SySc 512, Statistical Inference Session 17,

[NetLogo Demo | Central Limit Theorem] (sampling)

Define: Random number (sample mean)
$$\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i \quad , \quad w/ \ n = \# \ samples$$

where the sequence of random variables $X_i$ are drawn from distrib w/ mean $= \mu$ & variance $= \sigma^2$ (but distrib is unknown).

Define: Sample variance (random number)
$$S_n^2 = \frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X}_n)^2$$

Note: $S_n^2$ is a "biased estimator" of $\sigma^2$, because 2 free parameter of the data are being estimated: popul mean & variance. The "unbiased estimator" is $S_{n-1}^2$.

Prove: $\mathbb{E}\,\bar{X}_n = \mu$
$$= \mathbb{E}\left[\frac{1}{n}(X_1 + X_2 + \cdots + X_n)\right]$$
$$= \frac{1}{n}\left[\mathbb{E}X_1 + \mathbb{E}X_2 + \cdots + \mathbb{E}X_n\right]$$
$$= \frac{1}{n}\left[\mu + \mu + \cdots + \mu\right] = \frac{1}{n}\left[n\mu\right] = \mu$$

Prove: $\mathbb{D}\,\bar{X}_n = \frac{1}{n}\sigma^2$
$$= \mathbb{D}\left[\frac{1}{n}(X_1 + X_2 + \cdots + X_n)\right]$$
$$= \left(\frac{1}{n}\right)^2\left[\mathbb{D}X_1 + \mathbb{D}X_2 + \cdots + \mathbb{D}X_n\right]$$
$$= \left(\frac{1}{n}\right)^2\left\{\sigma^2 + \sigma^2 + \cdots \sigma^2\right\} = \frac{1}{n}\sigma^2$$

NetLogo Demo | Central Limit Theorem | sample dist.

Distrib of sample mean: $\longrightarrow$ Normal distrib
    For any underlying distrib.

Normal Distrib: $N(\mu,\sigma)$
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}}\exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$
(Gaussian) $w/ \ \mu = \mathbb{E}S, \quad \sigma^2 = \mathbb{D}S$

Setup

Create Random People

and/or

Create My Own People

Preset 1    Preset 2

Preset 3    Preset 4
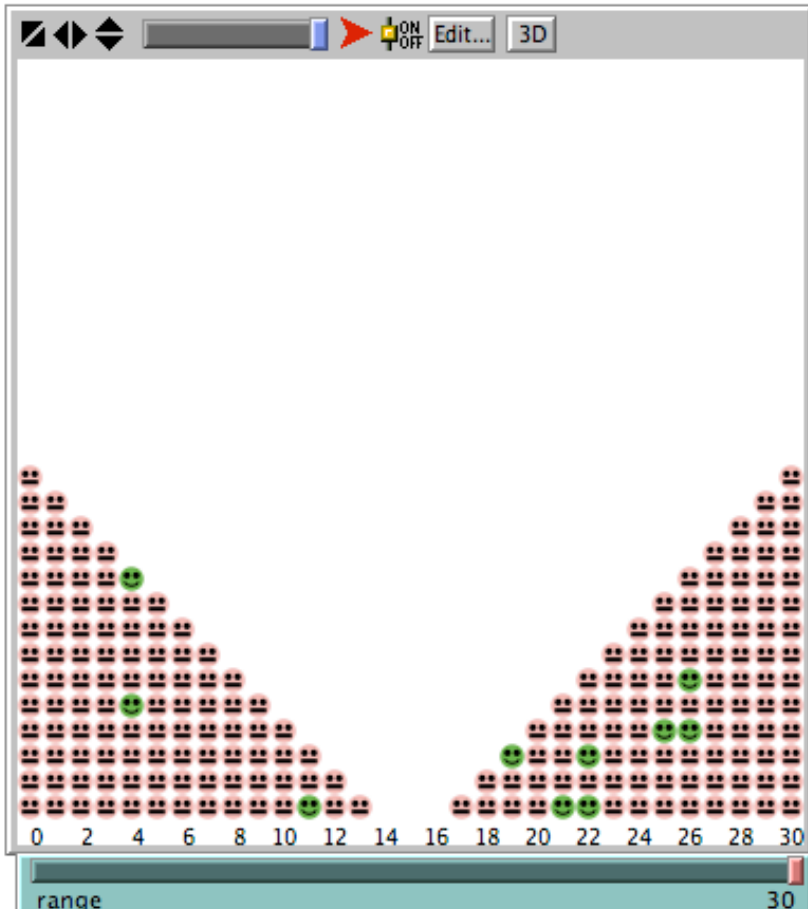
Preset 5    Preset 6

Sampling Commands:

sample-size          10

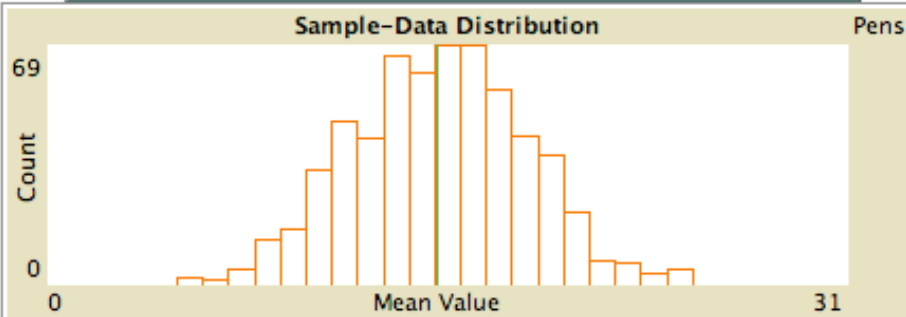On
Off    also-sums?

Go Once        Go

num-samples
600

std-dev-means
3.37

ON
OFF  Edit...  3D

0  2  4  6  8  10  12  14  16  18  20  22  24  26  28  30

range                                    30

Sample-Data Distribution                      Pens

69

Count

0

0            Mean Value                 31

Expected Value
15

On
Off    show-EV?

std-dev-sums
N/A

SySc 512, Session 17, Stat Inf (cont.)

Central Limit Theorem

("Distrib of $\bar{X}_n \to N(\mu, \sigma)$ as $n \to \infty$")

Or, more precisely:

Let $\bar{Z}_n = \dfrac{n \bar{X}_n - n\mu}{\sigma \sqrt{n}}$     (standardized)

If $\Phi(z)$ is the cumulative distrib of $N(0,1)$

$\Rightarrow \lim\limits_{n \to \infty} P\{\bar{Z}_n \le z\} = \Phi(z)$

Proof:

$1^{st}$ we need: Characteristic function

$$f_S(t) = \mathbb{E}\, e^{iSt} \quad , \quad t \in \mathbb{R}$$

Discrete: $f_S(t) = \sum\limits_{k=0}^{\infty} P_S(k) e^{ikt}$

Continuous: $f_S(t) = \int_{-\infty}^{\infty} P_S(x) e^{ixt} dx$

($f_S(t)$ is a Fourier transform of the distrib)

Taylor expand exponential:

$$f_S(t) = \mathbb{E}\, e^{iSt} = \mathbb{E}\left(1 + iSt - \tfrac{1}{2}S^2 t^2 + o(t^3)\right)$$
$$= 1 + it\mu - \tfrac{1}{2}t^2(\sigma^2 + \mu^2) + o(t^3)$$

$\underset{\uparrow}{}$ used: $\mathbb{D}S = \mathbb{E}S^2 - \mu^2 = \sigma^2$

(Aside: $\mu = -i \dfrac{d}{dt} f_S(t)\big|_{t=0}$ , $\sigma^2 = -\dfrac{d^2}{dt^2} f_S(t)\big|_{t=0} + \mu^2$)

For $N(0,1)$: $f_S(t) = 1 - \tfrac{1}{2}t^2 + o(t^3)$

Let $Y_i = \tfrac{1}{\sigma}(x_i - \mu)$

$\Rightarrow \bar{Z}_n = \sum\limits_{i=1}^{n} \dfrac{Y_i}{\sqrt{n}}$

$f_{\bar{Z}_n}(t) = \mathbb{E}\left(e^{i(Y_1 + Y_2 + \cdots + Y_n)t/\sqrt{n}}\right)$
$= \mathbb{E}\left(e^{iY_1 t/\sqrt{n}} e^{iY_2 t/\sqrt{n}} \cdots e^{iY_n t/\sqrt{n}}\right)$
$= \mathbb{E}\, e^{iY_1 t/\sqrt{n}} \cdots \mathbb{E}\, e^{iY_n t/\sqrt{n}}$
$= \left[f_Y\left(t/\sqrt{n}\right)\right]^n = \left[1 - \dfrac{t^2}{2n} + o(t^3/n)\right]^n \longrightarrow e^{-t^2/2}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ as $n \to \infty$

Same characteristic as $N(0,1)$

Weak Continuity Thm: Convergence of characteristic function $\Rightarrow$ convergence of distribution.

SySc 512, Session 17, Stat Inf (cont.)

Statistic: Value calculated from sample to characterize the population (mean, variance)

Purpose of statistics: Determine your confidence that your statistic ~~accurately~~ accurately characterizes the population.

Key: choose a statistic with a known distribution (indep. of population distrib)

Large-n:

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \text{ has a normal distrib: } N(0,1)$$

$$\frac{(n-1)\bar{S}_{n-1}^2}{\sigma^2} \text{ has a Chi-squared distrib: } \chi_{n-1}^2$$

$$\text{w/ } \chi_n^2(t) = \frac{1}{\Gamma(\frac{n}{2})}\left(\frac{1}{2}\right)^{\frac{n}{2}} t^{\frac{n}{2}-1} e^{-y/2}$$

$$\Gamma(n) = n! \, , \quad \Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$
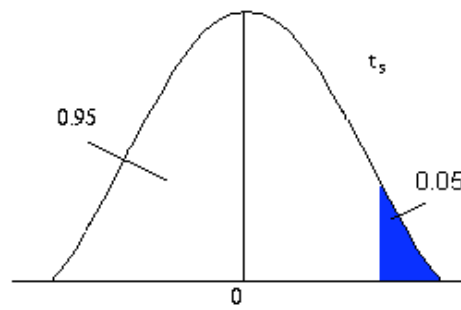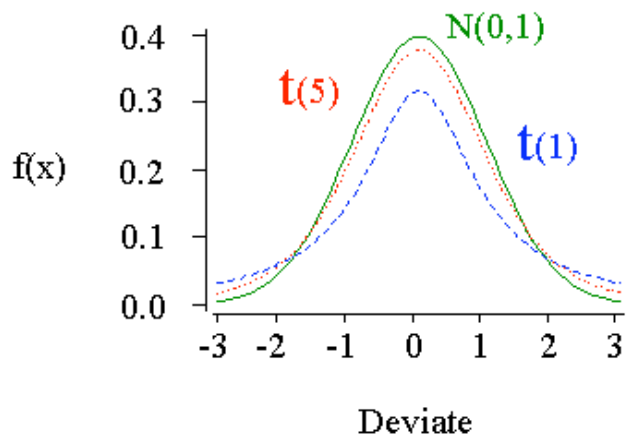
Small-n (n < 30)

Student's t-distr (Gossett)

$$\frac{\bar{X}_n - \mu}{\bar{S}_n/\sqrt{n}} = \frac{(\bar{X}_n - \mu)/(\sigma/\sqrt{n})}{\sqrt{(n-1)\bar{S}_n^2/(n-1)\sigma}} \, ; \text{ Distrib: } \frac{N(0,1)}{\sqrt{\frac{\chi_{n-1}^2}{n-1}}}$$

$$t\text{-distrib} \rightarrow N(0,1) \text{ for large } n$$

[ Figure: slide ]

## Student's t-distribution



P(t5 > 2.015) = 0.05
P(t20 > 1.725) = 0.05
P(t50 > 1.676) = 0.05
P(tinf > 1.645) = 0.05

Normal distr: $P(Z > 1.645) = 0.05$ where $Z \sim N(0, 1)$
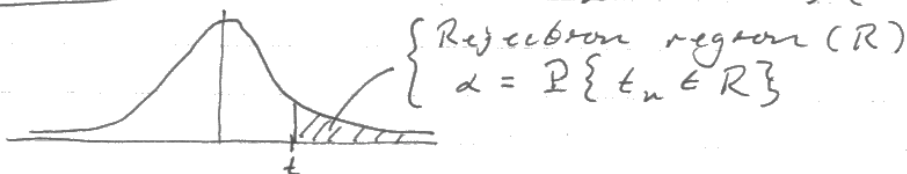
SySc 512, Session 17, Stat Inf (cont.)

Hypoth Testing

1) Choose null hypothesis: $H_0$
   ("not different") $\boxed{ex: \; \mathbb{E}\,\bar{x}_n = \mu}$

2) Choose alternative hypothesis: $H_1$
   ("something else") $\boxed{ex: \; \mathbb{E}\,\bar{x}_n \geq \mu}$

3) Choose significance level (confidence level)
   $\alpha$ = % risk of rejecting true $H_0$.
   $\boxed{ex: \; \alpha = 0.01}$

4) Calculate statistic:
   $\boxed{ex: \; \text{Student's } t\text{-test:} \; t_n = (\bar{x}_n - \mu)/(s_n/\sqrt{n})}$



$$\begin{cases} \text{Rejection region } (R) \\ \alpha = P\{t_n \in R\} \end{cases}$$

5) Get t-value from t-distrib w/ $n-1$ degrees of freedom:
   Let $\Phi_t(\nu) \equiv$ CDF of t-Distrib with $\nu$ deg. o.f.
   $\Rightarrow R = \{t_n \mid \Phi_t(\nu) > 1 - \alpha\}$

6) Reject $H_0$ if
   $$t_n > \Phi_t^{-1}(1 - \alpha, \, n-1)$$

$\boxed{\begin{aligned} &\underline{Reported}: p\text{-value (observed significance level)} \\ &\quad -\text{Smallest fixed level at which } H_0 \text{ is rejected} \\ &\quad p = P\{t_n > t\} = 1 - P\{t_n < t\} = 1 - \Phi_t(t_n, \nu) \end{aligned}}$

7) Interpretation: $\mathbb{E}\,\bar{x}_n$ is not significantly different from $\mu$ for signif level $\alpha$.

$\boxed{Matlab\ Demo}\; \boxed{hypothTest.m \quad w/ \; stixbox}$