

Statistical Methods

Nels Grevstad

Metropolitan State University of Denver
ngrevsta@msudenver.edu

September 26, 2019

Topics

- 1 Normal Probability Plots
- 2 One-Factor ANOVA for Population Means $\mu_1, \mu_2, \dots, \mu_I$

Objectives

Objectives:

- Use normal probability plots to assess whether a sample is from a normal population.
- Interpret sums of squares, degrees of freedom, and mean squares in a one-factor ANOVA context.
- State the ANOVA partition of the total variation in a data set.
- Carry out a one-factor ANOVA F test for population means $\mu_1, \mu_2, \dots, \mu_I$.

Normal Probability Plots

- Two ways to assess the **normality** of data:
 - A **histogram**. It should be roughly bell-shaped.
 - A **normal probability plot**. The points should hug the line.

Notes

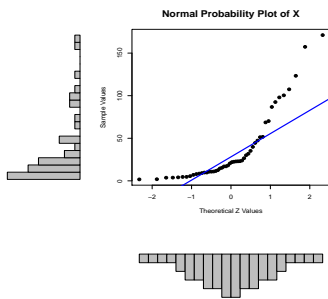
Notes

Notes

Notes

- If a sample is from a $N(\mu, \sigma)$ distribution, then
 - The points in a plot of $(\mu + z_i\sigma, X_{(i)})$, should fall close to the line $y = x$.
 - The points in a plot of $(z_i, X_{(i)})$, should fall close to the line $y = \mu + \sigma x$.
- A **normal probability plot** (or **quantile-quantile plot**) is a plot of the points $(z_i, X_{(i)})$.
- **Curved patterns indicate non-normality.**

Nels Grevstad



Nels Grevstad

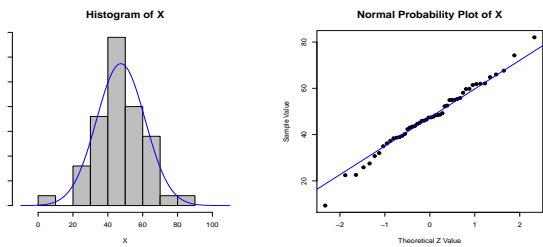


Figure: Histogram of symmetric, approximately normal data (left). Normal probability plot of the same data (right).

Nels Grevstad

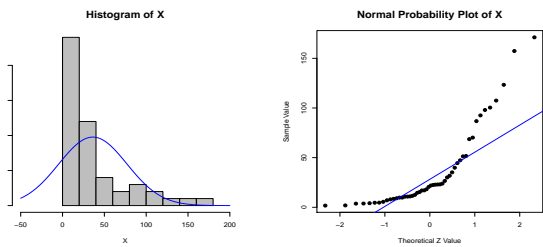


Figure: Histogram of non-normal, right skewed data (left). Normal probability plot of the same data (right).

Nels Grevstad

Notes

Notes

Notes

Notes

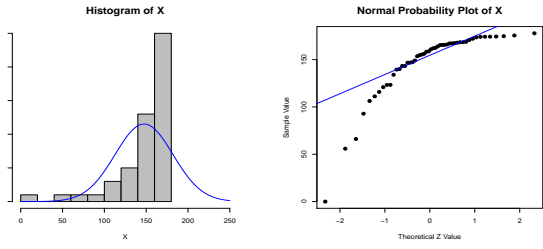


Figure: Histogram of non-normal, left skewed data (left). Normal probability plot of the same data (right).

Notes

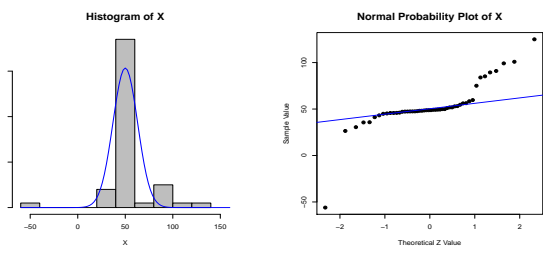


Figure: Histogram of non-normal, "heavy tailed" data (left). Normal probability plot of the same data (right).

Notes

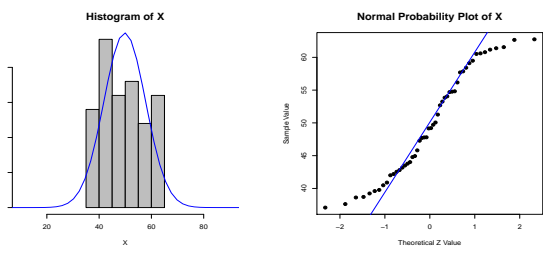


Figure: Histogram of non-normal, "light tailed" data (left). Normal probability plot of the same data (right).

Notes

One-Factor ANOVA for Population Means

$\mu_1, \mu_2, \dots, \mu_I$

Introduction

- Suppose we have independent random samples from I populations having **possibly different means** but **equal standard deviations**.

The populations might represent different **groups** or they might represent **treatments** in an experiment.

We want to decide if there are any differences among the population means.

Notes

Example

A quality assurance study was carried out to compare **lead measurements** made in water sent to $I = 5$ laboratories.

Differences among the five labs' results may signify improperly calibrated equipment or poorly trained technicians.

A vat of wastewater was split into **50** specimens randomized to the labs ($J = 10$ each) for analysis.

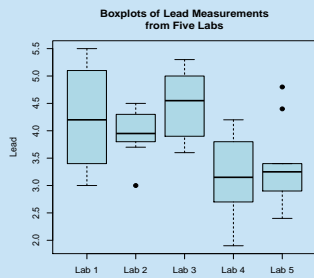
The **lead measurements** ($\mu\text{g/L}$) and their summary statistics are on the next slide.

Nels Grevstad

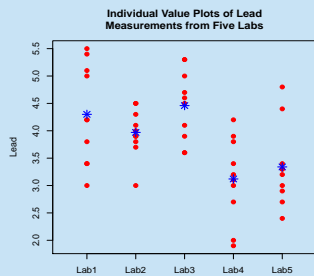
Measured Lead Concentrations

Lab 1	Lab 2	Lab 3	Lab 4	Lab 5
3.4	4.5	5.3	3.2	3.3
3.0	3.7	4.7	3.4	2.4
3.4	3.8	3.6	3.1	2.7
5.0	3.9	5.0	3.0	3.2
5.1	4.3	3.6	3.9	3.3
5.5	3.9	4.5	2.0	2.9
5.4	4.1	4.6	1.9	4.4
4.2	4.0	5.3	2.7	3.4
3.8	3.0	3.9	3.8	4.8
4.2	4.5	4.1	4.2	3.0
$\bar{X}_1 = 4.30$ $S_1 = 0.904$	$\bar{X}_2 = 3.97$ $S_2 = 0.440$	$\bar{X}_3 = 4.46$ $S_3 = 0.642$	$\bar{X}_4 = 3.12$ $S_4 = 0.764$	$\bar{X}_5 = 3.34$ $S_5 = 0.737$

Nels Grevstad



Nels Grevstad



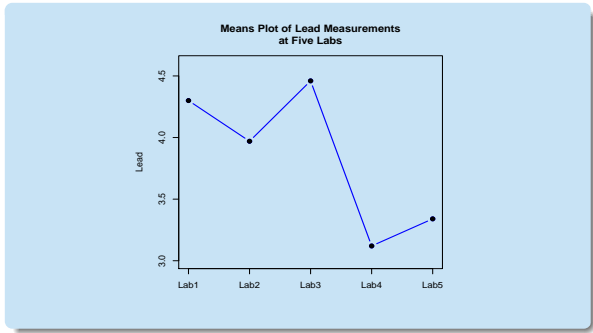
Nels Grevstad

Notes

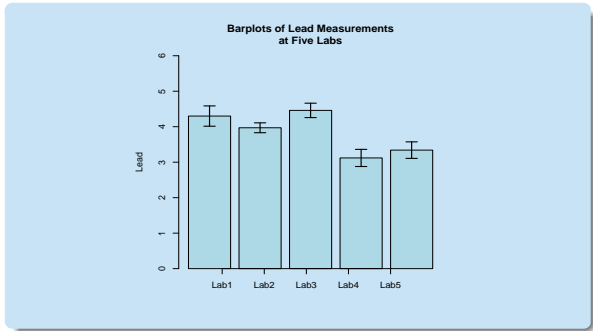
Notes

Notes

Notes



Nels Grevstad



Nels Grevstad

- Suppose we have random samples, each of size J , from I populations ($I \geq 2$),
 - We'll see how to use the samples to decide if there are differences among the **population means** μ_1, μ_2, \dots , and μ_I .
- The appropriate test is called the **one-factor ANOVA F test**.

Nels Grevstad

- **Comments:**
 - The sample sizes **don't** all have to be the same. But we'll only look at the equal-sample size case.
 - The data can be **samples** from populations or responses to treatments in a **randomized experiment**.

Nels Grevstad

Notes

Notes

Notes

Notes

- The **null hypothesis** is that there are no differences among the population means $\mu_1, \mu_2, \dots, \mu_I$:

Null Hypothesis:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_I$$

- The **alternative hypothesis** is that there's *at least one difference* among the set of means:

Alternative Hypothesis: The alternative hypothesis will be

$$H_a : \text{At least two of the } \mu_i \text{'s are different}$$

- Notation:**

- I = The number of treatment groups
- J = The common sample size for the I groups
- X_{ij} = The j th observation in the i th treatment group
- \bar{X}_i = The sample mean for the i th treatment group
- S_i = The sample standard deviation for the i th treatment group
- $\bar{X}_{..}$ = The **grand mean** of all IJ observations

Note:

$$\bar{X}_{..} = \frac{1}{I} \sum_{i=1}^I \bar{X}_i$$

(when the sample sizes are all the same).

Sums of Squares and the ANOVA Partition

- We can **partition** the **total variation** in the data into two parts:
 - One reflecting variation **between** the treatment groups.
 - The other reflecting variation **within** the groups.

The **ANOVA F test** is based on the amount of **between**-groups variation relative to the amount of **within**-groups variation.

Notes

Notes

Notes

Notes

- The **partition** will involve the following **sums of squares**:

- **SST** is the **total sum of squares**, defined as

$$SST = \sum_{i=1}^I \sum_{j=1}^J (X_{ij} - \bar{X}_{..})^2,$$

which measures the **total** variation in the X_{ij} 's.

- (cont'd):

- **SSTr** is the **treatment sum of squares**, defined as

$$SSTr = \sum_{i=1}^I \sum_{j=1}^J (\bar{X}_i - \bar{X}_{..})^2 = J \sum_{i=1}^I (\bar{X}_i - \bar{X}_{..})^2,$$

which measures variation **between** the treatment group means due to both **treatment effects** and **random error**.

- (cont'd):

- **SSE** is the **error sum of squares**, defined as

$$SSE = \sum_{i=1}^I \sum_{j=1}^J (X_{ij} - \bar{X}_i)^2,$$

which measures variation of the X_{ij} 's **within** treatment groups due to **random error**.

Proposition

ANOVA Partition of the Total Variation: It can be shown that

$$SST = SSTr + SSE.$$

Notes

Notes

Notes

Notes

- The **ANOVA partition** holds because we can write:

$$X_{ij} - \bar{X}_{..} = \bar{X}_i - \bar{X}_{..} + X_{ij} - \bar{X}_i$$

Upon squaring both sides and then summing over all i and j , the "cross product" terms on the right side sum to zero, and we get

$$\sum_i \sum_j (X_{ij} - \bar{X}_{..})^2 = \sum_i \sum_j (\bar{X}_i - \bar{X}_{..})^2 + \sum_i \sum_j (X_{ij} - \bar{X}_i)^2,$$

which is the **ANOVA partition**.

Nels Grevstad

Normal Probability Plots

One-Factor ANOVA for Population Means $\mu_1, \mu_2, \dots, \mu_I$

Example

For the data on lead measurements at five labs, software gives

$$\text{SST} = 36.758$$

$$\text{SSTr} = 13.813$$

$$\text{SSE} = 22.945$$

The **ANOVA partition** holds:

$$\begin{array}{rcc} 36.758 & = & 13.813 + 22.945 \\ \uparrow & & \uparrow \quad \uparrow \\ \text{Total} & \text{Between} & \text{Within} \\ \text{variation} & \text{groups} & \text{groups} \\ & \text{variation} & \text{variation} \end{array}$$

Nels Grevstad

Normal Probability Plots

One-Factor ANOVA for Population Means $\mu_1, \mu_2, \dots, \mu_I$

Degrees of Freedom

- Each sum of squares has an associated **degrees of freedom** (or **df**).

The **df** for a sum of squares is determined by how many deviations, among those used to compute the sum of squares, are "free to vary" (**unconstrained**).

Nels Grevstad

Normal Probability Plots

One-Factor ANOVA for Population Means $\mu_1, \mu_2, \dots, \mu_I$

Degrees of Freedom:

$$\text{SST has } IJ - 1 \text{ df}$$

$$\text{SSTr has } I - 1 \text{ df}$$

$$\text{SSE has } I(J - 1) = IJ - I \text{ df}$$

Nels Grevstad

Notes

Notes

Notes

Notes

Mean Squares

- The **ANOVA F test** is based on the amount of **between**-groups variation relative to the amount of **within**-groups variation.
But **SSTr** and **SSE** *aren't* directly comparable (they depend in different ways on I and J).
- A **mean square** a **sum of squares** divided by its **df**.
Example: A *sample variance* S^2 is a **mean square**.

Nels Grevstad

- (cont'd)
 - The **mean square for treatments**, denoted **MSTr**, is

$$MSTr = \frac{SSTr}{I - 1}.$$

Nels Grevstad

- (cont'd)
 - The **mean squared error**, denoted **MSE**, is

$$MSE = \frac{SSE}{I(J - 1)}.$$

It's easy to verify that

$$MSE = \frac{S_1^2 + S_2^2 + \dots + S_I^2}{I}$$

(when the sample sizes are all the same).

Thus **MSE** is the **average** (or **pooled**) **sample variance**.

Nels Grevstad

- MSTr** and **MSE** *are* directly comparable.

Nels Grevstad

Notes

Notes

Notes

Notes

Notes

The One-Factor ANOVA F Test

One-Factor ANOVA F Test Statistic:

$$F = \frac{MSTr}{MSE}$$

- F reflects **between**-groups variation (**MSTr**) relative to **within**-groups variation (**SSE**).
- **MSTr** will be **large** when there's substantial variation in $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_I$, which are estimates of the population means $\mu_1, \mu_2, \dots, \mu_I$.
It will be **large** when there are **differences** among $\mu_1, \mu_2, \dots, \mu_I$.

Notes

Large values of F provide evidence against H_0 in favor of H_a : At least two of the μ_i 's are different.

Notes

- Now suppose the I samples are from $N(\mu_1, \sigma), N(\mu_2, \sigma), \dots, N(\mu_I, \sigma)$ distributions and that they were drawn *independently* of each other.
Alternatively, the samples could be from **non-normal** populations as long as the common sample size J is **large**.

Notes

The ANOVA Table

- ANOVA results are summarized in an **ANOVA table**:

Source of Variation	df	Sum of Squares	Mean Square	f	P-value
Treatment	$I - 1$	SSTr	$MSTr = SSTr / (I - 1)$	$MSTr / MSE$	p
Error	$I(J - 1)$	SSE	$MSE = SSE / (I(J - 1))$		
Total	$IJ - 1$	SST			

Exercise

For lead measurements made at five labs, the **ANOVA table** is:

Source of Variation	df	Sum of Squares	Mean Square	f	P-value
Treatment	4	13.813	3.453	6.77	0.000
Error	45	22.945	0.510		
Total	49	36.758			

- Verify that **df for SSTr = $I - 1$** , that **df for SSE = $I(J - 1)$** , and that **df for SST = $IJ - 1$** .
- Verify that **SST = SSTr + SSE** and that the **df for SST = df for SSTr + df for SSE**.
- Verify that the **mean squares** are the **sums of squares** divided by their **df**.
- Verify that the **F statistic** is **MSTr** divided by **MSE**.

- State the **hypotheses**.
- Using $\alpha = 0.05$, is there statistically significant evidence for systematic differences in lead measurements among the five labs?
- If there are significant differences among the five labs, describe the nature of those differences (using the plots of the data given earlier in these slides).

- For comparing **two population means** μ_1 and μ_2 , the **ANOVA F test** and a **two-sided pooled two-sample t test** are **equivalent**.

The **square** of the t **statistic** is the F **statistic**, and the **p-values** will be the **same**.

Example

An example in a previous set of slides presented results of a computer simulation to compare the time (in seconds) to complete a semiconductor manufacturing process using one and two operators.

Here are the summary statistics:

One Operator	Two Operators
$m = 16$	$n = 16$
$\bar{X} = 373.6$	$\bar{Y} = 374.8$
$S_1 = 7.8$	$S_2 = 7.3$

If we carry out a **(pooled) two-sample t test** of

$$H_0: \mu_1 = \mu_2$$

$$H_a: \mu_1 \neq \mu_2$$

we get:

Pooled t	
Test Statistic	P-Value
$t = -0.445$	0.6596

If we carry out a **one-factor ANOVA**, we get:

Source of Variation	df	Sum of Squares	Mean Square	f	P-value
Treatment	1	11.3	11.3	0.198	0.6596
Error	30	1710.2	57.0		
Total	31	1721.5			

We see that $t^2 = F$ and the **p-values** for the two tests are the same.

Notes

Notes

Notes

Notes

- In general, the **square** of a t random variable is an F random variable.

Proposition

If

$$T \sim t(\nu)$$

then

$$T^2 \sim F(1, \nu).$$

Notes

Notes

Notes

Notes
