

Lecture - 20

Probability & Statistics

Use the tools of probability to solve statistics (problems).

The science of statistics deals with drawing conclusions from observed data.

Population: A large collection of items that have measurable values associated with them.
By suitably sampling from this collection, and then analyzing the sampled items, one hopes to draw some conclusions about the collection as a whole.

Sample/random sample: If x_1, x_2, \dots, x_n are independent random variables having a common distribution F , then we say that they constitute a sample or random sample from the distribution F .

① The population distribution F will not be completely specified and one will attempt to use the data to make inferences about F

— non parametric
~~interface for~~
inference problem.

② Sometimes, F is specified up to some unknown parameters.

For example, F is normal distribution having unknown mean and variance.

— parametric inference problem.

Big deal:

How to estimate the values of these parameters and how to verify that those are correct estimates.

statistic: A statistic is a random variable whose value is determined by the sample data.

Two important statistics are Sample mean and the sample variance.

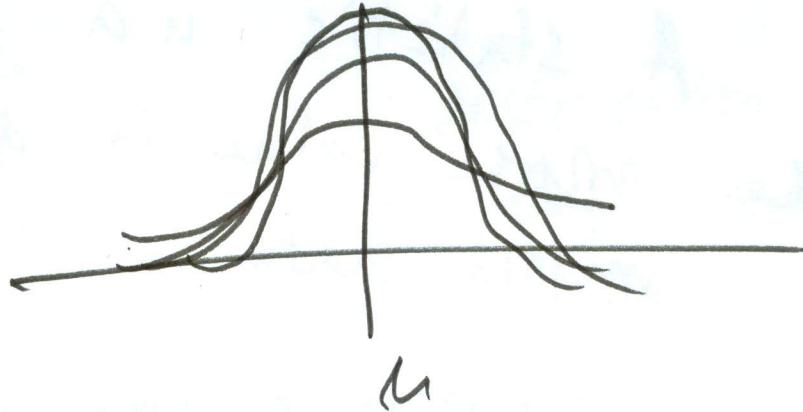
Sample mean: Let x_1, x_2, \dots, x_n be a sample of values from a population. Then the sample mean is defined by

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}.$$

Let x_i have mean μ and variance σ^2 , $i = 1, 2, \dots, n$. This μ is called the population mean and σ^2 is called the population variance.

Then $E[\bar{x}] = \frac{E[x_1] + E[x_2] + \dots + E[x_n]}{n} = \mu$

$$\text{var}(\bar{x}) = \frac{1}{n^2} \sum \text{var}(x_i) = \frac{\sigma^2}{n}$$



CLT

X_i 's are i.i.d

$$P \left\{ \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}} < x \right\} \approx \Phi(x)$$

as $n \rightarrow \infty$



$$P \left\{ \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < x \right\} \approx \Phi(x)$$



Here sample mean is μ and
variance is σ^2/n .

$\therefore \bar{X}$ will approximately be
normal when the sample size
 n is large.

Q. The weights of a population of workers have mean 167 and standard deviation 27.

If a sample of 36 workers is chosen, approximate the probability that the sample mean of their weights lie between 163 and 170.

Soln.

$$P\{163 < \bar{x} < 170\}$$

$$\begin{aligned} &\quad | \\ &\approx 2 P\{N(0,1) < 0.8889\} - 1 \end{aligned}$$

$$\approx 0.6259$$

Sample variance.

Let X_1, \dots, X_n be random sample from a distribution with mean μ and variance σ^2 .

Let \bar{X} be the sample mean. Then the statistic S^2 defined by

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

is called the sample variance as $S = \sqrt{S^2}$ is called the sample standard deviation.

$$\text{Notation} \quad \sum_{i=1}^n (x_i - \bar{x})^2, \quad \bar{x} = \frac{\sum x_i}{n}$$

$$= \sum_{i=1}^n \left(x_i^2 - 2\bar{x}x_i + \bar{x}^2 \right)$$

$$= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

Then it follows that

$$(n-1) s^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

$$\therefore (n-1) E[s^2] = E\left(\sum_{i=1}^n x_i^2\right) - n E[\bar{x}]^2$$

$$= n E(x_i^2) - n E[\bar{x}^2]$$

$$= n \text{Var}(x_i) + n(E[x_i]^2) - n \text{Var}(\bar{x}) - n(E[\bar{x}])^2$$

$$= n\sigma^2 + n\mu^2 - n(\sigma^2/n) - n\mu^2$$

$$= (n-1)\sigma^2$$

$$\therefore E[s^2] = \sigma^2$$

\therefore the expected value of the sample variance s^2 is equal to the population variance.