"around" the
Normal Distribution

Sums of normals:

If
$$Y_i \sim N(\mu_i, \sigma_i^2), \quad 1 \le i \le n,$$

then
$$S = \sum_{i=1}^{n} Y_i \sim N\left(\sum_{i=1}^{n} \mu_i, \sum_{i=1}^{n} \sigma_i^2\right).$$

Sum of normals is normal

## Log-Normal random variables

If $X \sim N(\mu, \sigma^2)$, then
$$Y = e^X \sim \ln N(\mu, \sigma^2)$$

and conversely if $Y \sim \ln N(\mu, \sigma^2)$, then
$$X = \ln(Y) \sim N(\mu, \sigma^2).$$

Note $-\infty < X < \infty$ and $0 < Y < \infty$.

If
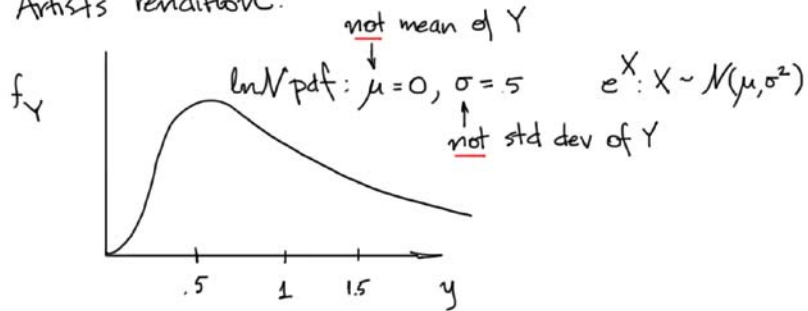$$Q = \prod_{i=1}^{m} Y_i, \quad Y_i \sim \ln N,$$

then $Q$ is also $\ln N$:

$$\ln Q = \sum_{i=1}^{m} \underbrace{\ln Y_i}_{normal} \implies \ln Q \sim N$$
$$\Downarrow$$
$$Q = e^{\ln Q} \sim \ln N$$

1

Artist's rendition:

$f_Y$

ln$N$ pdf: $\mu = 0$, $\sigma = 5$     $e^X$: $X \sim N(\mu, \sigma^2)$

not mean of Y

not std dev of Y

(axis labels: .5, 1, 1.5, y)

Note if $\mu \gg \sigma$ then ln$N \sim N$.

AT Patera      2.086 Prob_Stats      February 28, 2013    5

---

Pomelo estimation:

$$Q = \rho \cdot \frac{4\pi}{3} \cdot a \cdot b \cdot c$$

mass   density     principle axes of ellipsoid

if $\rho, a, b, c$ ln$N$

ln$N \longleftarrow$ then

$n = 53$

To "test": compare histogram of

$$\frac{\ln Q - \mu_{est}}{\sigma_{est}}$$

to histograms from standard normal $Z \sim N(0, 1)$.

AT Patera      2.086 Prob_Stats      February 28, 2013    6

---



AT Patera      2.086 Prob_Stats      February 28, 2013    7

---

Central Limit Theorem       (one version)

Let

$X_i$, $1 \le i \le n$, be i.i.d. r.v. $\sim f_X$
with mean $\mu$ and variance $\sigma^2$,

and

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ be the sample mean;

then as $n \to \infty$

standard normal cdf

$$P\left( \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \le z \right) \to \Phi(z)$$

AT Patera      2.086 Prob_Stats      February 28, 2013    8

**Slide 9:**

Example: Bernoulli                    binomial

Let

$X_i$, $1 \le i \le n$, be i.i.d. $\sim f_X^{Bernoulli}(x; \theta)$

$\Rightarrow$ mean $\theta$ and variance $\theta(1-\theta)$,
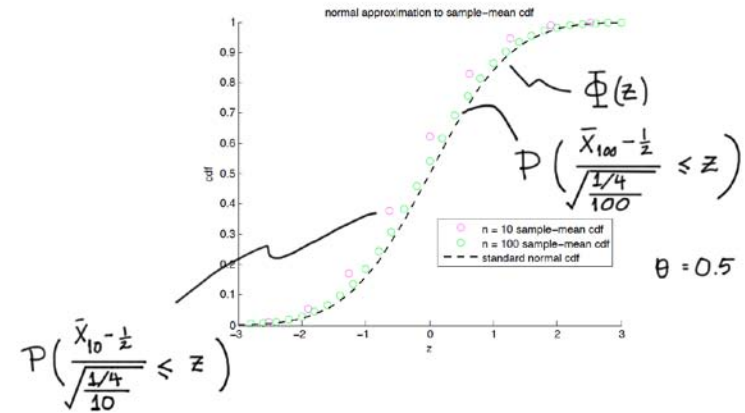
and

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ (fraction heads) be the sample mean,

then as $n \to \infty$

$$P\left( \frac{\bar{X}_n - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \le z \right) \to \Phi(z).$$

**Slide 10:**

"2.086" accuracy criterion: $n\theta > 5$ AND $n(1-\theta) > 5$



normal approximation to sample-mean cdf

$\Phi(z)$

$P\left( \dfrac{\bar{X}_{100} - \frac{1}{2}}{\sqrt{\frac{1/4}{100}}} \le z \right)$

n = 10 sample-mean cdf
n = 100 sample-mean cdf
standard normal cdf

$\theta = 0.5$

$P\left( \dfrac{\bar{X}_{10} - \frac{1}{2}}{\sqrt{\frac{1/4}{10}}} \le z \right)$

$\approx$

**Slide 11:**

Estimation: Bernoulli

**Slide 12:**

Estimator for $\theta$:

Recall if

$X_i$, $1 \le i \le n$, are i.i.d. $f_X^{Bernoulli}(x; \theta)$

and

$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ is the sample mean,

then

$$\mathbb{E}(\bar{X}_n) = \theta \quad \text{and} \quad \mathbb{E}\left((\bar{X}_n - \theta)^2\right) = \frac{\theta(1-\theta)}{n}.$$

↑ estimator for $\theta$      ↑ "good" estimator for $\theta$

Thus define

$$\hat{\Theta}_n \equiv \bar{X}_n \left( = \frac{1}{n} \sum_{i=1}^{n} X_i \right)$$

as our <u>estimator</u> for $\theta$, and

$$\hat{\theta}_n \equiv \bar{x}_n \left( = \frac{1}{n} \sum_{i=1}^{n} x_i \right)$$

as our <u>estimate</u> for $\theta$.

Note:

realization of

$\theta$ : parameter ;     $\hat{\Theta}_n$ : estimator ;     $\hat{\theta}_n$ : estimate .
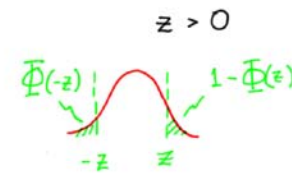deterministic     random variable     number

---

Confidence Interval :          2-sided, normal approximation

Recall for large $n$          $z > 0$

$$P \left( \frac{\hat{\Theta}_n - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} < z \right) \approx \Phi(z)$$

$\Phi(-z)$          $1-\Phi(z)$
$-z$     $z$

so

$$P \left( -z \leq \frac{\hat{\Theta}_n - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \leq z \right) \approx \Phi(z) - \Phi(-z)$$

$$= \Phi(z) - (1 - \Phi(z))$$

$$= 2\Phi(z) - 1 .$$

---

Choose

$$2\Phi(z_\gamma) - 1 = \gamma          \text{confidence level}$$

hence

$$\Phi(z_\gamma) = (1 + \gamma)/2$$

or

$$z_\gamma = \tilde{z}_{(1+\gamma)/2} \quad \left( \frac{1+\gamma}{2} \text{ quantile of } \Phi \right).$$

Note $\gamma = 0 \Rightarrow z_\gamma = 0$ (median), and
$\gamma \to 1 \Rightarrow z_\gamma \to \infty$, and
$\gamma = 0.95 \Rightarrow z_\gamma \approx 1.96$ .

---

Thus,

$$P \left( -z_\gamma \leq \frac{\hat{\Theta} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \leq z_\gamma \right) = \gamma$$

$\Updownarrow$

$$-z_\gamma \leq \frac{\hat{\Theta} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \Rightarrow \theta \leq \hat{\Theta} + z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}}$$

$$\frac{\hat{\Theta}_n - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \leq z_\gamma \Rightarrow \hat{\Theta}_n - z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}} \leq \theta$$

or $\hat{\Theta}_n - z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}} \leq \theta \leq \hat{\Theta}_n + z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}}$ .

with probability $\gamma$

Define
$$[CI]_n^0 \equiv \left[ \hat{\Theta}_n - z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}} \;,\; \hat{\Theta}_n + z_\gamma \sqrt{\frac{\theta(1-\theta)}{n}} \right]$$

and then (since $\theta$ unknown)                    $n$ large
$$[CI]_n \equiv \left[ \hat{\Theta}_n - z_\gamma \sqrt{\frac{\hat{\Theta}_n(1-\hat{\Theta}_n)}{n}} \;,\; \hat{\Theta}_n + z_\gamma \sqrt{\frac{\hat{\Theta}_n(1-\hat{\Theta}_n)}{n}} \right] .$$

Hence
$$P\left( \theta \text{ is inside } [CI]_n \right) = \gamma .    \approx$$

Note: $\theta$ is a <u>deterministic</u> parameter, whereas
$[CI]_n$ is a <u>random</u> interval.

---

In practice:  choose $n, \gamma$          $\Rightarrow z_\gamma$

sample $x_1, x_2, \ldots, x_n$ ;                    0 or 1

compute estimate for $\theta$,
$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^{n} x_i \;;$$

compute $[ci]_n$,
$$[ci]_n = \left[ \hat{\theta}_n - z_\gamma \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \;,\; \hat{\theta}_n + z_\gamma \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \right] ,$$

$$\text{HalfLength}_{\theta;n} = z_\gamma \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} ,$$

$$\text{RelErr}_{\theta;n} = z_\gamma \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \Big/ \hat{\theta}_n = z_\gamma \sqrt{\frac{1-\hat{\theta}_n}{\hat{\theta}_n \cdot n}} .$$

---

Note:

as $\gamma \to 1$, $z_\gamma \to \infty$, and $\text{HalfLength}_{\theta;n} \to \infty$

    more confidence $\Rightarrow$ less accuracy ;

as $n \to \infty$, $\text{HalfLength}_{\theta;n} \to 0$    but SLOWLY

    more samples $\Rightarrow$ more accuracy ,

as $\theta(\hat{\theta}_n) \to 0$, $\text{RelErr}_{\theta;n} \to \infty$    fixed $n$

    rare event   $\Rightarrow$  less accuracy . †

(Also require $n\hat{\theta} > 5$ AND $n(1-\hat{\theta})$ for normal approximation.)

5

---

Frequentist interpretation:



$\gamma = 0.8$

$\theta = .5$

$[ci]_n$

$\gamma = 0.95$

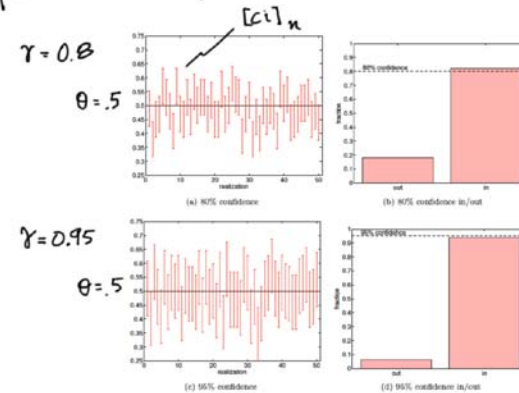$\theta = .5$

50 experiments
each
with $n = 100$
coin flips

Figure 10.3: An example of confidence intervals for estimating the mean of a Bernoulli random variable ($\theta = 0.5$) using 100 samples.

$\theta = 0.5$   $n = 100$

Some Applications
of
Bernoulli Estimation

The Birthmonth "Distribution"

Let

$$X = \begin{cases} 0 \equiv \text{birthmonth [Jan-June]} & \text{probability } 1-\theta \\ 1 \equiv \text{birthmonth [July-Dec]} & \text{probability } \theta \end{cases}$$

Choose

$n = 51$   (2086 class size)

$\gamma = 0.95$   (confidence level)

Collect data: $n = 51$

birthmonth_number $(i)$, $1 \le i \le 12$,

is #(occurrences of birthmonth $i$);

note sum(birthmonth_number) = 51.

Compute estimate for $\theta$:

$\hat{\theta}_{n=51}$ = sum(birthmonth_number([7:12]))/51

= 26/51 = .5098

$\left( = \sum_{i=1}^{n} x_i / n \right)$.     $n\hat{\theta}_n = 27.05 \; (>5) \; \checkmark$

$n(1-\hat{\theta}_n) = 25.98 \; (>5) \; \checkmark$

Calculate confidence interval for $\theta$:     $z_{0.95} = 1.96$

$$1.96 \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} = .1372$$

$\Downarrow$

$$[CI]_n = \left[ .5098 - .1372, \; .5098 + .1372 \right]$$

$$= \left[ .3726, \; .6470 \right]$$

Conclusion: If we are amongst the
lucky 9,500/10,000 parallel universes,

$$\boxed{.3726 \le \theta \le .6470}$$

(and no reason to reject hypothesis $\theta = \frac{1}{2}$).

6

Other applications (same game)

Quality control:

$$X = \begin{cases} 0 & \text{part not up to spec} \quad \text{probability } 1-\theta \\ 1 & \text{part up to spec} \quad \text{probability } \theta \end{cases}$$

$n \equiv \#(\text{parts})$ from large population inspected

$\hat{\theta}_n \equiv$ fraction of parts up to spec.

Failure:

$$X = \begin{cases} 0 & \text{strut } \sigma \le \sigma_{max} \quad \text{probability } 1-\theta \\ 1 & \text{strut } \sigma > \sigma_{max} \quad \text{probability } \theta \; (\ll 1) \end{cases}$$

$n \equiv \#(\text{struts})$ from large population inspected

$\hat{\theta}_n : \#(\text{cars in fleet}) \equiv \#(\text{repairs})$.

Preferences:

two choices: A, B
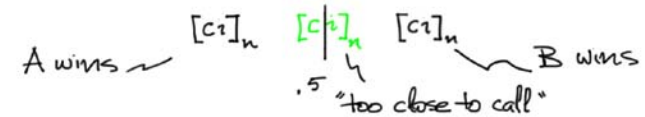└ competing candidates, or products, or...

$$X = \begin{cases} 0 & \text{prefer A} \quad \text{probability } 1-\theta \\ 1 & \text{prefer B} \quad \text{probability } \theta \end{cases}$$

$n$: number of voters in survey sample (focus group)

$\hat{\theta}_n$: fraction of voters who prefer B

Note if popular (simple majority) election,

$\gamma$

A wins $\sim$ $[c_i]_n$ $[c_i^+]_n$ $[c_i]_n$ $\sim$ B wins

.5 "too close to call"

2.086 Numerical Computation for Mechanical Engineers
Spring 2013