# Induction Course in Quantitative Methods for Finance
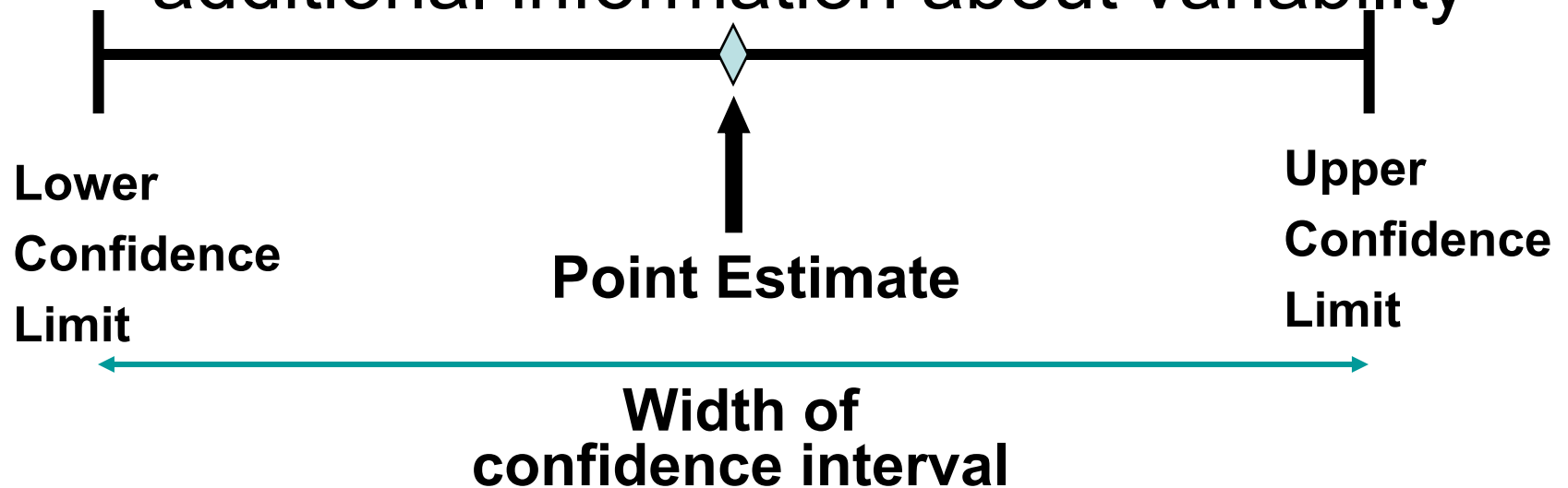
## Estimation: Single Population

# Definitions

- An estimator of a population parameter is
  - a random variable that depends on sample information . . .
  - whose value provides an approximation to this unknown parameter

- A specific value of that random variable is called an estimate

# Point and Interval Estimates

- A point estimate is a single number,

- a confidence interval provides additional information about variability

**Lower Confidence Limit**

**Point Estimate**

**Upper Confidence Limit**

**Width of confidence interval**

# Point Estimates

| We can estimate a Population Parameter … | | with a Sample Statistic (a Point Estimate) |
|:---:|:---:|:---:|
| Mean | $\mu$ | $\overline{x}$ |
| Proportion | $P$ | $\hat{p}$ |

# Unbiasedness

- A point estimator $\hat{\theta}$ is said to be an unbiased estimator of the parameter $\theta$ if the expected value, or mean, of the sampling distribution of $\hat{\theta}$ is $\theta$,
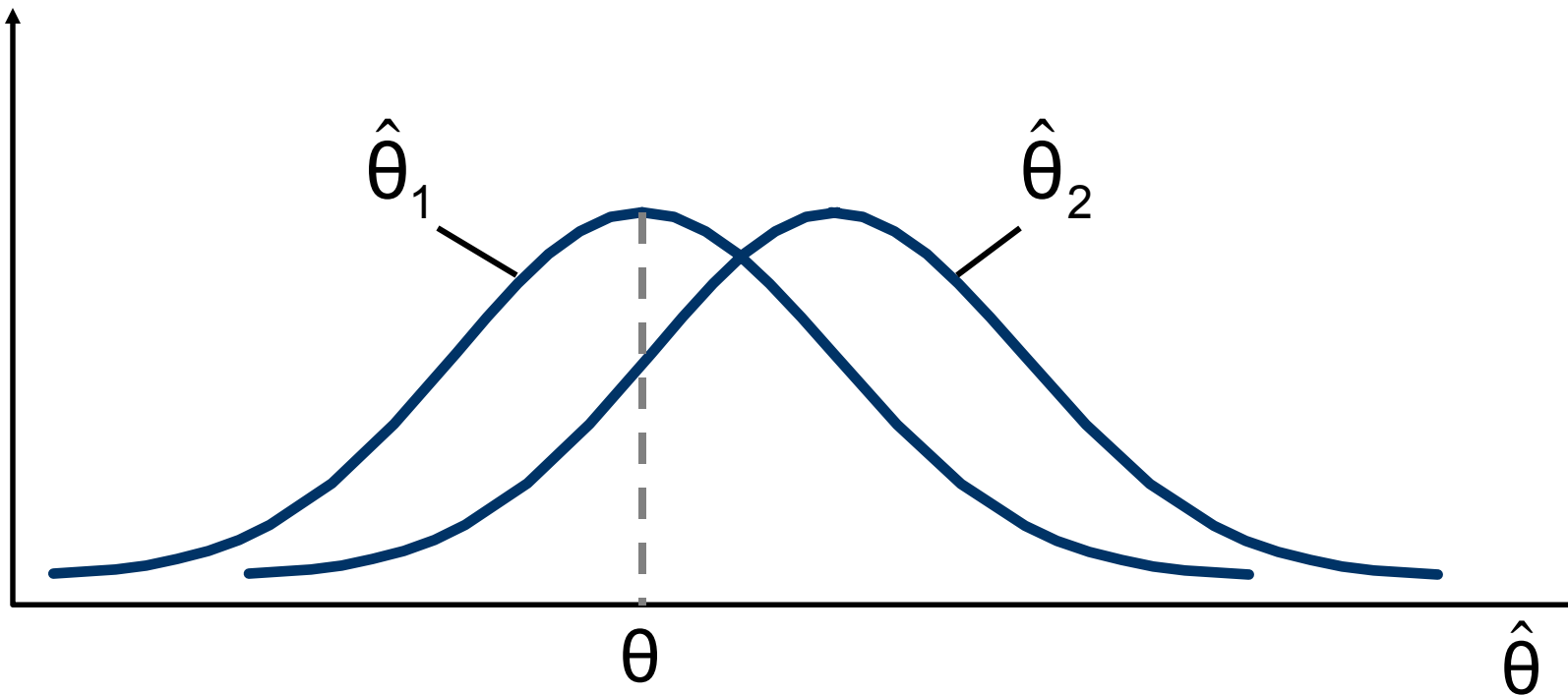
$$E(\hat{\theta}) = \theta$$

- Examples:
  - The sample mean is an unbiased estimator of $\mu$
  - The sample variance is an unbiased estimator of $\sigma^2$
  - The sample proportion is an unbiased estimator of P

# Unbiasedness

- $\hat{\theta}_1$ is an unbiased estimator, $\hat{\theta}_2$ is biased:



$\hat{\theta}_1$

$\hat{\theta}_2$

$\theta$

$\hat{\theta}$

K. Drakos, Quantitative Methods
for Finance

6

# Bias

- Let $\hat{\theta}$ be an estimator of $\theta$
- The bias in $\hat{\theta}$ is defined as the difference between its mean and $\theta$

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

- The bias of an unbiased estimator is 0

# Consistency

- Let $\hat{\theta}$ be an estimator of $\theta$

- $\hat{\theta}$ is a consistent estimator of $\theta$ if the difference between the expected value of and $\theta$ decreases as the sample size increases

- Consistency is desired when unbiased estimators cannot be obtained

# Most Efficient Estimator

- Suppose there are several unbiased estimators of $\theta$
- The most efficient estimator or the minimum variance unbiased estimator of $\theta$ is the unbiased estimator with the smallest variance

- Let $\hat{\theta}_1$ and $\hat{\theta}_2$ be two unbiased estimators of $\theta$, based on the same number of sample observations. Then,

  - $\hat{\theta}_1$ is said to be more efficient than $\hat{\theta}_2$ if $\boxed{Var(\hat{\theta}_1) < Var(\hat{\theta}_2)}$

  - The relative efficiency of $\hat{\theta}_1$ with respect to $\hat{\theta}_2$ is the ratio of their variances:

$$Relative\ Efficiency = \frac{Var(\hat{\theta}_2)}{Var(\hat{\theta}_1)}$$

# Confidence Intervals

- How much uncertainty is associated with a point estimate of a population parameter?

- An interval estimate provides more information about a population characteristic than does a point estimate

- Such interval estimates are called confidence intervals

# Confidence Interval Estimate

- An interval gives a range of values:

  - Takes into consideration variation in sample statistics from sample to sample

  - Based on observation from 1 sample

  - Gives information about closeness to unknown population parameters

  - Stated in terms of level of confidence

    - Can never be 100% confident

# Confidence Interval and Confidence Level

- If P(a < θ < b) = 1 - α then the interval from  a  to  b  is called a  100(1 - α)%  confidence interval of  θ.

- The quantity (1-α) is called the confidence level of the interval (α between 0 and 1)

  – In repeated samples of the population, the true value of the parameter θ would be contained in 100(1 - α)% of intervals calculated this way.
  – The confidence interval calculated in this manner is written as a < θ < b with 100(1 - α)% confidence

# Confidence Level, (1-$\alpha$)

*(continued)*

- Suppose confidence level = 95%

- Also written (1 - $\alpha$) = 0.95

- A relative frequency interpretation:

  – From repeated samples, 95% of all the confidence intervals that can be constructed will contain the unknown true parameter

- A specific interval either will contain or will not contain the true parameter

  – No probability involved in a specific interval

# General Formula

- The general formula for all confidence intervals is:

**Point Estimate $\pm$ (Reliability Factor)(Standard Error)**

- The value of the reliability factor depends on the desired level of confidence

# Confidence Interval for μ (σ² Known)

- Assumptions
  - Population variance $\sigma^2$ is known
  - Population is normally distributed
  - If population is not normal, use large sample

- Confidence interval estimate:

$$\overline{x} - z_{\alpha/2}\,\frac{\sigma}{\sqrt{n}} \;<\; \mu \;<\; \overline{x} + z_{\alpha/2}\,\frac{\sigma}{\sqrt{n}}$$

(where $z_{\alpha/2}$ is the normal distribution value for a probability of $\alpha/2$ in each tail)

# Margin of Error

- The confidence interval,

$$\bar{x} - z_{\alpha/2}\,\frac{\sigma}{\sqrt{n}} \;<\; \mu \;<\; \bar{x} + z_{\alpha/2}\,\frac{\sigma}{\sqrt{n}}$$

- Can also be written as $\bar{x} \pm ME$

  where ME is called the margin of error

$$ME = z_{\alpha/2}\,\frac{\sigma}{\sqrt{n}}$$

- The interval width, w, is equal to twice the margin of error
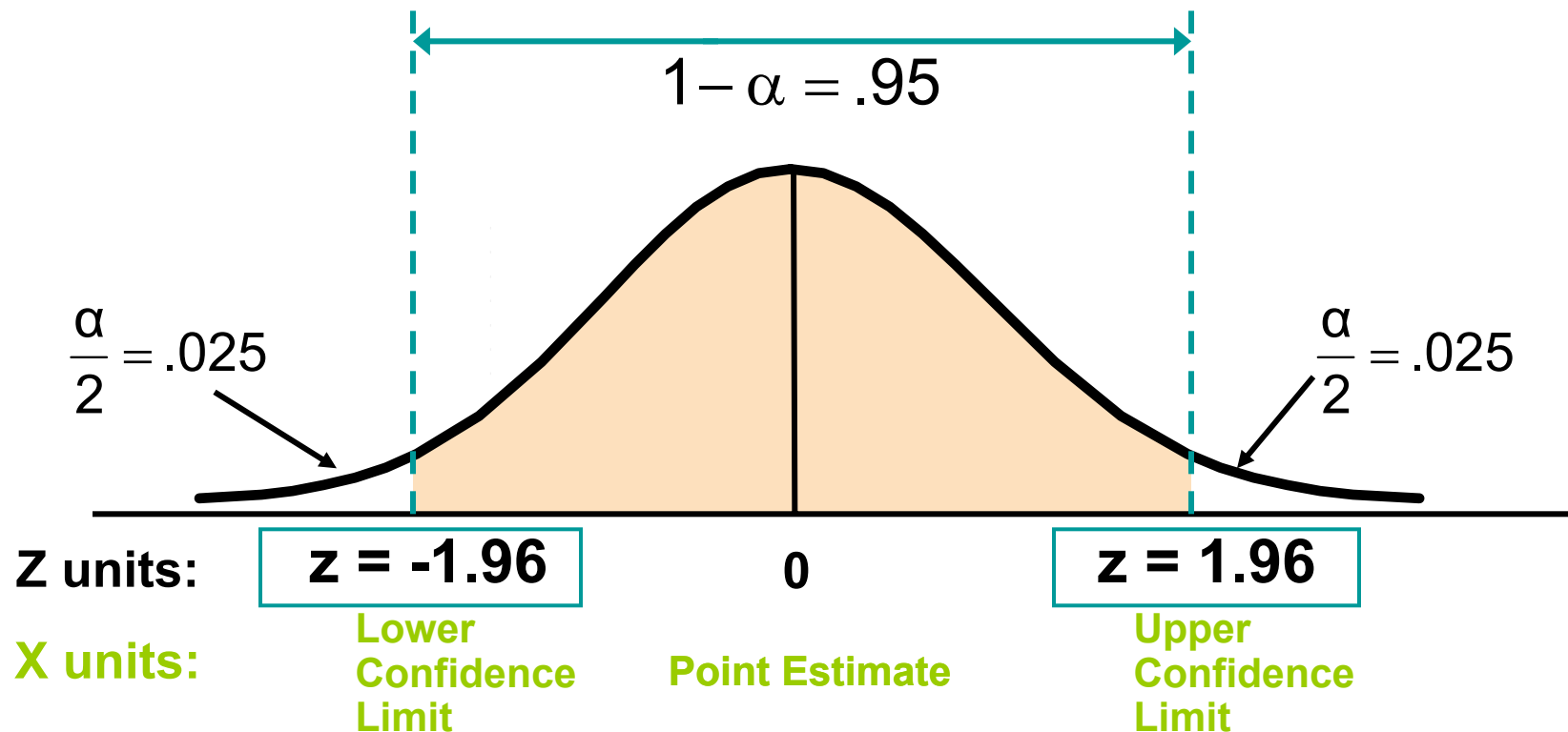
# Reducing the Margin of Error

$$ME = z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

The margin of error can be reduced if

- the population standard deviation can be reduced ($\sigma\downarrow$)
- The sample size is increased ($n\uparrow$)
- The confidence level is decreased, ($1-\alpha$) $\downarrow$

# Finding the Reliability Factor, $z_{\alpha/2}$

- Consider a 95% confidence interval:



$$1-\alpha = .95$$

$$\frac{\alpha}{2} = .025$$

$$\frac{\alpha}{2} = .025$$

Z units:   z = -1.96   0   z = 1.96

X units:   Lower Confidence Limit   Point Estimate   Upper Confidence Limit

- Find $z_{.025} = \pm1.96$ from the standard normal distribution table
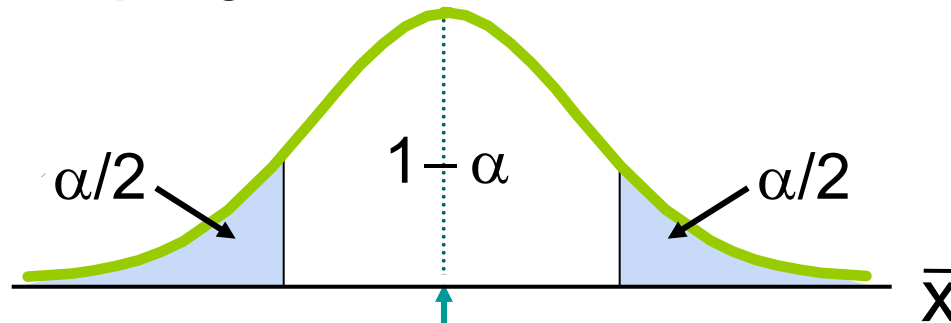
# Common Levels of Confidence

- Commonly used confidence levels are 90%, 95%, and 99%

| Confidence Level | Confidence Coefficient, $1-\alpha$ | $Z_{\alpha/2}$ value |
|---|---|---|
| 80% | .80 | 1.28 |
| 90% | .90 | 1.645 |
| 95% | .95 | 1.96 |
| 98% | .98 | 2.33 |
| 99% | .99 | 2.58 |
| 99.8% | .998 | 3.08 |
| 99.9% | .999 | 3.27 |

K. Drakos, Quantitative Methods for Finance

# Intervals and Level of Confidence

Sampling Distribution of the Mean

$$\alpha/2 \qquad 1-\alpha \qquad \alpha/2$$

$$\overline{x}$$

$$\mu_{\overline{x}} = \mu$$

$$\overline{x}_1$$

$$\overline{x}_2$$

Intervals extend from

$$\overline{x} - z\frac{\sigma}{\sqrt{n}}$$

to

$$\overline{x} + z\frac{\sigma}{\sqrt{n}}$$

$100(1-\alpha)\%$ of intervals constructed contain $\mu$;

$100(\alpha)\%$ do not.

Confidence Intervals

# Example

- A sample of 11 firms from a large normal population has a mean monthly return of 2.20%. We know from past testing that the population standard deviation is 0.35%.

- Determine a 95% confidence interval for the true mean return of the population.

# Example

$$\bar{x} \pm z \frac{\sigma}{\sqrt{n}}$$

$$= 2.20 \pm 1.96\,(.35/\sqrt{11})$$

$$= 2.20 \pm .2068$$

$$1.9932 < \mu < 2.4068$$

# Interpretation

- We are 95% confident that the true mean return is between 1.9932 and 2.4068 %

- Although the true mean may or may not be in this interval, 95% of intervals formed in this manner will contain the true mean

# Student's t Distribution

- Consider a random sample of n observations
  - with mean $\bar{x}$ and standard deviation s
  - from a normally distributed population with mean μ

- Then the variable $$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

- follows the Student's t distribution with (n - 1) degrees of freedom

# Confidence Interval for μ (σ² Unknown)

- If the population standard deviation σ is unknown, we can substitute the sample standard deviation, s

- This introduces extra uncertainty, since s is variable from sample to sample

- So we use the t distribution instead of the normal distribution

# Confidence Interval for μ (σ Unknown)

- Assumptions
  - Population standard deviation is unknown
  - Population is normally distributed
  - If population is not normal, use large sample

- Use Student's t  Distribution

- Confidence Interval Estimate:

$$\overline{x} - t_{n-1,\alpha/2}\,\frac{S}{\sqrt{n}} \;<\; \mu \;<\; \overline{x} + t_{n-1,\alpha/2}\,\frac{S}{\sqrt{n}}$$

where $t_{n-1,\alpha/2}$ is the critical value of the t distribution with n-1 d.f. and an area of α/2 in each tail:

$$P(t_{n-1} > t_{n-1,\alpha/2}) = \alpha/2$$
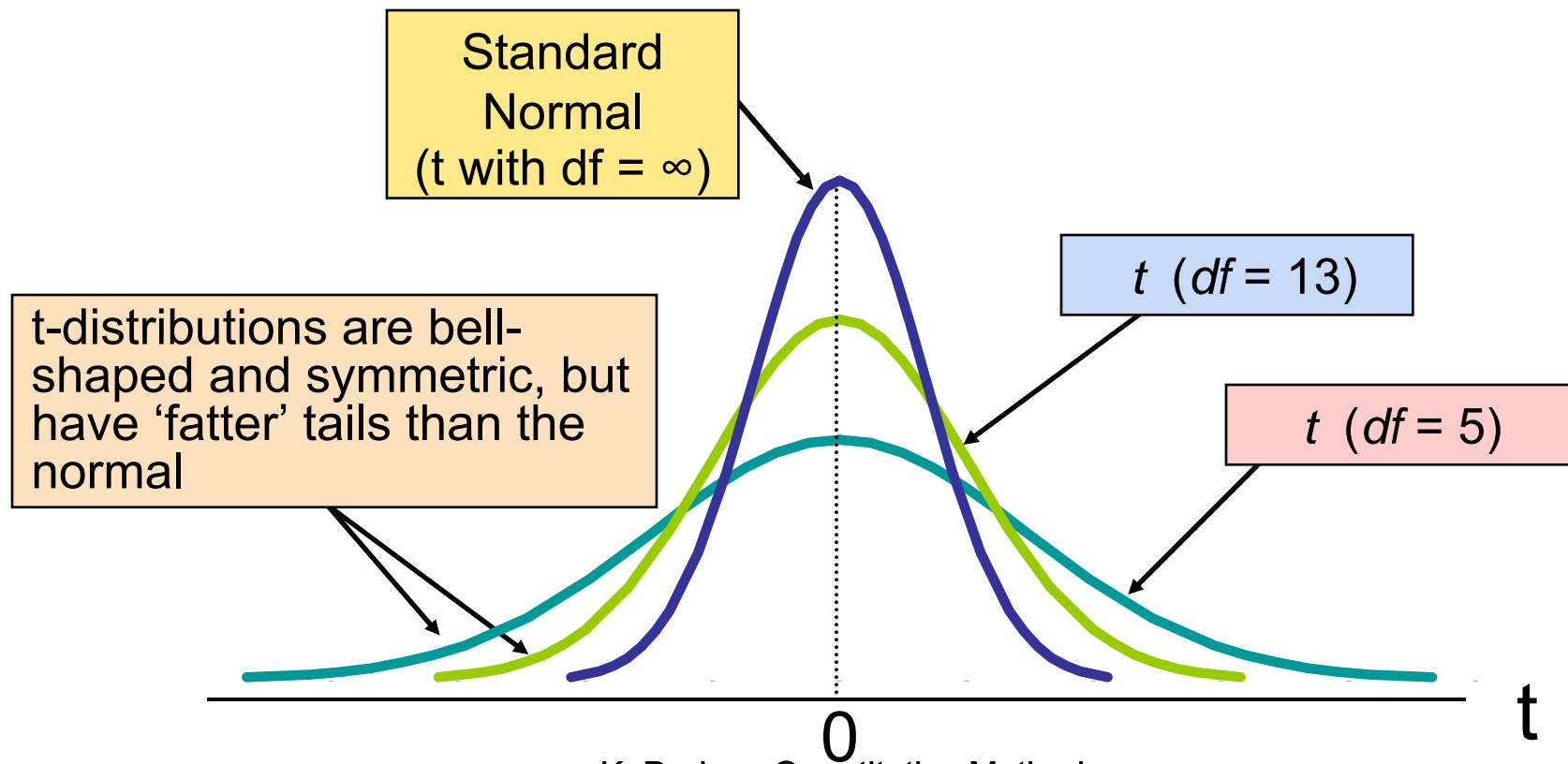
K. Drakos, for Finance

26

# Student's t Distribution

- The  t  is a family of distributions

- The  t value  depends on degrees of freedom (d.f.)

  - Number of observations that are free to vary after sample mean has been calculated
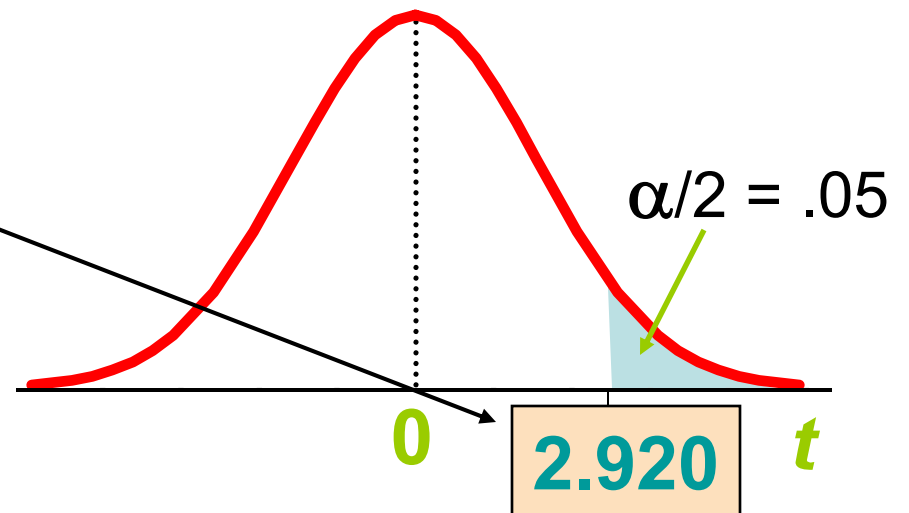
$$d.f. = n - 1$$

# Student's t Distribution

Note: t ⟶ Z as n increases

Standard
Normal
(t with df = ∞)

$t$ ($df$ = 13)

t-distributions are bell-
shaped and symmetric, but
have 'fatter' tails than the
normal

$t$ ($df$ = 5)

0

t

# Student's t Table

| df | Upper Tail Area | | |
|----|------|------|------|
|    | .10  | .05  | .025 |
| 1  | 3.078 | 6.314 | 12.706 |
| 2  | 1.886 | **2.920** | 4.303 |
| 3  | 1.638 | 2.353 | 3.182 |

Let: n = 3
df = $n$ - 1 = 2
$\alpha$ = .10
$\alpha$/2 =.05

The body of the table contains t values, not probabilities
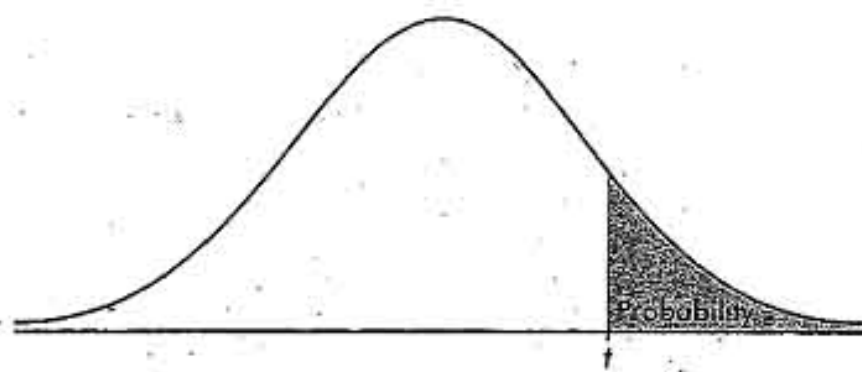
$\alpha$/2 = .05

0

**2.920**

$t$

# TABLE B: t-DISTRIBUTION CRITICAL VALUES

| df | \.25 | \.20 | \.15 | \.10 | \.05 | \.025 | \.02 | \.01 | \.005 | \.0025 | \.001 | \.0005 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Tail probability $p$ | | | | | | |
| 1 | 1.000 | 1.376 | 1.963 | 3.078 | 6.314 | 12.71 | 15.89 | 31.82 | 63.66 | 127.3 | 318.3 | 636.6 |
| 2 | .816 | 1.061 | 1.386 | 1.886 | 2.920 | 4.303 | 4.849 | 6.965 | 9.925 | 14.09 | 22.33 | 31.60 |
| 3 | .765 | .978 | 1.250 | 1.638 | 2.353 | 3.182 | 3.482 | 4.541 | 5.841 | 7.453 | 10.21 | 12.92 |
| 4 | .741 | .941 | 1.190 | 1.533 | 2.132 | 2.776 | 2.999 | 3.747 | 4.604 | 5.598 | 7.173 | 8.610 |
| 5 | .727 | .920 | 1.156 | 1.476 | 2.015 | 2.571 | 2.757 | 3.365 | 4.032 | 4.773 | 5.893 | 6.869 |
| 6 | .718 | .906 | 1.134 | 1.440 | 1.943 | 2.447 | 2.612 | 3.143 | 3.707 | 4.317 | 5.208 | 5.959 |
| 7 | .711 | .896 | 1.119 | 1.415 | 1.895 | 2.365 | 2.517 | 2.998 | 3.499 | 4.029 | 4.785 | 5.408 |
| 8 | .706 | .889 | 1.108 | 1.397 | 1.860 | 2.306 | 2.449 | 2.896 | 3.355 | 3.833 | 4.501 | 5.041 |
| 9 | .703 | .883 | 1.100 | 1.383 | 1.833 | 2.262 | 2.398 | 2.821 | 3.250 | 3.690 | 4.297 | 4.781 |
| 10 | .700 | .879 | 1.093 | 1.372 | 1.812 | 2.228 | 2.359 | 2.764 | 3.169 | 3.581 | 4.144 | 4.587 |
| 11 | .697 | .876 | 1.088 | 1.363 | 1.796 | 2.201 | 2.328 | 2.718 | 3.106 | 3.497 | 4.025 | 4.437 |
| 12 | .695 | .873 | 1.083 | 1.356 | 1.782 | 2.179 | 2.303 | 2.681 | 3.055 | 3.428 | 3.930 | 4.318 |
| 13 | .694 | .870 | 1.079 | 1.350 | 1.771 | 2.160 | 2.282 | 2.650 | 3.012 | 3.372 | 3.852 | 4.221 |
| 14 | .692 | .868 | 1.076 | 1.345 | 1.761 | 2.145 | 2.264 | 2.624 | 2.977 | 3.326 | 3.787 | 4.140 |
| 15 | .691 | .866 | 1.074 | 1.341 | 1.753 | 2.131 | 2.249 | 2.602 | 2.947 | 3.286 | 3.733 | 4.073 |
| 16 | .690 | .865 | 1.071 | 1.337 | 1.746 | 2.120 | 2.235 | 2.583 | 2.921 | 3.252- | 3.686 | 4.015 |
| 17 | .689 | .863 | 1.069 | 1.333 | 1.740 | 2.110 | 2.224 | 2.567 | 2.898 | 3.222 | 3.646 | 3.965 |
| 18 | .688 | .862 | 1.067 | 1.330 | 1.734 | 2.101 | 2.214 | 2.552 | 2.878 | 3.197 | 3.611 | 3.922 |
| 19 | .688 | .861 | 1.066 | 1.328 | 1.729 | 2.093 | 2.205 | 2.539 | 2.861 | 3.174 | 3.579 | 3.883 |
| 20 | .687 | .860 | 1.064 | 1.325 | 1.725 | 2.086 | 2.197 | 2.528 | 2.845 | 3.153 | 3.552 | 3.850 |
| 21 | .686 | .859 | 1.063 | 1.323 | 1.721 | 2.080 | 2.189 | 2.518 | 2.831 | 3.135 | 3.527 | 3.819 |
| 22 | .686 | .858 | 1.061 | 1.321 | 1.717 | 2.074 | 2.183 | 2.508 | 2.819 | 3.119 | 3.505 | 3.792 |
| 23 | .685 | .858 | 1.060 | 1.319 | 1.714 | 2.069 | 2.177 | 2.500 | 2.807 | 3.104 | 3.485 | 3.768 |
| 24 | .685 | .857 | 1.059 | 1.318 | 1.711 | 2.064 | 2.172 | 2.492 | 2.797 | 3.091 | 3.467 | 3.745 |
| 25 | .684 | .856 | 1.058 | 1.316 | 1.708 | 2.060 | 2.167 | 2.485 | 2.787 | 3.078 | 3.450 | 3.725 |
| 26 | .684 | .856 | 1.058 | 1.315 | 1.706 | 2.056 | 2.162 | 2.479 | 2.779 | 3.067 | 3.435 | 3.707 |
| 27 | .684 | .855 | 1.057 | 1.314 | 1.703 | 2.052 | 2.158 | 2.473 | 2.771 | 3.057 | 3.421 | 3.690 |
| 28 | .683 | .855 | 1.056 | 1.313 | 1.701 | 2.048 | 2.154 | 2.467 | 2.763 | 3.047 | 3.408 | 3.674 |
| 29 | .683 | .854 | 1.055 | 1.311 | 1.699 | 2.045 | 2.150 | 2.462 | 2.756 | 3.038 | 3.396 | 3.659 |
| 30 | .683 | .854 | 1.055 | 1.310 | 1.697 | 2.042 | 2.147 | 2.457 | 2.750 | 3.030 | 3.385 | 3.646 |
| 40 | .681 | .851 | 1.050 | 1.303 | 1.684 | 2.021 | 2.123 | 2.423 | 2.704 | 2.971 | 3.307 | 3.551 |
| 50 | .679 | .849 | 1.047 | 1.299 | 1.676 | 2.009 | 2.109 | 2.403 | 2.678 | 2.937 | 3.261 | 3.496 |
| 60 | .679 | .848 | 1.045 | 1.296 | 1.671 | 2.000 | 2.099 | 2.390 | 2.660 | 2.915 | 3.232 | 3.460 |
| 80 | .678 | .846 | 1.043 | 1.292 | 1.664 | 1.990 | 2.088 | 2.374 | 2.639 | 2.887 | 3.195 | 3.416 |
| 100 | .677 | .845 | 1.042 | 1.290 | 1.660 | 1.984 | 2.081 | 2.364 | 2.626 | 2.871 | 3.174 | 3.390 |
| 1000 | .675 | .842 | 1.037 | 1.282 | 1.646 | 1.962 | 2.056 | 2.330 | 2.581 | 2.813 | 3.098 | 3.300 |
| ∞ | .674 | .841 | 1.036 | 1.282 | 1.645 | 1.960 | 2.054 | 2.326 | 2.576 | 2.807 | 3.091 | 3.291 |
| | 50% | 60% | 70% | 80% | 90% | 95% | 96% | 98% | 99% | 99.5% | 99.8% | 99.9% |

Confidence level $C$

# t distribution values

## With comparison to the Z value

| Confidence Level | t (10 d.f.) | t (20 d.f.) | t (30 d.f.) | Z |
|---|---|---|---|---|
| .80 | 1.372 | 1.325 | 1.310 | 1.282 |
| .90 | 1.812 | 1.725 | 1.697 | 1.645 |
| .95 | 2.228 | 2.086 | 2.042 | 1.960 |
| .99 | 3.169 | 2.845 | 2.750 | 2.576 |

Note:  t ⟶ Z  as  n  increases

# Example

A random sample of $n = 25$ has $\bar{x} = 50$ and $s = 8$. Form a 95% confidence interval for $\mu$

– d.f. $= n - 1 = 24$, so $t_{n-1,\alpha/2} = t_{24,.025} = 2.0639$

The confidence interval is

$$\bar{x} - t_{n-1,\alpha/2}\,\frac{S}{\sqrt{n}} \;<\; \mu \;<\; \bar{x} + t_{n-1,\alpha/2}\,\frac{S}{\sqrt{n}}$$

$$50 - (2.0639)\,\frac{8}{\sqrt{25}} \;<\; \mu \;<\; 50 + (2.0639)\,\frac{8}{\sqrt{25}}$$

$$46.698 \;<\; \mu \;<\; 53.302$$

# Confidence Intervals for the Population Proportion, p

- An interval estimate for the population proportion ( P ) can be calculated by adding an allowance for uncertainty to the sample proportion ( $\hat{p}$ )

# Confidence Intervals for the Population Proportion, p

- Recall that the distribution of the sample proportion is approximately normal if the sample size is large, with standard deviation

$$\sigma_P = \sqrt{\frac{P(1-P)}{n}}$$

- We will estimate this with sample data:

$$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

# Confidence Interval Endpoints

- Upper and lower confidence limits for the population proportion are calculated with the formula

$$\hat{p} - z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < P < \hat{p} + z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

- where
    - $z_{\alpha/2}$ is the standard normal value for the level of confidence desired
    - $\hat{p}$ is the sample proportion
    - n is the sample size

# Example

- A random sample of 100 firms shows that 25 did not pay dividend

- Form a 95% confidence interval for the true proportion of non-paying firms

# Example

$$\hat{p} - z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \; < \; P \; < \; \hat{p} + z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\frac{25}{100} - 1.96\sqrt{\frac{.25(.75)}{100}} \; < \; P \; < \; \frac{25}{100} + 1.96\sqrt{\frac{.25(.75)}{100}}$$

$$0.1651 \; < \; P \; < \; 0.3349$$

# Interpretation

- We are 95% confident that the true percentage of non-paying firms in the population is between

  16.51% and 33.49%.

- Although the interval from 0.1651 to 0.3349 may or may not contain the true proportion, 95% of intervals formed from samples of size 100 in this manner will contain the true proportion.

# Dependent Samples

Tests Means of 2 Related Populations

- Paired or matched samples

- Repeated measures (before/after)

- Use difference between paired values:

$$d_i = x_i - y_i$$

- Assumptions:

  - Both Populations Are Normally Distributed

# Mean Difference

The i[th] paired difference is $d_i$, where

$$d_i = x_i - y_i$$

The point estimate for the population mean paired difference is $\bar{d}$:

$$\bar{d} = \frac{\sum\limits_{i=1}^{n} d_i}{n}$$

The sample standard deviation is:

$$S_d = \sqrt{\frac{\sum\limits_{i=1}^{n}(d_i - \bar{d})^2}{n-1}}$$

n is the number of matched pairs in the sample

# Confidence Interval for Mean Difference

The confidence interval for difference between population means, $\mu_d$, is

**Dependent samples**

$$\bar{d} - t_{n-1,\alpha/2}\frac{S_d}{\sqrt{n}} < \mu_d < \bar{d} + t_{n-1,\alpha/2}\frac{S_d}{\sqrt{n}}$$

Where
n = the sample size
(number of matched pairs in the paired sample)

# Confidence Interval for Mean Difference

*(continued)*

Dependent samples

- The margin of error is

$$ME = t_{n-1,\alpha/2}\,\frac{s_d}{\sqrt{n}}$$

- $t_{n-1,\alpha/2}$ is the value from the Student's t distribution with (n – 1) degrees of freedom for which

$$P(t_{n-1} > t_{n-1,\alpha/2}) = \frac{\alpha}{2}$$

# Paired Samples Example

- Six people sign up for a weight loss program. You collect the following data:

| Person | Weight: Before (x) | After (y) | Difference, $d_i$ |
|:---:|:---:|:---:|:---:|
| 1 | 136 | 125 | 11 |
| 2 | 205 | 195 | 10 |
| 3 | 157 | 150 | 7 |
| 4 | 138 | 140 | - 2 |
| 5 | 175 | 165 | 10 |
| 6 | 166 | 160 | 6 |
| | | | 42 |

$$\overline{d} = \frac{\sum d_i}{n}$$

$$= 7.0$$

$$S_d = \sqrt{\frac{\sum (d_i - \overline{d})^2}{n-1}}$$

$$= 4.82$$

# Paired Samples Example

- For a 95% confidence level, the appropriate t value is $t_{n-1,\alpha/2} = t_{5,.025} = 2.571$

- The 95% confidence interval for the difference between means, $\mu_d$, is

$$\overline{d} - t_{n-1,\alpha/2}\frac{S_d}{\sqrt{n}} < \mu_d < \overline{d} + t_{n-1,\alpha/2}\frac{S_d}{\sqrt{n}}$$

$$7 - (2.571)\frac{4.82}{\sqrt{6}} < \mu_d < 7 + (2.571)\frac{4.82}{\sqrt{6}}$$

$$-1.94 < \mu_d < 12.06$$

Since this interval contains zero, we cannot be 95% confident, given this limited data, that the weight loss program helps people lose weight

# Difference Between Two Means

**Population means, independent samples**

**Goal:** Form a confidence interval for the difference between two population means, $\mu_x - \mu_y$

- Different data sources
  - Unrelated
  - Independent
    - Sample selected from one population has no effect on the sample selected from the other population
- The point estimate is the difference between the two sample means:

$$\overline{x} - \overline{y}$$

# Difference Between Two Means

*(continued)*



Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known ⟶ Confidence interval uses $z_{\alpha/2}$

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal

Confidence interval uses a value from the Student's **t** distribution

# $\sigma_x^2$ and $\sigma_y^2$ Known

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known *

$\sigma_x^2$ and $\sigma_y^2$ unknown

Assumptions:

- Samples are randomly and independently drawn

- both population distributions are normal

- Population variances are known

# σ$_x^2$ and σ$_y^2$ Known

| Population means, independent samples |
|---|

| σ$_x^2$ and σ$_y^2$ known | * |
| σ$_x^2$ and σ$_y^2$ unknown |

When σ$_x$ and σ$_y$ are known and both populations are normal, the variance of $\overline{X} - \overline{Y}$ is

$$\sigma^2_{\overline{X}-\overline{Y}} = \frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}$$

…and the random variable

$$Z = \frac{(\overline{x}-\overline{y})-(\mu_X-\mu_Y)}{\sqrt{\dfrac{\sigma_x^2}{n_X} + \dfrac{\sigma_y^2}{n_Y}}}$$

has a standard normal distribution

# Confidence Interval, $\sigma_x^2$ and $\sigma_y^2$ Known

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

**\*** The confidence interval for $\mu_x - \mu_y$ is:

$$(\bar{x} - \bar{y}) - z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n_x} + \frac{\sigma_Y^2}{n_y}} < \mu_X - \mu_Y < (\bar{x} - \bar{y}) + z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n_x} + \frac{\sigma_Y^2}{n_y}}$$
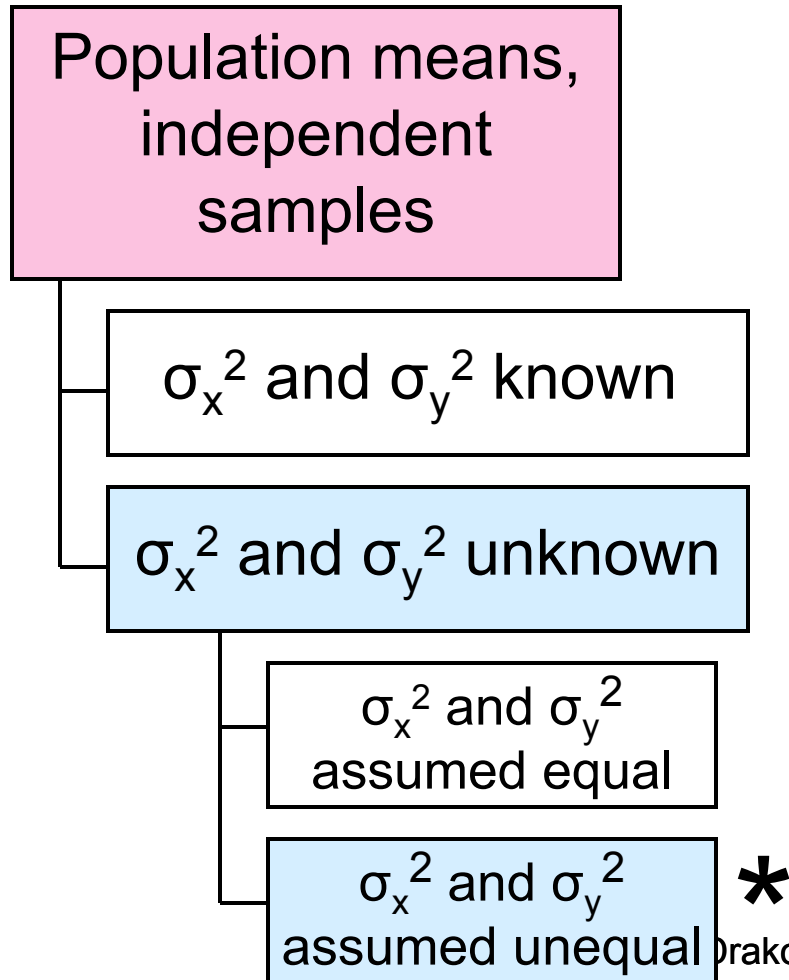
# $\sigma_x^2$ and $\sigma_y^2$ Unknown, Assumed Equal

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal **\***
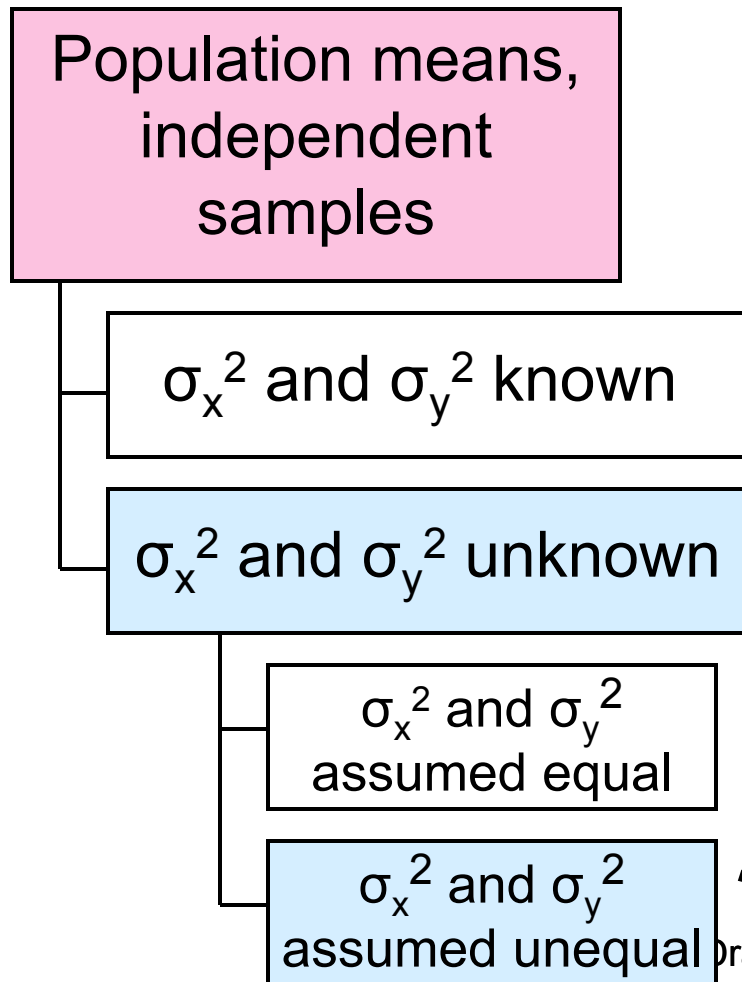
$\sigma_x^2$ and $\sigma_y^2$ assumed unequal

Assumptions:

- Samples are randomly and independently drawn

- Populations are normally distributed

- Population variances are unknown but assumed equal

# $\sigma_x 2$ and $\sigma_y^2$ Unknown, Assumed Equal

*(continued)*

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal **\***

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal

Forming interval estimates:

- The population variances are assumed equal, so use the two sample standard deviations and pool them to estimate $\sigma$

- use a t value with $(n_x + n_y - 2)$ degrees of freedom

# $\sigma_x^2$ and $\sigma_y^2$ Unknown, Assumed Equal

*(continued)*

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal **\***

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal

The pooled variance is

$$s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{n_x + n_y - 2}$$

# Confidence Interval, $\sigma_x^2$ and $\sigma_y^2$ Unknown, Equal

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal

**\*** The confidence interval for $\mu_1 - \mu_2$ is:

$$(\overline{x} - \overline{y}) - t_{n_x + n_y - 2, \alpha/2} \sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}} < \mu_X - \mu_Y < (\overline{x} - \overline{y}) + t_{n_x + n_y - 2, \alpha/2} \sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}}$$

Where $$s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{n_x + n_y - 2}$$

ve Methods

for Finance

52

# $\sigma_x^2$ and $\sigma_y^2$ Unknown, Assumed Unequal

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal $*$

Assumptions:

- Samples are randomly and independently drawn

- Populations are normally distributed

- Population variances are unknown and assumed unequal

# $\sigma_x^2$ and $\sigma_y^2$ Unknown, Assumed Unequal

*(continued)*

Population means, independent samples

$\sigma_x^2$ and $\sigma_y^2$ known

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal **✱**

Forming interval estimates:

- The population variances are assumed unequal, so a pooled variance is not appropriate

- use a t value with $\nu$ degrees of freedom, where

$$\nu = \frac{\left[ \left(\dfrac{s_x^2}{n_x}\right) + \left(\dfrac{s_y^2}{n_y}\right) \right]^2}{\left(\dfrac{s_x^2}{n_x}\right)^2 /(n_x - 1) + \left(\dfrac{s_y^2}{n_y}\right)^2 /(n_y - 1)}$$

Drakos, Quantit for Fina...

# Confidence Interval, $\sigma_x^2$ and $\sigma_y^2$ Unknown, Unequal

$\sigma_x^2$ and $\sigma_y^2$ unknown

$\sigma_x^2$ and $\sigma_y^2$ assumed equal

$\sigma_x^2$ and $\sigma_y^2$ assumed unequal   **\***

The confidence interval for $\mu_1 - \mu_2$ is:

$$(\bar{x} - \bar{y}) - t_{\nu,\alpha/2}\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}} < \mu_X - \mu_Y < (\bar{x} - \bar{y}) + t_{\nu,\alpha/2}\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

Where

$$\nu = \frac{\left[(\frac{s_x^2}{n_x}) + (\frac{s_y^2}{n_y})\right]^2}{\left(\frac{s_x^2}{n_x}\right)^2 /(n_x - 1) + \left(\frac{s_y^2}{n_y}\right)^2 /(n_y - 1)}$$

# Two Population Proportions

Population proportions

**Goal:** Form a confidence interval for the difference between two population proportions, $\boxed{P_x - P_y}$

**Assumptions:**

Both sample sizes are large (generally at least 40 observations in each sample)

The point estimate for the difference is $\hat{p}_x - \hat{p}_y$

# Two Population Proportions

**Population proportions**

- The random variable

$$Z = \frac{(\hat{p}_x - \hat{p}_y) - (p_x - p_y)}{\sqrt{\dfrac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \dfrac{\hat{p}_y(1-\hat{p}_y)}{n_y}}}$$

is approximately normally distributed

# Confidence Interval for Two Population Proportions

Population proportions

The confidence limits for $P_x - P_y$ are:

$$(\hat{p}_x - \hat{p}_y) \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}}$$

# Example:
# Two Population Proportions

Form a 90% confidence interval for the difference between the proportion of retail firms and the proportion of industrial firms who went bankrupt last year.

- In a random sample, 26 of 50 retail and 28 of 40 industrial firms had gone bankrupt

# Example:
# Two Population Proportions

*(continued)*

Retail:

$$\hat{p}_x = \frac{26}{50} = 0.52$$

Industrial:

$$\hat{p}_y = \frac{28}{40} = 0.70$$

$$\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}} = \sqrt{\frac{0.52(0.48)}{50} + \frac{0.70(0.30)}{40}} = 0.1012$$

For 90% confidence, $Z_{\alpha/2} = 1.645$

# Example:
## Two Population Proportions

*(continued)*

The confidence limits are:

$$(\hat{p}_x - \hat{p}_y) \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}}$$

$$= (.52 - .70) \pm 1.645(0.1012)$$

so the confidence interval is

$$-0.3465 < P_x - P_y < -0.0135$$

Since this interval does not contain zero we are 90% confident that the two proportions are not equal

# Confidence Intervals for the Population Variance

Population Variance

- Goal: Form a confidence interval for the population variance, $\sigma^2$

- The confidence interval is based on the sample variance, $s^2$

- Assumed: the population is normally distributed

# Confidence Intervals for the Population Variance

*(continued)*

| Population Variance |
|---|

The random variable

$$\chi_{n-1}^2 = \frac{(n-1)s^2}{\sigma^2}$$

follows a chi-square distribution with (n – 1) degrees of freedom

The chi-square value $\chi_{n-1,\alpha}^2$ denotes the number for which

$$P(\chi_{n-1}^2 > \chi_{n-1,\alpha}^2) = \alpha$$

# Confidence Intervals for the Population Variance

*(continued)*

Population Variance

The (1 - $\alpha$)% confidence interval for the population variance is

$$\frac{(n-1)s^2}{\chi^2_{n-1,\,\alpha/2}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{n-1,\,1-\alpha/2}}$$

# Hypothesis Testing

# What is a Hypothesis?

- A hypothesis is a claim (assumption) about a population parameter:
- population mean / population proportion

**Example: The mean monthly cell phone bill of this city is $\mu = \$42$**

**Example: The proportion of adults in this city with cell phones is $p = .68$**

# The Null Hypothesis, $H_0$

- States the assumption (numerical) to be tested

> Example: The average number of TV sets in U.S. Homes is equal to three $(H_0 : \mu = 3$ )

- Is always about a population parameter, not about a sample statistic

$$H_0 : \mu = 3 \qquad \cancel{H_0 : \overline{X} = 3}$$

# The Null Hypothesis, $H_0$

- Begin with the assumption that the null hypothesis is true

    – Similar to the notion of innocent until proven guilty

- Refers to the status quo

- Always contains "=" , "≤" or "≥" sign

- May or may not be rejected

# The Alternative Hypothesis, $H_1$

- Is the opposite of the null hypothesis
  - e.g., The average number of TV sets in U.S. homes is not equal to 3  ( $H_1$: $\mu \neq 3$ )
- Challenges the status quo
- May or may not be supported
- Is generally the hypothesis that the researcher is trying to support

# Hypothesis Testing Process

**Claim:** the population mean age is 50. (Null Hypothesis: $H_0$: $\mu = 50$ )



**Population**

Now select a random sample

**Sample**

Is $\overline{X} = 20$ likely if $\mu = 50$?

**If not likely, REJECT Null Hypothesis**

Suppose the sample mean age is 20: $\overline{X} = 20$

# Hypothesis Tests Design

- Is $\overline{X}$ likely given that μ = 50? If we believe that this not likely we will reject $H_0$.

- How can we determine if the event is likely to occur given that $H_0$ is true?

- We define a _rejection region_ of the sampling distribution, $\overline{X}$ < c.

$$P\left(\overline{X} < c \,\middle|\, \mu = 50\right) = P\left(\text{Re ject } H_0 \,\middle|\, H_0 \text{ true}\right) = \alpha$$

- So we want small values of α (_significance level_). According to α if we know the distribution of X we can determine c (_critical value_)

# Hypothesis Tests Design

- If we find that for $\alpha$ = 1% the c = 22 then since 20 < 22 *we will reject $H_0$ at the 1% significance level*.

- So for $H_0$ being true the 1% of the samples would have $\overline{X}$ < c. The rest 99% would have a sample mean $\overline{X}$ > c. So we are 99% confident that $H_0$ should be rejected.

K. Drakos, Quantitative Methods
for Finance

# Hypothesis Tests Design

*(continued)*

Sampling Distribution of $\overline{X}$



$\alpha$

20

c = 22

$\mu = 50$
If $H_0$ is true

If it is unlikely that we would get a sample mean of this value ...

... if in fact this were the population mean…

... then we reject the null hypothesis that $\mu = 50$.

# Level of Significance, $\alpha$

- Defines the unlikely values of the sample statistic if the null hypothesis is true

  - Defines rejection region of the sampling distribution

- Is designated by $\alpha$ , (level of significance)

  - Typical values are .01, .05, or .10

- Is selected by the researcher at the beginning

- Provides the critical value(s) of the test

# Level of Significance and the Rejection Region

Level of significance = $\alpha$

$H_0: \mu = 3$
$H_1: \mu \neq 3$

$\alpha/2$          $\alpha/2$

Two-tail test          0

$H_0: \mu \leq 3$
$H_1: \mu > 3$

$\alpha$

Upper-tail test          0

$H_0: \mu \geq 3$
$H_1: \mu < 3$

$\alpha$

Lower-tail test          0

Represents critical value

Rejection region is shaded

# Errors in Making Decisions

- **Type I Error**
  - Reject a true null hypothesis
  - Considered a serious type of error

The probability of Type I Error is $\alpha$

- Called level of significance of the test
- Set by researcher in advance

# Errors in Making Decisions

- **Type II Error**
    - Fail to reject a false null hypothesis

The probability of Type II Error is  $\beta$

# Outcomes and Probabilities

| | Possible Hypothesis Test Outcomes | |
|---|---|---|

| | **Actual Situation** | |
|---|---|---|
| **Decision** | $H_0$ True | $H_0$ False |
| Do Not Reject $H_0$ | **No error** $(1 - \alpha)$ | **Type II Error** $(\beta)$ |
| Reject $H_0$ | **Type I Error** $(\alpha)$ | **No Error** $(1 - \beta)$ |

**Key:**
**Outcome**
**(Probability)**

# Type I & II Error Relationship

- Type I and Type II errors can not happen at the same time

    - Type I error can only occur if $H_0$ is true

    - Type II error can only occur if $H_0$ is false

If Type I error probability ( $\alpha$ ) ⬆, then

Type II error probability ( $\beta$ ) ⬇

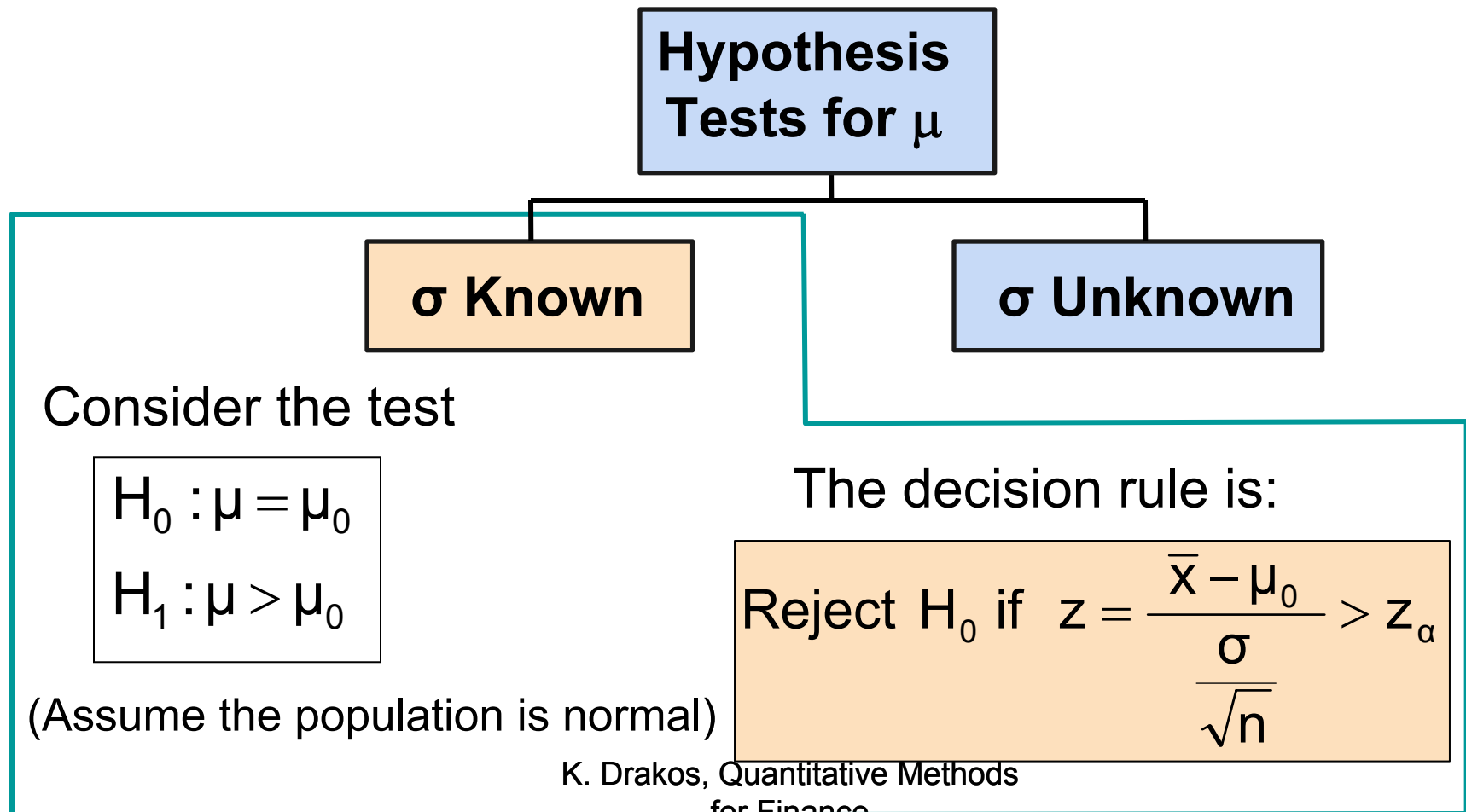# Factors Affecting Type II Error

- All else equal,
  - β ⬆ when the difference between hypothesized parameter and its true value ⬇

  - β ⬆ when α ⬇

  - β ⬆ when σ ⬆

  - β ⬆ when $n$ ⬇

# Power of the Test

- The power of a test is the probability of rejecting a null hypothesis that is false

- i.e.,     Power = P(Reject $H_0$ | $H_1$ is true)

  – Power of the test increases as the sample size increases

# Test of Hypothesis
# for the Mean (σ Known)
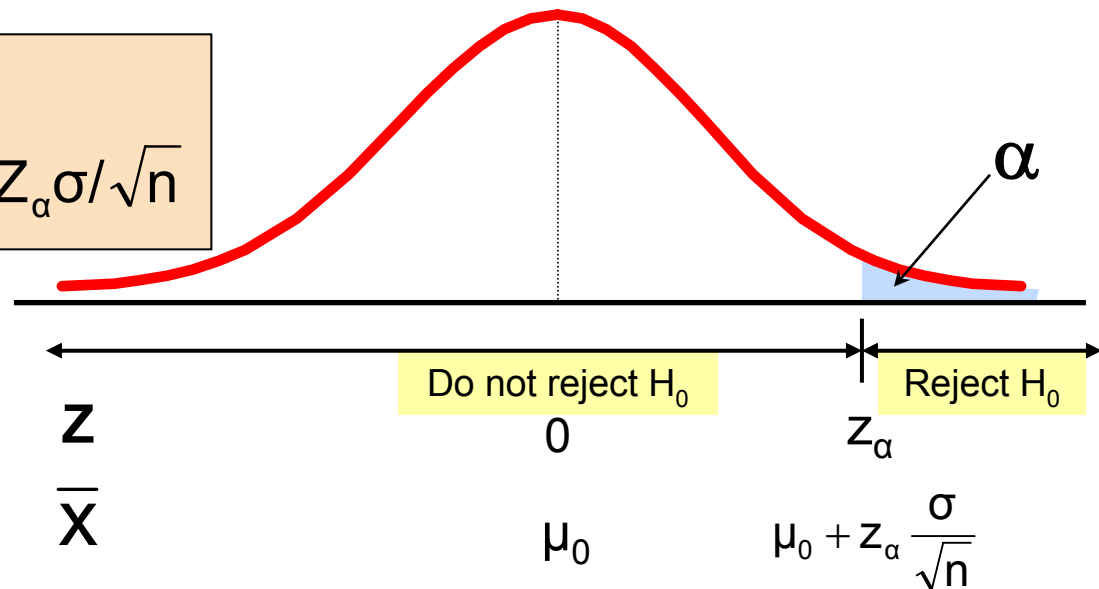
- Convert sample result ($\bar{x}$) to a z value



**Hypothesis Tests for** $\mu$

**σ Known**

**σ Unknown**

Consider the test

$$H_0 : \mu = \mu_0$$
$$H_1 : \mu > \mu_0$$

(Assume the population is normal)

The decision rule is:

Reject $H_0$ if $z = \dfrac{\bar{x} - \mu_0}{\dfrac{\sigma}{\sqrt{n}}} > z_\alpha$

K. Drakos, Quantitative Methods for Finance

82

# Decision Rule

$$H_0: \mu = \mu_0$$
$$H_1: \mu > \mu_0$$

Reject $H_0$ if $z = \dfrac{\bar{x} - \mu_0}{\dfrac{\sigma}{\sqrt{n}}} > z_\alpha$

Alternate rule:

Reject $H_0$ if $\bar{X} > \mu_0 + Z_\alpha \sigma / \sqrt{n}$

$\alpha$

| Do not reject $H_0$ | Reject $H_0$ |

Z     0     $z_\alpha$

$\bar{X}$     $\mu_0$     $\mu_0 + z_\alpha \dfrac{\sigma}{\sqrt{n}}$

Critical value

K. Drakos, Quantitative Methods
for Finance

# p-Value Approach to Testing

- p-value: Probability of obtaining a test statistic more extreme ( ≤ or ≥ ) than the observed sample value given $H_0$ is true

  – Also called observed level of significance

  – Smallest value of $\alpha$ for which $H_0$ can be rejected

# p-Value Approach to Testing

- Convert sample result (e.g., $\bar{x}$ ) to test statistic (e.g., z statistic )

- Obtain the p-value
  - For an upper tail test:

$$p\text{-value} = P(Z > \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}, \text{given that } H_0 \text{ is true})$$

$$= P(Z > \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \mid \mu = \mu_0)$$

- Decision rule: compare the p-value to $\alpha$

  - If  p-value $< \alpha$ ,  reject $H_0$

  - If  p-value $\geq \alpha$ ,  do not reject $H_0$

K. Drakos, Quantitative Methods for Finance

85

# Example: Upper-Tail Z Test for Mean ($\sigma$ Known)

A phone industry manager thinks that customer monthly cell phone bill have increased, and now average over $52 per month.  The company wishes to test this claim.  (Assume $\sigma = 10$ is known)
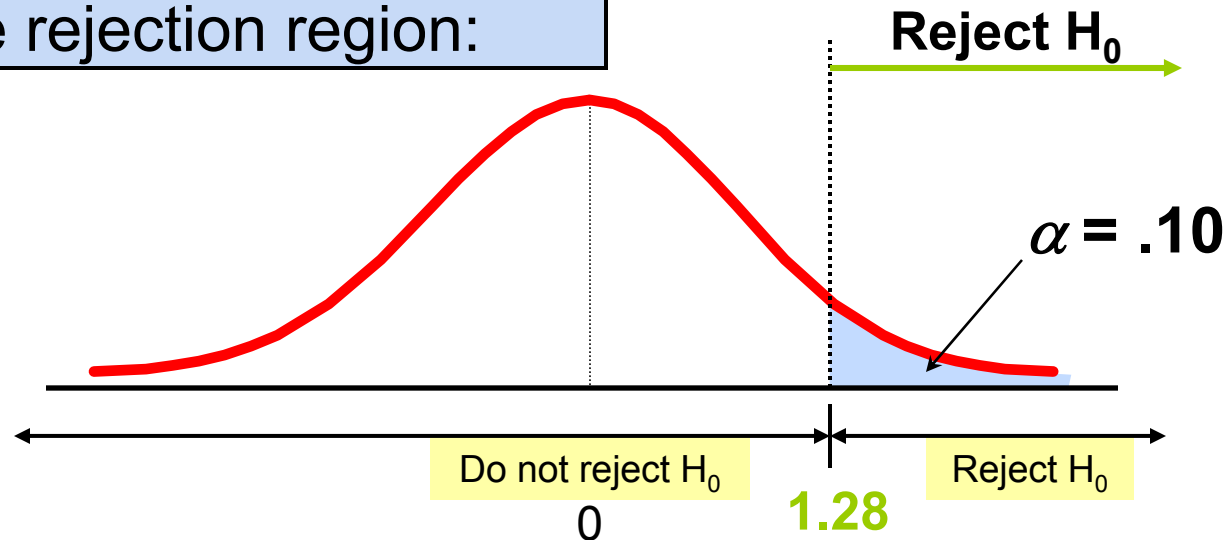
Form hypothesis test:

$H_0$: $\mu \leq 52$     the average is not over $52 per month

$H_1$: $\mu > 52$     the average **is** greater than $52 per month
(i.e., sufficient evidence exists to support the manager's claim)

# Example: Find Rejection Region

- Suppose that $\alpha = .10$ is chosen for this test

Find the rejection region:

**Reject $H_0$**

$\alpha = .10$

Do not reject $H_0$     Reject $H_0$

0     **1.28**

$$\text{Reject } H_0 \text{ if } z = \frac{\overline{x} - \mu_0}{\sigma/\sqrt{n}} > 1.28$$

K. Dral
for Finance

87

# Example: Sample Results

Obtain sample and compute the test statistic

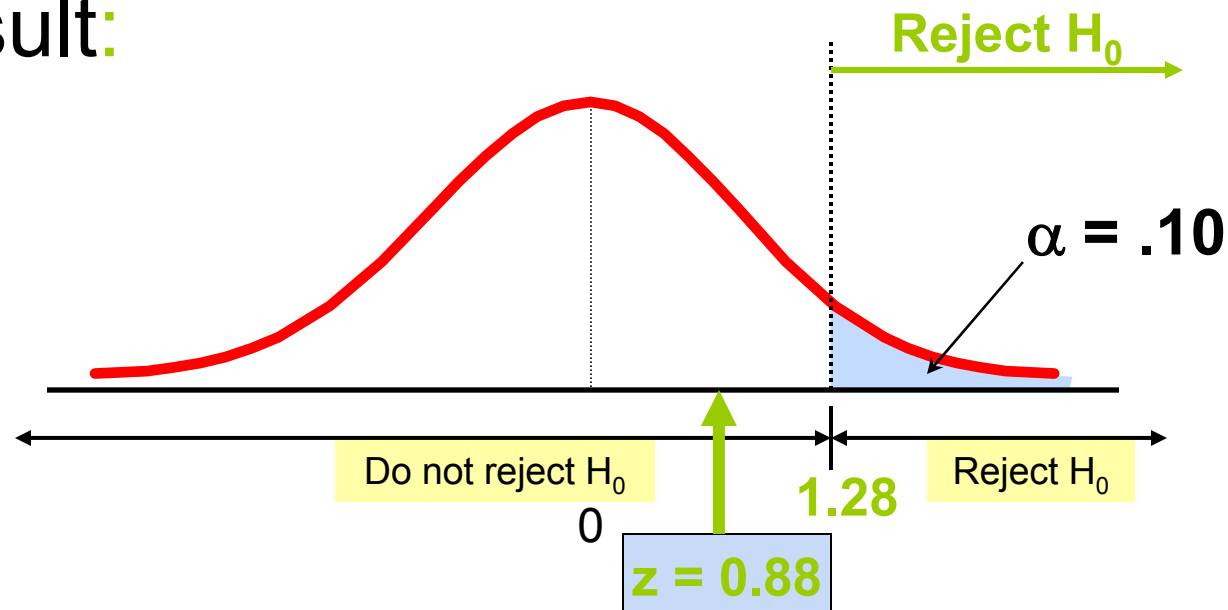Suppose a sample is taken with the following results: n = 64, $\bar{x}$ = 53.1 ($\sigma$=10 was assumed known)

Using the sample results**,**

$$z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{53.1 - 52}{\frac{10}{\sqrt{64}}} = 0.88$$

# Example: Decision

Reach a decision and interpret the result:

Reject $H_0$

$\alpha = .10$

Do not reject $H_0$

0

1.28

Reject $H_0$

z = 0.88

**Do not reject $H_0$ since z = 0.88 < 1.28**
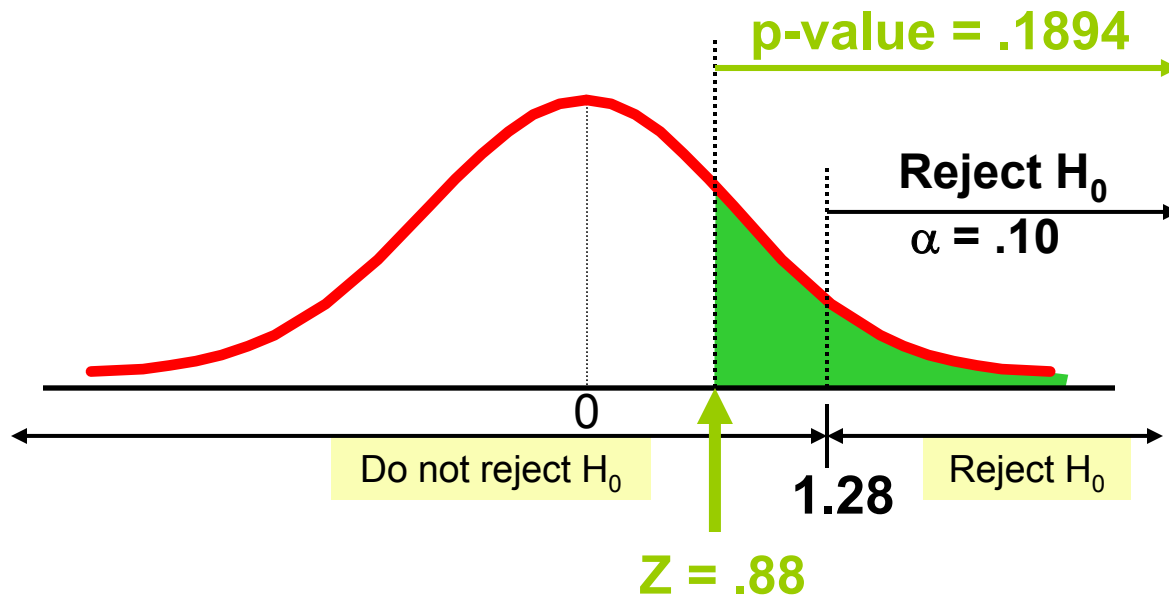
i.e.: there is not sufficient evidence that the mean bill is over $52

for Finance

# Example: p-Value Solution

Calculate the p-value and compare to $\alpha$

(assuming that $\mu = 52.0$)



**p-value = .1894**

**Reject $H_0$**
$\alpha = .10$

0

Do not reject $H_0$     **1.28**     Reject $H_0$

**Z = .88**

$$P(\overline{x} \geq 53.1 | \mu = 52.0)$$

$$= P\left( z \geq \frac{53.1 - 52.0}{10/\sqrt{64}} \right)$$

$$= P(z \geq 0.88) = 1 - .8106$$

$$= .1894$$

**Do not reject $H_0$ since p-value = .1894 > $\alpha$ = .10**

# One-Tail Tests

- In many cases, the alternative hypothesis focuses on one particular direction

$H_0: \mu \leq 3$

$H_1: \mu > 3$

⟶ This is an upper-tail test since the alternative hypothesis is focused on the upper tail above the mean of 3
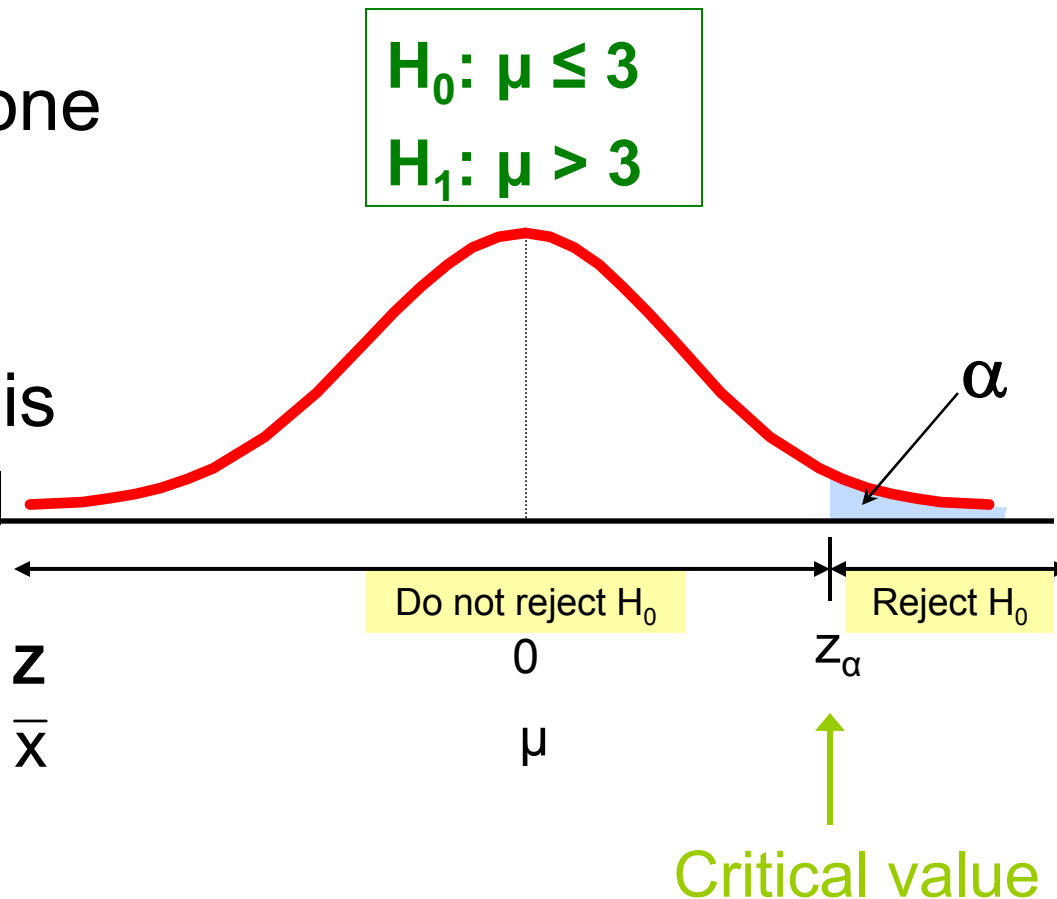
$H_0: \mu \geq 3$

$H_1: \mu < 3$

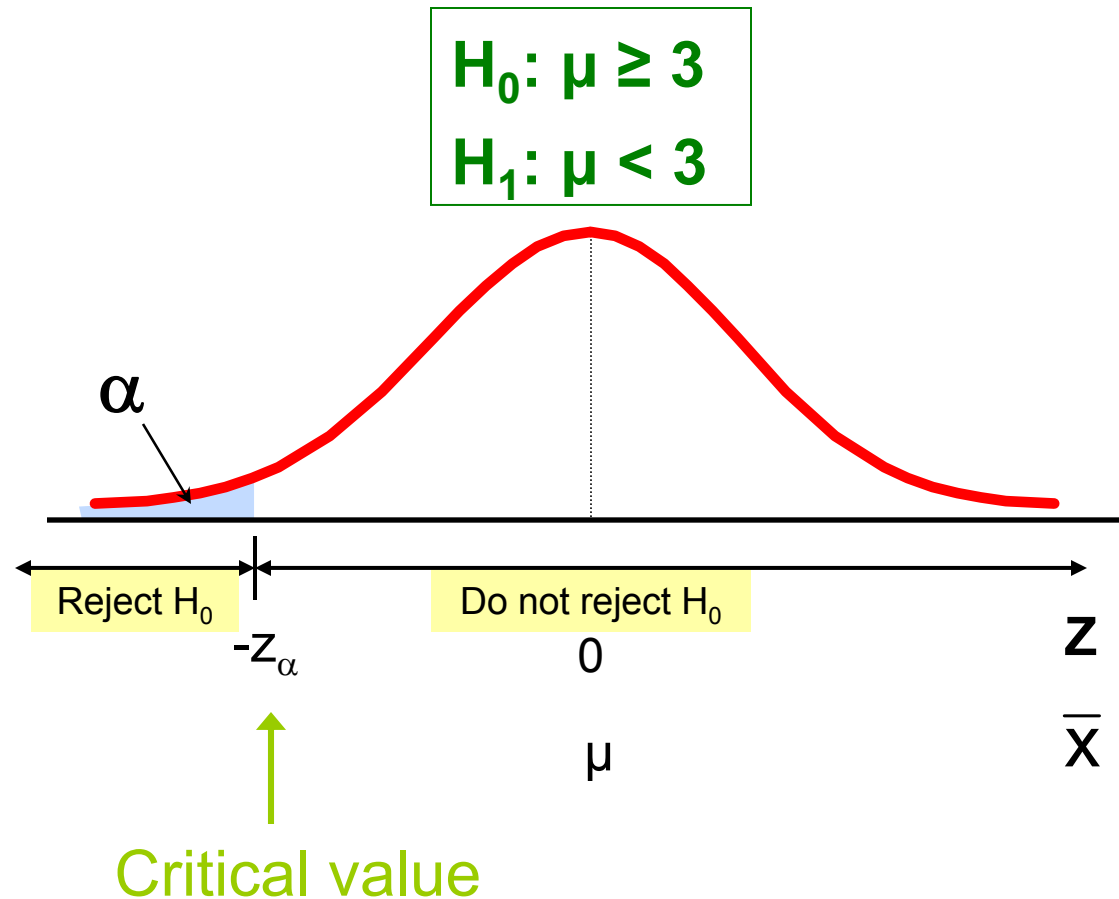⟶ This is a lower-tail test since the alternative hypothesis is focused on the lower tail below the mean of 3

# Upper-Tail Tests

- There is only one critical value, since the rejection area is in only one tail

$H_0$: $\mu \leq 3$

$H_1$: $\mu > 3$



$\alpha$

Do not reject $H_0$ | Reject $H_0$

z

0     $z_\alpha$

$\overline{x}$

$\mu$

Critical value

# Lower-Tail Tests

- There is only one critical value, since the rejection area is in only one tail

$H_0: \mu \geq 3$

$H_1: \mu < 3$

$\alpha$

Reject $H_0$

Do not reject $H_0$

$-z_\alpha$

0

$\mu$
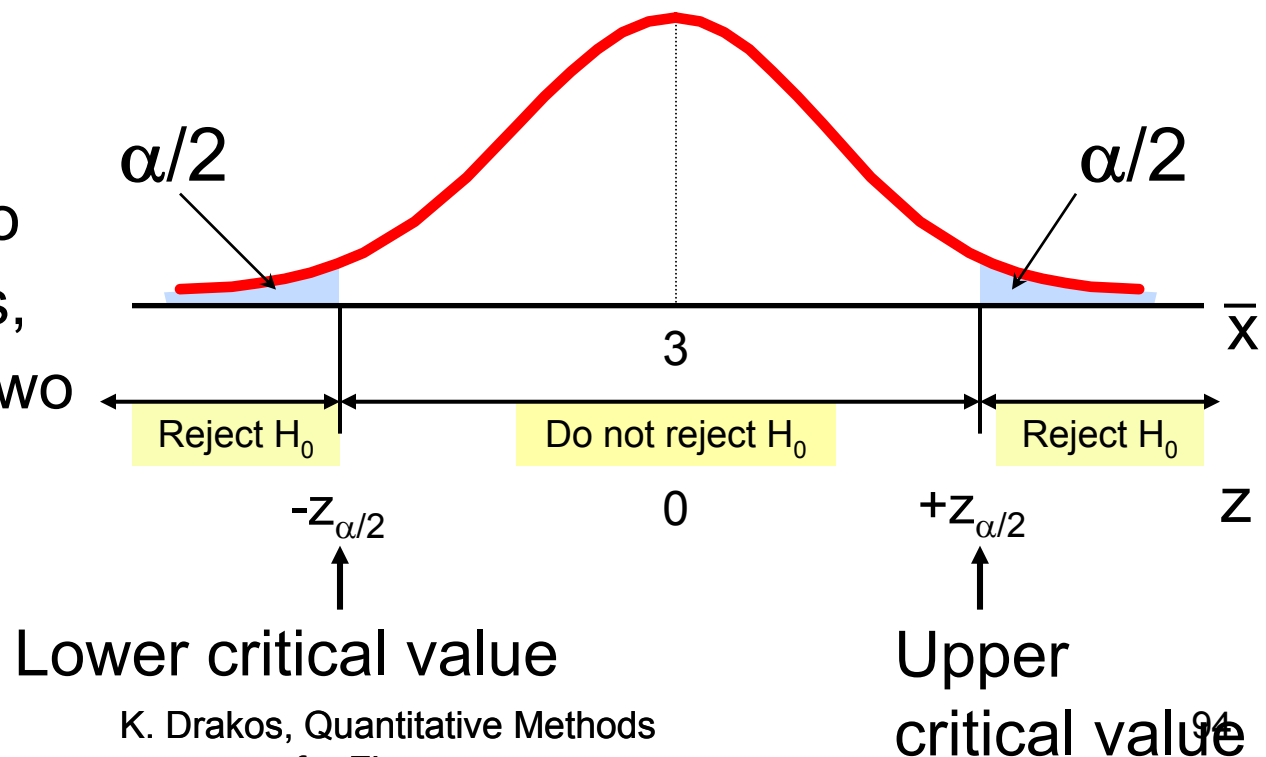
Z

$\overline{X}$

Critical value

# Two-Tail Tests

- In some settings, the alternative hypothesis does not specify a unique direction

$$H_0: \mu = 3$$
$$H_1: \mu \neq 3$$

- There are two critical values, defining the two regions of rejection

$\alpha/2$

$\alpha/2$

$\bar{X}$

3

Reject $H_0$

Do not reject $H_0$

Reject $H_0$

$-z_{\alpha/2}$

0

$+z_{\alpha/2}$

$z$

Lower critical value

Upper critical value

95

# Hypothesis Testing Example

**Test the claim that the true mean # of TV sets in US homes is equal to 3. (Assume σ = 0.8)**

- State the appropriate null and alternative hypotheses
  - $H_0: \mu = 3$ , $H_1: \mu \neq 3$    (This is a two tailed test)
- Specify the desired level of significance
  - Suppose that $\alpha$ = .05 is chosen for this test
- Choose a sample size
  - Suppose a sample of size n = 100 is selected

# Hypothesis Testing Example

- Determine the appropriate technique
  - $\sigma$ is known so this is a  z  test
- Set up the critical values
  - For $\alpha$ = .05 the critical  z  values are ±1.96
- Collect the data and compute the test statistic
  - Suppose the sample results are

    n = 100,  $\overline{x}$ = 2.84  ($\sigma$ = 0.8 is assumed known)
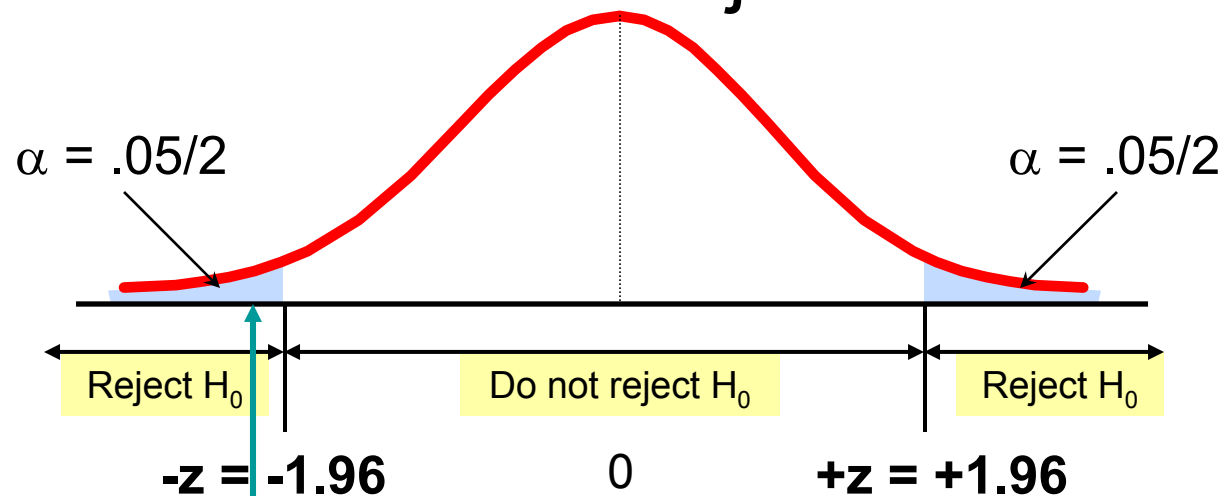
  So the test statistic is:

$$z = \frac{\overline{X} - \mu_0}{\dfrac{\sigma}{\sqrt{n}}} = \frac{2.84 - 3}{\dfrac{0.8}{\sqrt{100}}} = \frac{-.16}{.08} = -2.0$$

# Hypothesis Testing Example

- Is the test statistic in the rejection region?

Reject $H_0$ if $z < -1.96$ or $z > 1.96$; otherwise do not reject $H_0$

$\alpha = .05/2$

$\alpha = .05/2$

Reject $H_0$

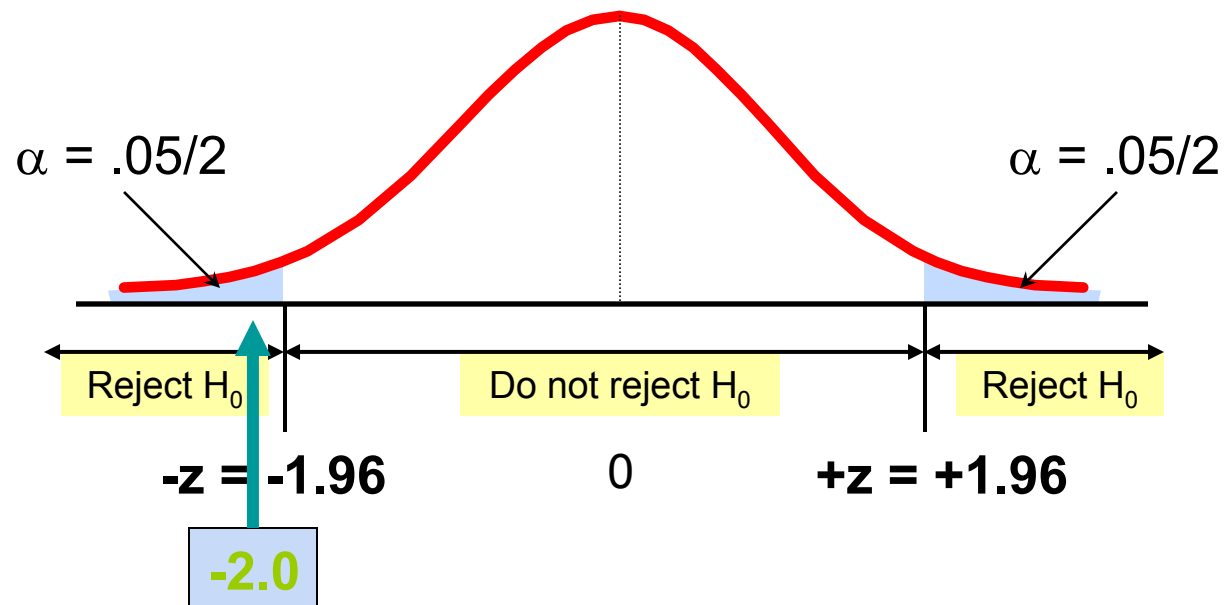Do not reject $H_0$

Reject $H_0$

**-z = -1.96**

0

**+z = +1.96**

Here, z = -2.0 < -1.96, so the test statistic is in the rejection region

# Hypothesis Testing Example

- Reach a decision and interpret the result

$\alpha = .05/2$                                            $\alpha = .05/2$

Reject $H_0$          Do not reject $H_0$          Reject $H_0$

$-z = -1.96$          0          $+z = +1.96$

**-2.0**

Since $z = -2.0 < -1.96$, we <u>reject the null hypothesis</u> and conclude that there is sufficient evidence that the mean number of TVs in US homes is not equal to 3

# Example: p-Value

- Example: How likely is it to see a sample mean of 2.84 (or something further from the mean, in either direction) if the true mean is $\mu = 3.0$?
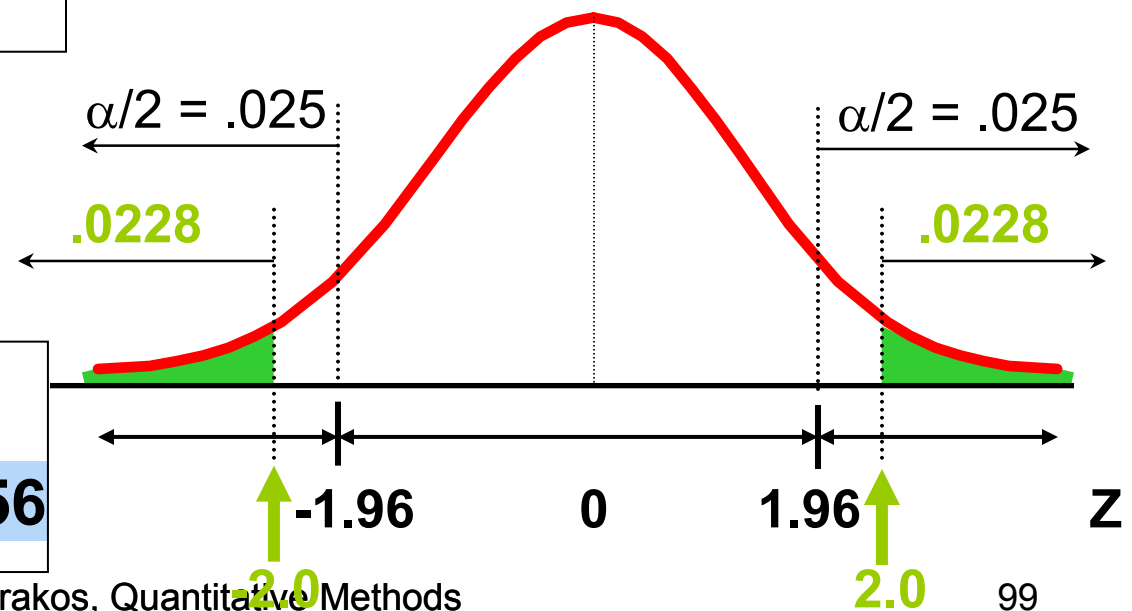
$\overline{x} = 2.84$ is translated to a z score of z = -2.0

$P(z < -2.0) = .0228$

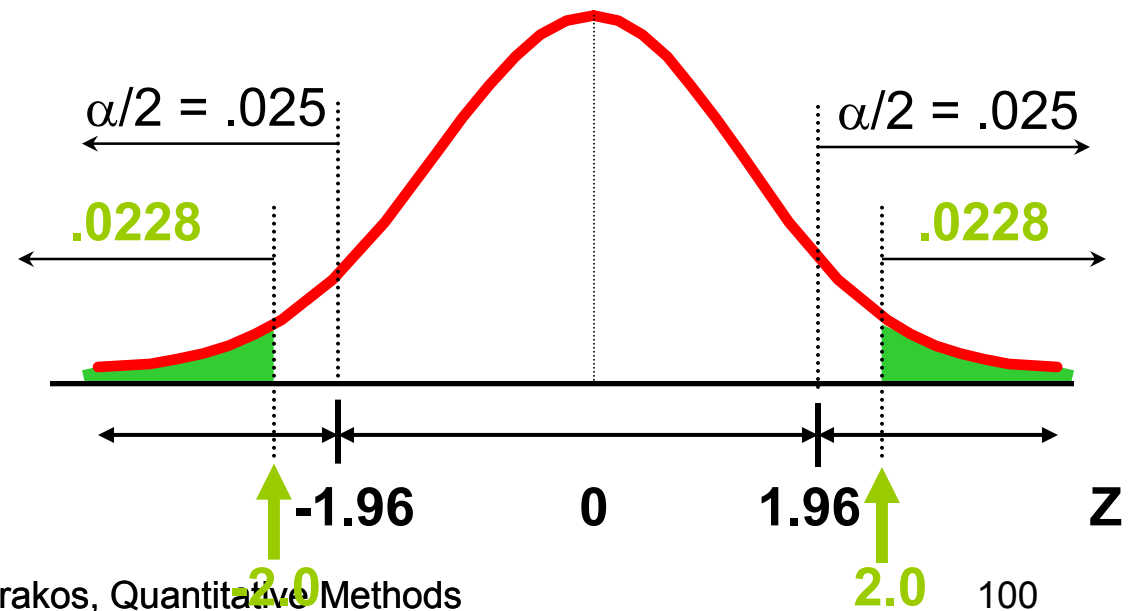$P(z > 2.0) = .0228$

**p-value**

**= .0228 + .0228 = .0456**

$\alpha/2 = .025$

$\alpha/2 = .025$

.0228

.0228

-1.96

0

1.96

Z

-2.0

2.0

# Example: p-Value

- Compare the p-value with $\alpha$

  – If p-value $<$ $\alpha$ , reject $H_0$

  – If p-value $\geq$ $\alpha$ , do not reject $H_0$

Here:  p-value = .0456
$\alpha$ = .05

**Since .0456 < .05, we reject the null hypothesis**

$\alpha/2 = .025$

$\alpha/2 = .025$

.0228

.0228

-1.96

0

1.96

Z

-2.0

2.0

# t Test of Hypothesis for the Mean (σ Unknown)

- Convert sample result ( $\bar{x}$ ) to a  t  test statistic

**Hypothesis Tests for $\mu$**

**σ Known**

**σ Unknown**

Consider the test

$$H_0 : \mu = \mu_0$$
$$H_1 : \mu > \mu_0$$

(Assume the population is normal)

The decision rule is:

Reject $H_0$ if $\quad t = \dfrac{\bar{x} - \mu_0}{\dfrac{s}{\sqrt{n}}} > t_{n-1, \alpha}$

# t Test of Hypothesis for the Mean (σ Unknown)

- ## For a two-tailed test:

Consider the test

$$H_0 : \mu = \mu_0$$
$$H_1 : \mu \neq \mu_0$$

(Assume the population is normal, and the population variance is unknown)

The decision rule is:

Reject $H_0$ if $\quad t = \dfrac{\bar{x} - \mu_0}{\dfrac{s}{\sqrt{n}}} < -t_{n-1,\,\alpha/2}\quad$ or if $\quad t = \dfrac{\bar{x} - \mu_0}{\dfrac{s}{\sqrt{n}}} > t_{n-1,\,\alpha/2}$

# Example: Two-Tail Test (σ Unknown)

The average cost of a 5-star hotel room in Athens is said to be 168 euros per night. A random sample of 25 hotels resulted in $\overline{x}$ = 172.50 euros and

s = 15.40 euros. Test at the

$\alpha$ = 0.05 level.

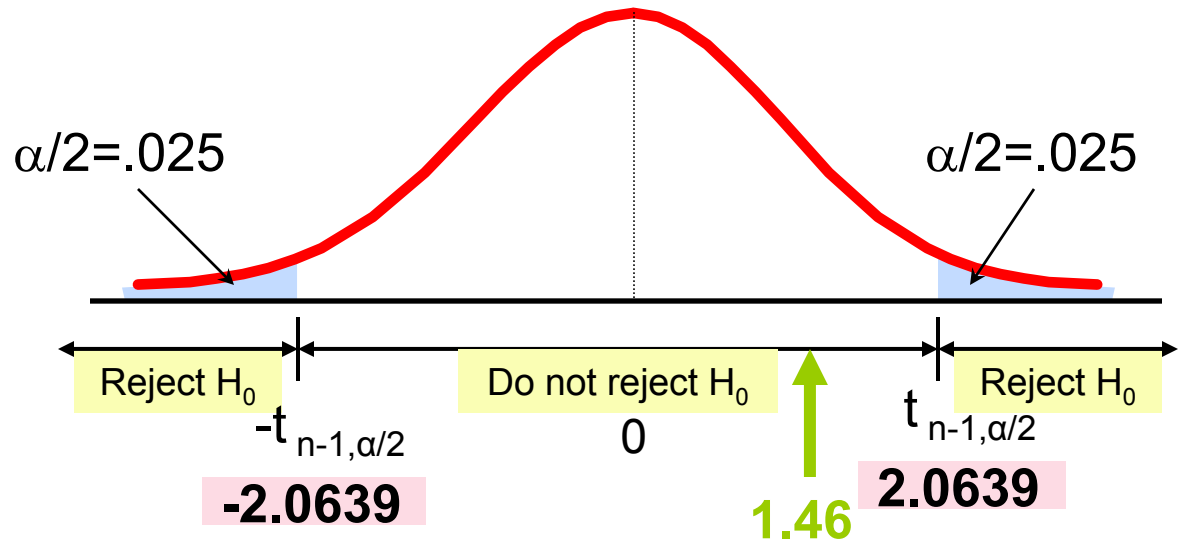(Assume the population distribution is normal)

$H_0$: μ = 168

$H_1$: μ ≠ 168

# Example Solution: Two-Tail Test

$H_0: \mu = 168$

$H_1: \mu \neq 168$

- n = 25
- a = 0.05
- σ is unknown, so use a t statistic
- Critical Value:

$t_{24,\,0.025} = \pm 2.0639$



$\alpha/2 = .025$

$\alpha/2 = .025$

Reject $H_0$

Do not reject $H_0$

Reject $H_0$

$-t_{n-1,\alpha/2}$

0

$t_{n-1,\alpha/2}$

-2.0639

1.46

2.0639

$$t_{n-1} = \frac{\overline{x} - \mu}{\dfrac{s}{\sqrt{n}}} = \frac{172.50 - 168}{\dfrac{15.40}{\sqrt{25}}} = 1.46$$

**Do not reject $H_0$:** not sufficient evidence that true mean cost is different than $168