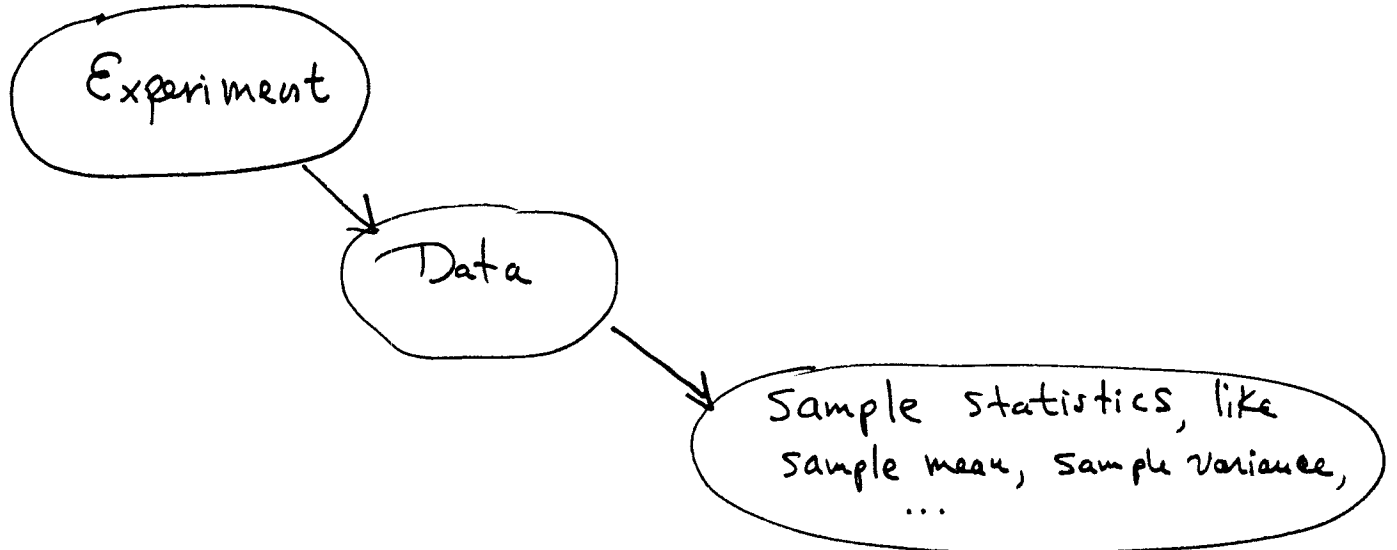
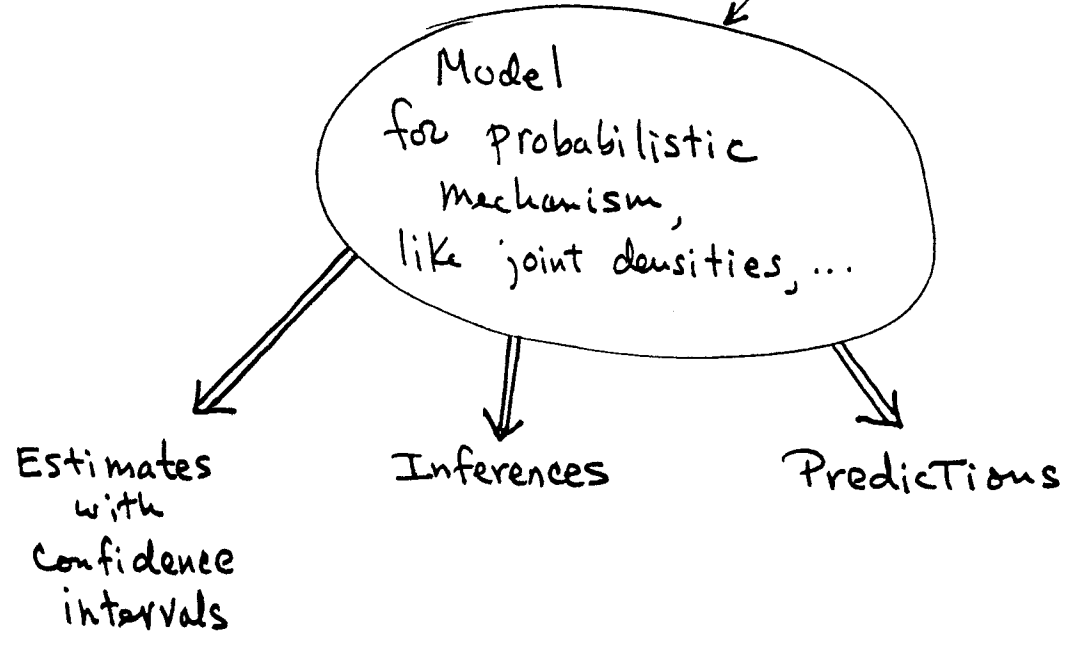


Descriptive Statistics



Inferential Statistics



Suppose x is a random variable with known prob. density fctn. $p(x)$.

MEAN

$$\mu = E(x) = \int_{-\infty}^{\infty} x p(x) dx \quad \left(\begin{array}{l} \text{discrete} \\ \sum x_i p_i \end{array} \right)$$

Variance

$$\begin{aligned} \sigma^2 &= E[(x-\mu)^2] \\ &= E[x^2] - 2E[x\mu] + E[\mu^2] \\ &= E[x^2] - \mu^2 \end{aligned}$$

σ is called standard deviation

Distinguish from sample mean, sample variance:

Suppose we have N independent observations of x , x_1, x_2, \dots, x_N

- Sample mean = $\frac{1}{N} \sum_{i=1}^N x_i = \bar{x}$ (a "Statistic"
a random variable)

$$E[\bar{x}] = \mu \quad \underline{\underline{\text{unbiased}}}$$

- Sample variance = $s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$
↑ sample mean

why divide by $N-1$? $\rightarrow E[s^2] = \sigma^2$ unbiased

A computational point:

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

Can compute by first finding \bar{x} , then direct evaluation.

But

$$\begin{aligned} s^2 &= \frac{1}{N-1} \left[\sum x_i^2 - 2\bar{x} \left[\sum x_i + \sum (\bar{x})^2 \right] \right] \\ &= \frac{1}{N-1} \left[\sum x_i^2 - 2(\bar{x})^2 N + N(\bar{x})^2 \right] \\ &= \frac{1}{N-1} \left[\sum x_i^2 - N(\bar{x})^2 \right] \end{aligned}$$

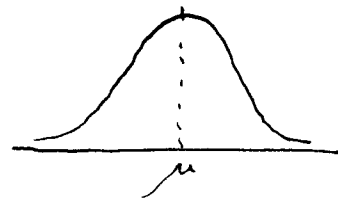
$$s^2 = \frac{\sum x_i^2 - N(\bar{x})^2}{N-1}$$

- fewer operations
- more accurate

Importance of Gaussian (Normal) Distribution

1.3.5

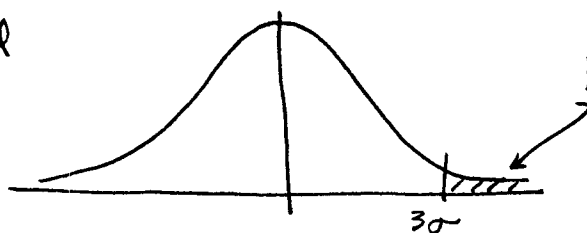
$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$



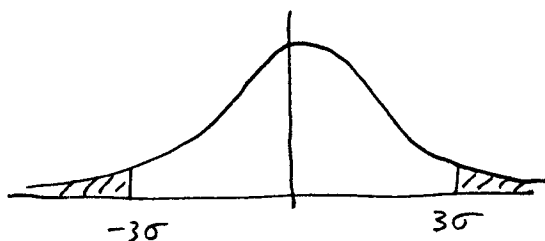
Often denoted $N(\mu, \sigma^2)$. $E(x) = \mu$ $\text{var}(x) = \sigma^2$

We often deal with $z =$ normalized Gaussian $N(0, 1)$

Tails are small

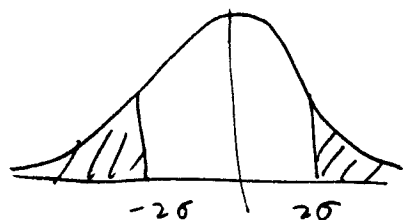


prob. = 0.00135



prob. of deviating
more than 3σ from
mean = 0.0027

prob. of not = 99.73%



& prob. of deviating less than $2\sigma = 95.45\%$

Why so Important?

Sum of independent observations converge to Gaussian under very general circumstances.

In nature, events that result from many small, independent effects tend to be Gaussian.

Central Limit Theorem

Suppose we sample x_1, \dots, x_n from a distribution with mean μ and variance σ^2 .

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

then

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \rightarrow N(0, 1) \quad \begin{array}{l} \text{Standard} \\ \text{Normal} \end{array}$$

Normalized
random variable

Notice this is true for any parent distribution.
(See probability books for technical details & conditions.)

Important Properties of Normal Distribution

1. Linear combination of Normals is also Normal
2. Normal has maximum entropy for given σ .
3. Least-squares becomes maximum likelihood
4. Many derived random variables have analytically known densities
5. Sample mean and variance of n identical, independent samples are independent; the sample mean is Normal

$$\bar{X}_n \sim N(\mu, \sigma/\sqrt{n})$$

Intuitively, the N differences $(x_i - \bar{x})$

are not independent, because $\sum_{i=1}^N (x_i - \bar{x}) = 0$

\therefore there are only $(N-1)$ degrees of freedom in statistic

A proof is straightforward but requires some algebra

Some details, then:

Lemma Suppose samples x_1, \dots, x_N come from distribution with mean μ & variance σ^2 , and are independent.

then $E(\bar{x} - \mu)^2 = \frac{\sigma^2}{N}$ \rightarrow shows that standard deviation decreases as $1/\sqrt{N}$

Proof algebra

$$\begin{aligned}
 \text{then } E[s^2] &= \cancel{\frac{1}{N}} E\left[\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2\right] \\
 &= \frac{1}{N-1} E\left[\underbrace{\sum x_i^2}_{\substack{\text{R1.32} \\ \downarrow}} - N \bar{x}^2\right] \\
 &= \frac{1}{N-1} \left[N(\sigma^2 + \mu^2) - N\left(\frac{\sigma^2}{N} + \mu^2\right) \right] \\
 &= \sigma^2 \quad \checkmark
 \end{aligned}$$

Summary of Distributions of random Variables derived from Normal

1.3.7

Let x_i be n indep. ident. distributed samples from $N(\mu, \sigma)$.

SAMPLE MEAN • $\bar{x}_n \triangleq \frac{1}{n} \sum_{i=1}^n x_i$ is distrib. as $N(\mu, \sigma/\sqrt{n})$

SAMPLE VARIANCE • $S_n^2 \triangleq \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$

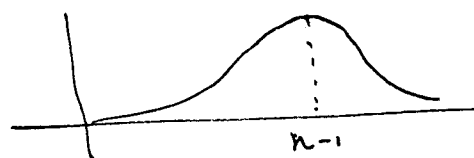
then $U = \frac{(n-1) S_n^2}{\sigma^2}$ (Normalized)

has a χ^2 -distribution with $(n-1)$ degrees of freedom:

$$p(\chi^2) = \left[2^{n/2} \Gamma\left(\frac{n}{2}\right) \right]^{-1} (\chi^2)^{\frac{n}{2}-1} e^{-\chi^2/2} \quad \chi^2 \geq 0$$

$$E[U] = n-1$$

$$\text{var}[U] = 2(n-1)$$

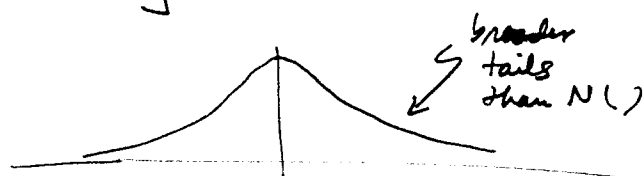


SAMPLE MEAN UNKNOWN VARIANCE • $\frac{\bar{x} - \mu}{S_n/\sqrt{n}}$ has a t -distribution

with $(n-1)$ degrees of freedom:

$$p(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \frac{1}{\sqrt{\pi n}} \frac{1}{\left[1 + t^2/n\right]^{\frac{n+1}{2}}}$$

"Student-t"
= W.S. Gosset



Confidence Intervals

1.3.8

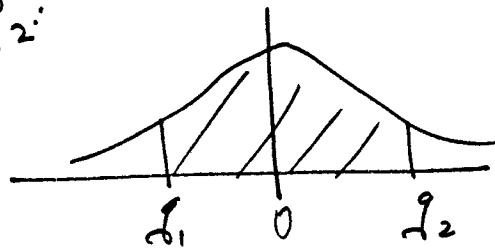
Assume again x_1, \dots, x_n i.i.d. Normal.

We want to know how far \bar{x}_n might be from μ .

We know

$\frac{\bar{x}_n - \mu}{S_n/\sqrt{n}}$ is student-t distributed.

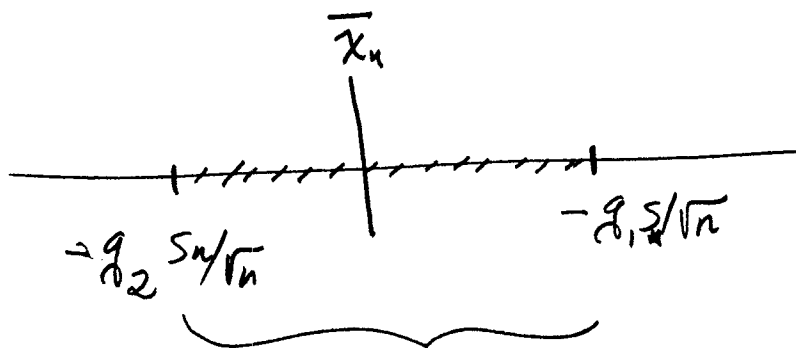
1. Pick $g_1, g_2 \ni$ student-t with $(n-1)$ degrees of freedom has ~~99%~~ 99% prob. (say) of lying between g_1, g_2 :



from TABLES

2. Prob. $\left\{ g_1 < \frac{\bar{x}_n - \mu}{S_n/\sqrt{n}} < g_2 \right\} = 0.99$

\Rightarrow Prob. $\left\{ \bar{x}_n - g_2 \frac{S_n}{\sqrt{n}} < \mu < \bar{x}_n - g_1 \frac{S_n}{\sqrt{n}} \right\} = 0.99$



"99% confidence interval"

\sim prob. μ is here is 99%

Confidence Interval for σ^2

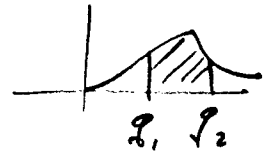
1.3.9

Again, assume X_1, \dots, X_n i.i.d. Normal.

We know $(n-1)S_n^2/\sigma^2$ is χ^2 -distr. with $n-1$ D.O.F.

1. Pick $g_1, g_2 \ni$

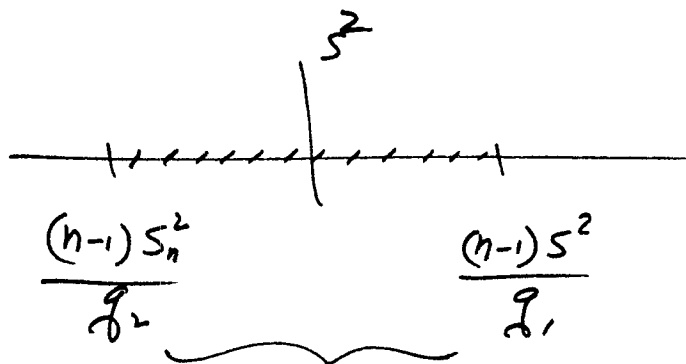
$$\text{prob. } \left\{ \chi^2 \text{ } n-1 \text{ D.O.F. } \right. \\ \left. \text{between } g_1, g_2 \right\} = 0.99$$



2. $\text{prob. } \left\{ g_1 < \frac{(n-1)S_n^2}{\sigma^2} < g_2 \right\} = 0.99$

TABLES

$$\Rightarrow \text{prob. } \left\{ \frac{(n-1)S_n^2}{g_2} < \sigma^2 < \frac{(n-1)S_n^2}{g_1} \right\} = 0.99$$



prob. that true σ is here is 99%

CAVEAT: Must be close to Normal to be valid

→ can often invoke Central Limit theorem to justify.