

Chapter 4 Continued - Test 2 Material

Inference on Two populations' variances

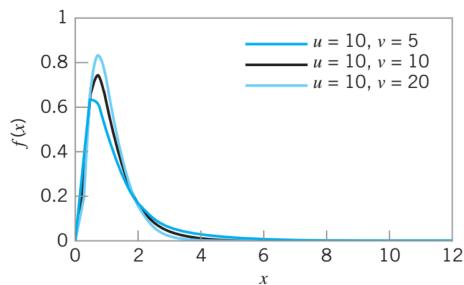
→ Two independent, normally distributed populations

$$\text{Null: } H_0: \sigma_1^2 = \sigma_2^2$$

$$\text{Alt: } H_a: \sigma_1^2 \neq \sigma_2^2 \quad \text{two-sided}$$

$$\text{Test Stat: } F_0 = \frac{s_1^2}{s_2^2} \leftarrow D.o.F_1 = n_1 - 1$$

Sensitive to normality – should check!



Relevance distribution: F-distribution

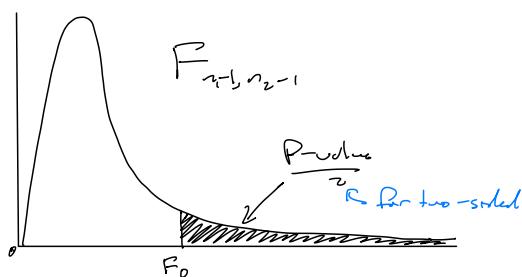
- (i) The curve is not symmetric (skewed right)
- (ii) D.o.F count curve for each set of D.o.Fs
- (iii) F statistic ≥ 0
- (iv) As D.o.F increase curve gets more normal

$$\text{Note: } F_{0.99, v_1, v_2} = 1 / F_{0.01, v_2, v_1}$$

Rejection Criterion: (two-sided)

P-value

$$P\text{-value} < \alpha$$



Summary

$$\text{Null hypothesis: } H_0: \sigma_1^2 = \sigma_2^2$$

$$\text{Test statistic: } F_0 = \frac{s_1^2}{s_2^2} \quad (10-31)$$

Alternative Hypotheses

Rejection Criterion

$$H_1: \sigma_1^2 \neq \sigma_2^2 \quad f_0 > f_{\alpha/2, n_1-1, n_2-1} \text{ OR } f_0 < f_{1-\alpha/2, n_1-1, n_2-1}$$

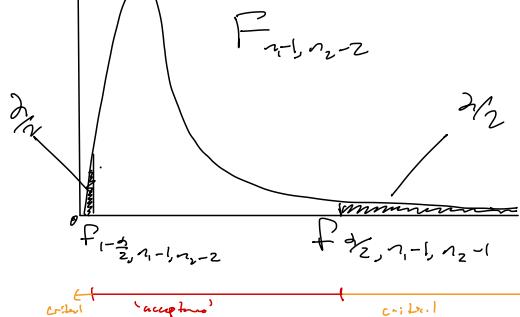
$$H_1: \sigma_1^2 > \sigma_2^2 \quad f_0 > f_{\alpha, n_1-1, n_2-1}$$

$$H_1: \sigma_1^2 < \sigma_2^2 \quad f_0 < f_{1-\alpha, n_1-1, n_2-1}$$

$$F_0 = F_{\alpha/2, n_1-1, n_2-1}$$

Fixed
Signif. F.
Level

$$\text{or } F_0 \leq F_{1-\alpha/2, n_1-1, n_2-1}$$



Assumptions

- The data are normally distributed
- Samples are independent

Sensitive to normality!

Example

Example 10-13 Semiconductor Etch Variability Oxide layers on semiconductor wafers are etched in a mixture of gases to achieve the proper thickness. The variability in the thickness of these oxide layers is a critical characteristic of the wafer, and low variability is desirable for subsequent processing steps. Two different mixtures of gases are being studied to determine whether one is superior in reducing the variability of the oxide thickness. Sixteen wafers are etched in each gas. The sample standard deviations of oxide thickness are $s_1 = 1.96$ angstroms and $s_2 = 2.13$ angstroms, respectively. Is there any evidence to indicate that either gas is preferable? Use a fixed-level test with $\alpha = 0.05$.

1) Parameter of interest: Variances oxide thickness for two mixtures of gases, unknown variances

• Two-sample test

$$2) H_0: \sigma_1^2 = \sigma_2^2$$

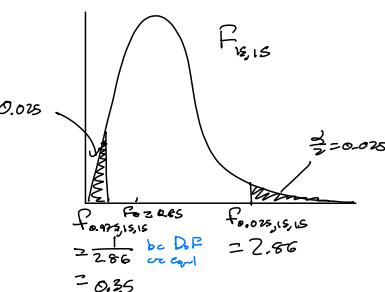
$$4) F_0 = \frac{s_1^2}{s_2^2} = \frac{1.96^2}{2.13^2} = 0.85$$

$$3) H_a: \sigma_1^2 \neq \sigma_2^2$$

$$5) \text{Rej. if: } F_0 > F_{0.025, 15, 15} \text{ or } F_0 < F_{0.975, 15, 15}$$

Minitab: Stat-> Basic Stat -> Two-sample Variance

Fail to reject H_0 , need to have enough evidence to reject the claim that the variances of the oxide layers produced by the different mixtures are different.



Single Factor Analysis of Variance

- Analysis of Variance: Used for comparing means when there are more than two levels in a single factor
 - Factor: Entity being tested
 - level: Settings within a factor

- Replicates: # of repeats taken at each level

- Randomization: Randomizing other data such that nuisance variables are eliminated

Factor levels	Tensile Strength of Paper (psi)						Totals	Averages
	Observations							
Hardwood Concentration (%)	1	2	3	4	5	6		
5	7	8	15	11	9	10	60	10.00
10	12	17	13	18	19	15	94	15.67
15	14	18	19	17	16	18	102	17.00
20	19	25	22	23	18	20	127	21.17
							383	15.96

HoR...
 Factors = 1
 Levels = 4
 Replicates = 6

Deriving and Development of ANOVA

Typical Data for a Single-Factor Experiment

Treatment	Observations						Totals	Averages
1	y_{11}	y_{12}	...	y_{1n}	$\bar{y}_{1\cdot}$	$\bar{y}_{\cdot 1}$		
2	y_{21}	y_{22}	...	y_{2n}	$\bar{y}_{2\cdot}$	$\bar{y}_{\cdot 2}$		
.		
.		
a	y_{a1}	y_{a2}	...	y_{an}	$\bar{y}_{a\cdot}$	$\bar{y}_{\cdot a}$		
					$\bar{y}_{\cdot \cdot}$	$\bar{y}_{\cdot \cdot}$		

$$a = \# \text{ of treatments} \quad (\text{8-min t+hi}) \quad (\text{end min})$$

y_{ij} = j^{th} observation of i^{th} treatment

n = # of observations at each level

Hypotheses Testing the Equality of Treatment Effects

Null: $H_0: \tau_1 = \tau_2 = \tau_3 = \dots = \tau_a = 0$

Alt: $H_a: \tau_i \neq 0$ for at least one i

Test Statistic: $F_0 = \frac{SS_{\text{Treatments}}/a}{SS_E/(a(n-1))}$ where $a = \# \text{ of treatments}$

$(n = \# \text{ of replicates @ each treatment level})$

Reference Distribution: F -dist $\sim F(a-1, n-a)$ [top] and $a(n-1)$ [bottom] D.F.

Rejection Criteria: $H_0: F_0 > F_{\alpha, a-1, n-a}$ (this is always one-sided)

Calculations for ANOVA

The Analysis of Variance for a Single-Factor Experiment					
Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F_0	
Treatments	$SS_{\text{Treatments}}$	$a - 1$	$MS_{\text{Treatments}}$	$\frac{SS_{\text{Treatments}}}{MS_E}$	
Error	SS_E	$a(n - 1)$	MS_E		
Total	SS_T	$an - 1$			

$$F_0 = \frac{MS_{\text{Treatments}}}{MS_E} = \frac{SS_{\text{Treatments}}/a}{SS_E/(a(n-1))}$$

$$SS_T = SS_{\text{Treatments}} + SS_E$$

$$SS_T = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{\cdot \cdot})^2 = \text{total sum of squares}$$

$$SS_{\text{Treatments}} = n \sum_{i=1}^a (\bar{y}_{i\cdot} - \bar{y}_{\cdot \cdot})^2 = \text{treatment sum of squares}$$

$n \times a$ squared terms

$a: \# \text{ of squared terms}$

$$\bar{y}_i = \sum_{j=1}^n y_{ij} \quad \bar{y}_{i\cdot} = y_i/n \quad i = 1, 2, \dots, a$$

$$\bar{y}_{\cdot \cdot} = \sum_{i=1}^a \sum_{j=1}^n y_{ij} / n = y_{\cdot \cdot}/N$$

Hypothesis tests the equality of treatment effects

Deriving & Development of ANOVA

a = number of treatments

y_{ij} = the j^{th} observation of the i^{th} treatment

n = number of observations for each level

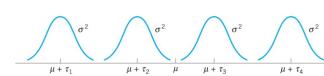
Linear statistical model: $y_{ij} = \mu + \tau_i + \epsilon_{ij}$ $\begin{cases} i = 1, 2, \dots, a \\ j = 1, 2, \dots, n \end{cases}$

Where: y_{ij} = observation

μ = overall mean (common to all treatments)

τ_i = the i^{th} treatment effect

ϵ_{ij} = random error



Null and Alternative Hypothesis:

$$H_0: \tau_1 = \tau_2 = \dots = \tau_a = 0 \quad H_1: \tau_i \neq 0 \text{ for at least one } i$$

Test Statistic: (reference distribution is the F-distribution)

$$F_0 = \frac{SS_{\text{Treatments}} / (a-1)}{SS_E / (a(n-1))} = \frac{MS_{\text{Treatments}}}{MS_E}$$

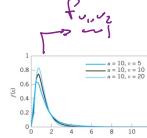
↑ DOF

↑ DOF

Rejection Criteria:

$$f_0 > f_{\alpha, a-1, n-a}$$

↑ Sign. 1
↑ noise



Calculations for ANOVA (same 1)

The Analysis of Variance for a Single-Factor Experiment

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F ₀
Treatments	SS _{Treatments}	a - 1	MS _{Treatments}	MS _{Treatments} / MS _E
Error	SS _E	a(n - 1)	MS _E	
Total	SS _T	an - 1		

$$SS_T = SS_{Treatments} + SS_E$$

$$SS_T = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2 \text{ total sum of squares}$$

$$SS_{Treatments} = n \sum_{i=1}^a (\bar{y}_i - \bar{y}_{..})^2 \text{ treatment sum of squares}$$

$$y_i = \sum_{j=1}^n y_{ij} \quad \bar{y}_i = y_i/n \quad i = 1, 2, \dots, a$$

$$y_{..} = \sum_{i=1}^a \sum_{j=1}^n y_{ij} \quad \bar{y}_{..} = y_{..}/N \quad \text{Total number of observations}$$

General

i = treatment

j = observation

b_{ij} = a value

dots = total

n = replicates

a = # of treatments

y_{ij} = jth observation of ith treatment

$\bar{y}_{..}$ = grand mean of all observations

$\bar{y}_{..}$ = total of all observations

\bar{y}_i = average of observations @ treatment i

What does that calculation look like?

Tensile Strength of Paper (psi)

Factor	Observations						Averages
	1	2	3	4	5	6	
Treatment Level	5	7	8	15	11	9	10.00
	10	12	17	13	18	19	15.67
	15	14	18	19	17	16	17.00
	20	19	25	22	23	18	20.17
							383 / 38.36

a=4

$$SS_T = SS_{Treatments} + SS_E$$

$$SS_T = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2 = (7 - 15.96)^2 + (8 - 15.96)^2 + \dots + (20 - 15.96)^2 = 513$$

$$SS_{Treatments} = n \sum_{i=1}^a (\bar{y}_i - \bar{y}_{..})^2 = 6 \left[(10 - 15.96)^2 + (15.67 - 15.96)^2 + (17 - 15.96)^2 + (20.17 - 15.96)^2 \right] = 382.8$$

$$\Rightarrow SS_E = 513 - 382.8 = 130.2$$

Hypothesis Test: Does hard wood concentration impact the tensile strength of paper? $\alpha = 0.01$

1) Parameter of F test: Multiple pop

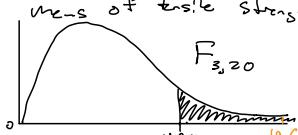
2) H₀: $\bar{y}_1 = \bar{y}_2 = \bar{y}_3 = \bar{y}_4 = 0$

3) $\exists i \in \{1, 4\}$ s.t. $\bar{y}_i \neq 0$

$$4) F_0 = \frac{\frac{SS_{Treatments}}{a(a-1)}}{\frac{SS_E}{a(n-1)}} = \frac{\frac{382.8}{6(6-1)}}{\frac{130.2}{20}} = 19.6$$

5) Reject H₀ if $F_0 > F_{0.01, 3, 20}$

$$F_{0.01, 3, 20} \approx 4.94 \text{ from charts}$$



Minitab: Stat -> ANOVA -> one-way

⇒ Conclusion: Reject H₀ as $F_0 > F_{0.01, 3, 20}$

At least one treatment has an impact on the tensile strength of paper. More generally, Hardwood Concentration impacts paper tensile strength!

What do you do when the null is rejected?
Multiple Comparison Methods (post-hoc tests) is performed to find which means are different after ANOVA has been performed and H₀ is rejected.

Fisher's Least Significant Difference (LSD) Test:

$$\text{Declare Significantly Different if: } |\bar{y}_i - \bar{y}_j| > LSD$$

$$\rightarrow LSD = t_{0.05, a(n-1)} \sqrt{\frac{2MS_E}{n}}$$

Example

Hardwood Concentration (%)	Observations						Totals	Averages
	1	2	3	4	5	6		
5	7	8	15	11	9	10	60	10.00
10	12	17	13	18	15	94	15.67	
15	14	18	19	17	16	18	102	17.00
20	19	25	22	23	18	20	127	21.17
							383	15.96

$$S_{vs. 10} : |10 - 15.67| = 5.67^*$$

$$S_{vs. 15} : |10 - 17| = 7^*$$

$$S_{vs. 20} : \dots = 11.17^*$$

$$10 \text{ vs } LS : \dots = 1.33$$

$$10 \text{ vs. } 20 : \dots = 5.5^*$$

$$LS \text{ vs. } 20 : \dots = 4.17$$

$$\begin{aligned} LSD &= t_{0.005, 20} \sqrt{\frac{2MS_E}{n}} \\ &= 2.845 \sqrt{2 \left(\frac{6.5}{6} \right)} \\ &= 4.19 \end{aligned}$$

$20 : C$ 10 vs 20 is
outlier
 $15 : C$ 10 vs 15 is
outlier
 $10 :$ B 10 vs 15 is
outlier
 $S :$ A Different
treatments

Summary

Declare significantly different if: $|\bar{y}_i - \bar{y}_j| > LSD$
where LSD, the least significant difference, is

$$LSD = t_{\alpha/2, a(n-1)} \sqrt{\frac{2MS_E}{n}}$$

Better notes on LSD test:

- Do this ONLY after ANOVA null hypothesis is rejected
- Find differences between all of the treatments
- Compare differences to the LSD to see if they are significantly different

Mini tab might produce an output like this:
If a pair has the same letter, they are not significantly different

Factor	N	Mean	Grouping
20%	6	21.17	A
15%	6	17.000	A
10%	6	15.67	B
5%	6	10.00	C

Assumptions

(i) Normal

(ii) No hidden variables

Fisher's LSD Test

Declare significantly different if: $|\bar{y}_i - \bar{y}_j| > LSD$
where LSD, the least significant difference, is

$$LSD = t_{\alpha/2, a(n-1)} \sqrt{\frac{2MS_E}{n}}$$

Better notes on LSD test:

- Do this after ANOVA null hypothesis is rejected
- Find differences between all of the treatments
- Compare differences to the LSD to see if they are significantly different

Mini tab might produce an output like this:

If a pair has the same letter, they are not significantly different

Sometimes funny overlap occurs – for example, 15% is not different than 20% or 10% but, 10% and 20% are different

Factor	N	Mean	Grouping
20%	6	21.17	A
15%	6	17.000	A
10%	6	15.67	B
5%	6	10.00	C

Confidence Intervals

A $100(1-\alpha)\%$ confidence interval on the mean of the i th treatment μ_i is

$$\bar{y}_i - t_{\alpha/2, a(n-1)} \sqrt{\frac{MS_E}{n}} \leq \mu_i \leq \bar{y}_i + t_{\alpha/2, a(n-1)} \sqrt{\frac{MS_E}{n}}$$

A $100(1-\alpha)$ percent confidence interval on the difference in two treatment means $\mu_i - \mu_j$ is

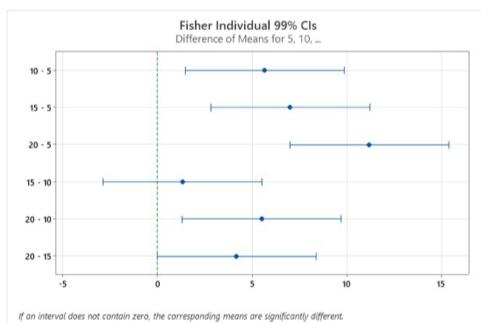
$$\bar{y}_i - \bar{y}_j - t_{\alpha/2, a(n-1)} \sqrt{\frac{2MS_E}{n}} \leq \mu_i - \mu_j \leq \bar{y}_i - \bar{y}_j + t_{\alpha/2, a(n-1)} \sqrt{\frac{2MS_E}{n}}$$

ANOVA Assumptions

(i) Errors (and observations) are normally distributed

(ii) Each treatment has the same variance (equal variance assumption)

(iii) Errors and observations are independently distributed



$$y_{ij0} = \mu_i + \epsilon_{ij0}$$

\uparrow obs.
 \uparrow estimated

$$E_{ij0} = y_{ij0} - \bar{y}_{i0}$$

\uparrow residual obs.
 \uparrow residual

Hardwood Concentration (%)	Observations						Totals	Averages
	1	2	3	4	5	6		
5	7	8	15	11	9	10	60	10.00
10	12	17	13	18	15	94	15.67	
15	14	18	19	17	16	18	102	17.00
20	19	25	22	23	18	20	127	21.17
							383	15.96

Hardwood Concentration (%)	Residuals					
	5	10	15	20	-3.00	-2.00
5					1.00	1.00
10					1.33	-2.67
15					1.00	2.33
20					0.83	-3.33
					-1.00	-0.67
					0.00	1.00
					-1.00	-1.00
					-3.17	-1.17

Residuals

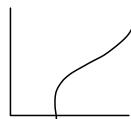
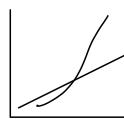
The difference between an observation (y_{ij0}) and its estimated (or fitted) value from the statistical model being studied

Normal Probability Plot of Residuals

- (i) Rank data
 - (ii) Calc Cumulative Frequency
 - (iii) Find Z_i 's
 - (iv) Plot residuals vs. Z_i 's
- Looking for straight line

Minitab

Can ask Minitab to plot a normal probability plot after pasting in residuals into one column, or use ANOVA function, and select graphs, then select – normal probability plot of residuals



Checking equal Variance at each Point Level

→ Minitab does not automatically generate

Checking Equal Variances as a function of predicted value

Residuals plotted as a function of predicted value

Minitab: use ANOVA function, and select graphs, then select – residuals vs. fits

Can also generate in Excel or Minitab: Residuals in one column, predicted value in the other, using scatter plot

Hardwood Concentration (%)	Residuals					
	-3.00	-2.00	5.00	1.00	-1.00	0.00
5						
10	-3.67	1.33	-2.67	2.33	3.33	-0.67
15	-3.00	1.00	2.00	0.00	-1.00	1.00
20	-2.17	3.83	0.83	1.83	-3.17	-1.17

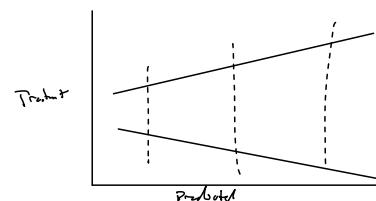
Hardwood Concentration (%)	Observations					
	1	2	3	4	5	Totals
5	7	8	15	11	9	60
10	12	17	13	18	19	94
15	14	18	19	17	16	102
20	19	25	22	23	18	127
						383
						15.96

Look for: even residuals, if not...

Rule of thumb we used for t-tests applies, or, you could check using your 2 variance test!

What if there is a pattern?

- Data transformation
- Relevel Analysis



Checking Independence

No formal plot... but:

1) Addressed via experimental design (randomize!)

2) Can plot residuals versus any suspected or nuisance variables

1) Example: plot residuals versus run order (or other)

When is my data not independent?

- (1) repeated measurements are taken on the same subject → paired
- (2) residuals or observations are correlated with a nuisance variable

Regression Modeling

• **Regression Modeling:** A mathematical model fit to a set of data relating a response variable (y) to an independent variable (x_1, x_2, \dots, x_n)

• **Empirical Models:** Models where the structure is determined by the observed relationship among experimental data, but not necessarily mechanistically relevant

• **Simple Linear Regression Model:** There is only one indep. or predictor variable (x) and the relationship with response (y) is assumed to be linear

$$y = \beta_0 + \beta_1 x + \epsilon$$

y = expected response

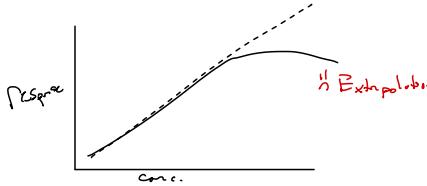
β_0 = intercept

β_1 = slope

x = predictor variable

ϵ = random error term

regression coeff

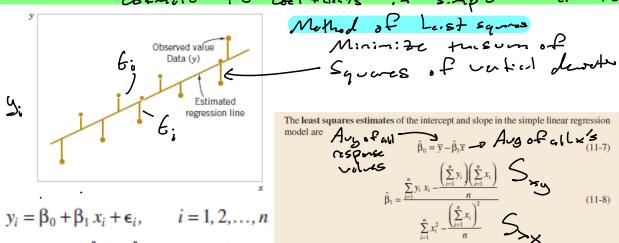


• **P+P HS:** Empirical M. Ls

→ Correlation does not equal causation!!!

→ Valid in the range of original data (be careful with extrapolation)

• How to estimate the coefficients in simple linear regression



Fitted or Estimated Regression Line

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

Predicted Value

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \epsilon_i$$

Each pair satisfies this

Simple Linear Regression

- Use simple linear regression for building an empirical model
- Understand how the method of least squares is used to estimate the parameters in a linear regression model
- Analyze residuals to determine whether the regression model is an adequate fit to the data or whether assumptions are violated
- Test statistical hypotheses and construct confidence intervals on regression model parameters
- Apply the correlation model
- Use simple transformations to achieve a linear regression model

Example: Fit a simple linear regression line to these data

TABLE 11-1 Oxygen and Hydrocarbon Levels

Observation Number	Hydrocarbon Level ($x\%$)	Purity ($y\%$)
1	0.99	90.01
2	1.02	89.05
3	1.15	91.43
4	1.29	93.74
5	1.46	96.73
6	1.36	94.45
7	0.87	87.59
8	1.23	91.77
9	1.55	99.42
10	1.40	93.65
11	1.19	93.54
12	1.15	92.52
13	0.98	90.56
14	1.01	89.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	87.33

$$n=20$$

$$\sum x_i = 23.42 \quad \sum y_i = 1848.21 \quad \sum x_i y_i = 2214$$

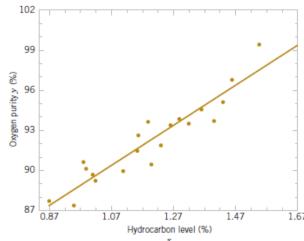
$$\sum x_i^2 = 29.28 \quad \sum y_i^2 = 170044$$

The least squares estimates of the intercept and slope in the simple linear regression model are

$$\hat{y}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (11-7)$$

$$\hat{\beta}_1 = \frac{\left(\sum_{i=1}^n y_i \right) \left(\sum_{i=1}^n x_i \right)}{\sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} S_{xy} \quad (11-8)$$

where $\bar{y} = (1/n) \sum_{i=1}^n y_i$ and $\bar{x} = (1/n) \sum_{i=1}^n x_i$.



$$\hat{y} = 74.28 + 14.95x$$

Minitab: Stat > Regression > Regression > Fit Regression Model

Minitab: Stat > Regression > Fitted Line -> choose x, y and linear

ANOVA For Linear Regression

Null & Alternative Hypotheses:

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Test Statistic

$$F_0 = \frac{SS_R / 1}{SS_E / n-2} = \frac{MS_R}{MS_E} \quad \begin{matrix} \text{Signal} \\ \text{noise} \end{matrix}$$

Refined Distribution: F-distribution

\hookrightarrow Want to vs: $F_{1, n-2}$

Rejection Criteria:

$$\text{Reject } H_0: F_0 > F_{\alpha/2, n-2}$$

or $p\text{-value} < \alpha$

Calculating the Error

$$\text{Recall: } SS_T = SS_E + SS_R$$

$$\text{Residual: } e_i = y_i - \hat{y}_i \quad \text{Actual vs predicted}$$

$$\text{Error Sum of Squares: } SS_E = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{Mean Squared Error: } SS_E / n-2$$

$$\text{Sum of Squares Total: } SS_T = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\text{Regression Sum of Squares: } SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Checking Assumptions

- 1) Errors (and observations) are normally distributed
- 2) Errors (and observations) are independently distributed
- 3) Variance is constant

$$e_i = y_i - \hat{y}_i \quad i=1, 2, \dots, n$$

residuals
check for
constant variance



Linear?

Normality

Plot residuals
on Normal
Probability Plot

Pathogenic?

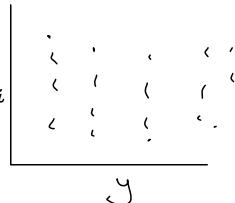
Independent Variance

Plot residuals vs.
Run order or
nonsystematic variables

Pathogenic or Spurious

Constant Variance

① Plot residuals
vs. predicted (constant)
② Plot vs. independent
variable (x)



Example

TABLE 11-1 Oxygen and Hydrocarbon Levels		
Observation Number	Hydrocarbon Level x (%)	Purity y (%)
1	0.99	90.05
2	1.02	91.43
3	1.15	93.74
4	1.29	96.73
5	1.46	99.42
6	1.58	97.59
7	0.87	91.77
8	1.23	99.42
9	1.55	93.25
10	1.40	93.68
11	1.19	95.54
12	1.15	92.52
13	0.98	90.56
14	1.00	91.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	97.33

1) Parameter of Interest: Slope of Linear Model; one regressor \rightarrow simple linear regression

2) $H_0: \beta_1 = 0$

$H_a: \beta_1 \neq 0$

$$4) F_0 = \frac{SS_E / 1}{SS_T / n - 2}$$

5) Reject if $F_0 > F_{0.05, 18, 18}$

$$\bar{y} = 92.16$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i \Rightarrow \hat{y} = 74.28 + 14.947x$$

$$n=1 \Rightarrow \hat{y} = 74.28 + 14.947(0.99) = 89.08$$

↓ all data points

$$n=20 \Rightarrow \hat{y} = 74.28 + 14.947(0.95) = 88.48$$

$$SS_E = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = (89.08 - 92.16)^2 + \dots + (88.48 - 92.16)^2 = 152.17$$

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2 = (90.01 - 92.16)^2 + \dots + (82.33 - 92.16)^2 = 173.38$$

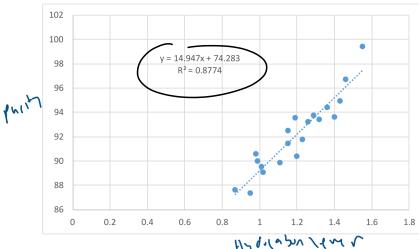
$$SS_B = SS_T - SS_E = 21.2$$

$$\Rightarrow F_0 = \frac{152.17 / 1}{21.2 / (20-2)} = 129 \quad | \quad 129 > 4.41 \rightarrow \text{Reject } H_0$$

There is a significant linear relationship between hydrocarbon level and purity

$$F_{0.05, 18, 18} = 4.41 \text{ from t-tables}$$

$$R^2 = \frac{SS_B}{SS_T} = 1 - \frac{SS_E}{SS_T} = \frac{152}{173} = 0.88$$



R^2 represents how much error is accounted for by the model

Confidence Intervals around Coefficients

Under the assumption that the observations are normally and independently distributed, a $100(1 - \alpha)\%$ confidence interval on the slope β_1 in simple linear regression is

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \quad (11-29)$$

Similarly, a $100(1 - \alpha)\%$ confidence interval on the intercept β_0 is

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \quad (11-30)$$

How to write and interpret:

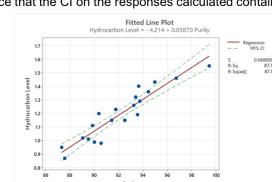
In repeated sampling the intervals calculated capture the true slope or intercept $(1-\alpha)\%$ of the time

Notes on CIs around Response

- CIs (confidence intervals) around a response: give a range of values associated with the true response for a given x
- You can construct confidence limits of your regression line by repeating the CIs for different values of x

How to Write?

CIs: In repeated experiments, there 95% chance that the CI on the responses calculated contain the true responses



How to generate?

To get the bounds on a graph: Minitab: Stat > Regression -> Fitted Line -> choose x, y and linear (Options tab: display CIs)

You can also use the "predict" function to have Minitab give you the CIs for any given x

Confidence Intervals around the Response

A $100(1 - \alpha)\%$ confidence interval on the mean response at the value of $x = x_0$, say $\mu_{Y|x_0}$, is given by

$$\hat{\mu}_{Y|x_0} - t\alpha/2, n-2 \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]} \leq \mu_{Y|x_0} \leq \hat{\mu}_{Y|x_0} + t\alpha/2, n-2 \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]} \quad (11-31)$$

where $\hat{\mu}_{Y|x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$ is computed from the fitted regression model.

(Nicely worth) ANOVA for Linear Regression

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$
 $H_1: \beta_j \neq 0$ for at least one j

Test Statistic:

$$F_0 = \frac{SS_E/1}{SS_E/(n-2)} = \frac{MS_E}{MS_E}$$

Reference Distribution:

F-distribution
 Write as: $F_{k,n-2}$

Rejection Criteria: Reject H_0 if $f_0 > f_{k,n-2}$.

Regression on Transformed Variables

Intrinsically Linear Models: Models that can be transformed into linear equations

Transformation is useful when: Relationship isn't linear and/or when residuals show a pattern

Examples

$$Y = \beta_0 e^{\beta_1 x} \rightarrow \ln Y = \ln \beta_0 + \beta_1 x + \ln \epsilon$$

$$Y = \beta_0 + \beta_1 \left(\frac{1}{x} \right) + \epsilon \rightarrow Y = \beta_0 + \beta_1 z + \epsilon$$

Example

Find a statistically satisfactory ($\alpha = 0.05$) model for:

Observation Number	Wind Velocity (m/sec)	DC Output	Observation Number	Wind Velocity (m/sec)	DC Output
1	5.00	1.502	14	5.00	1.727
2	6.00	1.822	15	5.40	1.988
3	7.00	2.007	16	5.60	2.137
4	2.70	0.500	17	3.85	2.179
5	9.70	2.256	18	8.00	2.112
6	7.90	2.395	19	7.00	2.089
7	8.55	2.294	20	5.45	1.961
8	6.50	1.656	21	6.00	1.876
9	8.15	2.166	22	10.20	2.310
10	5.20	1.366	23	4.00	1.394
11	7.20	1.853	24	3.50	1.441
12	6.50	1.930	25	2.45	0.123
13	—	—	—	—	—

$$\text{Final Model: } y = \beta_0 + \beta_1 \left(\frac{1}{x} \right) + \epsilon$$

Steps

- Plot data out
- Fit data with a model we think could work (try linear)
- Analyze the residuals (ANOVA, R² – plot residuals, normality)
- Transform the data if there is a problem and regress data again
- Check the results again (ANOVA, R², residuals, normality)
- Determine if model is better
- Write out final model: If ANOVA is significant, and assumptions are met

We will rely on Minitab a lot more!

1) Scatter Plot doesn't look linear

2-8) Try linear → Residuals vs Fit results → clear patterns

↳ Equal variances assumption violated, try a different model

4) Try reciprocal of data

5-6) Results look good now! If they don't just try a different model

7) $y = \text{DC output}$ $\left\{ \begin{array}{l} y = 2.924 - 6.488 \left(\frac{1}{x} \right) \\ x = \text{wind velocity} \end{array} \right.$ Always specify variables

Conclusions: The transformed model for DC output and wind velocity is significant.

Common Transformations

Method	Math Operation	Good for	Bad for:
Log	$\ln(x)$ $\log_b(x)$	Right-skewed data Might not be really good at handling higher order powers of 10 (e.g. 1000, 10000)	Zero values Negative values
Square root	\sqrt{x}	Right-skewed data	Negative values
Square	x^2	Left-skewed data	Negative values
Cube root	$x^{1/3}$	Right-skewed data Negative values	Not as effective at normalizing as log transform
Reciprocal	$1/x$	Making small values trigger and big values smaller	Zero values Negative values



Multiple Linear Regression

Multiple Regression Models: A linear regression that contains more than one regressor (K)

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$$

but contains more than one regressor (K)

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

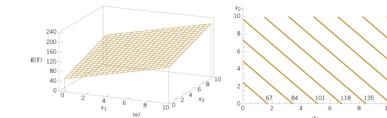


FIGURE 12-1 (a) The regression plane for the model $E(Y) = 50 + 10x_1 + 7x_2$. (b) The contour plot.

Cubic polynomial

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$$

Interaction

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

Second order with interaction

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 x_1 x_2 + \epsilon$$

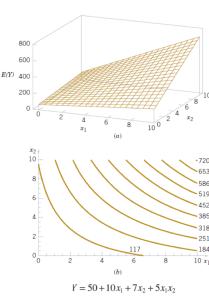


FIGURE 12-3 (a) Three-dimensional plot of the regression model $E(Y) = 800 + 10x_1 + 8.5x_2 - 5x_1^2 + 4x_1x_2$. (b) The contour plot.

Estimated Coefficients and Variance

Coefficients: Found by method of least squares

$$\Rightarrow \text{Residuals: } e_i = y_i - \hat{y}_i \quad \text{sum of these minimized}$$

Variance from different sources

• Variance from Random Error: SS_E from calculation of residuals

$$SS_T = SS_R + SS_E \quad \Rightarrow \quad SS_E = \sum e_i^2$$

• Variance from Regression: SS_R = variance from regression

• Variance from individual regression terms: $SS_{\beta_j}(\beta_j)$ = variance coming from specific regression term

• Partial F-test: Compares variance from specific regression terms to random error.

Some concepts simple
just a lot more terms

Hypothesis Testing

Recall

$n = \# \text{ of data points}$
 $k = \# \text{ of regression terms}$

Test for significance of regression:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_j \neq 0 \text{ for at least one } j$$

Test Statistic:

$$F_0 = \frac{SS_R/k}{SS_E/(n-k-1)}$$

Partial F-test:

$$H_0: \beta_i = 0$$

$$H_1: \beta_i \neq 0$$

$$\text{Test Statistic: } F_0 = \frac{\frac{SS_{R,i}(\beta_i)}{1}}{\frac{SS_E}{n-k-1}}$$

Reference Distribution:

$$F_{-dist}$$

$$F_{k, n-k-1}$$

Rejection criteria: (F-test)

$$F_0 > F_{\alpha, k, n-k-1}$$

Reference Distribution:

$$F_{-dist}$$

$$F_{1, n-k-1}$$

Rejection criteria:

$$F_0 > F_{\alpha, 1, n-k-1}$$

Nicely Typed Notes

Test for significance of regression:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_j \neq 0 \text{ for at least one } j$$

Reference Distribution:
F-distribution: F k, n-k-1

$$F_0 = \frac{SS_R/k}{SS_E/(n-k-1)}$$

Rejection criteria:
Reject H0 if F0 exceeds F alpha, k, n-k-1
Or if p-value is < alpha

Partial F-test:

$$H_0: \beta_j = 0$$

$$H_1: \beta_j \neq 0$$

Reference Distribution:
F-distribution: F 1, n-k-1

$$F_0 = \frac{SS_{R,j}/1}{SS_E/(n-k-1)}$$

Rejection criteria:
Reject H0 if F0 exceeds F alpha, 1, n-k-1
Or if p-value is < alpha

Adjusted R^2

$$SS_T = SS_R + SS_E \Rightarrow R^2 = \frac{SS_R}{SS_T} \quad \text{Good for simple linear regression}$$

$$\text{Use Adj. } R^2 \text{ for multiple!} \quad R^2 = \frac{SS_R}{k} / \frac{SS_T}{n-1} \quad \text{Adjusts for fitting}$$

Nonetheless DoF avoids innundated R^2
arbitrarily!

Example

Minimizing

Find a statistically satisfactory ($\alpha = 0.05$) model for:

Brake Horsepower (y) rpm	Road Octane Number	Compression
2000	90	100
212	90	95
229	94	95
222	90	110
219	90	98
278	96	110
246	94	98
237	90	100
233	98	105
224	86	97
223	90	100
230	89	104

Plot: Graph > Scatter plot > plot response against all X's

Fit Model: Stat > Regression > Fit Regression Model (add Response, and Predictors by double click) (Determine model by button "Model" to add or subtract terms)

Contour Plots: Stat > regression > regression > contour or surface plot

Final Result

Regression Equation

$$\text{Brake Horsepower} = -266.0 + 0.01071 \text{ rpm} + 3.135 \text{ Road Octane Number}$$

$$+ 1.867 \text{ Compression}$$

As rms low \rightarrow significant !!

1) Plot the data

2) Fit a multiple linear regression model. Check the ANOVA results, R^2_{adj} , residuals, normality

3) Based on results, consider adding or removing regressors

4) If model is changed – repeat 2-3

5) Generate final model and contour plots

Goal: simplest model which significantly describes the data

Control Charts

- Statistical Process Control:** A collection of problem-solving tools useful in achieving process stability and reducing variability.

The Magnificent Seven

1. Histogram or stem-and-leaf plot
2. Check sheet
3. Pareto chart
4. Cause-and-effect diagram
5. Defect concentration diagram
6. Scatter diagram
7. Control chart

Shewhart's Control Charts & Theory of Variability

- Chance causes of variation:** Inherent variability in a process or background noise - a process operating with only chance causes is in statistical control

- Assignable Causes of Variation:** Sources of variability that are not due to chance - a process operating outside of assignable ranges is an out-of-control process

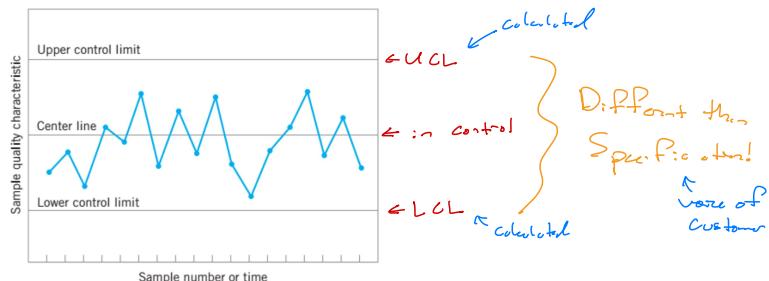
- Sources:**
 - Improperly adjusted/controlled machines
 - Operator errors
 - Different raw material *

→ **Control Charts:** Online process monitoring technique

Control Chart Uses

- Monitoring:** Using the control chart to determine if the process is in control, and quickly detect any assignable causes
- Estimating:** Using the control chart to estimate process parameters (mean or variance)
- Improving:** Using the control chart to improve process parameters

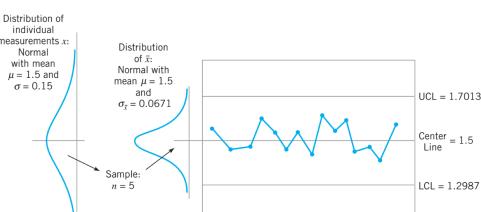
Features of a Control Chart



Relationship to Hypothesis Testing

Hypothesis: $M = M_0$ in statistical control

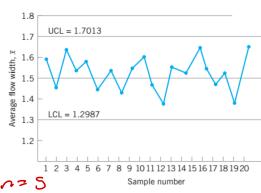
Alternative: $M \neq M_0$



Type I error: Rejecting H_0 when H_0 is true; concluding the process is out of control when it is not

Type II error: Not rejecting H_0 when it is false; Assuming a process is in control when it really is out of control

Example



→ **Control chart:** uses sample averages to monitor a process mean

$100(1-\alpha)\%$ of sample means should fall between control limits

Three Sigma Control Limits: Define $Z_{\alpha/2} = 3 \Rightarrow \alpha = 0.0027$

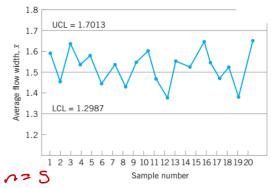
Control limits
(Known values)

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{0.15}{\sqrt{5}} = 0.0671$$

$$UCL = \bar{x} + 3(\sigma_{\bar{x}}) = 1.7013$$

$$LCL = \bar{x} - 3(\sigma_{\bar{x}}) = 1.2987$$

Example



Control chart uses sample means to monitor a process mean

100(1-2)% of sample means should fall between control limits

Three Sigma Control limits: Define $Z_{\alpha/2} = 3 \Rightarrow \alpha = 0.0027$

$$\text{Control limits} : \bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

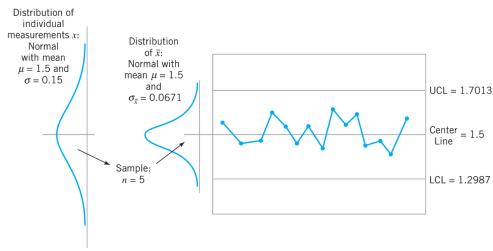
(Known values)

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{0.15}{\sqrt{5}} = 0.0671$$

$$UCL = \bar{x} + 3(0.0671) = 1.7013$$

$$LCL = \bar{x} - 3(0.0671) = 1.2987$$

General Model



w = quality characteristic of interest

μ_w = mean of w

σ_w = std dev of w

L = distance from centerline expressed in std dev

$$UCL = \mu_w + L \sigma_w$$

$$\text{Center} = \mu_w$$

$$LCL = \mu_w - L \sigma_w$$

Design of a Control Chart

Selecting Control Limits:

- Set alpha (Europe)

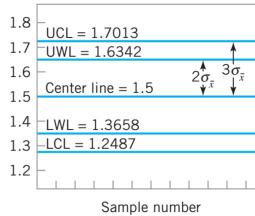
- Set Sigma Level (US)

3 is pretty standard compromise
balancing type I and II error

Two Limit Control Charts

- Action limits - outer limits that trigger search for assignable cause

- Warning limits - inner limits which raise alarms



Average Run Length (ARL) - the average # of points that must be plotted before a point will indicate an out-of-control condition

$$ARL = \frac{1}{p} ; p = \text{probability the point exceeds control limit}$$

Average Run Signal (ARS) - the average amount of time between out-of-control signals

$$ARS = ARL \times (\text{time between samples})$$

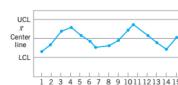
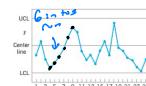
Analysis of Patterns A process can be out of control even if all points are within limits!

Run: A series of observations of the same type

Run Length: How many points are in the series

Pattern: Non-random behavior

Zone Rules: Criteria to help identify patterns



Sinusoidal

Dotsy Patterns and Sensitizing Rules

Standard Action Signal:	1. One or more points outside of the control limits 2. Two or three consecutive points outside the two-sigma warning limits but still inside the control limits 3. Four or five consecutive points beyond the one-sigma limits 4. A run of eight consecutive points on one side of the center line 5. Six points in a row steadily increasing or decreasing 6. Fifteen points in a row in zone C (both above and below the center line) 7. Fourteen points in a row alternating up and down 8. Eight points in a row on both sides of the center line with none in zone C 9. An unusual or nonrandom pattern in the data 10. One or more points near a warning or control limit	Western Electric Rules Out of control if yes
Sensitizing Rules:		
If yes, investigate more, get more data Raises alarm		

Sensitizing Rules: Criteria to Increase Sensitivity



\Rightarrow Out of Control automatically.

\bar{x} chart: Average of samples

R-Charts: the difference between min and max value of a sample

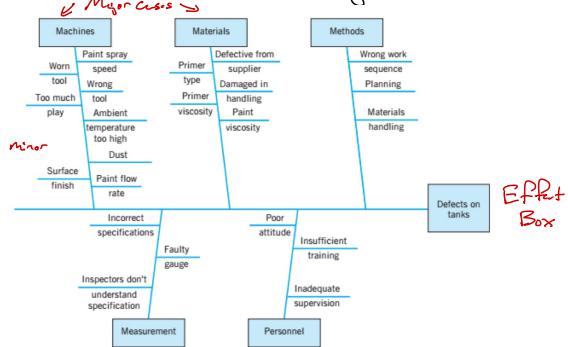
Octo of Control Action Plans

OCAP: A Plan chart or text-based description of the sequence of actions that take place following an actuality event

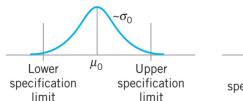
Champions: Potential assignable causes

Commands: Action tasks to Resolve

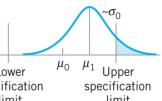
Cause and Effect Diagram



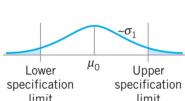
Controlling Mean & Variability



In control



Mean shifts



sigma shifts

Known Population and Varience

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\mu + Z_{\alpha/2} \sigma_{\bar{x}} = \mu + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

LSL

$$\mu - Z_{\alpha/2} \sigma_{\bar{x}} = \mu - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

USL

Why use both R and S charts?

Mean Sharts: R shows change R chart won't

Variance Sharts: R chart shows change, R won't

Parameters vs. Unknowns

Estimated from "in control" samples of at least 20-25 (m) each with at least 4-5 replicates (n)

where m is the total # of samples of size n

$$M = \bar{\bar{x}} = \frac{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_m}{m}$$

$$\bar{R} = \frac{R_1 + R_2 + \dots + R_m}{m} \rightarrow \bar{R}_m = \bar{x}_{max} - \bar{x}_{min}$$

for any sample

Designed to give limits 3 standard deviations away from mean:

Control Limits for the \bar{x} Chart

$$\begin{aligned} \text{UCL} &= \bar{x} + A_2 \bar{R} \\ \text{Center line} &= \bar{x} \\ \text{LCL} &= \bar{x} - A_2 \bar{R} \end{aligned} \quad (6.4)$$

The constant A_2 is tabulated for various sample sizes in Appendix Table VI.

Control Limits for the R Chart

$$\begin{aligned} \text{UCL} &= D_4 \bar{R} \\ \text{Center line} &= \bar{R} \\ \text{LCL} &= D_3 \bar{R} \end{aligned} \quad (6.5)$$

The constants D_3 and D_4 are tabulated for various values of n in Appendix Table VI.

Example

EXAMPLE 6.1 \bar{x} and R Charts for a Manufacturing Process

A hard-bake process (see Section 5.3.1) is used in conjunction with photolithography in semiconductor manufacturing. We wish to establish statistical control of the flow width of the resist in this process using \bar{x} and R charts. Twenty-five samples,

each of size five wafers, have been taken when we think the process is in control. The interval of time between samples or subgroups is one hour. The flow width measurement data (in microns) from these samples are shown in Table 6.1.

TABLE 6.1 Flow Width Measurements (microns) for the Hard-Bake Process

Sample Number	Wafers						
	1	2	3	4	5	\bar{x}_i	R_i
1	1.3235	1.4128	1.6744	1.4573	1.6914	1.5119	0.3679
2	1.4314	1.3592	1.6075	1.4666	1.6109	1.4951	0.2517
3	1.4284	1.4871	1.4932	1.4324	1.5674	1.4817	0.1390
4	1.5028	1.6352	1.3841	1.2831	1.5507	1.4712	0.3521
5	1.5603	1.2735	1.5261	1.4363	1.6441	1.4882	0.3706
6	1.5955	1.3451	1.3574	1.3281	1.4198	1.4492	0.2674
7	1.4747	1.5544	1.5606	1.5265	1.5717	1.5260	0.2657
8	1.4190	1.4203	1.6037	1.6097	1.5519	1.5343	0.2471
9	1.3884	1.7277	1.5355	1.5176	1.3608	1.5076	0.3589
10	1.4039	1.6697	1.5089	1.4627	1.5229	1.5134	0.2658
11	1.4158	1.7667	1.4278	1.5928	1.4181	1.5242	0.3598
12	1.5821	1.3355	1.5777	1.3908	1.7559	1.5284	0.4204
13	1.2856	1.4106	1.4447	1.6309	1.1928	1.3947	0.4707
14	1.4951	1.4036	1.5893	1.6458	1.4969	1.5261	0.2423
15	1.3589	1.2863	1.5996	1.2497	1.5471	1.4083	0.3499
16	1.5747	1.5301	1.5171	1.1839	1.8662	1.5344	0.6823
17	1.3680	1.7269	1.3957	1.5014	1.4449	1.4874	0.3589
18	1.4163	1.3864	1.3057	1.6210	1.5573	1.4573	0.3153
19	1.5798	1.4185	1.6541	1.5116	1.7247	1.5777	0.3062
20	1.7106	1.4412	1.2361	1.3820	1.7601	1.5060	0.5240
21	1.4371	1.5051	1.3485	1.5670	1.4880	1.4691	0.2185
22	1.4738	1.5936	1.6583	1.4973	1.4720	1.5390	0.1863
23	1.5917	1.4333	1.5551	1.5295	1.6866	1.5952	0.2533
24	1.6399	1.5243	1.5705	1.5563	1.5530	1.5688	0.1156
25	1.5797	1.3663	1.6240	1.3732	1.6887	1.5264	0.3224
						$\Sigma x_i = 37.4400$	$\Sigma R_i = 8.1302$
						$\bar{x} = 1.5056$	$\bar{R} = 0.32521$

$$\begin{aligned} \bar{R} &= \frac{0.3679 + 0.2812 + \dots + 0.2251}{25} \\ &= 0.2821 \end{aligned}$$

$$LCL = \bar{R} - D_3 \bar{R} = 0$$

$$\begin{aligned} UCL &= \bar{R} + D_4 \bar{R} = 0.32821(2.114) \\ &= 0.68249 \end{aligned}$$

$$\begin{aligned} \bar{x}_{\text{Cl},+} &= \frac{1.5119 + 1.4951 + \dots + 1.5261}{25} \\ &= 1.5086 \end{aligned}$$

$$\begin{aligned} UCL &= \bar{x} + A_2 \bar{R} = 1.5086 + 0.5773(0.2821) \\ &= 1.69325 \end{aligned}$$

$$\begin{aligned} LCL &= \bar{x} - A_2 \bar{R} = 1.5086 - 0.5773(0.2821) \\ &= 1.31795 \end{aligned}$$

Minutab

Stat -> Control Charts -> Variables Charts for Subgroups -> Select Xbar or R -> dropdown to observations from subgroup are in one row of columns -> select all columns with data (n=number of selections)

Estimates from Control Charts

$$\begin{aligned} \text{Estimate process standard deviation:} \\ \hat{\sigma} &= \frac{\bar{R}}{d_2} = \frac{0.2821}{2.114} = 0.1398 \end{aligned}$$

$$\begin{aligned} \text{Estimate of process mean:} \\ \bar{M} &= \bar{x} \end{aligned}$$

Note: Can chart versus subgroups, but be careful!

Topics

General:

- how to correctly phrase conclusions in a hypothesis test
- how to correctly phrase what the p-value is
- how to correctly phrase what confidence intervals are representing
- how to sketch hypothesis tests on a reference (sampling) distribution

Inferring the variance of two populations:

- how to conduct the hypothesis test
- how to draw conclusions

ANOVA:

- how to conduct the test (by hand and interpret Minitab) and draw conclusions
- what situations ANOVA is used in (when to pick it as the statistical test)
- what, generally, the test statistic represents
- what the reference distribution is as related to the test statistic and when the test statistic follows the reference distribution
- know how to check assumptions and what assumptions are made
- how and when to perform a posthoc test
- how and when to transform data

Linear Regression:

- know how to perform simple linear regression by hand (how to perform the hypothesis test - won't ask the regression terms by hand)
- how to set up equations and hypothesis for multiple and simple linear regression
- know how to interpret results from Minitab output on simple and multiple linear regression
- explain what R^2 is (+adjusted), what it can and cannot be used for, and how to calculate by hand for simple and multiple linear regression
- know how to check assumptions and what assumptions are made
- how and when to transform data
- describe what the confidence intervals on the response and coefficients are

Control Charts:

- Be able to explain what impact wide or narrow control limits have on type I and type II error, and be able to define type I and type II errors for control charts
- Be able to explain in words and sketches why it is important to use control charts for both process mean and variability
- Calculate control limits and warning limits for xbar and R charts
- Be able to identify in and out of control process
- ARL calculation and concept

Specs \neq Cont. Lims

ECHE 313 Exam 2

Write your name: _____

Write the following statement, and sign your name acknowledging your agreement:

"I neither received nor gave external assistance on this exam"

The exam layout is such that all problems add up to 100 points. This point total is meant to guide you in time allocation for this test, will be adjusted such that this exam is 25% of your total class grade. You may use your class materials, class notes, homework, book, and Minitab to aid in the completion of the exam, but you are not allowed to use any outside resources. The exam is designed to be completed in 1 hr and 15 min, but you will have between 10 AM 4/9 to 10 AM 4/10 to complete the exam. Note: Once you start the exam, you have 4 hours to complete the test. All of the necessary tables are provided in the back of your book. Please show your work and all steps necessary to arrive at your answer. It is expected you do this entire test by hand, using a calculator and the tables found in the back of your book. So while you are allowed to use Minitab if you wish to check your answers, it should not replace the work you have done by hand using the tables for your answers. Where applicable, it is best to write out the general formulas you want to use first, then write in your subbed in values so it is clear. Also, when reporting numbers, keep the same number of decimals as the data you are working with unless there are specific directions on how many sig figs to report. Good luck!

Premise: You work as a process engineer for a company that makes eclipse glasses. What a great time to be in this business! These glasses work by filtering light, reducing light exposure to the eye. Unless otherwise stated, your company uses $\alpha=0.05$.

Problem 1 (67 points total) Your team wants to change the concentration of materials that are used to block light in your product to save money. Material B is particularly expensive and is typically set at 5.00 wt% in the material. You want to know if the concentration of material B has an impact on the % transmittance of light, so you vary material B and test the amount of light that is transmitted through the material, and get the following data:

$\alpha = 0.05$

Material B (wt%)	% Transmittance			Ave.	Standard Deviation
	Replicate 1	Replicate 2	Replicate 3		
(typical) 5.00	0.0023	0.0023	0.0024	0.0023	0.000058
4.00	0.0023	0.0025	0.0028	0.0025	0.000252
3.00	0.0027	0.0024	0.0021	0.0024	0.000300
2.00	0.0032	0.0038	0.0035	0.0035	0.000300
1.00	0.0041	0.0039	0.0045	0.0042	0.000306

$$\bar{y}_{..} = \frac{0.0023 + 0.0025 + 0.0024 + 0.0035 + 0.0048}{5} = 0.00298$$

- a. (1 points) List the *name* of the hypothesis test you would use

ANOVA

- b. (2 points) Name the parameter of interest in the context of the problem statement

Multiple means, specifically $\mu_1, \mu_2, \mu_3, \mu_4$ at different levels of material B.

- c. (2 points) Write out the correct null hypothesis for this test

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4 = 0$$

- d. (2 points) Write out the correct alternative hypothesis for this test

$$H_a: \mu_i \neq 0 \text{ for some } i \in \{1, 2, 3, 4\}$$

- e. (6 points) Write the formula of the test statistic, and other formulas needed to fully calculate the test statistic from the data above, and calculate its value given that SS_E is 0.00000068

$$F_0 = \frac{\sum_{i=1}^n (\bar{y}_{i..} - \bar{y}_{..})^2 / (a-1)}{\sum_{i=1}^n s_{i..}^2 / (a(n-1))}$$

$$SS_{\text{treatments}} = n \sum_{i=1}^n (\bar{y}_{i..} - \bar{y}_{..})^2$$

$$= 3 \left[(0.0028 - 0.00298)^2 + (0.0025 - 0.00298)^2 + (0.0024 - 0.00298)^2 + (0.0035 - 0.00298)^2 + (0.0042 - 0.00298)^2 \right]$$

$$= 8.364 \times 10^{-6} \quad \text{If you use calc/excel, you get to } 0.0000079$$

$$F_0 = \frac{(8.364 \times 10^{-6}) / (5-1)}{(0.00000068) / (5(5-1))} = 30.75$$

$$F_0 = 30.75 \quad F_0 = 29 \text{ if using exact value } 8$$

- f. (4 points) Describe what the test statistic represents (specifically the numerator and denominator of the test statistic as they related to sources of variation) and what it means if the test statistic is very large versus very small

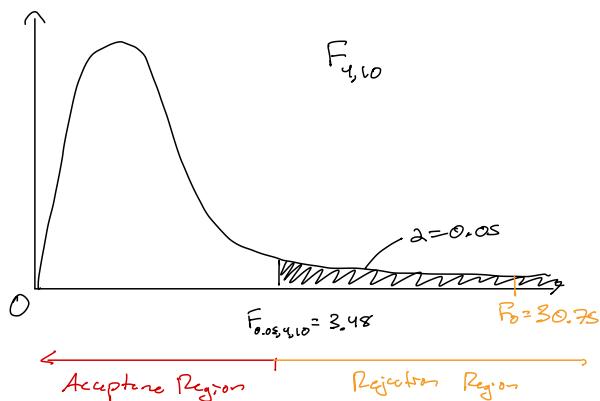
F_0 is effectively $\frac{\text{Var}_{\text{treatments}}}{\text{Var}_{\text{error}}}$. The numerator describes variance from treatments and the denominator describes variance from random error. If F_0 is very large variance is coming from treatments (but hypothesis test determines significance). If F_0 is small, variance is comparable to random error. $D.F$ has an impact on the value.

- g. (4 points) Define your general rejection criteria using fixed significance level testing, sub in values and write out the critical value

$$\text{Reject } H_0 : F_0 > F_{\alpha, n-1, n(n-1)}$$

$$\Rightarrow F_0 > F_{0.05, 4, 10} = 3.48$$

- h. (10 points) Sketch out the reference distribution for this problem and label: the type of reference distribution used and any defining characteristics (shape, DOF if applicable, and where zero is), the critical value(s), the test statistic, alpha, the acceptance region, and the rejection region.



- i. (3 points) Perform any remaining calculations and state your conclusions. Remember to state them in the context of the problem statement, and phrase them appropriately

$$F_0 = 30.75 > F_{\text{crit}} = 3.48 \rightarrow \text{Reject } H_0$$

Reject the null and accept the alternative. We have significant evidence to indicate that the wt% of material B impacts % transmission.

- j. (7 points) Describe when it is appropriate to do post-hoc testing, and assuming it is appropriate in this case, conduct a Fisher's post hoc test and identify which groups are significantly different than the normal process setting (Material B at 5%).

A post-hoc test is only appropriate once the ANOVA H_0 is rejected!

Determine significantly different: $f: |\bar{y}_{i\cdot} - \bar{y}_{j\cdot}| > LSD$

$$S_{us. 4} = |0.0023 - 0.0028| = 0.0005$$

$$LSD = t_{\frac{\alpha}{2}, n-1} \sqrt{\frac{2MSE}{n}}$$

$$S_{us. 3} = |0.0023 - 0.0021| = 0.0002$$

$$= t_{0.025, 10} \sqrt{\frac{2(0.0000063)}{3}}$$

$$S_{us. 2} = |0.0023 - 0.0028| = 0.0005$$

$$= 2.228 (2.129 \times 10^{-4})$$

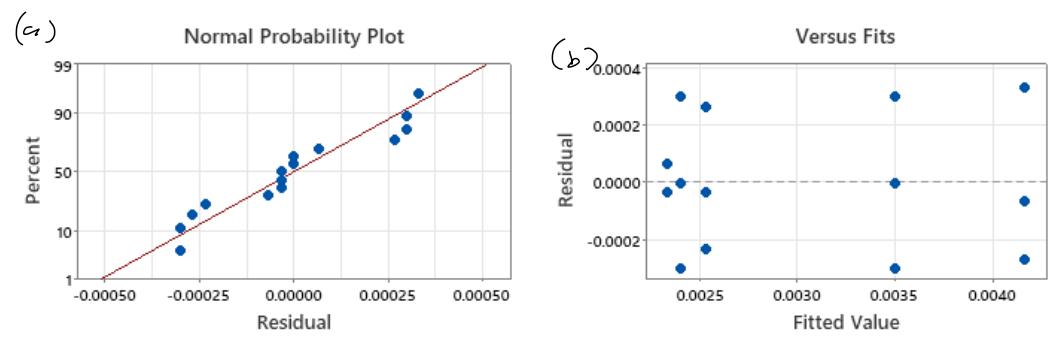
$$S_{us. 1} = |0.0023 - 0.0042| = 0.0019$$

$$= 4.74 \times 10^{-4}$$

$$= 0.000474$$

1 and 2 wtr% are significantly different than S-wtr% (norm)

- k. (4 points) If these are the residual plots:



1. State what assumptions are being checked for each plot

(a) Errors and observations are normally distributed

(b) Constant Variance

2. Describe what you should do if the assumptions are violated

Transform the data and perform ANOVA again.

- I. You are concerned that the variance of the data for 5 wt% treatment is different than the other treatments.
- (1 points) What is the name of the hypothesis test I should perform if I want to know if the variance from 5% treatment and 4% treatment are different?

Z -variance test!

- (2 points) Write out the correct null and alternative hypotheses for this test

$$H_0: \sigma_1^2 = \sigma_2^2 \quad \text{where } 1 = 5\text{wt\%}, \text{ and}$$

$$H_a: \sigma_1^2 \neq \sigma_2^2 \quad 2 = 4\text{wt\%}$$

- (5 points) Calculate the test statistic (report 3 significant figures)

Write out the general formula:

Sub in values and calculate:

$$F_0 = \frac{s_1^2}{s_2^2}$$

$$F_0 = \frac{0.000058^2}{0.000252^2} = 0.052973$$

- (8 points) 1. Write out the rejection critical using p-value method and 2. determine the p-value from your test statistic (report 2 sig figs in your p-value). 3. Show the p-value graphically. In your sketch, you should have the reference distribution labeled with any defining characteristics (shape, zero, and degrees of freedom if applicable), the test statistic, and the p-value).

(i) $R_{\text{cut}}: F_{p-\text{val}} < 2 = 0.05$

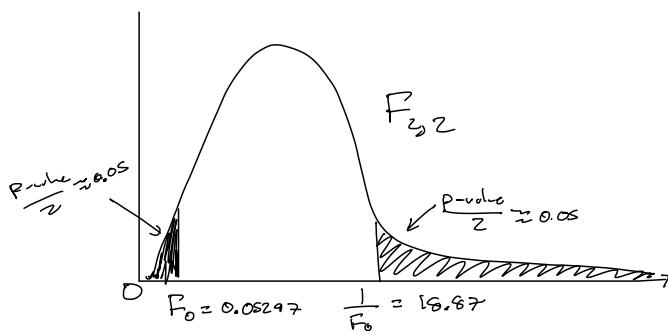
(ii) $F_{z_1 z_2} (F_{n_1-1, n_2-1})$

$$19.00 = F_{0.05, 2, 2}$$

$$\therefore p\text{-val} \approx 2(0.05) = 0.10$$

$$\Rightarrow p\text{-val} \approx 0.10$$

(iii)



5. (3 points) State your conclusions. Remember to state them in the context of the problem statement, and phrase them appropriately. Also use your conclusion to answer: is there cause for concern?

$P-value = 0.1 > 0.05 \Rightarrow$ Fail to reject H_0 . We do not have enough evidence to say that the variance at 5 wt% is different than 4 wt%. No cause for concern as the implies constant variance is still valid.

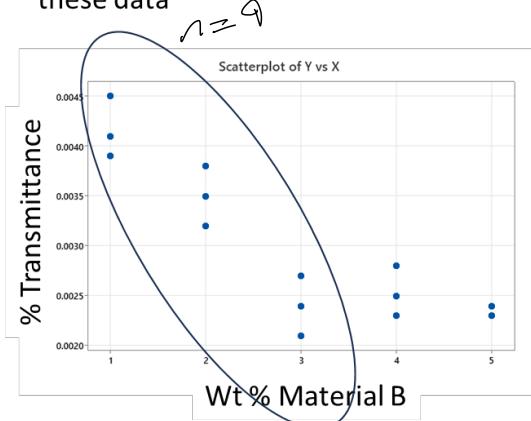
- m. (3 points) Provide and answer to your boss who was wondering: "Why is the rejection criteria for ANOVA one sided, whereas sometimes an F-test can be one- or two-sided?"

ANOVA is looking to see if a "signal" is larger than "noise". We therefore only care if the variance of the "signal" is larger than random error, thus it is one sided. If using an F-test we are determining if two variances are different, thus we need a two-sided test.

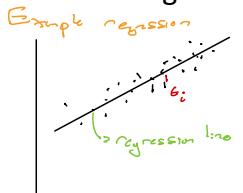
- Problem 2.** (22 points total) You suspect that the Material B wt% (from 1 to 3 wt% - predictor variable) is linearly correlated to % Transmittance (response variable). You want to perform a hypothesis test to see if the linear correlation is significant.

Performing linear regression on these data

$k = \# \text{ of regression lines}$
 $n = \# \text{ of data points}$



- a. (3 points) You perform the liner regression and your estimated coefficient for the slope is -0.00088 and the intercept is 0.0051. What is the name of the method that is used to estimate these coefficients, and describe generally how that method works. (Hint: I am NOT looking for how to calculate S_{xx} or S_{xy} here). You can use a diagram and general formulas for your answer if needed.



The method of least squares is used to estimate the coefficients in linear regression. In this method, the vertical deviations from data and the regression line (e_i 's) are minimized. These are the residuals, and the model aims to minimize $\sum_{i=1}^n e_i^2$.

- b. This is your Minitab output from the simple linear regression:

Analysis of Variance

Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	1	0.00000468167	0.000468167	51.16	0.000
Error	6	0.00000064056	0.0000010676		
Total	7				

- (1) (2 points) State the degrees of freedom for the $SS_{\text{Regression}}$ and the SS_{Error}

$$SS_{\text{Regression}} : \nu_C = 1$$

$$SS_{\text{Error}} : n - \nu_C - 1 = n - 2 = 7$$

- (2) (2 points) Write out the formula used and calculate the SS_{Total} using information given in the table above

$$\begin{aligned} SS_{\text{Total}} &= SS_R + SS_E \\ &= 0.00000468167 + 0.00000064056 \\ &= 0.0000053 \end{aligned}$$

- (3) (2 points) Write out the formula for and calculate R^2

$$\frac{SS_R}{SS_T} = \frac{0.00000468167}{0.0000053} = 0.88$$

- (4) (3 points) State your conclusions. Remember to state them in the context of the problem statement, and phrase them appropriately. Assume that no issues were noted with the residual plots.

$P_{\text{value}} = 0.2 > 0.05 \rightarrow P_{\text{reject }} H_0 \text{ and accept } H_a$. We have significant evidence to reject the claim that there is no linear relationship between material B from 1-3% and % transmittance. The wt% of material B from 1-3% is linearly correlated with % transmittance.

- c. (10 points) You think there might be another variable that could be useful in your model (wt% of another material, Material A). You add in this variable and run multiple linear regression to get the following output:

Analysis of Variance

Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	2	0.000005	0.000002	26.21	0.001
Wt% B	1	0.000002	0.000002	19.92	0.004
Wt% A	1	0.000000	0.000000	1.03	0.349
Error	6	0.000001	0.000000		
Total	8	0.000005			

Model Summary

S	R-sq	R-sq(adj)	R-sq(pred)
0.0003018	89.73%	86.30%	76.89%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value	VIF
Constant	0.005844	0.000760	7.69	0.000	
X	-0.001100	0.000246	-4.46	0.004	4.00
X2	-0.000289	0.000285	-1.02	0.349	4.00

- (1) (3 points) The R^2 increased after adding the new variable. Should we keep the variable in the model? Why or why not?

No. R^2 is not a valid tool for determining the significance of a variable in regression. In addition, the partial F-test shows P-value > 0.05 = α for wt% A, so there isn't enough evidence to say it has a significant relationship with % transmittance.

- (2) (3 points) R^2 adjusted before adding Material A as a predictor was 86%. Explain why R^2 adjusted did not go up when adding the new variable. Use the formula in your answer.

$$R_{adj}^2 = \frac{SS_R/k}{SS_T/(n-1)}$$

R_{adj}^2 includes degrees of freedom which prevents k^2 from increasing simply by adding more regression terms.

- (3) (4 points) State your conclusions appropriately in the context of the problem statement and write out the final model

Fail to reject H_0 . We don't have enough evidence to say wt% A has a significant relationship with % transmittance.

The final model is: $y = 0.0051 - 0.00088x$

$\begin{cases} y = \% \text{ transmittance} \\ x = \text{wt\% material B} \end{cases}$

Problem 3. (11 points total) Your company wants to establish an Xbar control chart on the process for component C wt%. There are 22 preliminary samples, each of size 8, on the internal diameter of the seal. The summary data (in mm) are as follows:

Sample #	Range (R)	Average (Xbar)
1	5	15
2	9	11
3	8	11
4	4	13
5	9	13
6	5	10
7	10	14
8	5	14
9	9	12
10	8	13
11	8	15
12	8	11
13	8	15
14	8	12
15	5	13
16	8	14
17	7	14
18	4	11
19	9	12
20	8	15
21	8	14
22	5	11
Sum	158	283

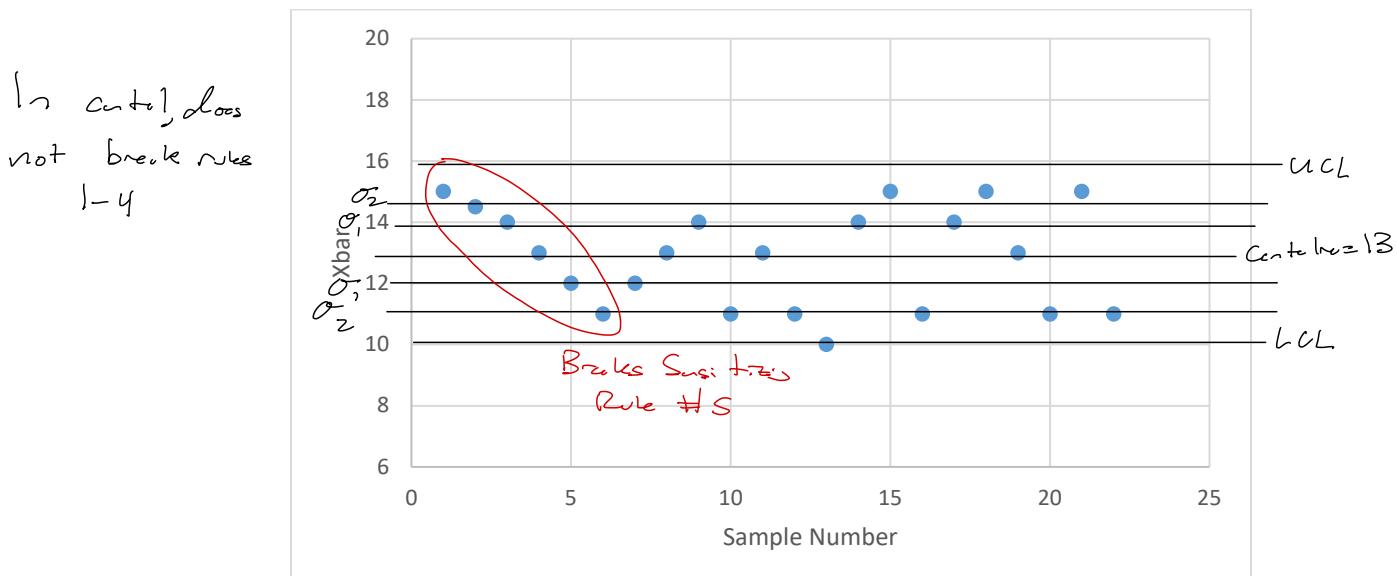
$$\bar{R} = \frac{158}{22} = 7.18 \quad \bar{\bar{x}} = \frac{283}{22} = 13.0$$

- a. (8 Points) Write out the general formulas for and calculate the centerline as well as the three sigma control limits that should be used on the Xbar and R control chart

$$X_{bar} \left\{ \begin{array}{l} UCL = \bar{x} + A_2 \bar{R} = 13.0 + 0.373(7.18) = 15.7 \approx 16 \\ Cen-tre = \bar{x} = 13.0 \approx 13 \\ LCL = \bar{x} - A_2 \bar{R} = 13.0 - 0.373(7.18) = 10.3 \approx 10 \end{array} \right.$$

$$R \left\{ \begin{array}{l} UCL = D_4 \bar{R} = 1.864(7.18) = 13.4 \approx 13 \\ Cen-tre = \bar{R} = 7.18 \\ LCL = D_3 \bar{R} = 0.136(7.18) = 0.98 \approx 1 \end{array} \right.$$

- b. (3 points) In the next month, you take an additional 20 samples in a similar manner as before. Sketch in your sigma control limits (UCL and LCL) and centerline from part A, as well as the 1 sigma and 2 sigma warning limits on this graph of your data. After analyzing this plot, what is your recommendation?



Process should be investigated as sampling #5 is broken

Homework 6

Due 3/20/25

Reminder: Reading Chapter 4 of the textbook will help you answer these questions. In lectures 12 and 13 and some of 14 you will have seen relevant background and examples. *Unless otherwise noted, write out by hand how to calculate the desired values (you can abbreviate but sub in some values). Also please pay attention to sig figs, and remember to sketch out the reference distribution when you do the hypothesis tests – each problem should have a reference distribution sketch.*

- 1) (a) Do problem 4.13 (a) (4.23 in 7th edition) in your book doing the step-by-step hypothesis method by hand (either p-value or critical value is fine). Don't do part b from the book but instead do the following: (b) check your answer in part a in Minitab, print the output and report the p-value you obtained
- 2) You wish to know if changing the catalyst in a process causes the yield to have more variability. Data are collected on the yield for each catalyst:

Cat. 1 Cat. 2

57.9 66.4

66.2 71.7

65.4 70.3

65.4 69.3

65.2 64.8

62.6 69.6

67.6 68.6

63.7 69.4

67.2 65.3

71.0 68.8

- a) Cut and paste this data into Minitab and test the hypothesis that the variances are equal. Assume normality (there is a box to check in this test that assumes normality). Report the output and your conclusions – answer the question of whether you could assume equal variance for these samples if you were to conduct a t-test. Use $\alpha=0.05$
- b) Check the normality assumption using Minitab – does it seem appropriate?

3.) Do problem 4.23 from your book (4.39 in the 7th edition) and use Excel to process the data for parts a-c. Provide the excel sheets and write out the equations used to get your answers. For part b, put the calculated residuals you got from Excel into Minitab to produce your normal probability plot. Use $\alpha=0.05$.

For the final part of this problem (part d) check your answers in Minitab (including the residual plot via the ANOVA function), click assume equal variance box. Report the P-value.

4.) Do problem 4.25 (4.41 in 7th edition) from the book. Use Minitab to conduct the hypothesis testing (assuming equal variances).

a) Conduct the analysis and state your conclusions and sketch out the situation on a reference distribution. Also answer the following questions:

a-1) how many factors are in this experiment?

a-2) how many levels?

a-3) What is the critical value of F if you were to compare the test statistic to it instead of using the P-value?

b) Analyze residuals according to the following directions

b-1) Calculate the residuals in Excel. (Write out the equations needed and report the sheet you used to get your answers).

b-2) Analyze the residuals by plotting them on a normal probability plot in Minitab (cutting and pasting from your table you made in Excel). State if you think normality is an OK assumption.

b-3) Plot the residuals you calculated in Excel, versus the factor levels. (Can use Minitab or Excel) Answer what assumption this graph checks, and comment if the assumption is still valid.

b-4) Plot the residuals vs. the fitted value (\hat{y}_i). Can use Minitab or Excel. Do the residuals depend on the response? State what you would do if they did appear to have a pattern.

C) Do this next extra part (not in the book)

c-1) Conduct the Fisher's LSD post-hoc test using Excel given the $MSE = 6.625$. Compare: 0.37 vs. 0.51, 0.37 vs. 0.71, and 0.51 vs. 0.71 by hand and conduct the full analysis using Minitab

c-2) Construct a plot that represent your data with an appropriate caption by plotting the average \pm the standard deviation at each treatment mean, and above the bars, listing the letter assignments from the post hoc testing. For the caption, write something like this

example: "All data are represented by the average \pm standard error. ANOVA was performed to confirm statistical significance of (insert what factor you are investigating), followed by the (name here) *post hoc* test to determine differences between treatment means using $\alpha=X$. Bars that do not share a letter within a graph are statistically different."

Of course variations of this can be written, but you should include all of the above information.

Remember for sketching reference distributions for a hypothesis test:

P-value method: Sketch out the reference distribution, labeling the test statistic, and the p-value. In addition, fully define the distribution you are sketching on with the correct shape, any defining characteristics and where 0 is. It is also good to draw in alpha – remembering that alpha is not the same as your p-value.

Fixed Significance Level Testing: Sketch out the reference distribution, labeling the acceptance region, the rejection region, α , the critical value(s) and your test statistic. In addition, fully define the distribution you are sketching on with the correct shape, any defining characteristics and where 0 is.

ECHE313 Homework 6 - Due 03/20/25

Trevor Swan (tcs94)

U.13: ~
book

1. $\bar{x}_1 = 9.85$	$\bar{x}_2 = 8.08$
$s_1^2 = 6.29$	$s_2^2 = 6.18$
$n_1 = 10$	$n_2 = 8$

a) Inference on Two chemical processes variance - independent & normality assumed.
 ↳ Two sample variance test

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$F_0 = \frac{s_1^2}{s_2^2} = \frac{6.29}{6.18} = 1.10$$

$$H_a: \sigma_1^2 \neq \sigma_2^2$$

$$D.F_1 = 10 - 1 = 9$$

$$D.F_2 = 8 - 1 = 7$$

$$\alpha = 0.05$$

$$\frac{\alpha}{2} = 0.025$$

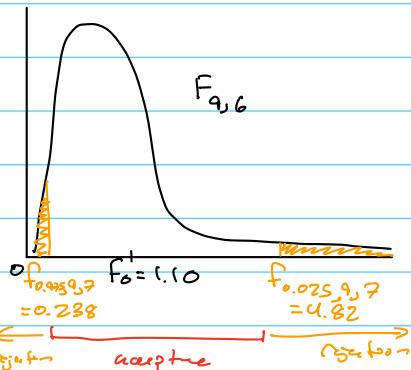
Rejection Criteria:

$$\begin{cases} F_0 > F_{0.025, 9, 7} \\ F_0 < F_{0.975, 9, 7} = \frac{1}{F_{0.025, 9, 7}} \end{cases}$$

$$F_{0.025, 9, 7} = 4.82 \checkmark$$

$$F_{0.975, 9, 7} = \frac{1}{4.20} = 0.238$$

$$F_0 < F_{0.025, 9, 7} \text{ & } F_0 > F_{0.975, 9, 7}$$



Fail to reject H_0 . We have insufficient evidence to reject the claim that the variances of the two chemical processes are different. We cannot conclude that the variance of the new production unit differs from the old production units' variances.

b) Minitab P-value: 0.922 → Fail to reject as $p\text{-value} = 0.922 > \alpha = 0.05$

Test

Null hypothesis $H_0: \sigma_1^2 / \sigma_2^2 = 1$

Alternative hypothesis $H_1: \sigma_1^2 / \sigma_2^2 \neq 1$

Significance level $\alpha = 0.05$

Test

Method Statistic DF1 DF2 P-Value

F	1.10	9	7	0.922
---	------	---	---	-------

b) Residuals

b-1) $E_{ij} = y_{ij} - \bar{y}_{j\cdot}$

$$\bar{y}_{1\cdot} = \frac{80 + 83 + 83 + 85}{4} = 82.75$$

$$\bar{y}_{2\cdot} = \frac{75 + 78 + 79 + 79}{4} = 77$$

$$\bar{y}_{3\cdot} = \frac{74 + 73 + 76 + 72}{4} = 75$$

$$\bar{y}_{4\cdot} = \frac{67 + 72 + 74 + 74}{4} = 71.75$$

$$\bar{y}_{5\cdot} = \frac{62 + 62 + 67 + 69}{4} = 65$$

$$\bar{y}_{6\cdot} = \frac{60 + 61 + 64 + 66}{4} = 62.75$$

$$E_{11} = y_{11} - \bar{y}_{1\cdot} = 80 - 82.75 = -2.75$$

$$E_{12} = y_{12} - \bar{y}_{1\cdot} = 83 - 82.75 = 0.25$$

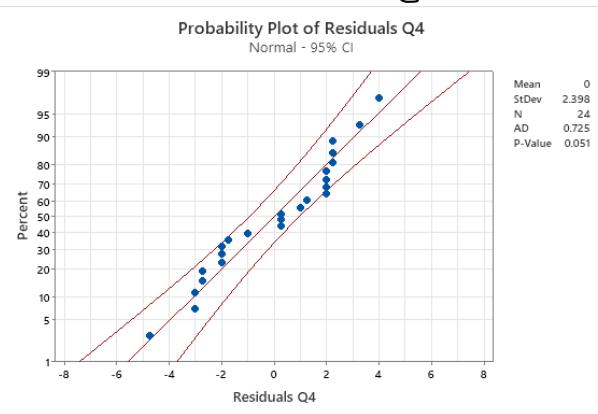
$$E_{13} = y_{13} - \bar{y}_{1\cdot} = 83 - 82.75 = 0.25$$

⋮

$$E_{64} = y_{64} - \bar{y}_{6\cdot} = 66 - 62.75 = 3.25$$

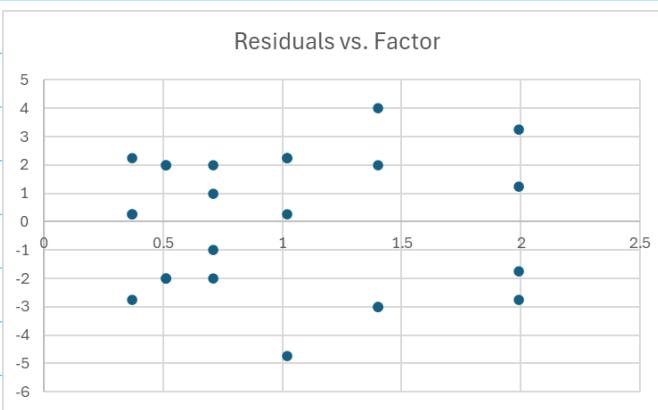
Residuals on 'Q1stn4' Sheet from
Excel ↴

b-2) Normal Probability Plot of Results



$p\text{-Value} = 0.051 < 0.05 = \alpha$, so we fail to reject the normality assumption. With heavily clear visible and practice being so high, normality is an OK assumption. This p-value was close to 2, so I would proceed with caution...

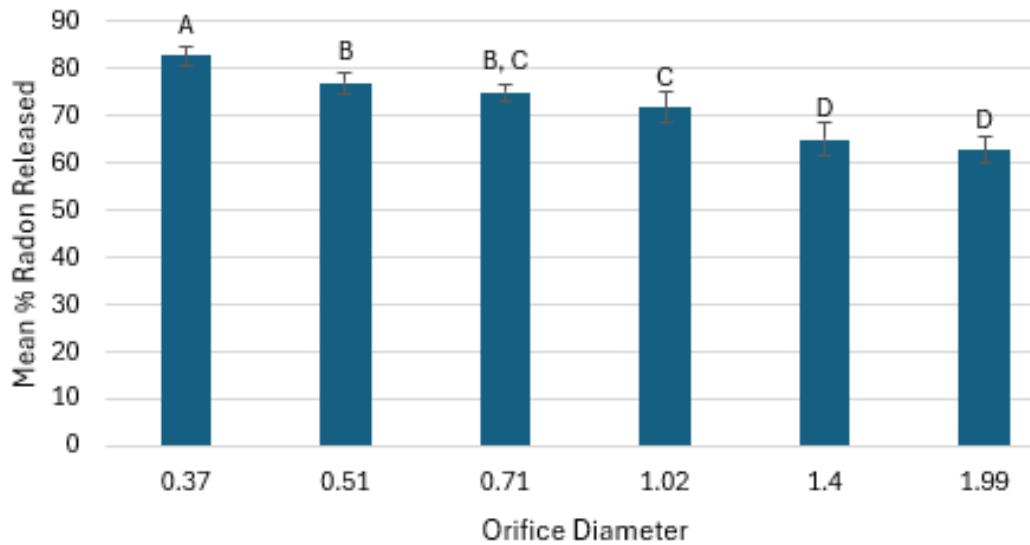
b-3) Residuals vs. Factor Levels



This graph checks equal variance at each factor level. The graph shows the residuals at each factor level being randomly, but evenly, scattered along each factor, so the equal variance assumption is ok here.

C-Z)

Average % Randon Released by Orifice Diameter



All data are represented by the average \pm standard error. ANOVA was performed to confirm statistical significance of Orifice Diameter difference on the mean radon percentage released, followed by Fischer's LSD post hoc test to determine differences between treatment means using $\alpha=0.05$. Bars that do not share a letter within a graph are statistically different.

Homework 7

Due 3/27/25 by end of the day

Directions: Reading the supplemental reading posted as Lecture 14 reading and the last part of Chapter 4 will help you answer these questions. In lectures 14, 15 and 16 you will have seen relevant background and examples. Do problem 1-2 by hand using Excel (unless otherwise noted). For the rest of the homework, use Minitab. Try to be aware of the significant figures when reporting model coefficients.

- 1) Do problem 4.26 (4.43 in 7th edition) in your book, and answer parts a and b by hand (using Excel to aid in the calculations) but use Minitab to check parts a and b and calculate the CI in part c. Keep 3 significant figures and use $\alpha=0.05$. Also, do extra part (not in book):
 - d) calculate the R^2 by hand and explain what this tells you
- 2) Do problem 4.28 (4.45 in 7th edition) by plotting: (1) residuals on a normal probability plot, (2) residuals versus predicted value, (3) residuals vs. x. To complete these, calculate the residuals in Excel and paste them as part of your answer. You can produce the plots using Excel or Minitab. Remember to comment on if you think the assumptions are valid.
- 3) Do problem 4.27 (4.44 in the 7th edition) in your book – you can use Minitab for the entire problem (no need to calculate by hand). For part a) also generate a scatter plot of the data and convince yourself a linear model is appropriate and be sure to highlight or report your regression equation. For part b) be sure to state the null and alternative hypothesis and your conclusion. For part c) report the CI and describe what it means. Keep 3 significant figures and use $\alpha=0.05$. Also add these parts to this problem (not in your book):
 - d) use Minitab to generate a plot of the linear regression lines with 95% confidence intervals and explain what the intervals mean
- 4) An electric utility is interested in developing a model relating peak-hour demand (y in kilowatts) to total monthly energy usage during the month (x , kilowatt hours). Data for 50 residential customers are shown in the table below (last page). You can use Minitab to solve this entire problem.
 - a. Plot the scatter diagram of y vs x .
 - b. Fit the simple linear regression model (report the equation)
 - c. Test for significance of the simple linear regression using $\alpha=0.05$. State the null and alternative hypothesis.
 - d. Plot the residuals versus predicted, and comment on the underlying assumptions, specifically, is the equality of variance assumption valid?
 - e. Find a simple linear regression model using \sqrt{y} as the response. Does the transformation stabilize the inequality found in part d

- 5) Do problem 4.30 – 4.33 (4.49-4.52 in the 7th ed) using Minitab, alpha = 0.05. Create a contour plot using Minitab and explain why adding age to the model might be helpful. Using the contour plot, report the range of satisfaction for a 40 year old with a disease severity of 30 and the range for a 40 year old with severity of 60 according to your model.

Data for Problem 4 (cut and paste into Minitab):

	x	y
1	679	0.79
2	292	0.44
3	1012	0.56
4	493	0.79
5	582	2.7
6	1156	3.64
7	997	4.73
8	2189	9.5
9	1097	5.34
10	2078	6.85
11	1818	5.84
12	1700	5.21
13	747	3.25
14	2030	4.43
15	1643	3.16
16	414	0.5
17	354	0.17
18	1276	1.88
19	745	0.77
20	795	3.7
21	540	0.56
22	874	1.56
23	1534	5.28
24	1029	0.64
25	710	4
26	1434	0.31
27	837	4.2
28	1748	4.88
29	1381	3.48
30	1428	7.58
31	1255	2.63
32	1777	4.99
33	370	0.59
34	2316	8.19
35	1130	4.79
36	436	0.51
37	770	1.74
38	724	4.1
39	808	3.94
40	790	0.96
41	783	3.29
42	406	0.44
43	1242	3.24
44	658	2.14
45	1746	5.71
46	895	4.12
47	1114	1.9
48	413	0.51
49	1787	8.33
50	3560	14.94

ECHE313 Homework 7 - Due 03/27/25

Trevor Swan (tcs94)

4.26 in book

Table 4E.6 Tensile Strength Data for Exercise 4.26

Strength	Percentage Hardwood	Strength	Percentage Hardwood
160	10	181	20
171	15	188	25
175	15	193	25
182	20	195	28
184	20	200	30

a) $n=10$

$$\sum x_i = 160 + 171 + \dots + 200 = 1829$$

$$\sum x_i^2 = 160^2 + 171^2 + \dots + 200^2 = 335825$$

$$\sum y_i = 10 + 15 + \dots + 30 = 208$$

$$\sum y_i^2 = 10^2 + 15^2 + \dots + 30^2 = 4684$$

$$\sum x_i y_i = 160(10) + 171(15) + \dots + 200(30) = 38715$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 38715 - \frac{(208)(1829)}{10} = 671.8$$

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 335825 - \frac{(1829)^2}{10} = 1300.9$$

$$\tilde{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{671.8}{1300.9} = 0.5164 ; \quad \bar{x} = \frac{1}{10}(\sum x_i) = \frac{1}{10}(1829) = 182.9 \\ \bar{y} = \frac{1}{10}(\sum y_i) = \frac{1}{10}(208) = 20.8$$

$$\beta_0 = \bar{y} - \tilde{\beta}_1 \cdot \bar{x} = 20.8 - 0.5164(182.9) = -73.65$$

$$\hat{y} = 0.516x - 73.65 \quad \text{where } \begin{cases} x = \text{Strength} \\ y = \text{Percentage Hardwood} \end{cases}$$

correct procedure

b) ANOVA for linear regression

$$H_0: \beta_1 = 0$$

$$H_a: \beta_1 \neq 0$$

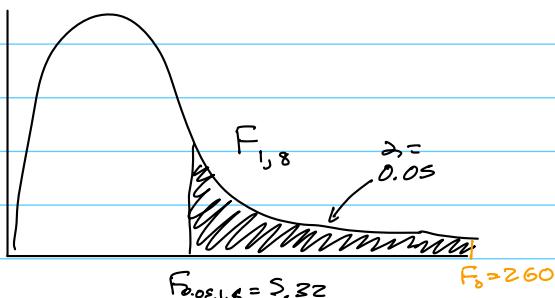
Test Statistic:

$$SS_E = \sum (\hat{y}_i - \bar{y})^2 = (8.97 - 20.8)^2 + (16.65 - 20.8)^2 + \dots + (29.63 - 20.8)^2 = 346.93$$

$$SS_B = \sum (y_i - \bar{y})^2 = (10 - 20.8)^2 + \dots + (30 - 20.8)^2 = 10.67$$

$$F_0 = \frac{SS_B/1}{SS_E/(n-2)} = \frac{346.93}{10.67/(8)} = 2.60$$

$$\text{Reject } F_0: F_0 > F_{0.05, 1, 8}$$



$$2.60 > 5.32 \rightarrow \text{Reject } H_0$$

We have sufficient evidence to support the claim that there is a significant linear relationship between percentage hardwood (y) and strength (x).

c) Regression Equation

$$\text{Percentage Hardwood} = -73.65 + 0.5164 \text{ Strength}$$

Coefficients

Term	Coeff	SE Coef	95% CI	T-Value	P-Value	VIF
Constant	-73.65	5.87	(-87.19, -60.12)	-12.55	0.000	
Strength	0.5164	0.0320	(0.4426, 0.5903)	16.12	0.000	1.00

Model Summary

S	R-sq	R-sq(adj)	PRESS	R-sq(pred)	AICc	BIC
1.15513	97.01%	96.64%	17.4802	95.11%	39.03	35.94

Analysis of Variance

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	1	346.93	97.01%	346.93	346.925	260.00	0.000
Strength	1	346.93	97.01%	346.93	346.925	260.00	0.000
Error	8	10.67	2.99%	10.67	1.334		
Total	9	357.60	100.00%				

$$d) R^2 = \frac{SS_R}{SS_T} = \frac{SS_R}{SS_R + SS_E} = \frac{346.93}{346.93 + 10.67} = 97.0\%$$

97% of the variability in the target variable is accounted for by the regression model.

95% CI on β_1 (slope) \rightarrow (0.4426, 0.5903) \rightarrow (1.61, 2.15) \rightarrow Proprietary errors

In repeated sampling, the interval calculated captures the true slope β_1 95% of the time.

Also, calculations from previous page (Excel assist) match the results output

4.28 in book

2.

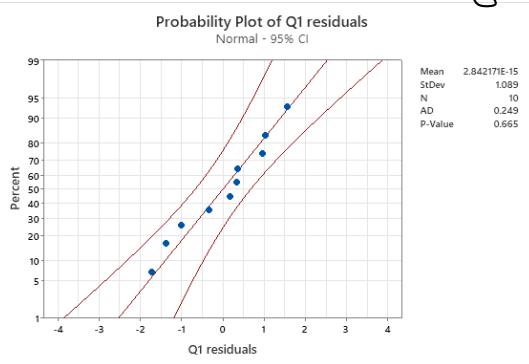
Table 4E.6 Tensile Strength Data for Exercise 4.26

Strength	Percentage Hardwood	Strength	Percentage Hardwood
160	10	181	20
171	15	188	25
175	15	193	25
182	20	195	28
184	20	200	30

Data from
problem ①

$$x = \text{strength}; y = \% \text{ Hardwood}$$

a) Residual Normal Probability Plot



Residuals

$$e_1 = 10 - 8.97 = 1.03$$

$$e_2 = 15 - 8.97 = 6.345$$

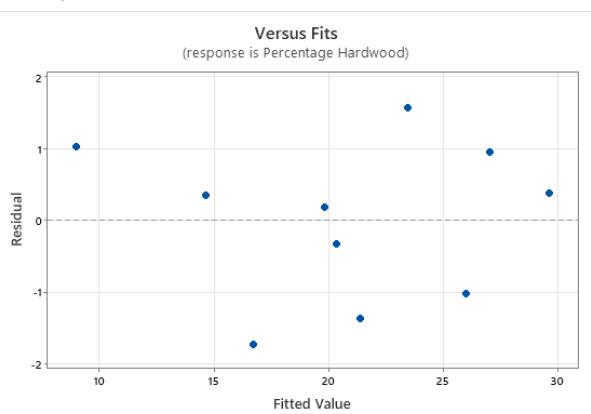
:

$$e_{10} = 30 - 29.62 = 0.369$$

Insignificant evidence to reject the claim that the data is not normal.

As $p\text{-value} = 0.665 > 0.05 = \alpha$, we can support the normality assumption of the model.

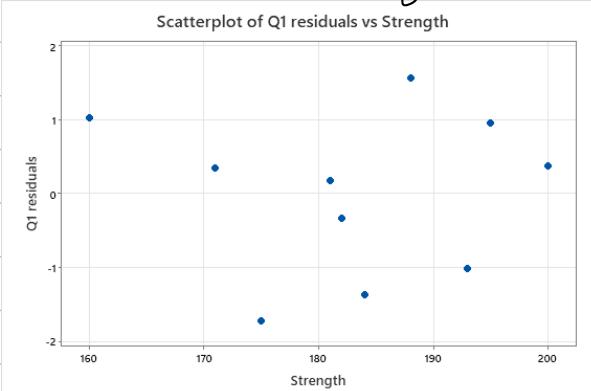
b) Residuals vs. Fitted Value



The residuals are evenly distributed about $y=0$. There is no clear pattern, and the equal variance assumption is valid here.

If there was a pattern, I would apply a transformation

c) Residuals vs. Strength (x)



The residuals show no patterns and are evenly distributed about 0. The equal variance assumption is again valid.

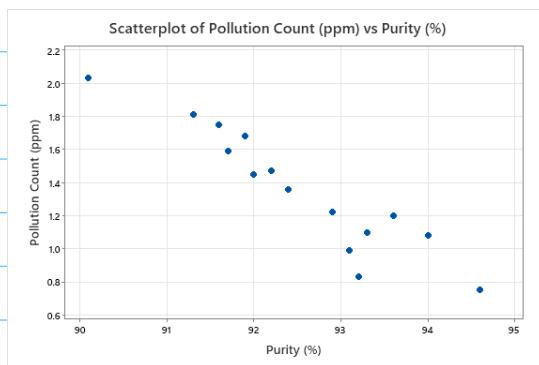
4.27 in book

3.

Purity (%)	93.3	92.0	92.4	91.7	94.0	94.6	93.6	
Pollution count (ppm)	1.10	1.45	1.36	1.59	1.08	0.75	1.20	
Purity (%)	93.1	93.2	92.9	92.2	91.3	90.1	91.6	91.9
Pollution count (ppm)	0.99	0.83	1.22	1.47	1.81	2.03	1.75	1.68

MSS
outlier
not lookin'

a)



With that a fit function a clear negative

linear relationship between Pollution Count (y)
and Purity Percent (x). I'm convinced!

Regression:

$$\hat{y} = 29.23 - 3.03x$$

$$\begin{cases} \hat{y} = \text{Pollution Count (ppm)} \\ x = \text{Purity Percent} \end{cases}$$

b) ANOVA test for linear regression

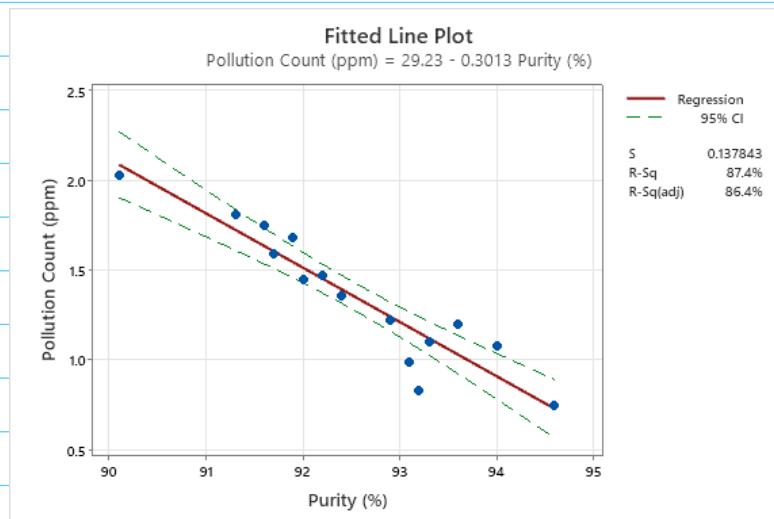
$$\begin{aligned} H_0: \beta_1 &= 0 \\ H_a: \beta_1 &\neq 0 \end{aligned} \quad \left. \begin{array}{l} \text{Min. t.b. p-value} \approx 0 < \alpha = 0.05 \end{array} \right.$$

p-value < α , so we have significant evidence to reject the claim that there is no linear relationship. We can accept there is a significant linear relationship between Pollution Count (y) and Purity Percent (x).

c) 95% CI for Purity %: $(-0.3698, -0.2827)$ ← slope
 $\text{Pollut. Count (Const.)} = (22.89, 35.57)$ ← intercept

In repeated samples, the intervals calculated capture the true slope or intercept 95% of the time.

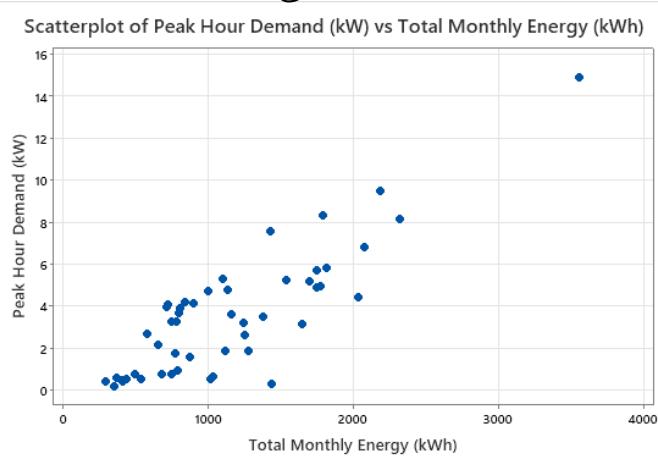
d)



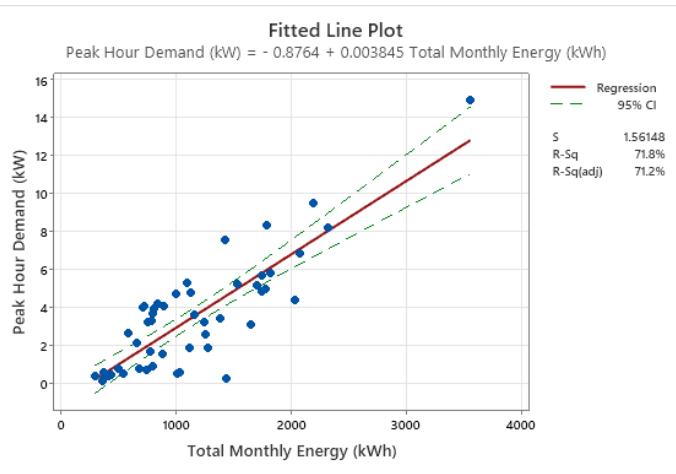
In repeated experiments
there is a 95% chance
that the CI on the
responses calculated contains
the true responses.

$$4. \begin{cases} x = \text{Total Monthly Energy Used (kWh)} \\ y = \text{Peak Hour Demand (kW)} \end{cases}$$

a) Scatter plot of y vs x



b) Linear Fit



c) $H_0: \beta_1 = 0$ $H_a: \beta_1 \neq 0$

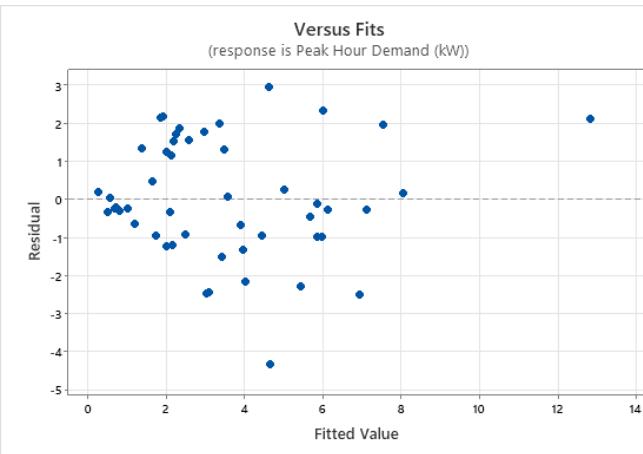
Analysis of Variance

Source	DF	Seq SS	Contribution	Adj SS	Adj MS	F-Value	P-Value
Regression	1	297.7	71.78%	297.7	297.728	122.11	0.000
Total Monthly Energy (kWh)	1	297.7	71.78%	297.7	297.728	122.11	0.000
Error	48	117.0	28.22%	117.0	2.438		
Total	49	414.8	100.00%				

p-value $< 0 < \alpha = 0.05$

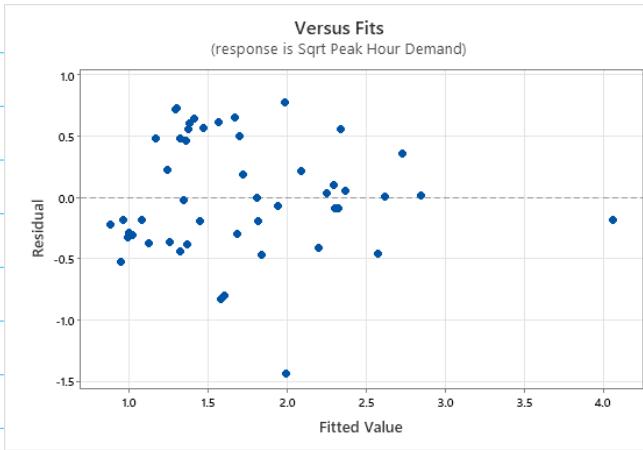
Since $p < \alpha$, we have significant evidence to reject the claim that there is no relationship. We can accept there is a significant linear relationship between the total monthly energy used (x), and the peak hour demand (y).

d)



There seems to be a pattern like pattern despite being roughly dispersed about 0. At low fitted values, residuals are closer together. Thus, the equality of variance may not be a good assumption and we should apply a data transformation to remedy this.

e)



Yes, the d.t. is more evenly distributed about zero and there is not a pattern corresponding to the fitted values. With this transformation, the equality of variance assumption appears to be a valid assumption.

S. a) Problem 4.30

$$(i) \text{Satisfaction } (y) = 131.10 - 1.290 \cdot \text{Age } (x_1)$$

$$\begin{cases} y = \text{Satisfaction} \\ x_1 = \text{Age} \end{cases}$$

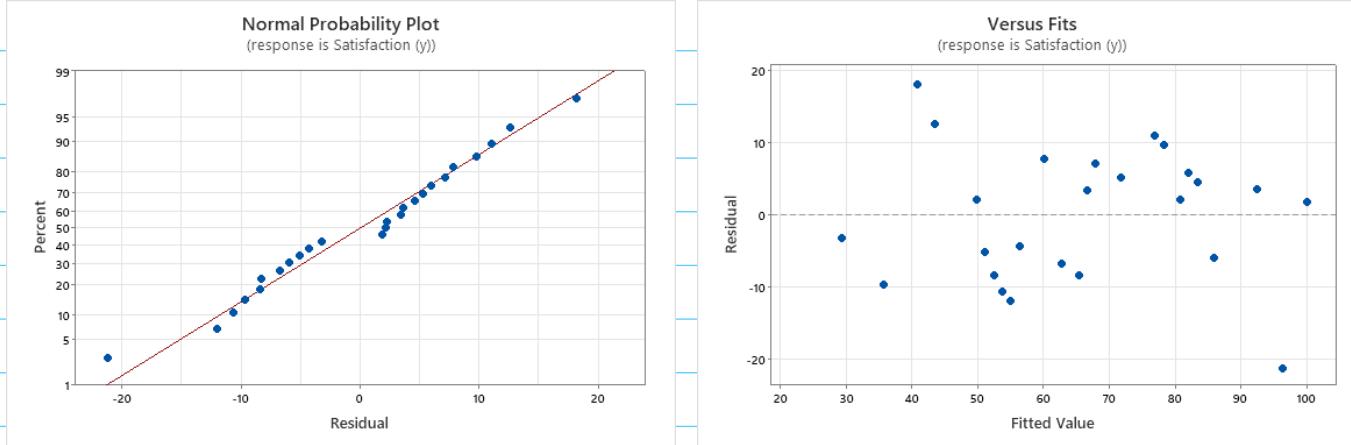
$$(ii) H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

$$\text{Min. t.b} \rightarrow p\text{-value} = 0 < \alpha = 0.05$$

We have significant evidence to reject the claim that there is no relationship. We can accept that there is a significant linear relationship between the Age of the patient (x_1), and their Satisfaction (y).

(iii) 81.24% of the variability is accounted for by the regressor variable age.

b) Problem 4.31



The normal probability plot of the residuals is linear with no outstanding variability, so the normality assumption is valid. The residuals vs. Fitted value is evenly distributed about 0, and no clear pattern is seen. Thus the equality of variance assumption is valid. The model seems to be a good fit.

c) Problem 4.32

$$(i) \text{Substitution } (y) = 143.47 - 1.03 \cdot \text{Age}(x_1) - 0.556 \cdot \text{Severity}(x_2)$$

$$\begin{cases} y = \text{Substitution} \\ x_1 = \text{Age} \\ x_2 = \text{Severity} \end{cases}$$

$$H_0: B_1 = B_2 = 0$$

$$H_a: B_i \neq 0 \text{ for at least one } i$$

$$\text{Min. t.b.} \rightarrow p\text{-value} = 0 < \alpha = 0.05$$

We have significant evidence to reject the claim that there is no relationship. We can accept that there is a significant linear relationship between the Age of the patient (x_1), Severity of the illness (x_2), and their satisfaction (y).

$$(ii) R^2(x_1) = 81.24\% \leftarrow \text{Age}$$

$$R^2(x_2) = 8.42\% \leftarrow \text{Severity}$$

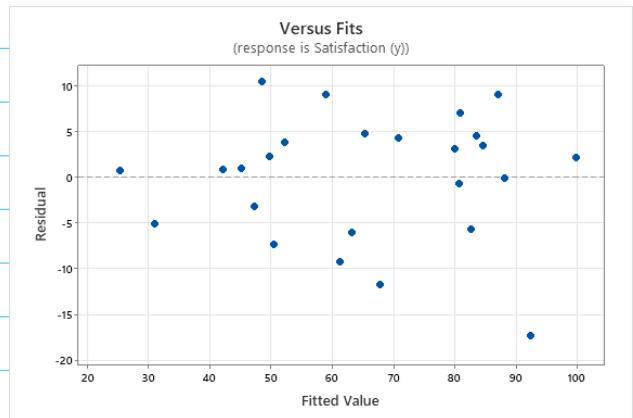
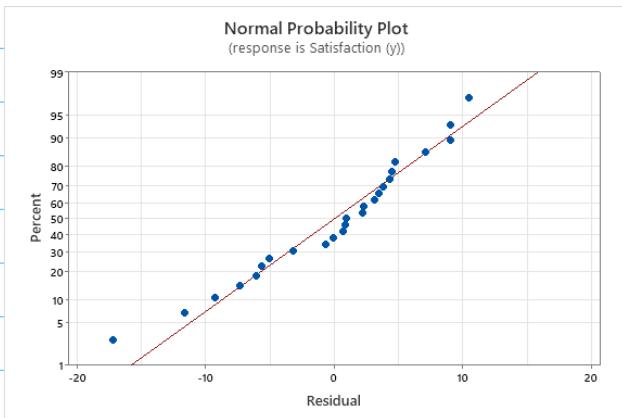
The portion of variability is significant ($> 2 \sigma^2$) in both cases. Both answers are needed.

$$(iii) R^2_{\text{mult.}} = 88.72\% \quad (\text{adj.})$$

$$R^2_{\text{single}} = 80.43\% \quad (\text{adj.})$$

The adjusted R^2 , which accounts for overfitting effects, improved with the full fit expression. Thus we can conclude added severity improved the model's quality.

d) Problem 4.33



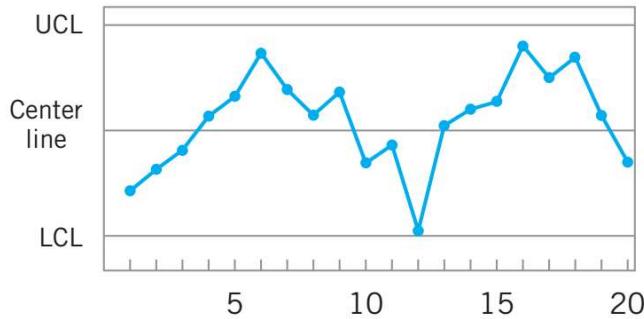
The normal probability plot of the residuals is linear with no outliers, so the normality assumption is valid. The residuals w.r.t. fitted value is evenly distributed about 0 and no clear pattern is seen. Thus the equality of variance assumption is valid. The model seems to be a good fit.

Homework 8

Due 4/3/25 by end of day

Directions: Reading the chapter 5 and 6 will help solve these problems along with content seen in class.

- 1) Most control charts choose 3 sigma levels for their control limits because it seems to be a good balance between type I and type II errors. Discuss the following:
 - a. If narrower limits are chosen, what happens to the magnitude of type I and II error?
 - b. What effect does the sigma level (high or low) have on alpha?
- 2) Laboratory glassware shipped from the manufacturer to Dr. Renner's Lab via an overnight package service has arrived damaged. Develop a cause-and-effect diagram that identifies and outlines the possible causes of this event. You won't necessarily be graded on the details, but include the major components of the "fishbone" diagram.
- 3) Sketch out diagrams and explain why it is important to control both process mean and variability (see Figure 6.1 in the book)
- 4) Problem 5.16, 5.17, 5.18 (5.19, 5.22, and 5.23 in the 7th edition) – use the below control chart (you can cut and paste into a separate file and print it off as part of your answer). For 5.17 use Sensitizing Rules 5-10, and for 5.18 use the Western Electric Rules (1-4).



- 5) Problem 6.1
- 6) A hospital emergency department is monitoring the time required to admit a patient using \bar{x} and R charts. The table below presents summary data for 20 subgroups of two patients each (time is in minutes)

Subgroup	Xbar	R
1	8.3	2
2	8.1	3
3	7.9	1
4	6.3	5
5	8.5	3
6	7.5	4
7	8	3
8	7.4	2
9	6.4	2
10	7.5	4
11	8.8	3
12	9.1	5
13	5.9	3
14	9	6
15	6.4	3
16	7.3	3
17	5.3	2
18	7.6	4
19	8.1	3
20	8	2

- a) Use these data to determine the control limits for the \bar{x} and R control charts for this patient admitting process.
- b) Plot the preliminary data from the first 20 samples on the control charts that you set up in part (a). Is the process in statistical control?
- 7) A high-voltage power supply should have a nominal output voltage of 350V. A sample of four units is selected each day and tested for process-control purposes. The data shown in the table below give the difference between the observed reading on each unit and the nominal voltage times 10; that is $x_i = (\text{observed voltage unit } i - 350) * 10$. Use Minitab to set up the \bar{x} and R charts on this process. Is the process in statistical control?

Sample #	X1	X2	X3	X4
1	6	9	10	15
2	10	4	6	11
3	7	8	10	5
4	8	9	6	13
5	9	10	7	13
6	12	11	10	10

7	16	10	8	9
8	7	5	10	4
9	9	7	8	12
10	15	16	10	13
11	8	12	14	16
12	6	13	9	11
13	16	9	13	15
14	7	13	10	12
15	11	7	10	16
16	15	10	11	14
17	9	8	12	10
18	15	7	10	11
19	8	6	9	12
20	13	14	11	15

ECHE313 Homework 8 - Due 4/3/25

Trevor Swan (tcs94)

1. Control Charts

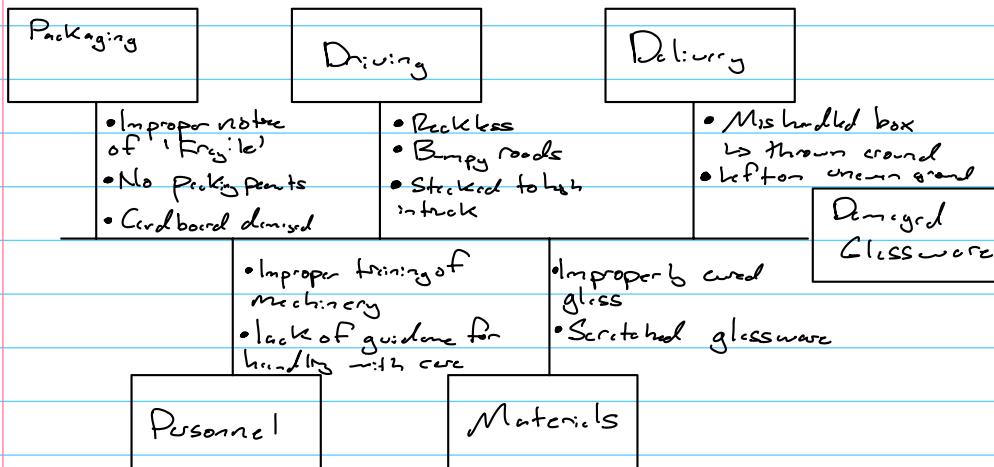
a) Effect of Narrower limits on mag. on type I & II Errors?

If we choose narrower limits, Type I error magnitude (false alarm) increases. Narrower limits \rightarrow more points fall outside acceptable range \rightarrow increased chance of false alarms. Type II error decreases b/c, as narrow limits make it more likely to detect small shifts, reduces B_2 .

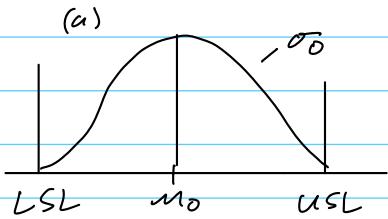
b) Effect of Sigma level on Alpha

A high Sigma level increases the control limit \Rightarrow harder for points to fall outside control limit range \rightarrow reduces α_1 , but increases B_2 . Lower Sigma level results: - tighter limits \rightarrow more points fall outside \rightarrow increasing α_1 but reduces B_2 .

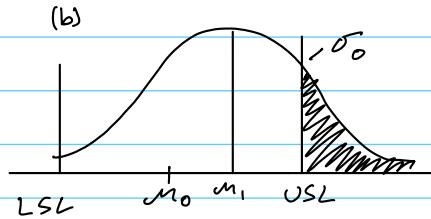
2. Cause & Effect Diagram for Glassware Shipping Error



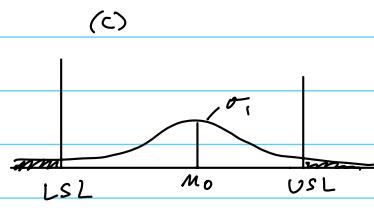
3. Sketches to explain why controlling mean & variability is needed



Normal Operation



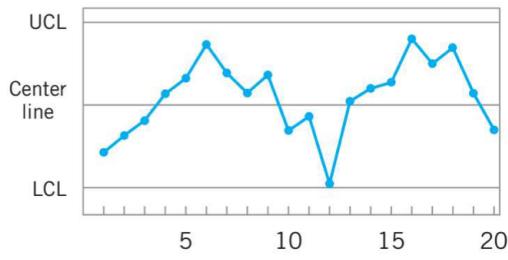
Process mean $M_1 > M_0$



Process std dev $\sigma_1 > \sigma_0$

Controlling mean ensures process lies within specs as in (a), while preventing systematic errors which cause consistent deviation like in (b). Controlling variability reduces out-of-spec variability as in (a), prevents situations like in (c). Controlling both ensures high quality, consistent output, within specs.

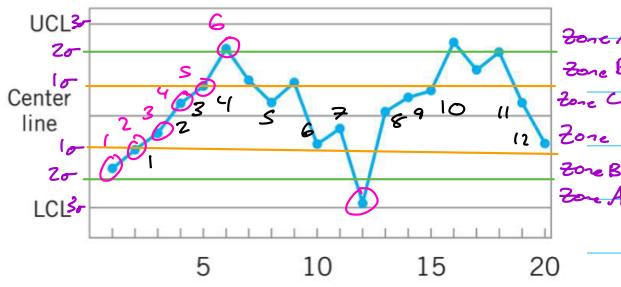
4.



a) Problem S-16

No, it does not appear random. The sharp dip between two roughly random regions suggests assignable cause. The plot could also be viewed as sinusoidal, suggesting a lack of randomness.

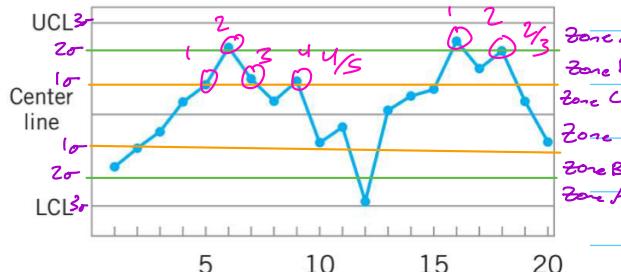
b) Problem S-17 - Sensitizing Rules (S-10)



- 5. Six points in a row steadily increasing or decreasing ✓
- 6. Fifteen points in a row in zone C (both above and below the center line) ✗
- 7. Fourteen points in a row alternating up and down ✗
- 8. Eight points in a row on both sides of the center line with none in zone C ✗
- 9. An unusual or nonrandom pattern in the data ✓
- 10. One or more points near a warning or control limit ✓

Rules 5, 9, and 10 are violated on this control chart, as noted directly on the chart. This indicates potentially out-of-control conditions which should be investigated.

c) Problem S-18 - Western Electric Rules (1-4)



- 1. One or more points outside of the control limits ✗
- 2. Two of three consecutive points outside the two-sigma warning limits but still inside the control limits ✓
- 3. Four of five consecutive points beyond the one-sigma limits ✓
- 4. A run of eight consecutive points on one side of the center line ✗

Rules 2 and 3 are violated on this chart, noted directly on the chart. This indicates a potential out-of-control process or conditions which should be investigated.

S. Problem 6.1

A manufacturer of components for automobile transmissions wants to use control charts to monitor a process producing a shaft. The resulting data from 20 samples of 4 shaft diameters that have been measured are:

$$\sum_{i=1}^{20} \bar{x}_i = 10.275, \quad \sum_{i=1}^{20} R_i = 1.012$$

- a. Find the control limits that should be used on the \bar{x} and R control charts.
- b. Assume that the 20 preliminary samples plot in control on both charts. Estimate the process mean and standard deviation.

a) $\bar{x} \text{ Ch.} +$
 $UCL = \bar{x} + A_2 \bar{R}$
 $LCL = \bar{x} - A_2 \bar{R}$

$R \text{ Ch.} +$
 $UCL = D_4 \bar{R}$
 $LCL = D_3 \bar{R}$

Appendix $n=20$

$$A_2 = 0.180$$

$$D_3 = 0.415$$

$$D_4 = 1.585$$

$$\bar{x} = \frac{\bar{x}}{n} = \frac{10.275}{20} = 0.51375$$

$$\bar{R} = \frac{R}{n} = \frac{1.012}{20} = 0.0506$$

\bar{x} $UCL = 0.51375 + (0.180)(0.0506) = 0.523$
 $LCL = 0.51375 - (0.180)(0.0506) = 0.505$

R $UCL = (1.585) 0.0506 = 0.0802$
 $LCL = (0.415) 0.0506 = 0.0210$

\bar{x} $UCL = 0.523$
 $LCL = 0.505$

R $UCL = 0.0802$
 $LCL = 0.0210$

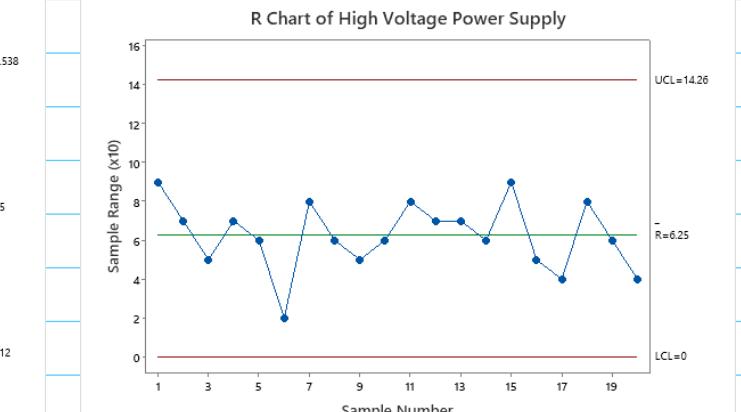
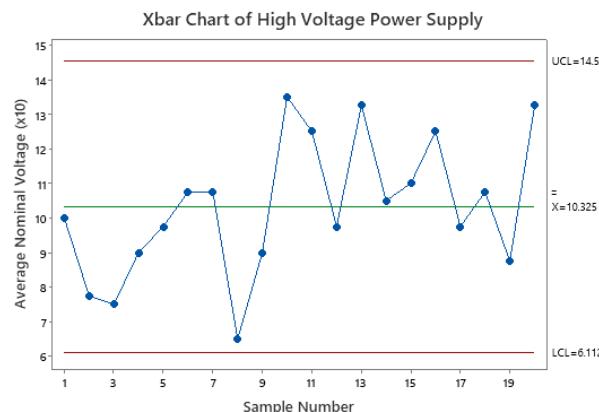
b) $m \approx \bar{x} = 0.51375$

Appendix $n=20$

$$d_2 = 3.735$$

$$\hat{\sigma} \approx \frac{\bar{R}}{d_2} = \frac{0.0506}{3.735} = 0.0135$$

7. All points on the \bar{x} and R charts are within the LCL and UCL shown below.



6.

Subgroup	Xbar	R
1	8.3	
2	8.1	
3	7.9	
4	6.3	
5	8.5	
6	7.5	
7	8	
8	7.4	
9	6.4	
10	7.5	
11	8.8	
12	9.1	
13	5.9	
14	9	
15	6.4	
16	7.3	
17	5.3	
18	7.6	
19	8.1	
20	8	

a) $n = 20$

$$\bar{X} = \frac{\sum_{i=1}^{20} \bar{x}_i}{20} = \frac{8.3 + 8.1 + 7.9 + \dots + 8.1 + 8}{20} = 7.57$$

$$\bar{R} = \frac{\sum_{i=1}^{20} R_i}{20} = \frac{2+3+1+\dots+3+2}{20} = 3.15$$

Appendix: $A_2 = 0.180$

$D_3 = 0.415$

$D_4 = 1.585$

$$UCL = 7.57 + 0.180(3.15) = 8.14$$

$$LCL = 7.57 - 0.180(3.15) = 7$$

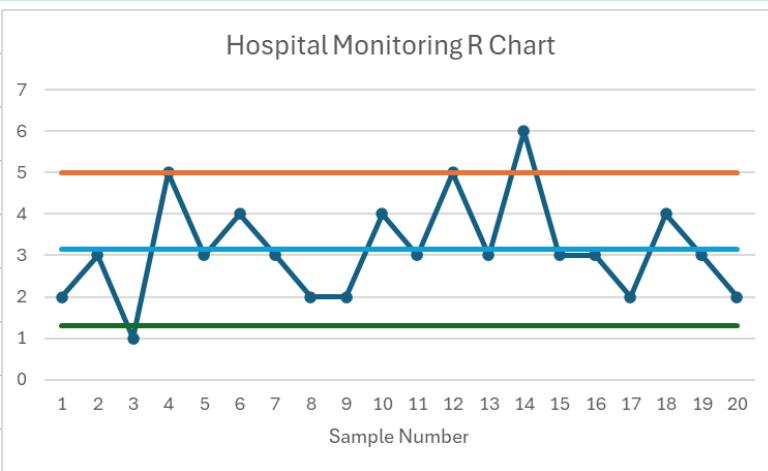
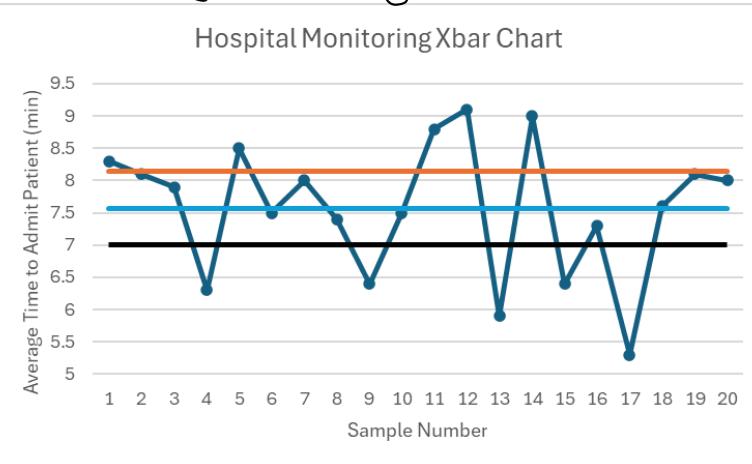
$$UCL = 1.585(3.15) = 4.99$$

$$LCL = 0.415(3.15) = 1.31$$

$\bar{X} \begin{cases} UCL = 8.14 \\ LCL = 7 \end{cases}$

$R \begin{cases} UCL = 4.99 \\ LCL = 1.31 \end{cases}$

b) Plots generated using excel



No, both the Xbar and R chart have points that lie outside the Upper and Lower control limits, indicating that the process is not in Statistical control.