# ECON 57 - Lecture 1
## A Brief History of Probability and Statistics

Augusto Gonzalez-Bonorino

Pomona College

Fall 2023

# Table of Contents

## Introduction

- Studying the history of any subject helps reveal details hidden by time, biases, assumptions, and set the context for the topics to be studied. It can help us avoid misinterpretations, identify long-term trends or interdisciplinary connections, and find inspiration to innovate.
- Logic and reasoning can then be employed to assess if:
  - The assumptions are still valid.
  - The biases are not strong enough to invalidate the methods.

# Collecting Experiences

1. First evidence of data collection is from 19000 BC
2. Data analysis as a concept has existed since 1663
3. Think of data as a collection of experiences or observations
4. Statistics emerged from previous efforts to make sense of these experiences.

# Eugenics and Early Statistics

- Francis Galton's pioneering role
- Birth of eugenics and its implications
- Statistical techniques developed by Galton:
  - Standard deviation
  - Correlation
  - Linear Regression
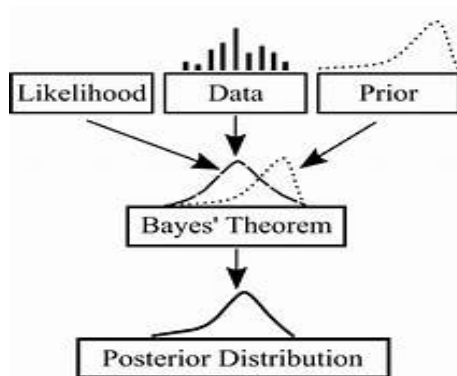- Ethical considerations and lessons learned

# Ronald Fisher, Karl Pearson, and the Frequentist School

- Ronald Fisher and Karl Pearson's contributions
- Introduction of statistical significance and unbiased authority
- Correlation as a measure of causation
- Limitations of inferring causation from correlation

# Emergence of Bayesian School of Thought

- Introduction to the Bayesian school of thought
- Thomas Bayes' theorem and incorporating prior beliefs
- Computational challenges and adoption of Bayesian methods
- Balancing uncertainty and complex analysis

# Wright's Opposition

- American geneticist and statistician from the '20s.
- Proposed a new methodology for studying causal relationships
    - He accounted for the direction of the relationship, unlike correlation which as a general measure of association.
- His methodology culminated in the development of path diagrams

# Path Diagrams

- Path diagrams are a reasoning tool. To help you break down relationships into relevant factors and map the direction of these relationshis.
- Moving beyond simplistic correlations.
- Help guide the formulation and testing of causal hypothesis.
- Fundamental tool in Structural Equation Modeling (SEM) and other causal inference techniques widely used by econometricians.
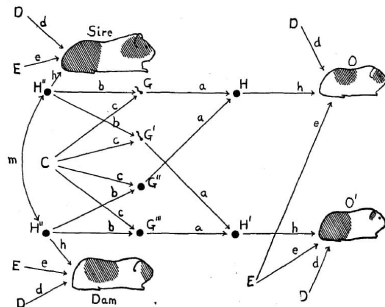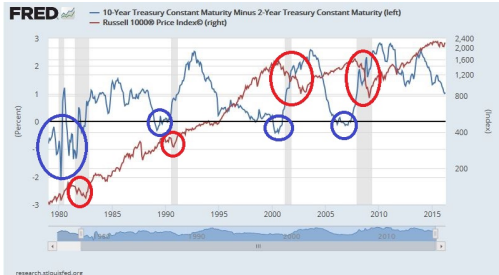


FIGURE 2.—A diagram illustrating the relations between two mated individuals and their progeny. $H$, $H'$, $H''$ and $H'''$ are the genetic constitutions of the four individuals. $G$, $G'$, $G''$ and $G'''$ are four germ-cells. $E$ and $D$ represent tangible external conditions and chance irregularities as factors in development. $C$ represents chance at segregation as a factor in determining the composition of the germ-cells. Path coefficients are represented by small letters.
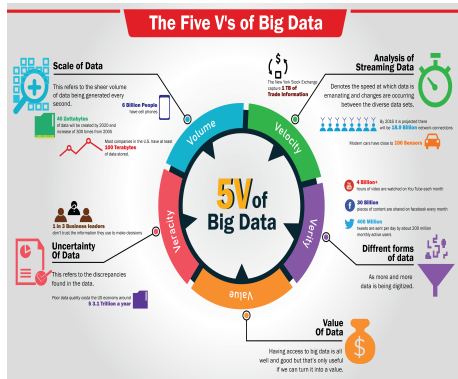
# Correlation is not Causation

- Shifting from relying solely on correlation
- Highlighting confounding variables and chance
- Recognizing the complexity of cause-and-effect relationships
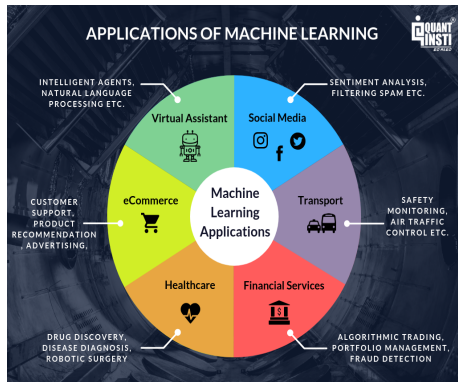- Nevertheless, correlations have taken us a long way.

- Transformation with computing power and Big Data
- Rise of data science as an interdisciplinary field
- Integration of machine learning algorithms
- Extracting insights from large datasets

# Machine Learning and Neural Networks

- ML combines statistics with algorithms with the goal of creating reasoning engines.
- Machine learning's impact on diverse domains
- Machine learning vs Neural networks vs Deep learning
- Sparks of AGI: computer vision and NLP
- Python and R



APPLICATIONS OF MACHINE LEARNING

# Judea Pearl and Bayesian Networks

- Judea Pearl's contribution to causal inference
- Introducing Bayesian networks for complex systems
- Incorporating probability theory and causal reasoning
- Applications in artificial intelligence and epidemiology

**Basic Concepts of Statistics**

## Populations and Samples

In statistics, we often deal with two main concepts: populations and samples.

**Population:** The entire group of individuals or items that we are interested in studying.

**Sample:** A subset of the population, carefully selected to represent the entire population.

For example, if we want to study the average income of all households in a country, the *population* would be all households in that country, and a *sample* would be a smaller group of households from different regions.

# Parameter vs Statistic

In statistics, we use both parameters and statistics to describe data:

- **Parameter:** A numerical measure that describes a characteristic of a population. For instance, the average income of all households in a country is a population parameter.
- **Statistic:** A numerical measure that describes a characteristic of a sample. The average income of a sample of households from different regions is a sample statistic.

In practice, we often use statistics to estimate parameters. This involves collecting data from a sample and using sample statistics to make educated guesses about population parameters.

# Random Sampling

Random sampling helps us representative samples:

- **Random Sample:** A sample selected from a population in such a way that each individual or item has an equal chance of being included in the sample. This helps to minimize bias and ensures that the sample is representative of the population.

- **Simple Random Sampling:** A type of random sampling where each possible sample of a given size has an equal chance of being selected. This is typically achieved using random number generators.

- **Stratified Sampling:** Dividing the population into subgroups (strata) based on some characteristic, and then selecting samples from each stratum. This ensures that each subgroup is represented in the sample.

- **Cluster Sampling:** Dividing the population into clusters, and then randomly selecting entire clusters to include in the sample. This can be more practical when the population is large.

## Importance of Sampling in Economics

- **Resource Efficiency:** Sampling yields accurate insights using fewer resources than surveying the entire population.
- **Representativeness:** Well-designed samples mirror population characteristics, enhancing result reliability.
- **Inference:** Accurate inferences about the entire population are drawn from well-conducted samples, vital for policymaking.
- **Reducing Bias:** Sampling mitigates bias risks from over/underrepresentation in analysis.
- **Feasibility:** Sampling offers a practical approach when surveying the entire population is impossible.
- **Time Sensitivity:** Sampling enables frequent data collection for swift response to changing economic trends.
- **Statistical Analysis:** Sampling provides manageable datasets for complex economic modeling.

# Statistical Data Types

Data can be classified into two main types:

1. **Categorical:** Data points that can be grouped into categories. For instance, food groups or gender groups or economic indicators (a special category of general statistical indicators that are related to economics). Remember the XOR problem from last class?

2. **Numerical:** Data points that are solely represented by numbers and do not belong to a group or category.
    1. **Discrete:** Data points that take distinct, separate values. For example, the number of students in a class, or the count of cars passing through a toll booth in a given time period.
    2. **Continuous:** Data points that can take any value within a range. These values are not limited to specific, separate points. Examples include height, weight, or temperature measurements.

# R Data Types

1. **character/string:** The character data type in R is used to represent text, such as words, sentences, or any sequence of characters. In R, strings are enclosed within single or double quotes. For example, "Hello, World!" and 'R Programming' are both examples of character data.

2. **integer:** Integers are used to represent whole numbers in R without any decimal points. They can be either positive or negative.

3. **float:** Float numbers are used to represent decimals in R. They differ from integers in memory requirements, which must be left "floating" to account for the fractional parts.

4. **logical/boolean:** Logical data types in R represent binary values, either TRUE or FALSE, and are used for making decisions and logical operations. They are commonly used in conditional statements and filtering data.

# R Data Structures
Unidimensional data structures

1. **vector:** A vector is a one-dimensional data structure in R that can hold elements of the same data type. It can be created using the 'c()' function, and elements can be accessed using numeric indices. Vectors are fundamental in R and are widely used for data manipulation and calculations.

2. **factor:** Factors are used to represent categorical data in R. They are particularly useful when dealing with data that have predefined categories or levels. Factors are created using the 'factor()' function, and each level represents a distinct category. Factors play a significant role in statistical modeling and are used for tasks such as creating bar charts and performing categorical data analysis.

## More on numerical data
Qualitative vs Quantitative

We will be focusing primarily on numerical data in this course. In general, numerical data and variables can be of two broader categories: Qualitative and Quantitative

1. **Qualitative:** Qualitative data are non-numeric in nature and represent categories or attributes. They convey qualities or characteristics that cannot be measured using numbers alone, they simply represent different labels.
   - There is no inherent meaning in the difference of numbers. This means that comparing two qualitative numbers yields little to no actionable information. For instance, economic regions.

2. **Quantitative:** There is a difference, thus quantitative variables carry information. These values can be measured, compared, and subjected to mathematical operations. For instance, the height of the players in a basketball team (i.e., x is taller than y) or inflation and GDP growth rates.

# More on numerical data
## Scales of Measurement

- **Qualitative:**
  1. **Nominal Data:** Categories without any inherent order, such as colors or names.
  2. **Ordinal Data:** Categories with a natural order, but the differences between them are not well-defined, such as rankings.
- **Quantitative:**
  1. **Interval Data:** Numerical data with a consistent interval between values, but no true zero point, such as temperature in Celsius.
  2. **Ratio Data:** Numerical data with a consistent interval between values and a true zero point, such as height or weight.

## Conceptual Exercises

1. You are analyzing the data related to countries' GDP (Gross Domestic Product) values. Discuss whether GDP values are nominal, ordinal, interval, or ratio data. Explain your reasoning and provide examples to support your answer.

2. Consider a dataset of student exam scores. The scores are recorded as "A," "B," "C," "D," and "F." Determine the level of measurement for these scores (nominal, ordinal, interval, or ratio). Also, discuss whether arithmetic operations like addition and subtraction are meaningful for this dataset.

3. A dataset contains the heights of students in centimeters. Classify the heights as interval or ratio data. Explain the difference between these two levels of measurement and provide an example for each.

# R programming

**Installing RStudio**

Install R
R studio