

# Confidence Intervals

# Midterm exam results

	p1 (10)	p2 (10)	P3 (15)	P4(10)	P5 (10)	P6 (10)
Average	12	9.96	14.4	9.9	9.8	7.58
St.Dev	3.1	0.2	1.21	0.3	0.7	2.56
	<b>3.04</b>	0.04	<b>0.58</b>	0.1	0.2	<b>2.42</b>

**(15 points)** If the letters of ILLINI are randomly ordered, all orderings being equally likely, what is the probability that not a single position has the same letter as in the original order? Hint: 3 Is (and 2 Ls) are identical.

**Answer:** Three letters I must go in places of L, L, and N. Once I pick where N goes, the rest is determined. There are 3 places to put N. There are 3 solutions. The total number of ways to order these 6 letters is  $6!/(3!*2!*1!)=6*5*4/2=60$ . Hence the probability is  $3/60=1/20=0.05$

6. **(10 points)** In a data communication system, several messages that arrive at a node are bundled into a packet before they are transmitted over the network. Assume the messages arrive according to a Poisson process with the mean rate equal to one message per two minutes. Five messages are required to form a packet and the packet is formed immediately after the last message has arrived.

**(a) (5 points)** What is the probability that a time interval between two consecutive messages is longer than 4 minutes?

**Answer:**  $\lambda = 1 \text{ message} / 2 \text{ minutes} = 0.5 \text{ messages/minute}$ .  
Exponential distribution  $P(X > 4) = \exp(-0.5 * 4) = \exp(-2) = 0.1353$

**(b) (5 points)** What is the mean time until a packet is formed, that is, until exactly five messages have arrived at the node?

**Answer:** Using Erlang distribution with  $r=5$ ,  $\lambda=0.5$  one gets  $(5/0.5) \text{ minutes} = 10 \text{ minutes}$

3. **(15 points)** Sequencing technologies can “read” many short fragments (simply called reads) from a genome. Given that the process through which the read sequences are generated is random, it is possible that certain parts of the genome will remain uncovered unless an impractical amount of sequences are generated. Human genome is  $3 \times 10^9$  bp long. A patient’s genome has been sequenced and it is randomly covered by  $300 \times 10^5$  reads (each read is 100 bp long). We assume that the number of times a base in the human genome is covered follows a Poisson distribution.

**(a) (5 points)** What is the probability that a particular base is not covered by any reads?

**Answer:**  $\lambda = \frac{300 \times 10^5 \times 100}{3 \times 10^9} = 1$        $P(X = 0) = e^{-1} = 0.3679$

**(b) (5 points)** What is the probability that a particular base is covered by at least two reads?

**Answer**  $P(X \geq 2) = 1 - P(X = 0) - P(X = 1) = 1 - e^{-1} - \lambda e^{-1} = 1 - 2e^{-1} = 0.2642$

**(c) (5 points)** We now start randomly selecting bases in our genome. What is the expected number of bases we have to look at before exactly 4 such uncovered bases are identified?

**Answer:**  $n_{bases} = 4 / 0.3679 = 10.8725$

# Two-sided confidence intervals

- Want to make a two-sided confidence interval of population average  $\mu$  based on the sample  $x_1, x_2, \dots, x_n$  and its sample mean  $\bar{x}$
- Assume population standard deviation  $\sigma$  is known
- Characterized by:
  - lower- and upper- confidence limits  $L$  and  $R$
  - the confidence coefficient  $1-\alpha$
- Find  $L$  and  $R$  such that
  - $\text{Prob}(\mu > R) = \alpha/2$
  - $\text{Prob}(\mu < L) = \alpha/2$
  - Therefore,  $\text{Prob}(L < \mu < R) = 1-\alpha$
- For a one-sided confidence interval, say, upper bound of  $\mu$ , find  $R$  that  $\text{Prob}(\mu > R) = \alpha$

Based on confidence level  
 $1-\alpha$  select  $z_{\alpha/2}$  such that

From CLT

$1-\alpha =$

$= \text{Prob} \left( \mu - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \bar{x} < \mu + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right) =$   
subtract  $-\bar{x} - \mu$

$= \text{Prob} \left( -\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right) =$   
invert sign

$= \text{Prob} \left( \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} > \mu > \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$

For  $1-\alpha = 0.95$  Confidence Interval  
 $z_{\alpha/2} = 1.96$  on  $\mu$

# Matlab exercise

- 1000 labs measured average P53 gene expression using  $n=20$  samples drawn from the Gaussian distribution with  $\mu=3$ ;  $\sigma=2$ ;
- Each lab found 95% confidence estimates of the population mean  $\mu$  **based on its sample only**
- Count the number of labs, where the population mean lies **outside their bounds**
- You should get  $\sim 50$  labs out of 1000 labs

# How I did it

- `n=20; k_labs=1000;`
- `rand_table=2.*randn(n,k_labs)+3;`
- `sample_mean=mean(rand_table,1);`
- `CI_low=sample_mean-1.96.*2./sqrt(n);`
- `CI_high=sample_mean+1.96.*2./sqrt(n);`
- `k_above=sum(3>CI_high)`
- `k_below=sum(3<CI_low)`
- `figure; ndisp=100; errorbar(1:ndisp,  
sample_mean(1:ndisp),  
ones(ndisp,1).*1.96.*2./sqrt(n),'ko');`
- `hold on; plot(1:ndisp, 3.*ones(ndisp,1),'r-');`

## 8-2 Confidence Interval on the Mean of a Normal Distribution, Variance Known

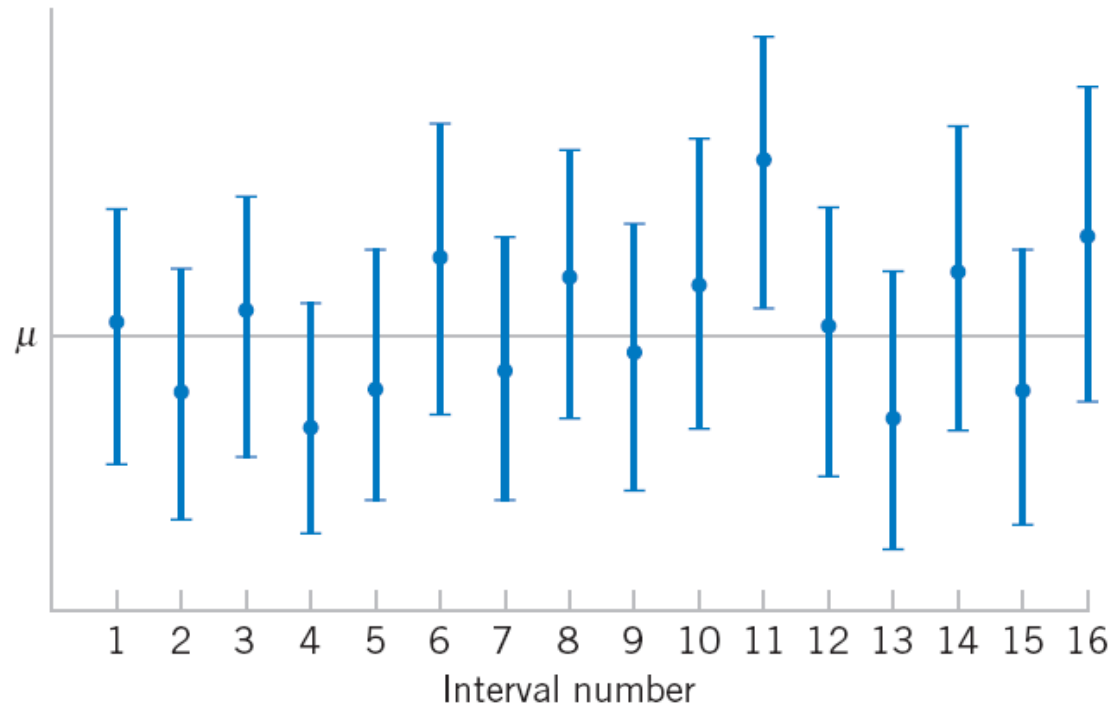


Figure 8-1 Repeated construction of a confidence interval for  $\mu$ .

Figure 8-1 Repeated construction of a confidence interval for  $\mu$ .

So far in estimating  
confidence intervals for population mean  $\mu$   
we assumed that the population variance  $\sigma^2$   
**is known**

Then (or when  $n \gg 1$ , say 20 and above)  
**one can use the Normal Distribution**  
to calculate confidence intervals

Q: What to do if the sample is small and population variance is not known?

A: Use the sample variance

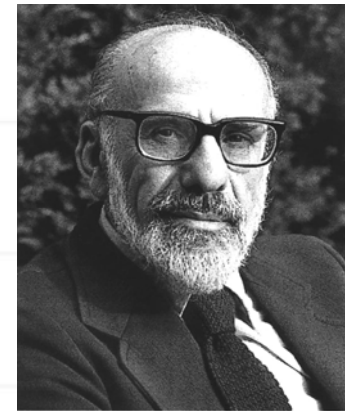
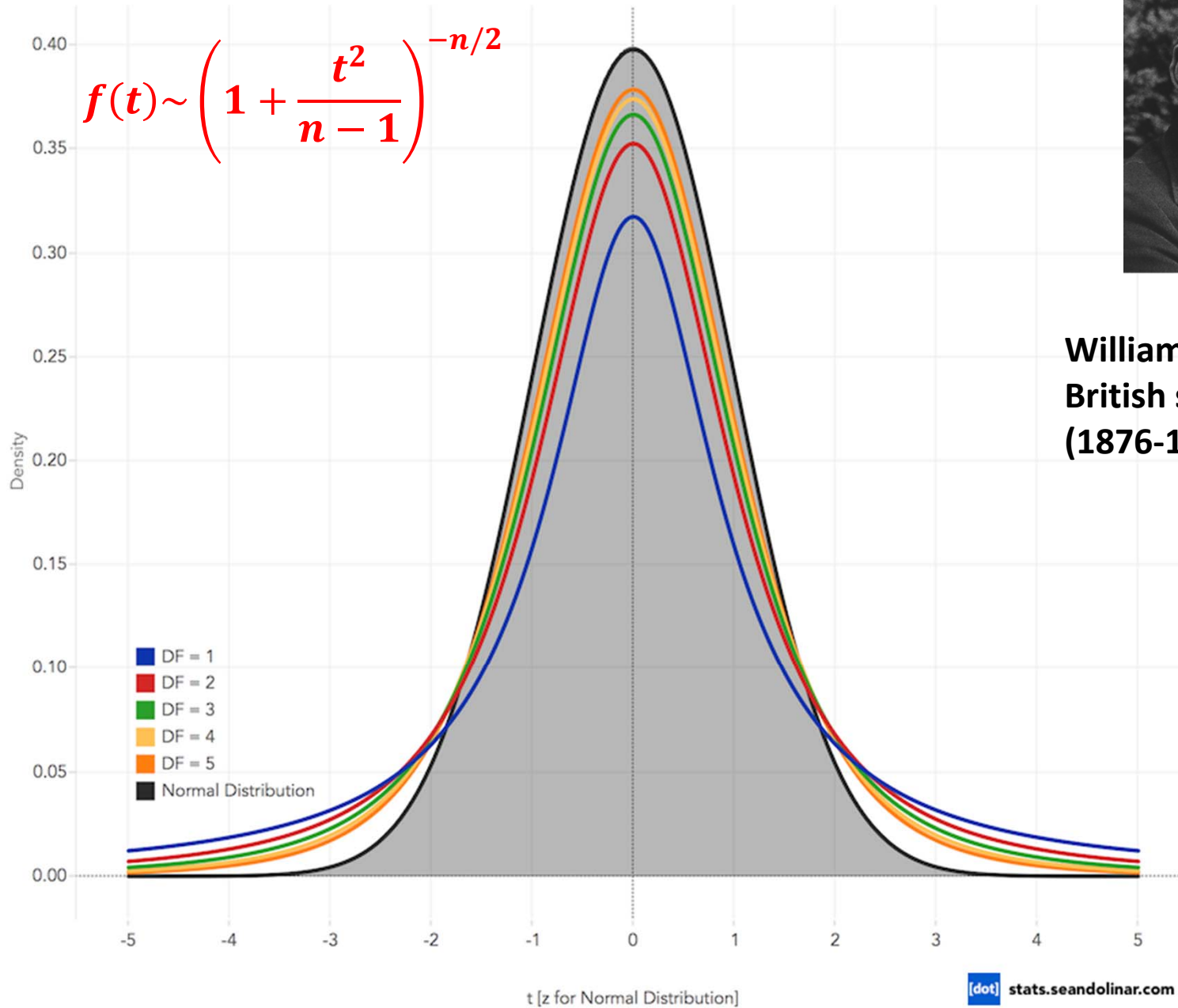
$$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 \text{ but carefully}$$

for small samples:

- Variable  $X$  has to be normally distributed
- t-distribution has to be used instead of the normal distribution.

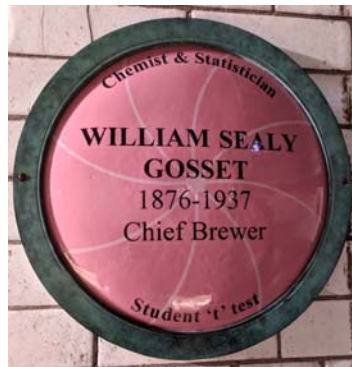
# Student's t-distribution

t-Distribution vs. Normal Distribution



**William Sealy Gosset**  
British statistician  
(1876-1937)

Another researcher at Guinness had previously published a paper containing trade secrets of the Guinness brewery. To prevent further disclosure of confidential information, Guinness prohibited its employees from publishing any papers regardless of the contained information. However, after pleading with the brewery and explaining that his mathematical and philosophical conclusions were of no possible practical use to competing brewers, he was allowed to publish them, but under a pseudonym ("Student"), to avoid difficulties with the rest of the staff. Thus his most noteworthy achievement is now called Student's, rather than Gosset's, t-distribution.



Gosset had almost all his papers including “The probable error of a mean” published in Pearson's journal *Biometrika* under the pseudonym Student

# Play with Mathematica notebook

<http://demonstrations.wolfram.com/ComparingNormalAndStudentsTDistributions/>

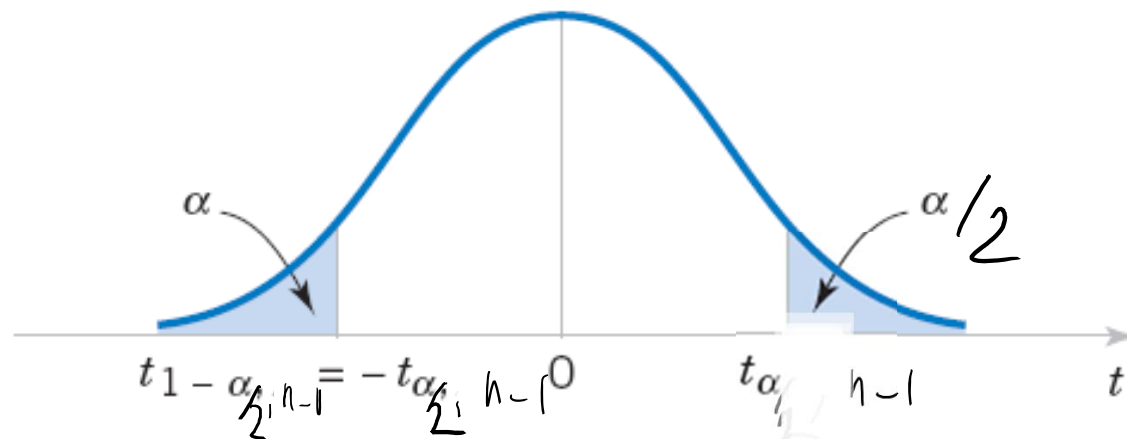
By Gary McClelland

# 8-3 Confidence Interval on the Mean of a Normal Distribution, Variance Unknown

---

## 8-3.1 Student's $t$ distribution

$$f(t) \sim \left(1 + \frac{t^2}{n-1}\right)^{-n/2}$$



**Figure 8-5** Percentage points of the  $t$  distribution.

# 8-3 Confidence Interval on the Mean of a Normal Distribution, Variance Unknown

---

## 8-3.2 The $t$ Confidence Interval on $\mu$

(Eq. 8-16)

If  $\bar{x}$  and  $s$  are the mean and standard deviation of a random sample from a normal distribution with unknown variance  $\sigma^2$ , a **100(1 -  $\alpha$ )% confidence interval on  $\mu$**  is given by

$$\bar{x} - t_{\alpha/2, n-1}s/\sqrt{n} \leq \mu \leq \bar{x} + t_{\alpha/2, n-1}s/\sqrt{n} \quad (8-16)$$

where  $t_{\alpha/2, n-1}$  is the upper 100 $\alpha$ /2 percentage point of the  $t$  distribution with  $n - 1$  degrees of freedom.

**One-sided confidence bounds** on the mean are found by replacing  $t_{\alpha/2, n-1}$  in Equation 8-16 with  $t_{\alpha, n-1}$ .

## Confidence intervals for population variance $\sigma^2$

Should be around

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

## 8-4 Confidence Interval on the Variance and Standard Deviation of a Normal Distribution

### Definition

(Eq. 8-17)

Let  $X_1, X_2, \dots, X_n$  be a random sample from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and let  $S^2$  be the sample variance. Then the random variable

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2} \quad (8-17)$$

has a chi-square ( $\chi^2$ ) distribution with  $n - 1$  degrees of freedom.

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

## 8-4 Confidence Interval on the Variance and Standard Deviation of a Normal Distribution

$$X = (n-1)S^2/\sigma^2$$

We know  $n, S^2$

want to estimate  $\sigma^2$

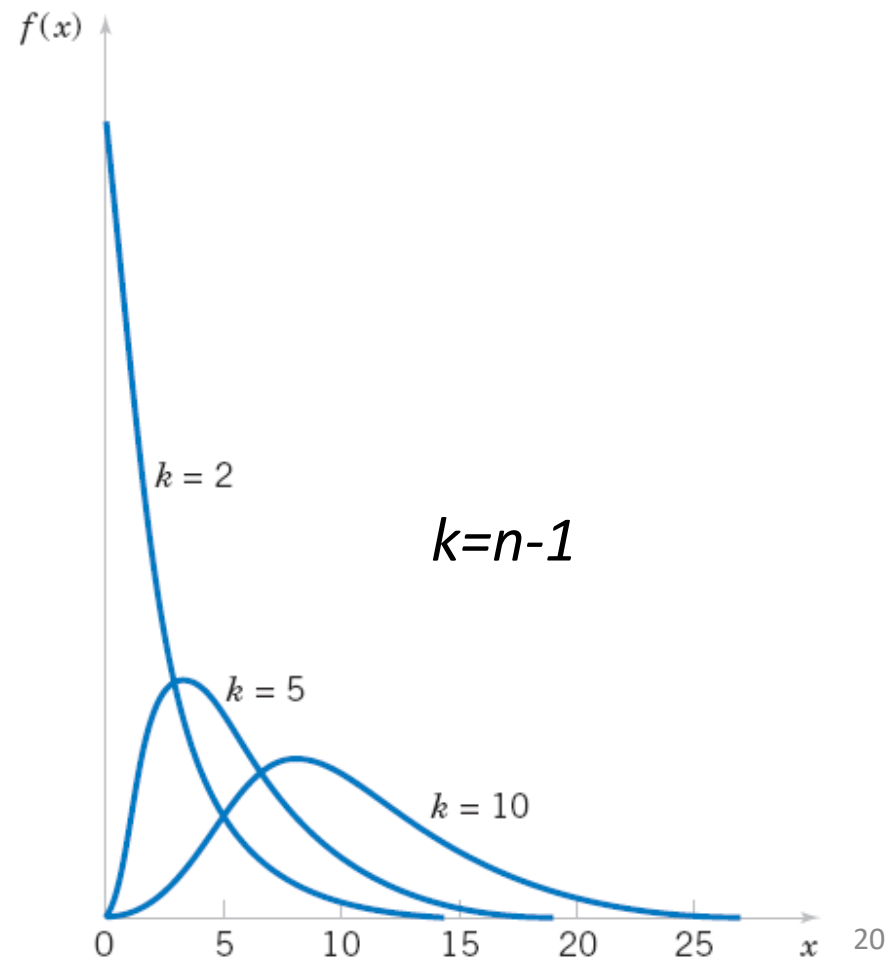
$$f(x, n) \sim x^{(n-1)/2-1} \exp(-x/2)$$

It is just Gamma PDF  
with  $r = (n-1)/2$ , and  $\lambda = 1/2$

Mean value:  
 $n-1$

Standard deviation:

$$\sqrt{2(n-1)}$$



# Play with Mathematica notebook

<http://demonstrations.wolfram.com/ChiSquaredDistributionAndTheCentralLimitTheorem/>

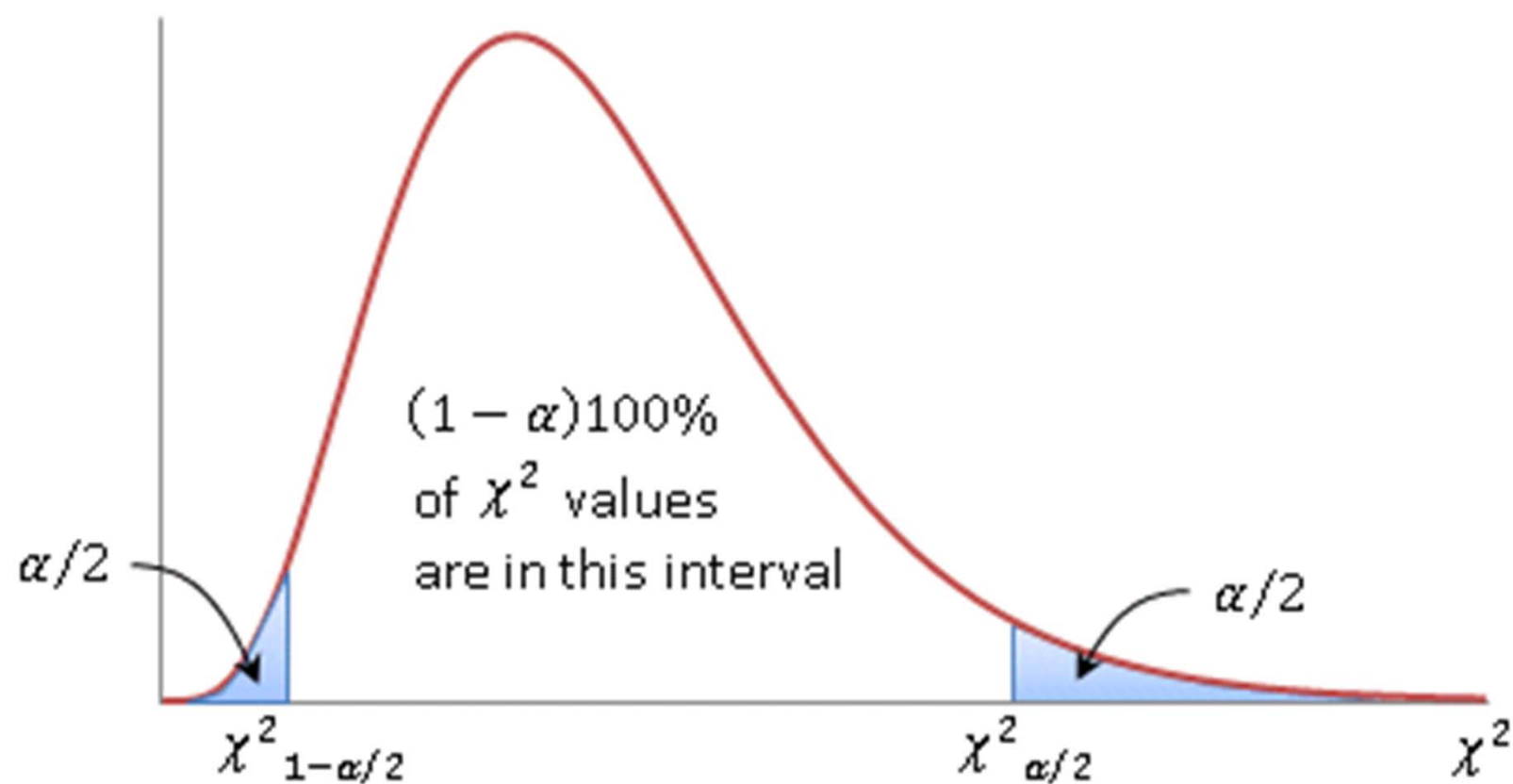
By Peter Falloon

# Matlab exercise

- Generate 100,000 experiments.  
Each experiment generates a sample with  $n=8$ , made out of normal variable with  $\sigma=5$ .
- For each sample calculate sample variance:  $s^2$
- Plot PDF-histogram of  $(n-1) s^2 / \sigma^2$  for 100,000 experiments
- Compare with Matlab function `chi2pdf(x,n-1)`

# Matlab exercise: solution

- **Stats=100000; n = 8;**
- **X = 5 \* randn([n, Stats]);**
- **ch2 = (n-1) \* var(X)/25;**
- **histogram(ch2,0:0.1:30,'Normalization','pdf')**
- **hold on**
- **plot( (0:0.1:30), chi2pdf((0:0.1:30), n-1),'r-')**



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

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

## 8-4 Confidence Interval on the Variance and Standard Deviation of a Normal Distribution

---

### Definition

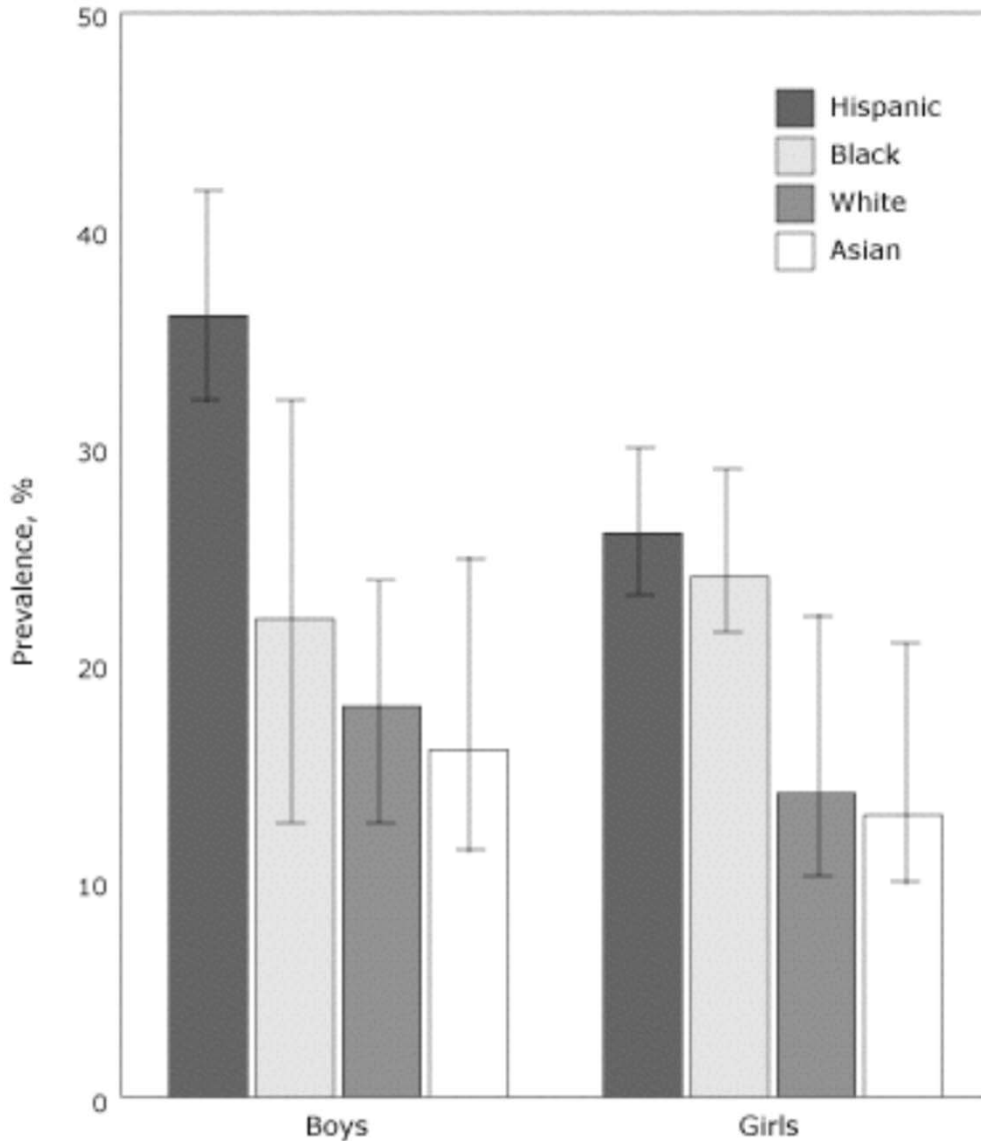
(Eq. 8-19)

If  $s^2$  is the sample variance from a random sample of  $n$  observations from a normal distribution with unknown variance  $\sigma^2$ , then a **100(1 -  $\alpha$ )% confidence interval on  $\sigma^2$**  is

$$\frac{(n - 1)s^2}{\chi_{\alpha/2, n-1}^2} \leq \sigma^2 \leq \frac{(n - 1)s^2}{\chi_{1-\alpha/2, n-1}^2} \quad (8-19)$$

where  $\chi_{\alpha/2, n-1}^2$  and  $\chi_{1-\alpha/2, n-1}^2$  are the upper and lower 100 $\alpha$ /2 percentage points of the chi-square distribution with  $n - 1$  degrees of freedom, respectively. A **confidence interval for  $\sigma$**  has lower and upper limits that are the square roots of the corresponding limits in Equation 8-19.

# Confidence estimates of the population proportion



Prevalence (with 95% CI bars) of obesity among New York City public elementary schoolchildren, by sex and race/ethnicity, 2003.

(source: CDC.GOV)

Collect a sample of BMI values  
Obese means  $BMI > 30$

What do those bars actually mean?

# Large sample confidence estimate of population proportion

- Want to know the **fraction  $p$  of the population** that belongs to a class, e.g. the class “obese” kids defined by  $BMI > 30$ .
- Each variable is a Bernoulli trial with one parameter  $p$ . We can use **moments** or **MLE estimator** to estimate  $p$
- Both give the same estimate: **sample fraction  $\hat{P} = (\# \text{ of obese kids in the sample}) / (\text{sample size } n)$**
- How to put confidence bounds on  $p$  based on  $\hat{P}$
- # of obese kids in the sample follows the binomial distribution: “success” = sampled kid is obese : -(  
 $p$  – probability of success,  $1-p$  – failure
- Expected # of successes is  $np$   $\rightarrow$  Expected fraction of successes is  $p$
- Standard deviation of # of successes is  $\sqrt{np(1-p)}$   $\rightarrow$   
Standard deviation of fraction of successes is  $\sqrt{p(1-p)/n}$

# 8-5 A Large-Sample Confidence Interval For a Population Proportion

---

## Normal Approximation for Binomial Proportion

If  $n$  is large, the distribution of

$$Z = \frac{X - np}{\sqrt{np(1-p)}} = \frac{\hat{P} - p}{\sqrt{\frac{p(1-p)}{n}}} \approx \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}$$

is approximately standard normal.

The quantity  $\sqrt{\hat{p}(1-\hat{p})/n}$  is the standard error of the point estimator  $\hat{p}$ .

## 8-5 A Large-Sample Confidence Interval For a Population Proportion (Eq. 8-23)

---

If  $\hat{p}$  is the proportion of observations in a random sample of size  $n$  that belongs to a class of interest, an approximate  $100(1 - \alpha)\%$  confidence interval on the proportion  $p$  of the population that belongs to this class is

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \quad (8-23)$$

where  $z_{\alpha/2}$  is the upper  $\alpha/2$  percentage point of the standard normal distribution.

This interval is known as the Wald interval (Wald and Wolfowitz, 1939).

Did you know that M&M's<sup>®</sup> Milk Chocolate Candies are supposed to come in the following percentages: 24% blue, 20% orange, 16% green, 14% yellow, 13% red, 13% brown?

<http://www.scientificameriken.com/candy5.asp>

"To our surprise M&Ms met our demand to review their procedures in determining candy ratios. It is, however, noted that the figures presented in their email differ from the information provided from their website (<http://us.mms.com/us/about/products/milkchocolate/>). An email was sent back informing them of this fact. To which M&Ms corrected themselves with one last email:

In response to your email regarding M&M'S CHOCOLATE CANDIES

Thank you for your email.

On average, our new mix of colors for M&M'S<sup>®</sup> Chocolate Candies is:

M&M'S<sup>®</sup> Milk Chocolate: 24% blue, 20% orange, 16% green, 14% yellow, 13% red, 13% brown.

M&M'S<sup>®</sup> Peanut: 23% blue, 23% orange, 15% green, 15% yellow, 12% red, 12% brown.

M&M'S<sup>®</sup> Kids MINIS<sup>®</sup>: 25% blue, 25% orange, 12% green, 13% yellow, 12% red, 13% brown.

M&M'S<sup>®</sup> Crispy: 17% blue, 16% orange, 16% green, 17% yellow, 17% red, 17% brown.

M&M'S<sup>®</sup> Peanut Butter and Almond: 20% blue, 20% orange, 20% green, 20% yellow, 10% red, 10% brown.

Have a great day!

Your Friends at Masterfoods USA  
A Division of Mars, Incorporated



How to estimate these probabilities from a finite sample and how to set confidence interval on these estimates?

Did you know that M&M's® Milk Chocolate Candies are supposed to come in the following percentages: 24% blue, 20% orange, 16% green, 14% yellow, 13% red, 13% brown?

How large is a sample needed for 95% CI on the percentage of blue M&Ms to be less than +/- 4%  
Same question for red M&Ms?



Did you know that M&M's® Milk Chocolate Candies are supposed to come in the following percentages: 24% blue, 20% orange, 16% green, 14% yellow, 13% red, 13% brown?



How large is a sample needed for 95% CI on the percentage of blue M&Ms to be less than +/- 4%  
Same question for red M&Ms?

$$\text{For blue M\&Ms } p = 0.24$$
$$1.96 \sqrt{\frac{0.24(1-0.24)}{n}} < 0.04$$

$$n > \left(\frac{1.96}{0.04}\right)^2 0.24 \times (1-0.24) = 438 \text{ M\&Ms or}$$

~ 2 x 7oz bags with 210 candies each

$$\text{For red M\&Ms } p = 0.13$$

$$n > \left(\frac{1.96}{0.04}\right)^2 \times 0.13 \times (1-0.13) \approx 271 \text{ M\&Ms or}$$

~ 1 x 7oz bag