

**ECE 3075A**  
**Random Signals**

**Lecture 18**  
**Introduction to Statistics & Sampling - I**

School of Electrical and Computer Engineering  
Georgia Institute of Technology  
Fall, 2003

**Probability Theory & Statistics**

- Probability Theory
  - A mathematical framework for dealing with uncertain aspects of the physical reality, the truth of which may otherwise take infinite effort to characterize without it. (Think about the example of radio noise that impairs reception, or how long two light bulbs connected in series or in parallel will continue to produce light.)
  - The framework is built upon set theory and set operations (Borel field and  $\sigma$ -algebra), and when extended to random variables it is practically supported by the Reimann integral (from the more general Lebesgue integral).

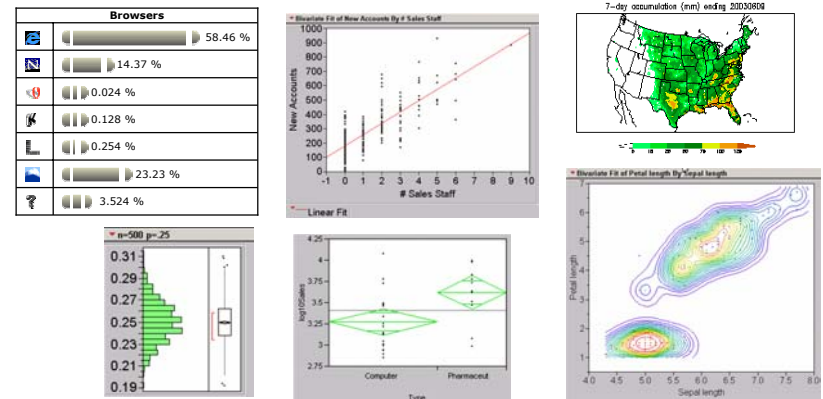
$$\int_{\Lambda(x)} g(x) d\mu(x) \rightarrow \int_{\Lambda(x)} g(x) f(x) dx$$

**Probability Theory & Statistics**

- Statistics
  - The science of assembling, sorting, tabulating, and analyzing **data** or a set of facts, in order to find the relationship or implicit regularity among them for other use.
    - Descriptive statistics: collecting, grouping, and presenting data for easy understanding and assimilation – tie to instinct and perception (or contrast to them)
    - Inductive Statistics: inferring some aspect of the truth about the origin or the environment where the data came from.

**Graphic Methods of Everyday Statistics**

- In statistics, we use data and some derived results to convey a message; graphic representation of the data or results often makes interpretation easy – a picture is worth a thousand words or millions of data items. Examples:



## Scope of Statistics Here

- **Sampling theory**: deals with problems in selecting data set that is manageably small but enough to allow meaningful inference
- **Estimation theory**: concerned with figuring out the value of key parameters from data
- **Hypothesis testing**: to verify if an assertion is true based on the evidence provided in the data set
- **Regression or curve fitting**: to find mathematical expressions that best represent or characterize the data
- **Analysis of variance**: to assess the significance of variation in data and how they relate to reality

## Sampling

- A **population** is a general collection of data whose statistically behavior is being studied. A population has size denoted by  $N$ ;  $N$  may be infinitely large.
- A **sample**, or random sample, is part of the population that has been selected at random for the study; its size, i.e., the sample size, denoted by  $n$ ,  $n \ll N$ , is the number of items or pieces of data selected for study. A sample is denoted by  $\{x_i\}_{i=1}^n$ .

The average of the data in the sample is called the **sample mean**.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

We can also view the sampling as a realization of an experiment involving random variables,  $X_1, X_2, X_3, \dots, X_n$ , which leads to a new r.v.

$$\hat{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad E[\hat{X}] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \bar{X}$$

## Sample Mean

- Therefore, the mean of the sample mean is the true mean of the population.
- We can use the sample mean as an estimate of the mean of the population; and since its expected value is identical to the true mean, it is an **unbiased estimator** – that is, it does not contain statistical deviation (for any sample, the deviation may be there, but the statistical average of the deviation is zero) from the true expected value.
- To determine if the sample mean is a good estimator, we need to evaluate its variance – to know how much the sample mean value will fluctuate. An estimator that fluctuate more is not as good as one that fluctuates less, provided that they're all unbiased.

## Variance of Sample Mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad E[\hat{X}] = \bar{X}$$

$$\text{var}(\hat{X}) = E\left[\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n X_i X_j\right] - (\bar{X})^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E[X_i X_j] - (\bar{X})^2$$

$$E[X_i X_j] = \begin{cases} \overline{X^2}, & i = j \\ \overline{X^2}, & i \neq j \end{cases} \quad \text{var}(\hat{X}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E[X_i X_j] - (\bar{X})^2$$

$$= \frac{1}{n^2} \left[ n \overline{X^2} + (n^2 - n) (\bar{X})^2 \right] - (\bar{X})^2 = \frac{\overline{X^2} - (\bar{X})^2}{n} = \frac{\sigma_x^2}{n}$$

That is, the fluctuation of the value of the sample mean is inverse proportional to the sample size; the more items included in the sample, the better the sample mean as an estimate of the mean of the population.

## Distribution of Sample Mean

$$\hat{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad E[\hat{X}] = \bar{X}, \quad \text{var}(\hat{X}) = \frac{\sigma_X^2}{n}$$

If  $n$  is large,  $f_{\hat{X}}(x) \approx N(\bar{X}, \frac{\sigma_X^2}{n})$  from the central limit theorem

Example: Let  $\bar{X} = 10$ , and  $\sigma_X^2 = 9$

We select a sample of size  $n$  to estimate  $\bar{X}$  and hope that the estimate has a standard deviation that is only 1% of the true mean.

$$\text{var}(\hat{X}) = \frac{\sigma_X^2}{n} = \frac{9}{n} = (0.01 \times 10)^2 = 0.01. \quad \text{Therefore, } n = 900$$

We can also ask the question: What is the probability that the sample mean is within 1% range of the true mean?

$$\begin{aligned} \Pr\{9.9 \leq \hat{X} \leq 10.1\} &= F_{\hat{X}}(10.1) - F_{\hat{X}}(9.9) = \Phi\left(\frac{10.1 - 10}{10 * 0.01}\right) - \Phi\left(\frac{9.9 - 10}{10 * 0.01}\right) \\ &= \Phi(1) - \Phi(-1) = 2\Phi(1) - 1 = 2 \times 0.8413 - 1 = 0.6826 \end{aligned}$$

## Weak Law & Strong Law of Large Numbers

Random variables  $X_1, X_2, X_3, \dots, X_n$  are identically distributed (with same mean and same finite variance, of course) and are at least pair-wise statistically independent, then, the sample mean satisfies

### Weak Law

$$\lim_{n \rightarrow \infty} \Pr\left\{ \left| \hat{X}_n - \bar{X} \right| < \varepsilon \right\} = 1 \quad \text{for any } \varepsilon > 0$$

### Strong Law

$$\Pr\left\{ \lim_{n \rightarrow \infty} \hat{X}_n = \bar{X} \right\} = 1$$

## Sampling without Replacement

- When population is large,  $N \rightarrow \infty$ , selecting sample with and without replacing the item back into the population does not seriously affect the statistical results.
- When population is not as large, sampling with and without replacement would lead to different statistical result. In particular, the variance of the sample mean is now

$$\text{var}(\hat{X}) = \frac{\sigma^2}{n} \left( \frac{N-n}{N-1} \right) \leftarrow \text{sampling without replacement}$$

as opposed to

$$\text{var}(\hat{X}) = \frac{\sigma^2}{n} \leftarrow \text{sampling with replacement or } N \rightarrow \infty$$