# Lecture 19

# Introduction to ANOVA

**STAT 512**

**Spring 2011**

**Background Reading**

**KNNL:  15.1-15.3, 16.1-16.2**

# Topic Overview

- Categorical Variables
- Analysis of Variance
- Lots of Terminology
- An ANOVA example

# Categorical Variables

- To this point, with the exception of the last lecture, all explanatory variables have been quantitative; e.g. comparing $X = 3$ to $X = 5$ makes sense numerically

- For *categorical* or *qualitative* variables there is no 'numerical' labeling; or if there is, it isn't meaningful.
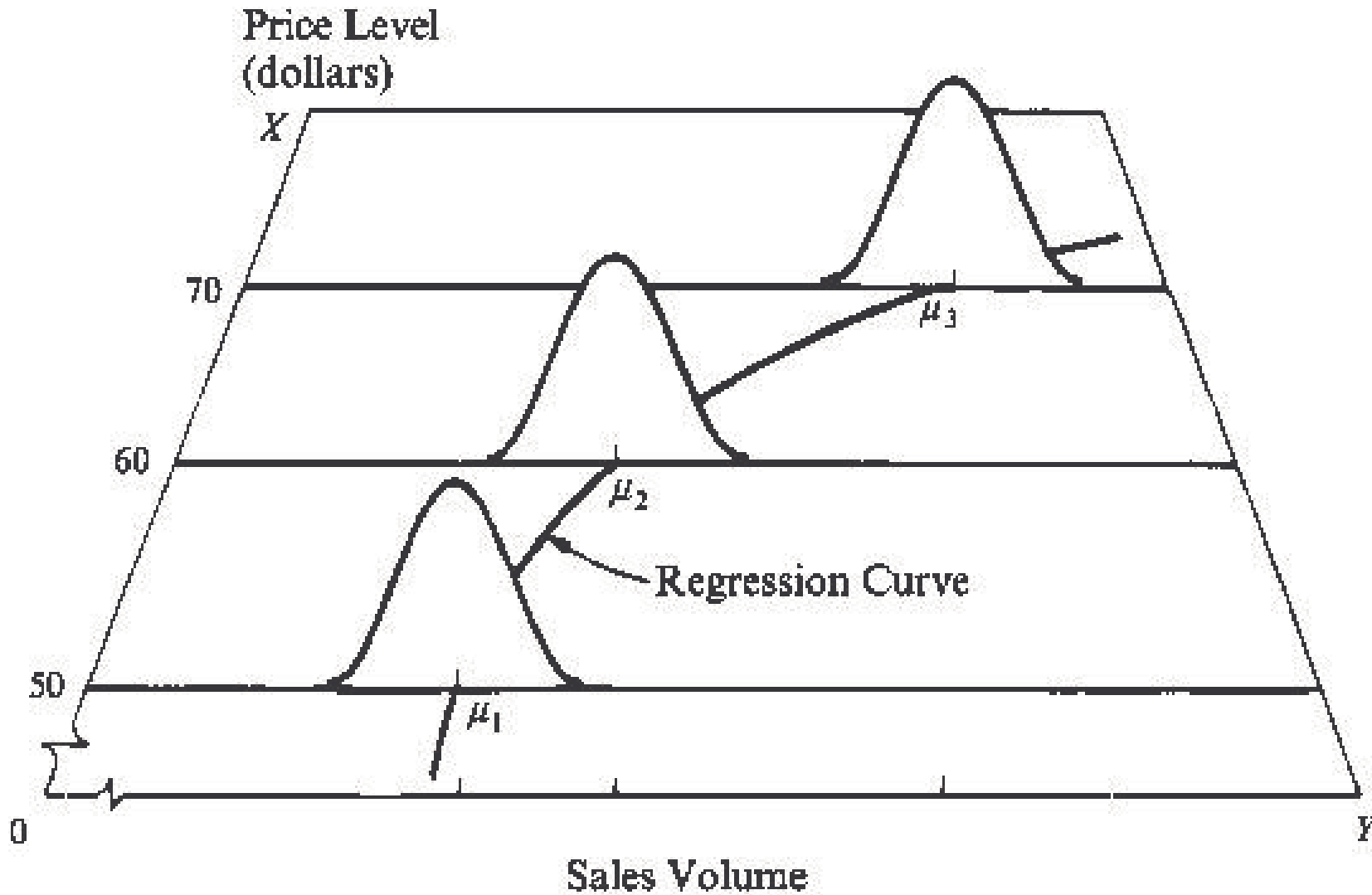
# Example

- Five medical treatments – ten subjects on each treatment.

- Goal: Compare the treatments in terms of their effectiveness
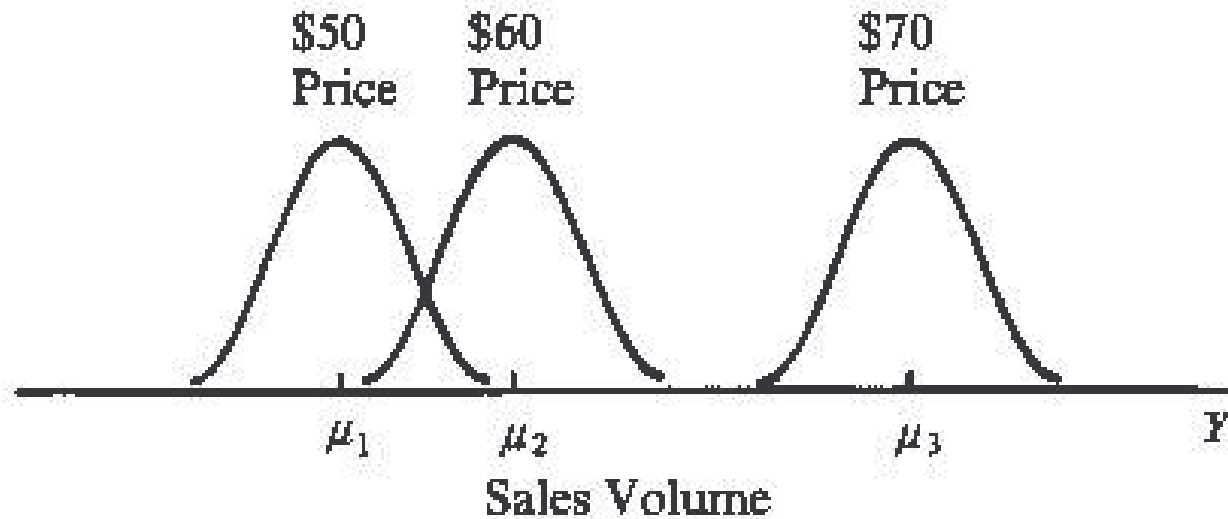  - If there were two treatments, what would we use?

# ANOVA

- *ANOVA = Analysis of Variance*

- Compare means among treatment groups, without assuming any parametric relationships (regression does assume such a relationship).

- Example: Price vs. Sales Volume

# Regression Model

# ANOVA Model



KEY DIFFERENCE: No assumption is made about the manner in which Price and Sales Volume are related.

# Similarities to Regression

- Assumptions on errors identical as to regression

- We assume each population is normal and the variances are identical. We also assume independence.

- Can get "predicted values" for each group, as well as CI's.

# Differences

- No specific relationship is assumed.

- Goal becomes: look for differences among the groups.

# Terminology

- We may refer to any qualitative predictor variable as a *factor*.

- Each factor has a certain number of *levels*.

- *Experimental factors* are "set" or "assigned" to the experimental units; *observational factors* are characteristics of the experimental units that cannot be assigned.

# Terminology (2)

- Factors are *qualitative* if they represent traits that could not be placed in some logical numerical order.
  - GENDER, BRAND, DRUG
- Factors are *quantitative* if levels are described by numerical quantities on an equal interval scale.
  - AGE, TEMPERATURE

# Terminology (3)

- A ***Treatment*** is a specific experimental condition (determined by factors and levels of each factor).

- The ***Experimental Unit*** (Basic Unit of Study) is the smallest unit to which a treatment can be assigned.

- A design is called ***balanced*** if each treatment is replicated the same number of times (i.e. same number of EU's per treatment).

# Examples

Five medications – each used for 10 subjects

- Medication is an experimental factor; EU is the subject (person) receiving the medication.

- There are five treatments, which may or may not have any logical "ordering"

- Design is balanced (generally) since we are able to assign the treatments.

Ten age groups – 50 subjects

- Age is an observational, quantitative factor; subject is again the EU; Design is probably not balanced

# Examples (2)

Blood Type

- Observational factor
- Qualitative factor
- Again design probably not balanced

Brand of Product

- Observational, qualitative factor
- Design likely balanced by arrangement

# Multiple Factors

- With two or more factors, each combination of levels is generally called a ***treatment combination***

- Can treat as single variable if desired

- Example:  Blood Type * Medication
    - 4 blood types
    - 5 medications
    - 20 treatment combinations

# **Crossed Factors**

- Two factors are *crossed* if all factor combinations are represented.

- Example:  Blood Type * Medication

|     | 1  | 2  | 3  | 4  | 5  |
|-----|----|----|----|----|----|
| A   | xx | xx | xx | xx | xx |
| B   | xx | xx | xx | xx | xx |
| AB  | xx | xx | xx | xx | xx |
| O   | xx | xx | xx | xx | xx |

Note:  This type of table is called a *design chart*.

# Nested Factors

- One factor has levels that are unique to a given level of another factor

- Example:  Plant * Operator

| Plant #1 | Plant #2 | Plant #3 |
|----------|----------|----------|
| Op #1    | Op #4    | Op #7    |
| Op #2    | Op #5    | Op #8    |
| Op #3    | Op #6    | Op #9    |

- We say:  Operators are nested within manufacturing plants.

# Control Groups

- Often a ***control*** or ***placebo*** treatment is used. This treatment is more of a "standard" than a treatment, as it is the case of no treatment at all.

- Comparing treatments to controls can be a very effective way of showing that a treatment is effective.

# Fixed vs. Random Factors

- For the most part, we will consider only *fixed effect models* in this class. A factor is called *fixed* because the levels are chosen in advance of the experiment and we were interested in differences in response among <u>those specific levels</u>.

- Note: *Random* factors will need to be treated differently, since their levels are chosen randomly from a large population of possible levels.

# Randomization

- Completely separate concept from random effects.

- In an experimental study, generally want to <u>avoid any potential bias</u> in the design by ***randomizing treatments to experimental units*** whenever possible.

- Randomization may be ***constrained***. Example: Have 100 people, 50 men and 50 women. Randomly assign each of the 5 treatments to 10 men and 10 women.

# Experimental Designs

- Completely Randomized Design

- Factorial Experiments

- Randomized Complete Block Designs

- Nested Designs

- Repeated Measures Designs

- Incomplete Block Designs

- We'll discuss some of these. More thorough experimental design course: STAT 514.
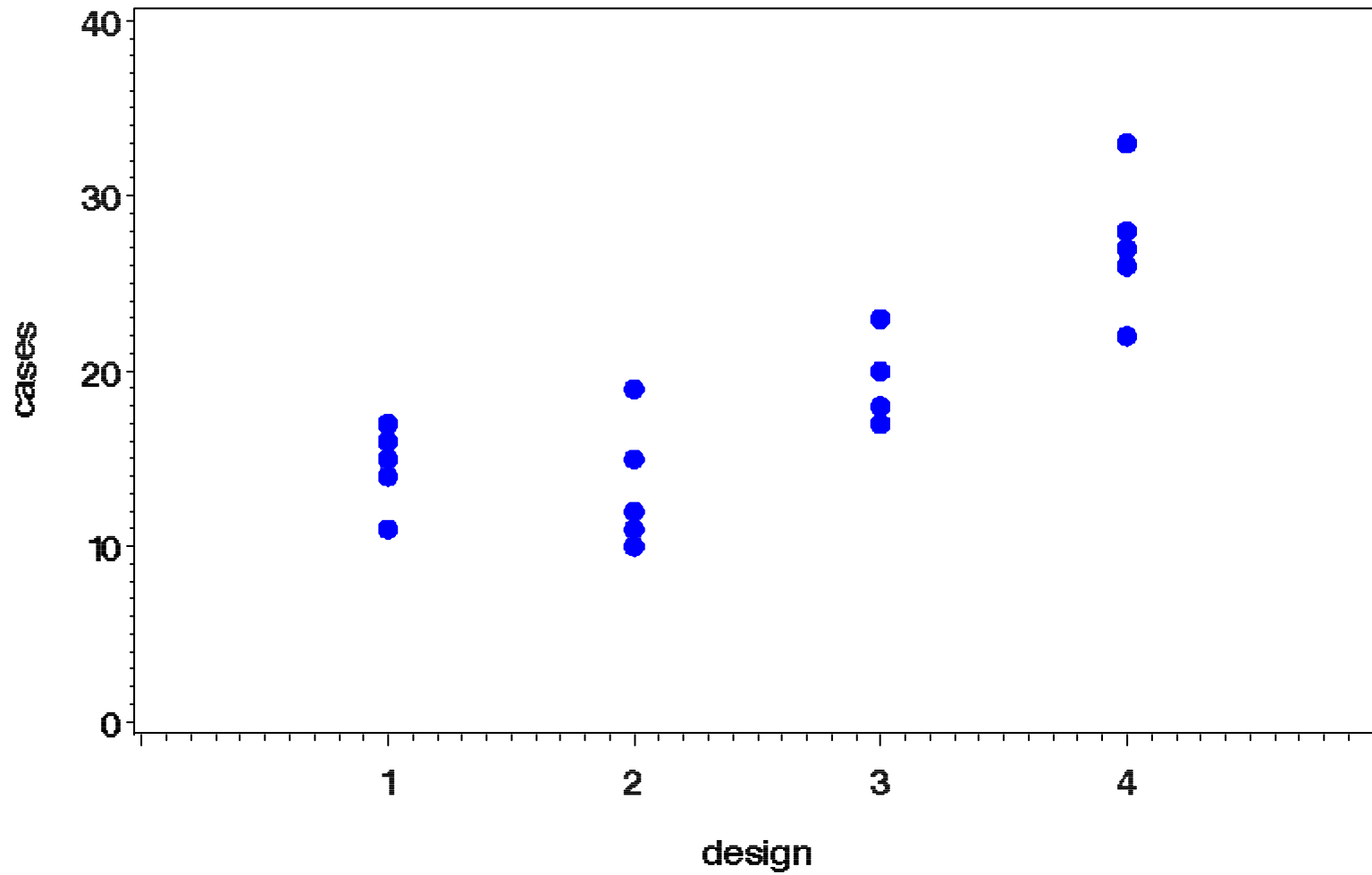
# Example

- Kenton Food Company Example (p685)
- Compare four different package designs (numbered 1, 2, 3, 4 in no particular order)
- Response:  # of cases sold
- 20 stores, but one was destroyed by fire during the study; 19 observations
- SAS file:  kenton.sas

# Data

| Design 1 | Design 2 | Design 3 | Design 4 |
|----------|----------|----------|----------|
| 11 | 12 | 23 | 27 |
| 17 | 10 | 20 | 33 |
| 16 | 15 | 18 | 22 |
| 14 | 19 | 17 | 26 |
| 15 | 11 |    | 28 |

# Scatter Plot

# ANOVA Code (SAS)

```
proc glm data=kenton;
   class design;
   model cases=design;
   means design /bon lines cldiff;
```

- Class statement identifies ALL categorical variables (separate by spaces as in model)
- Means statement requests comparisons of the group means (lots of options)

# Output

| Source | DF | SS | MS | F Value | Pr > F |
|--------|-----|-----|------|---------|--------|
| Model | 3 | 588 | 196 | 18.59 | <.0001 |
| Error | 15 | 158 | 10.5 | | |
| Total | 18 | 746 | | | |

| R-Square | Coeff Var | Root MSE | cases Mean |
|----------|-----------|----------|------------|
| 0.788055 | 17.43042 | 3.247563 | 18.63158 |

# Output (2)

Bonferroni (Dunn) t Tests for cases
NOTE: This test controls the Type I
experimentwise error rate, but it generally
has a higher Type II error rate than Tukey's
for all pairwise comparisons.

| | |
|---|---|
| Alpha | 0.05 |
| Error Degrees of Freedom | 15 |
| Error Mean Square | 10.54667 |
| Critical Value of t | 3.03628 |

Comparisons significant at the 0.05 level
are indicated by ***.

# Output (3)

```
design          Difference   Simultaneous 95%
Comparison        Means          Confidence Limits
4 - 3       7.700        1.085    14.315    ***
4 - 1      12.600        6.364    18.836    ***
4 - 2      13.800        7.564    20.036    ***
3 - 4      -7.700      -14.315    -1.085    ***
3 - 1       4.900       -1.715    11.515
3 - 2       6.100       -0.515    12.715
1 - 4     -12.600      -18.836    -6.364    ***
1 - 3      -4.900      -11.515     1.715
1 - 2       1.200       -5.036     7.436
2 - 4     -13.800      -20.036    -7.564    ***
2 - 3      -6.100      -12.715     0.515
2 - 1      -1.200       -7.436     5.036
```
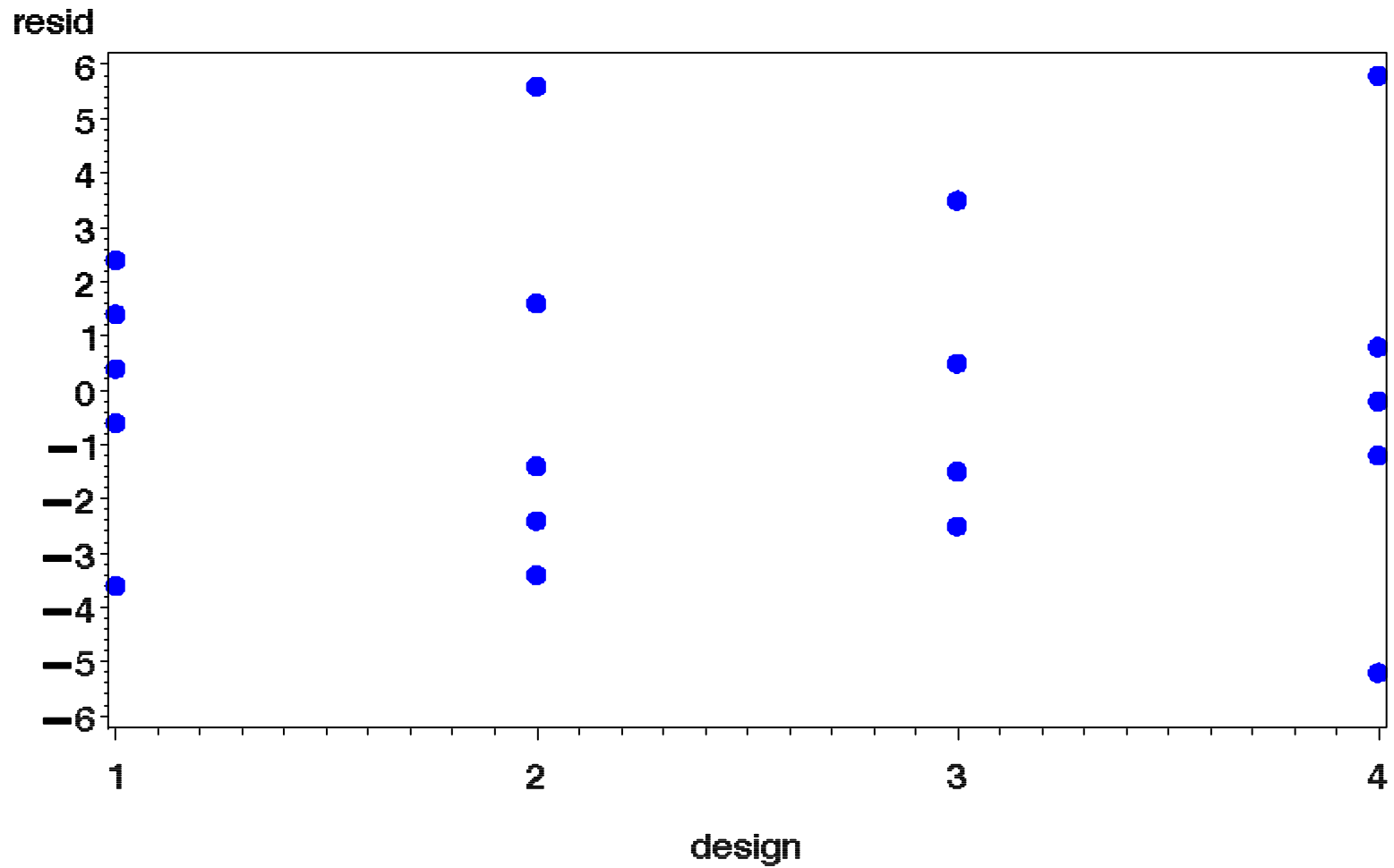
# Output (4)

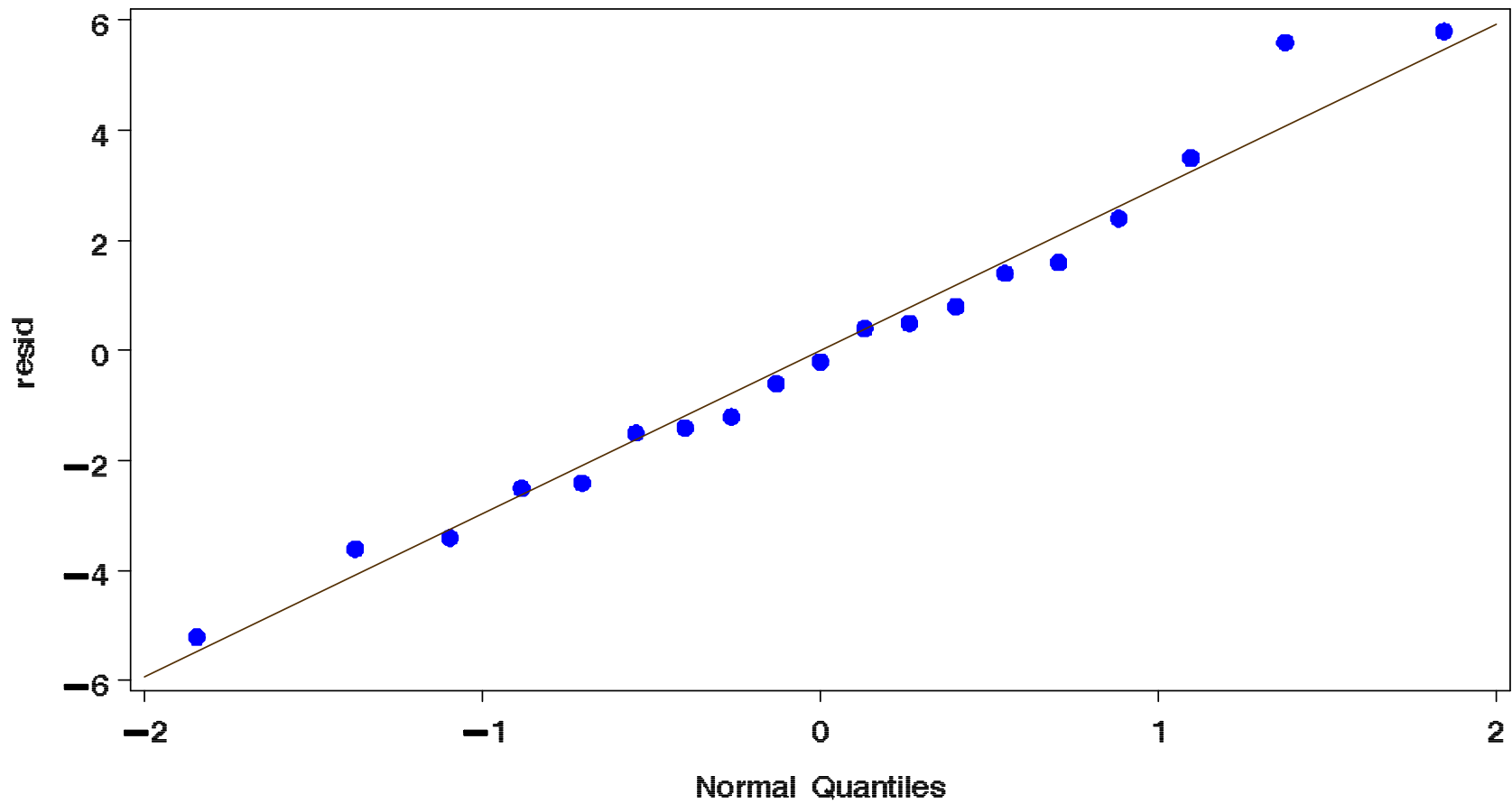| Group | Mean | N | design |
|-------|--------|---|--------|
| A | 27.200 | 5 | 4 |
| B | 19.500 | 4 | 3 |
| B | | | |
| B | 14.600 | 5 | 1 |
| B | | | |
| B | 13.400 | 5 | 2 |

# Assumptions

- Should always check normality, constancy of variance assumptions

- Plots to check these are as before

- No obvious problems for this dataset

# Residual Plot

# Normal QQ Plot

# Upcoming in Lecture 20...

- ANOVA Model I (Cell Means)
- Sections 16.3 – 16.6