

# One Factor ANOVA

# What is Experimental Design?

- A designed experiment is a test in which purposeful changes are made to the input variables ( $x$ ) so that we may observe and identify the reasons for change in the output response ( $y$ ).
- Objectives:
  - Which variables are most influential on the response  $y$ ?
  - Where to set the influential  $x$ 's so  $y$  is near a desired value?
  - Where to set the influential  $x$ 's so variability in  $y$  is small?

# Principals of design

- Replication: repetition of the basic experiment
  - Allows the estimation experimental error
  - Obtain a more precise estimate of the effect of a factor
- Randomization: both the allocation of the experimental material and the order in which the individual runs are performed are randomly determined
- Blocking: a block is a homogeneous portion of the experimental material
  - Used to increase the precision of an experiment

# Business Example

- A manufacturer is interested in **maximizing** the tensile strength of a new synthetic fiber that will be used to make cloth for men's shirt.
- From previous tests, the manufacturer knows:
  - the strength is affected by the **percentage of cotton** in the fiber
  - The **range** of the percentage is 10% to 40%
- Experiment with a Single Factor

# Example – continued

- Maximizing the tensile strength of a new synthetic fiber used to make shirts\*
- Strength is affected by the % of cotton in the fiber
  - Test five levels of cotton percentage
    - 15%, 20%, 25%, 30%, 35%
  - Test five specimens at each cotton %
  - So,  $a = 5$  levels of the factor, and  $n = 5$  replicates
  - All 25 runs are made in random order

\*DCM pp 50

# Example – continued

## Experimental Runs

- This randomized test sequence, known as a Completely Randomized Design (CRD), is necessary to prevent the effects of unknown nuisance variables from contaminating the results

Test Sequences	Run Numbers	Percentage of Cotton
1	8	20
2	18	30
3	10	20
...		
...		
24	19	30
25	3	15

# Data

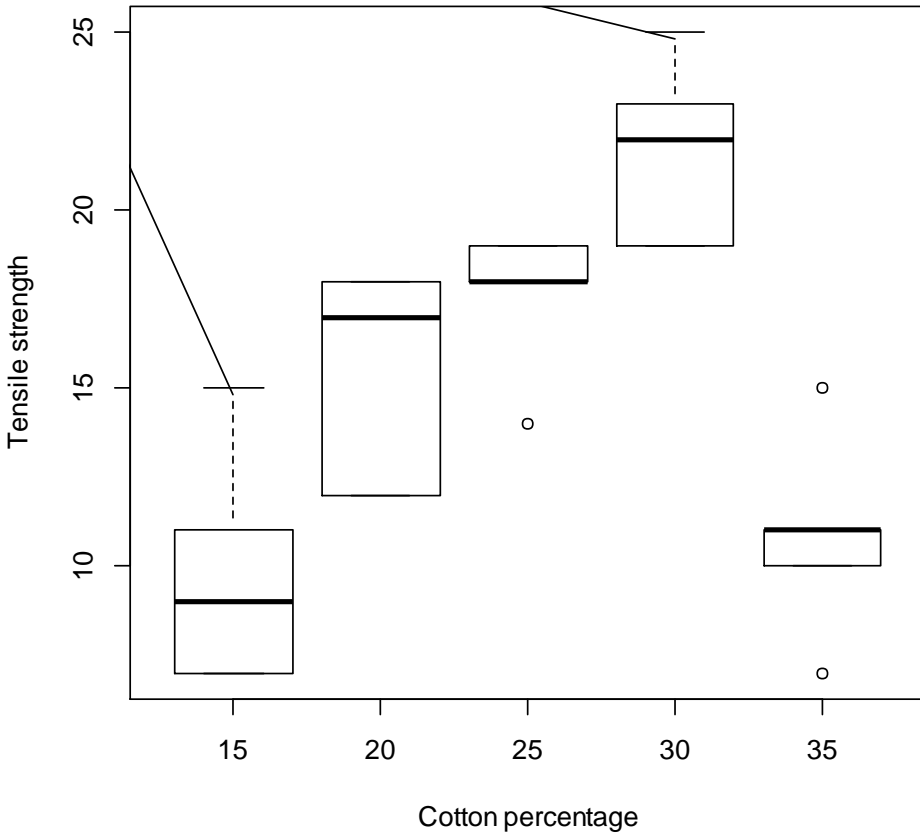
for the Cotton Percentage Example

---

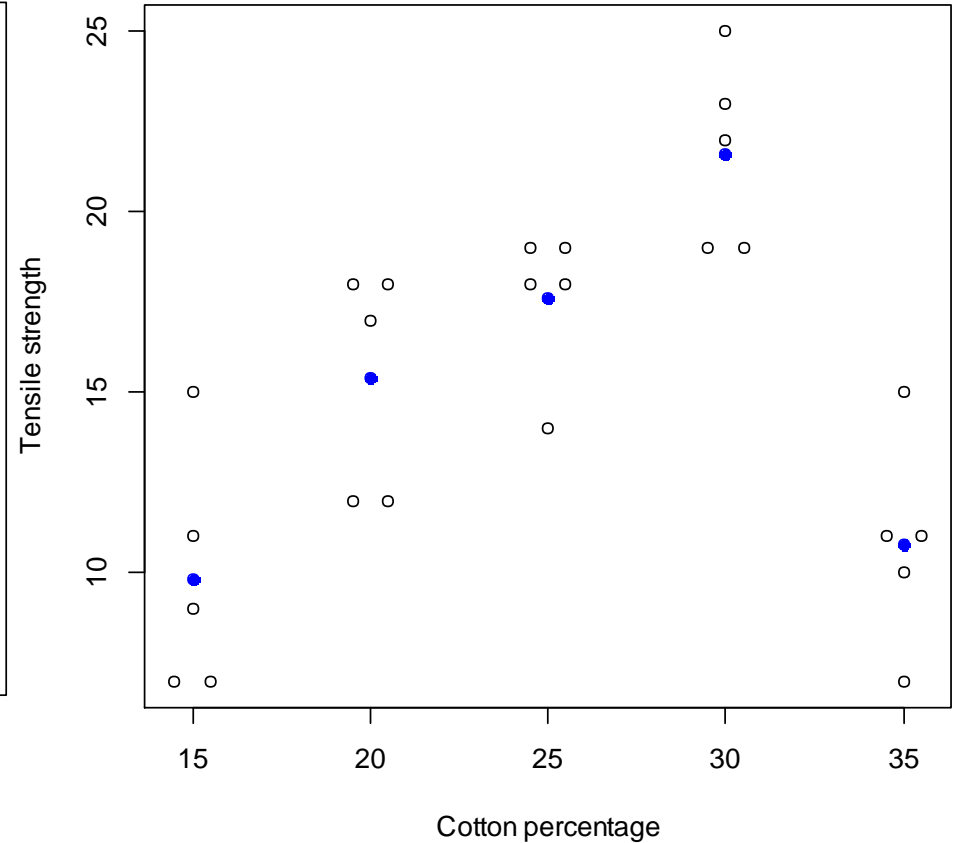
Cotton Percentage	Observations					Total	Average
	1	2	3	4	5		
15	7	7	15	11	9	49	9.8
20	12	17	12	18	18	77	15.4
25	14	18	18	19	19	88	17.6
30	19	25	22	19	23	108	21.6
35	7	10	11	15	11	<u>54</u>	<u>10.8</u>
						376	15.04

---

# Tensile Strength Vs Cotton Percentage



Box plots



Scatter diagram



# Exploratory Results

- We strongly suspect that:
  - Cotton content affects tensile strength
  - Around 30% cotton would result in maximum strength
- A more objective analysis:
  - Should we perform a *t*-test on all possible pairs of means?
  - No, results in a substantial increase in the type I error
- The appropriate procedure for testing the equality of several means is the analysis of variance (**ANOVA**)

# ANOVA

- The name ANOVA is derived from a partitioning of total variability into its component parts: ANalysisOfVAriance
- The total sum of squares is a measure of overall variability in the data:

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

- The total variability in the data can be partitioned into a sum of squares of the differences *between* the treatment averages and the grand average, PLUS a sum of squares of the differences of observations *within* treatments from the treatment average:

$$- SS_{\text{total}} = SS_{\text{treatment}} + SS_{\text{within}}$$

# ANOVA

## Degrees of freedom

- There are  $N = an$  total observations, so  $SS_{\text{total}}$  has  $N - 1$  degrees of freedom.
- There are  $a$  levels of the factor, so  $SS_{\text{treatment}}$  has  $a - 1$  degrees of freedom.
- Thus, we have  $N - a$  degrees of freedom for  $SS_{\text{within}}$

# ANOVA

- The analysis of variance identity provides us with two estimates of  $\sigma^2$  – one based on the inherent variability within treatments and one based on the variability between treatments
  - $SS_{\text{total}} = SS_{\text{treatment}} + SS_{\text{within}}$
- If there are no differences in the treatment means, then these two estimates should be very similar.
- If they are not, we suspect that the observed differences must be caused by differences in treatment means.

# ANOVA

## for the Cotton Percentage Data

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	$F_0$
Cotton Percentage	$SS_{\text{treatment}}$ → 475.76	4	118.94	<b>Test Statistics</b> $F_0=14.76$
Error	$SS_{\text{within}}$ → 161.20	20	8.06	
Total	636.96	24		

# ANOVA

## Additional Concepts

- A factor is *fixed* if the levels of a factor are predetermined and the experimenter is interested only in those particular levels (e.g. 250, 300, or 350°F).
- A factor is classified as *random* if the levels are selected at random from a population of levels (e.g. 254, 287, and 326°F).
- When there is only one factor, the classification does not have any effect on how the data are analyzed.

# Fixed Effects Model

$$y_{ij} = \mu + \tau_i + \epsilon_{ij}$$

- $y_{ij}$ : the  $ij^{\text{th}}$  observation
- $\mu$ : grand mean, a parameter common to all treatments
- $\tau_i$ : treatment effect, a parameter unique to the  $i^{\text{th}}$  treatment,
- $\epsilon_{ij}$ : random error

# Least Squares Estimate

The Sum of squares of error :

$$L = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \mu - \tau_i)^2$$

Choose  $\hat{\mu}$  and  $\hat{\tau}$  that minimize L:

$$\left. \frac{\partial L}{\partial \mu} \right|_{\hat{\mu}, \hat{\tau}_i} = 0$$

$$\left. \frac{\partial L}{\partial \tau} \right|_{\hat{\mu}, \hat{\tau}_i} = 0$$



# Least Squares Estimate

Adding a constraint:

$$\sum_{i=1}^a \hat{\tau}_i = 0$$

The least squares estimate for  $\mu, \tau_i$  :

$$\hat{\mu} = \bar{y}_{..}$$
$$\hat{\tau}_i = \bar{y}_{i.} - \bar{y}_{..}, i = 1, 2, \dots, a$$

# Cotton Percentage example

-continued-

Estimate of the overall mean:

$$\hat{\mu} = \hat{y}_{..} = \frac{376}{25} = 15.04$$

# Cotton Percentage example

-continued-

Estimate of treatment effects:

- $\hat{\tau}_1 = \overline{y_{1.}} - \overline{y_{..}} = 9.80 - 15.04 = -5.24$
  - $\hat{\tau}_2 = \overline{y_{2.}} - \overline{y_{..}} = 15.40 - 15.04 = 0.36$
  - $\hat{\tau}_3 = \overline{y_{3.}} - \overline{y_{..}} = 17.60 - 15.04 = -2.56$
  - $\hat{\tau}_4 = \overline{y_{4.}} - \overline{y_{..}} = 21.60 - 15.04 = 6.56$
  - $\hat{\tau}_5 = \overline{y_{5.}} - \overline{y_{..}} = 10.80 - 15.04 = -4.24$
- 
- \*see slide 7 for the data

# Confidence Interval (CI)

100(1 -  $\alpha$ ) percent CI on the  $i^{\text{th}}$  treatment mean:

$$\bar{y}_{i.} \pm t_{\alpha/2, N-a} \sqrt{MS_E/n}$$

100(1 -  $\alpha$ ) percent CI on the difference of two treatments mean:

$$\bar{y}_{i.} - \bar{y}_{j.} \pm t_{\alpha/2, N-a} \sqrt{2MS_E/n}$$

# Cotton Percentage example

-continued-

95% CI on the mean of treatment 4:

$$\left[ \underbrace{21.60}_{\bar{y}_{4.}} \pm \underbrace{(2.086)}_{t_{0.025,20}} \sqrt{\underbrace{8.06}_{MS_E} / \underbrace{5}_n} \right]$$

# Temperature Example

- Determine the effects of temperature on process yields
  - Case I: Two levels of temperature setting
  - Case II: Three levels of temperature setting

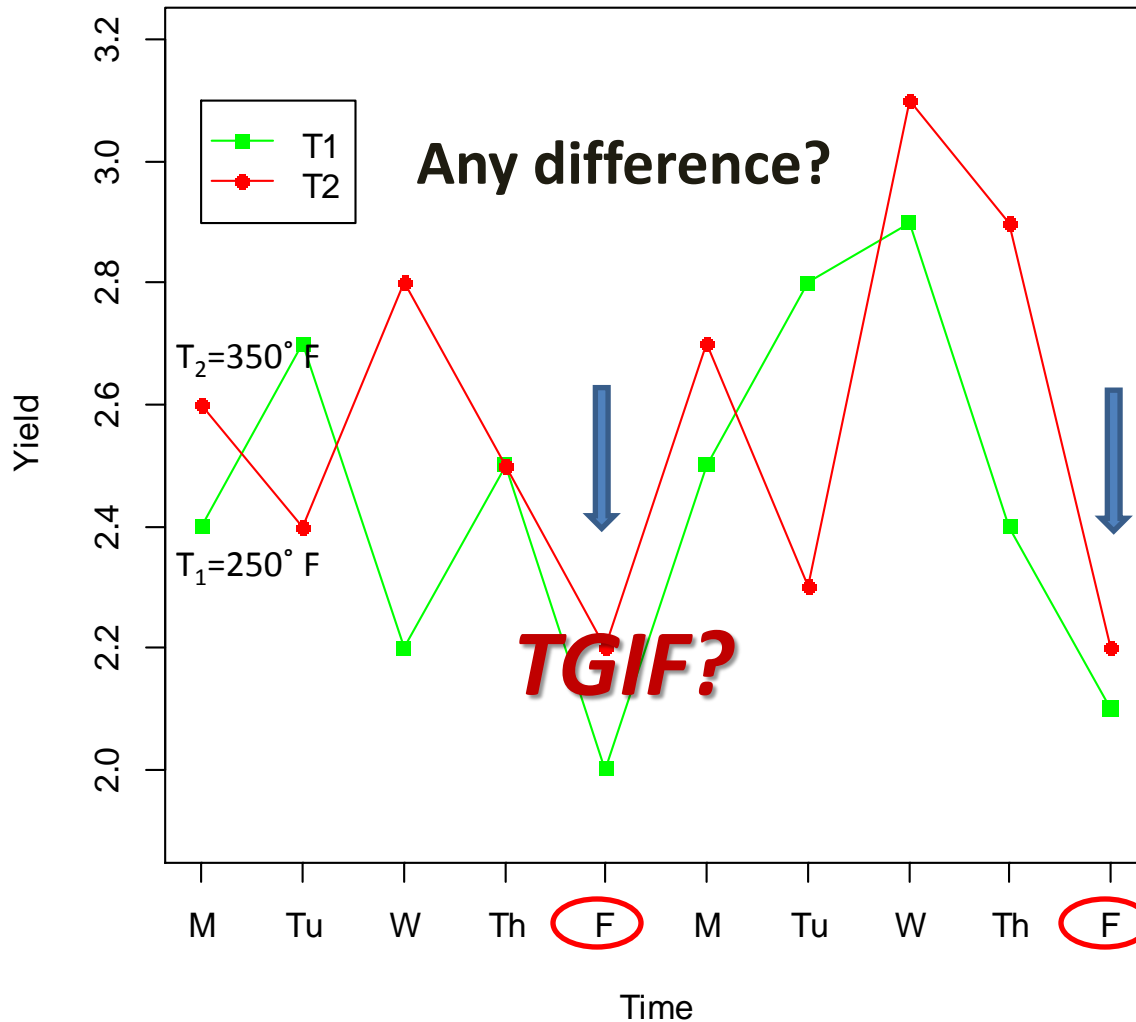
# Temperature Vs Process yields

---

		Temperature			
		250 °F	300 °F		
Week # 1	M	2.4	2.6	Week #3	
	Tu	2.7	2.4		
	W	2.2	2.8		
	Th	2.5	2.5		
	F	2.0	2.2		
Week # 2	M	2.5	2.7	Week # 4	
	Tu	2.8	2.3		
	W	2.9	3.1		
	Th	2.4	2.9		
	F	2.1	2.2		

---

# Exploratory analysis: Time Sequence plot





# Analysis

## t-test

- General form of t Statistics for one population

$$t = \frac{\hat{\theta} - \theta}{s_{\hat{\theta}}}$$

$\theta$ : the parameter to be estimated

$\hat{\theta}$ : the sample statistic that is the estimate of  $\theta$

$s_{\hat{\theta}}$ : the estimator of the standard deviation of  $\hat{\theta}$

# Temperature Data I

## t-test

Our interest : the difference  $\mu_1 - \mu_2$

Let

$$\hat{\theta} = \bar{X}_1 - \bar{X}_2$$
$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- Generally,  $\sigma_1, \sigma_2$  unknown, use estimates:
  - Two methods of estimating  $\sigma_{\bar{X}_1 - \bar{X}_2}$ , yielding *exact* and *approximate* t-test.

# Exact t-test

$$\hat{\theta} = \bar{X}_1 - \bar{X}_2$$

Exact t-test:

zero or a constant

$$t = \frac{\bar{X}_1 - \bar{X}_2 - d}{s_p \sqrt{1/n_1 + 1/n_2}}$$

Estimate of  $\sigma_{\bar{X}_1 - \bar{X}_2}$  in the form of pooled variance

where  $n_1$  and  $n_2$  are sample sizes,  $s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$ ,

$$s_a^2 = \frac{1}{n_a-1} \sum_{i=1}^{n_a} (x_{ai} - \bar{x}_a)^2, a = 1, 2$$

# Temperature Data I

## exact t-test

For the temperature data:

$$T = \frac{(2.45 - 2.57) - 0}{0.3005 \sqrt{1/10 + 1/10}} = -0.893$$

$$P(t \leq -0.893) = 0.191$$

p-value = 0.191, at 5% significance level, *fail to reject  $H_0$*

Conclusion -- two scenarios:

- $\mu_1 \cong \mu_2$  OR
- Highly variable data in each sample

# Exact t-test assumptions

- Normality of population
- Independence of samples
- Independence of observations
- Small ( $n < 30$ ) and equal or similar sample sizes ( $n_1 \cong n_2$ )
- *Equal variance*  $\sigma_1^2 = \sigma_2^2 = \sigma^2$

# Alternative approach: Confidence Interval for difference

- One sided test, such as  $\mu_1 < \mu_2$   
$$(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha, d.f.} S_{\bar{x}_1 - \bar{x}_2}$$
- Two Sided test, such as  $\mu_1 \neq \mu_2$   
$$(\bar{X}_1 - \bar{X}_2) \pm t_{\frac{\alpha}{2}, d.f.} S_{\bar{x}_1 - \bar{x}_2}$$

**Bounds can be calculated and the decision can be made based on interval information**

# ANOVA for Temperature Data I

Source of Variations	d.f.	SS	MS	F
Temperature	<b>a-1= 1</b>	0.072	0.072	<b>0.797</b>
Within	<b>an-a= 18</b>	1.626	0.090	Test Statistics Under $H_0$ , $F \sim F(a-1, n-1)$ , $p=0.3838$
Total	<b>an-1= 19</b>	1.698		

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

$$SS_{\text{within}} = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 \leftarrow \text{Numerator of } s_p \text{ in the exact t test!}$$

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

**a:** levels of treatment (temperature), **a=2**  
**n:** replication of each treatment, **n=10**

**ANOVA uses SAME assumptions as t test!**

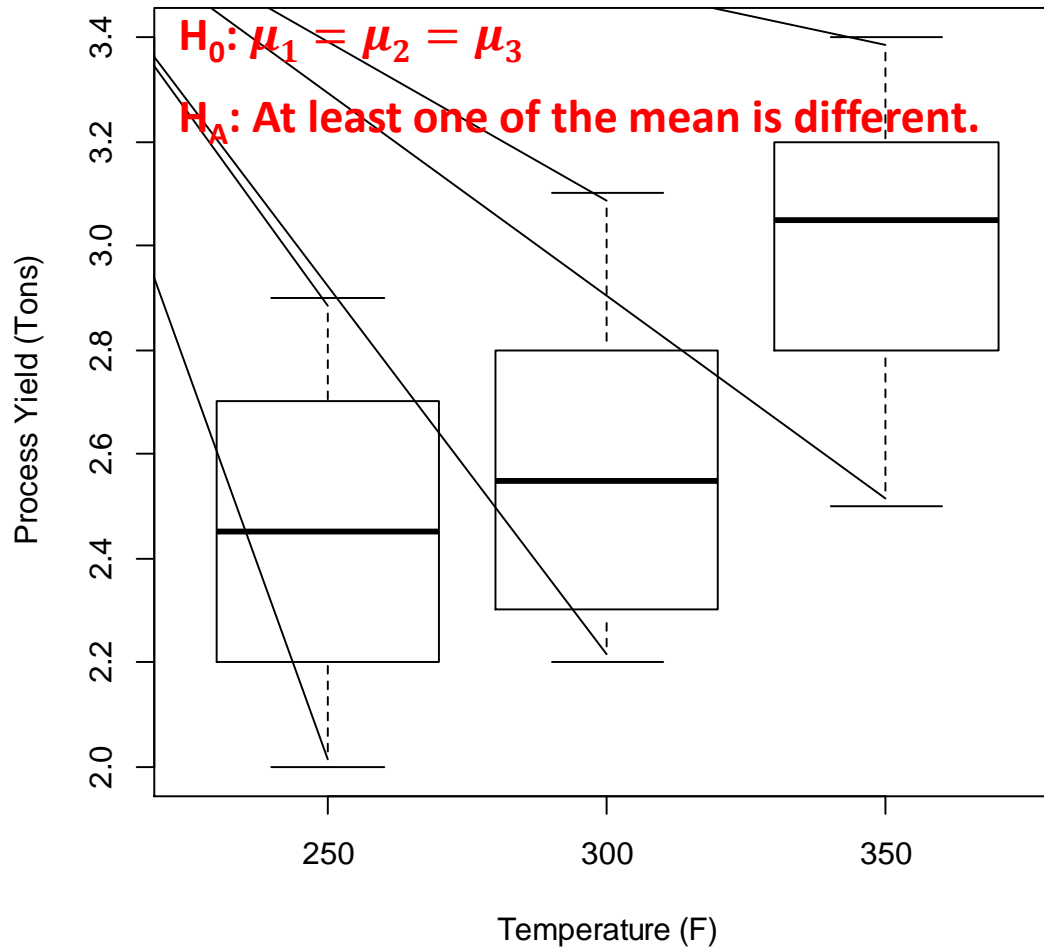
# Temperature Data II

Temperature				
Day	250 °F	300 °F	350°F	
M	2.4	2.6	3.2	
Tu	2.7	2.4	3.0	
W	2.2	2.8	3.1	
Th	2.5	2.5	2.8	
F	2.0	2.2	2.5	
M	2.5	2.7	2.9	
Tu	2.8	2.3	3.1	
W	2.9	3.1	3.4	
Th	2.4	2.9	3.2	
F	2.1	2.2	2.6	

**Variations:  
Week to Week?  
Day to Day?**



# Exploratory Analysis: Temperature Data II- box plots



**variability within each setting about the same**

# ANOVA for Temperature Data(3 levels)

Source of Variations	d.f.	SS	MS	F
Temperature	2	1.545	0.7725	8.91
Within	27	2.342	0.0867	<b>p-value=.001</b>
Total	29	3.887		<b>Reject H<sub>0</sub></b>

$$SS_{\text{temp}} = \sum_{i=1}^3 \frac{(\sum_{j=1}^{10} y_{ij})^2}{n} - \frac{(\sum_{i=1}^3 \sum_{j=1}^{10} y_{ij})^2}{an}$$

$$SS_{\text{total}} = \sum_{i=1}^3 \sum_{j=1}^{10} y_{ij}^2 - (\sum_{i=1}^3 \sum_{j=1}^{10} y_{ij})^2 / an$$

$$SS_{\text{within}} = SS_{\text{total}} - SS_{\text{temp}}$$

# Typical Data for a Single Factor Experiments

	Treatment (level)	Observations				Totals	Averages
a	1	$Y_{11}$	$Y_{12}$	...	$Y_{1n}$	$Y_{1.}$	$\bar{Y}_{1.}$
	2	$Y_{21}$	$Y_{22}$	...	$Y_{2n}$	$Y_{2.}$	$\bar{Y}_{2.}$
	.	.	.	...	.	.	.
	i	$Y_{i1}$	$Y_{i2}$	...	$Y_{in}$	$Y_{i.}$	.
	.	.	.	...	.	.	.
	a	$Y_{a1}$	$Y_{a2}$	...	$Y_{an}$	$Y_{a.}$	$\bar{Y}_{a.}$
		n				$Y_{..}$	$\bar{Y}$

# Clinical Trial Example

- A clinical trial is designed to estimate the efficacy of an experimental drug (D) compared to placebo (P) in congestive heart failure (CHF)
- Efficacy measure: the rate of change per week in distance walked after being administered therapy (D or P)
  - Baseline measure of left ventricular ejection fraction (LVEF) --- the lower the LVEF, the more serious CHF
  - 2 investigators
- Design and Analysis Tandem

# Rate of Change in Distance Walked per Week by Drug, Investigator, and Baseline LVEF

Investigator	LVEF≤15		15<LVEF≤25		25<LVEF≤30		30<LVEF	
	D	P	D	P	D	P	D	P
1	0	69	842	-451	319	-178	165	-141
	56	1276	107	-132	-59	8	-342	-533
	327	-20	-131		1075	191	173	-244
	-255	-109			1173	-181		
		-107			-5	-220		
2	-144	-248	228	-84	168	-286	-59	-383
	604	-71	-657	-316	885	155	291	-631
	1221	144	-211	-90	745	30	0	-397
	497		75	37		-123	540	10
	168		-96					-27

# Analysis of CHF data

- Means, standard deviations and sample sizes:
  - $\bar{X}_D = 240.6, S_D = 451.2, n_D = 32,$
  - $\bar{X}_P = -98.5, S_P = 322.4, n_P = 31,$
- The two-sample t with d.f. = 61 is

$$T = \frac{\bar{X}_D - \bar{X}_P}{\sqrt{S_D^2/n_D + S_P^2/n_P}} = 3.43$$

corresponds to a p-value = 0.001 : *reject  $H_0$ , drug is significantly more efficacious than placebo **overall***

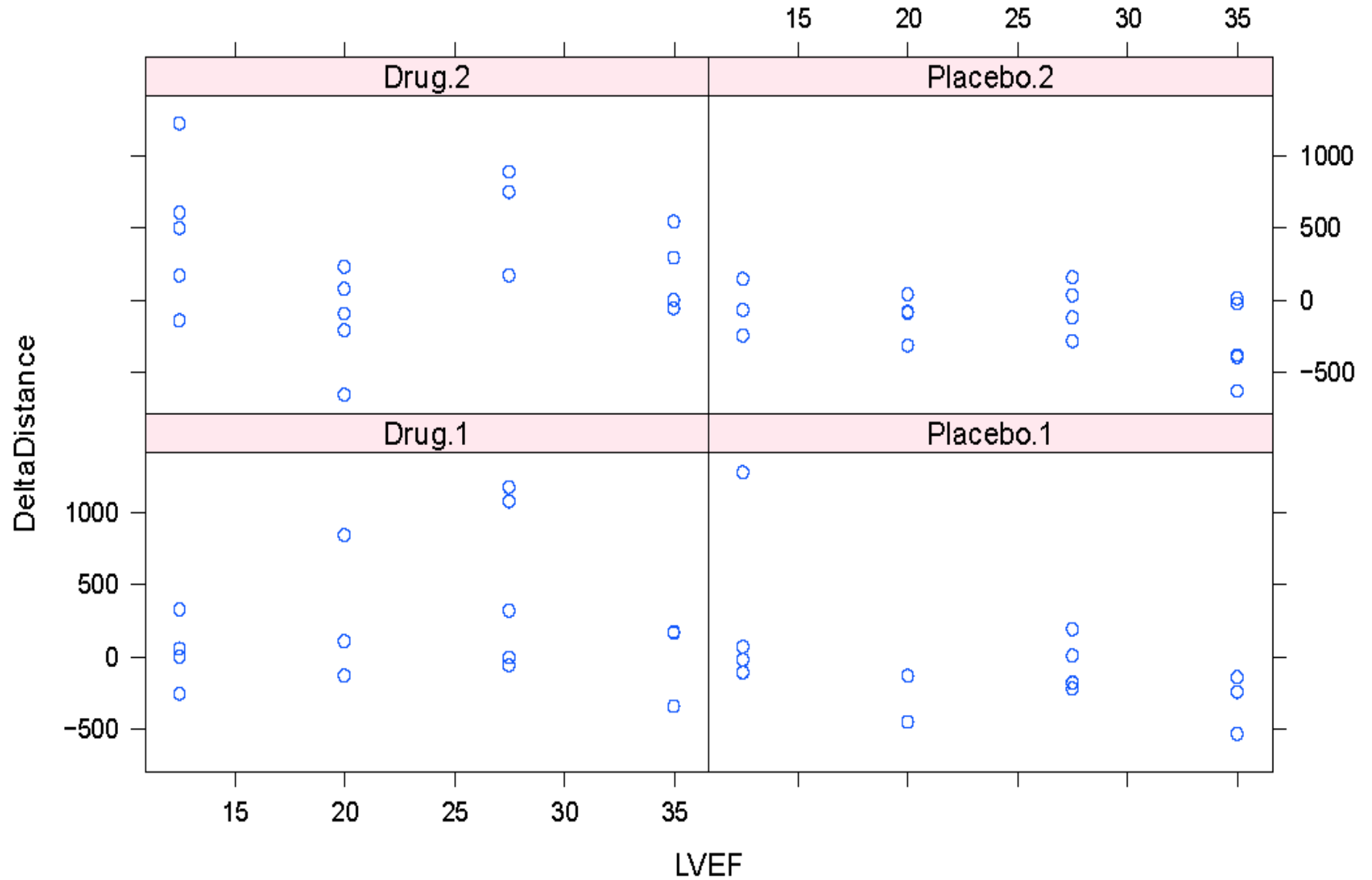
# ANOVA for CHF data

---

Factor	d.f.	Original data		Rank transform	
		F	P	F	P
Drug	1	11.66	0.001	15.96	0.0002
Investigator	1	0.09	0.77	0.02	0.9
LVEF	1	1.48	0.23	1.18	0.28
Inv. * drug	1	0.21	0.65	0.06	0.81
LVEF*drug	1	1.76	0.19	2.63	0.11

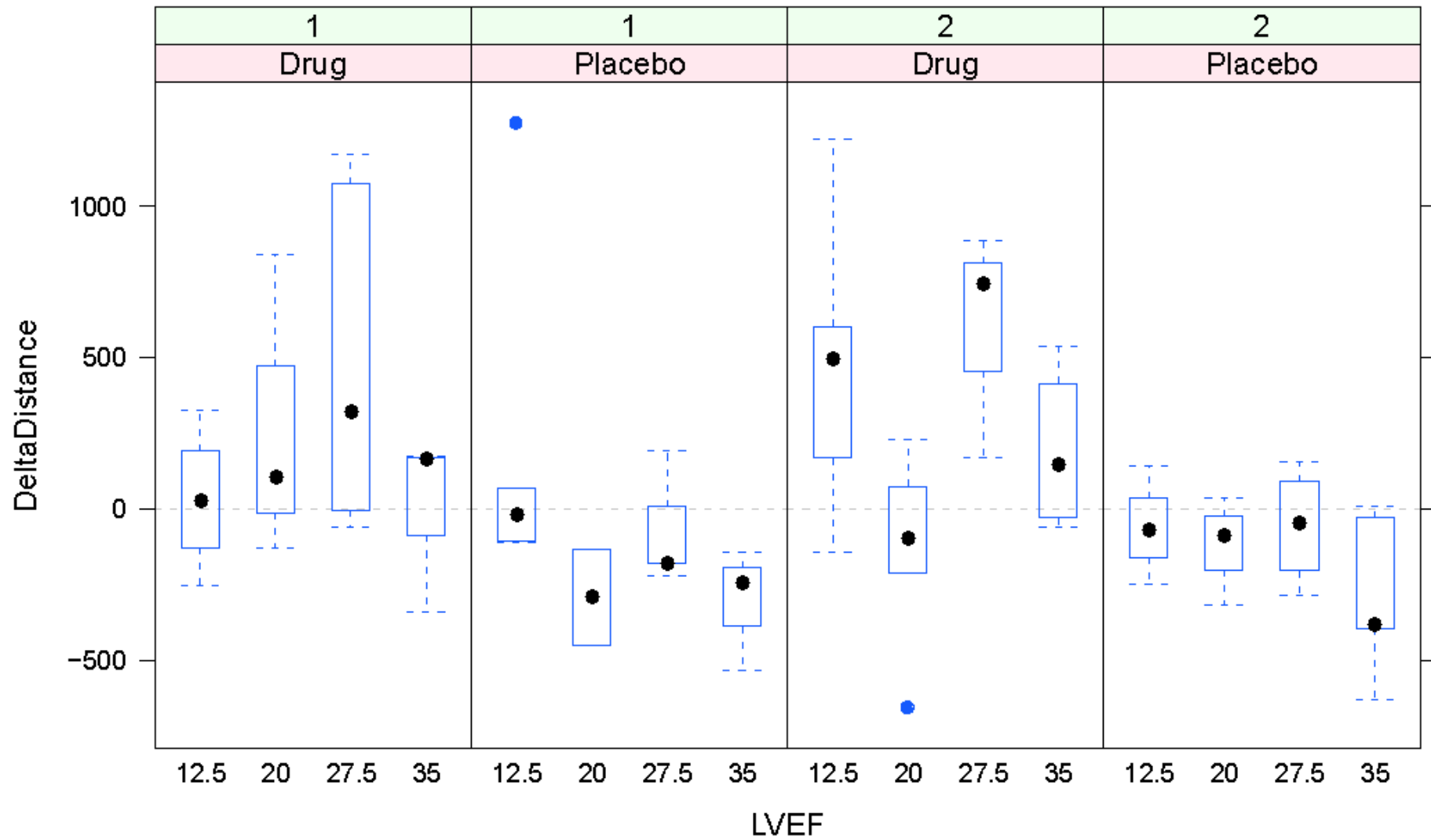
---

# Exploratory Data Analysis

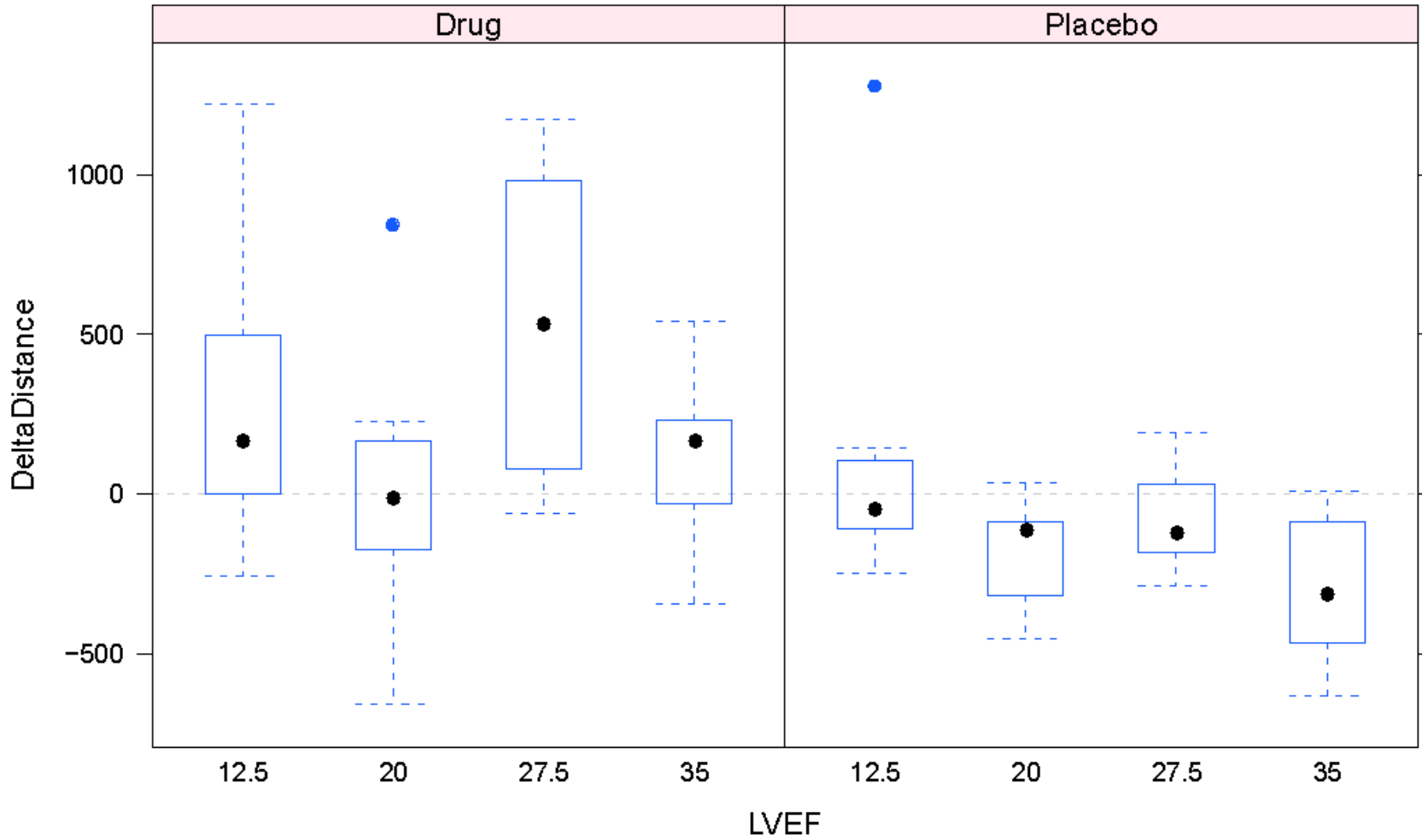




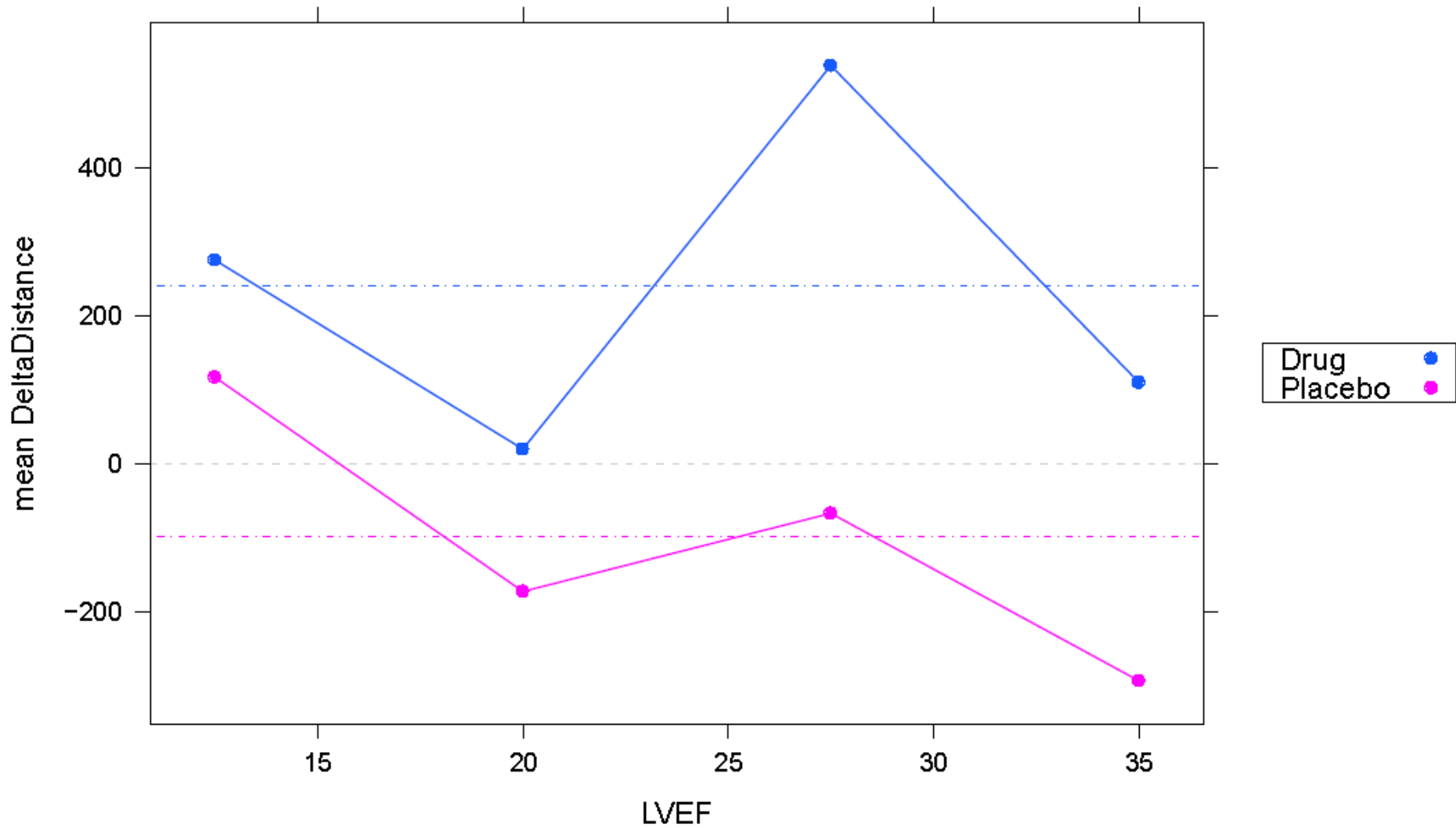
# Box plots by Investigator



# Box plots – Investigator effect pooled



# Mean Distance by LVEF category



# Mean Distance (D-P) by LVEF category

