# STAC 58 - Statistical Inference

Evaluations - Midterm   40%
           - final   60%.


Textbook: Probability and statistics - Evans & Rosenthal, Chapter 5.9

Measuring Statistical Evidence using relative belief

- M. Evans

website: http://www.utstat.utoronto.ca/mikevans/stac58/stac58.html

## Statistical Inference

January 6, 2019     9:08 PM

### 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 or 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 (STAC58 - 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?
  eg. 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 strenght.
- when can statistical inference play a role?

Theory tells Ja 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 $\underline{\Omega}$ = a finite set of objects of interests.

Eg. $\underline{\Omega}$ = set of all students enrolled at UofT on Jan 7, 2019

- $\#(\underline{\Omega}) < \infty$

cardinality / # of items in the set

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

$$X: \underline{\Omega} \to \mathcal{X} \qquad \omega \in \Omega \quad X(\omega)$$

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

Define

$X_1(w)$ = height of $w$ in cm (interval)

$X_2(w)$ = weight of $w$ in kg (interval)

$X_3(w)$ = gender of $w$ (categorical)

$$X(w) = \begin{pmatrix} X_1(w) \\ X_2(w) \\ X_3(w) \end{pmatrix}$$
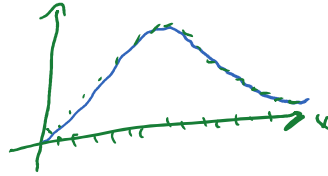
$X(w)$ is 3d measurements

$$X: \Omega \to \mathbb{R}^1 \times \mathbb{R}^1 \times \{M, F\}$$
$\underbrace{\qquad}_{\mathcal{X}}$

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

- $\Omega$ and $X$ define relative frequency distribution over $\mathcal{X}$.

$$f_X(x) = \frac{\# \{w: X(w) = x\}}{\#(\Omega)}$$

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

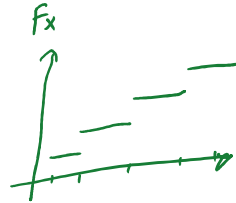# simplify by introducing continuous approximation

# discrete is too hard for us! approx it.

- note: i) $0 \leq f_X(x) \leq 1 \quad \forall x \in \mathcal{X}$

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)

$$F_X(x) = \frac{\# \{w: X(w) \leq x\}}{\#(\Omega)} = \text{cumulative distributive function of } X \text{ (CDF of } X\text{)}$$
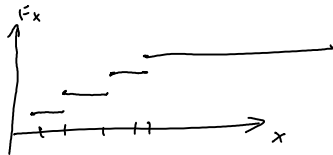
$$= \sum_{z \leq x} f_X(z)$$

Step function

$$f_X(x) = F_X(x) - F_X(x-\epsilon) \quad , \text{ where } F_X(x-\epsilon) = \lim_{z \uparrow x} F_X(z)$$
$\qquad\qquad\qquad\underbrace{\quad}_{\text{epsilon}}$

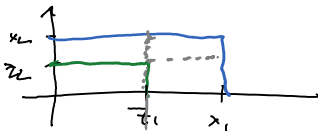- So $F_X$ and $f_X$ are two equivalent ways of presenting a frequency distribution.

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

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

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

$$f_X(x_1, x_2) = \lim_{z_1 \uparrow x_1} \left[ F_X(x_1, x_2) - F_X(x_1, z_2) - F_X(z_1, x_2) + F_X(z_1, z_2) \right]$$

So, $F_X \iff f_X$

- The whole point of any statistical analysis is to learn something about $F_X$.

- how do we do this?

- If possible we do a census, namely compute $X(w) \quad \forall w \in \Omega$ of the form $f_X$.

- Typically count (return to this in a moment)

- why do we want to know $F_X$?

eg 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.

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

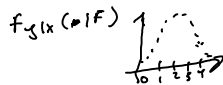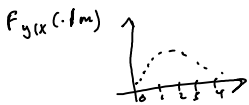$$= \frac{f_{(x,y)}(x, y_2)}{f_x(z)}$$

Definition: $x$ 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 though its mean as $x$ changes, $E(y|x)(v)$

$$\sum y(v)$$

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

$$\frac{Ex}{} = \sum_{y} y \, f_{y|x}(y|x)$$