rightarrow \mathcal{X}$$

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

Define

$X_1(\omega)$ = height of $\omega$ in cm (interval)

$X_2(\omega)$ = weight of $\omega$ in kg (interval)

$X_3(\omega)$ = gender of $\omega$ (categorical)

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

- $\Omega$ and $X$ define relative frequency function over $X$.

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

= proportion of individuals in $\Omega$ whose $X$ measurements 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 $f_x(x) > 0$.

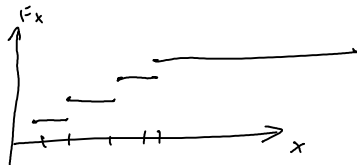- when $\mathcal{X} = R$ (or an interval)

$$F_x(x) = \frac{\#\{\omega : X(\omega) \leq x\}}{\#(\Omega)} = \text{Cumulative distributive function of } X \text{ (CDF of } X)$$

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

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

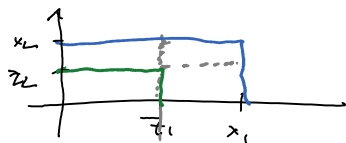- So $F_x$ and $f_x$ are two equivalent ways of presenting a frequency distribution.



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

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

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

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



So, $F_x \Longleftrightarrow f_x$

— The whole point of any statistical analysis is to learn something about $F_x$.

— how do we do this?

— If possible we do a <u>census</u>, namely compute $x(\omega) \ \forall \omega \in \Omega$ of the form $F_x$.

— Typically count (return to this in a moment)

— why do we want to know $F_x$?

<u>eg</u> relationships among variables.

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

— form the conditional relative frequency distribution.

$$f_{y|x}(y|x) = \frac{\# \{ \omega | x(\omega) = x, \ y(\omega) = y \}}{\# \{ \omega | x(\omega) = x \}}$$
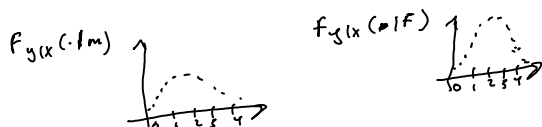
$$= \frac{f_{(x,y)}(x, y_2)}{f_x(z)}$$

<u>Definition:</u> $x$ 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 $f_{y|x}(\cdot | z)$
changes as $x_2$ changes.

eg. $\Omega = 1^{st}$ year students at UofT
$y = $ GPA as of Dec 31, 2015.
$x = $ gender

$f_{y|x}(\cdot | m)$



$f_{y|x}(\cdot | F)$



— often simplifying assumptions are introduced.

— regression assumption: $f_{y|x}(\cdot | x)$ changes at most through its mean as
$x$ changes, $E(y|x)(v)$

$$\underset{y}{\Sigma} y(\omega)$$

$$= \frac{1}{\# \{ \omega: x(\omega) = x \} \{ \omega: x(\omega) = x \}}$$

$$\overset{Ex}{=} \underset{y}{\Sigma} y \, f_{y|x}(y|x)$$