

ANOVA and pair-wise comparisons

REEU

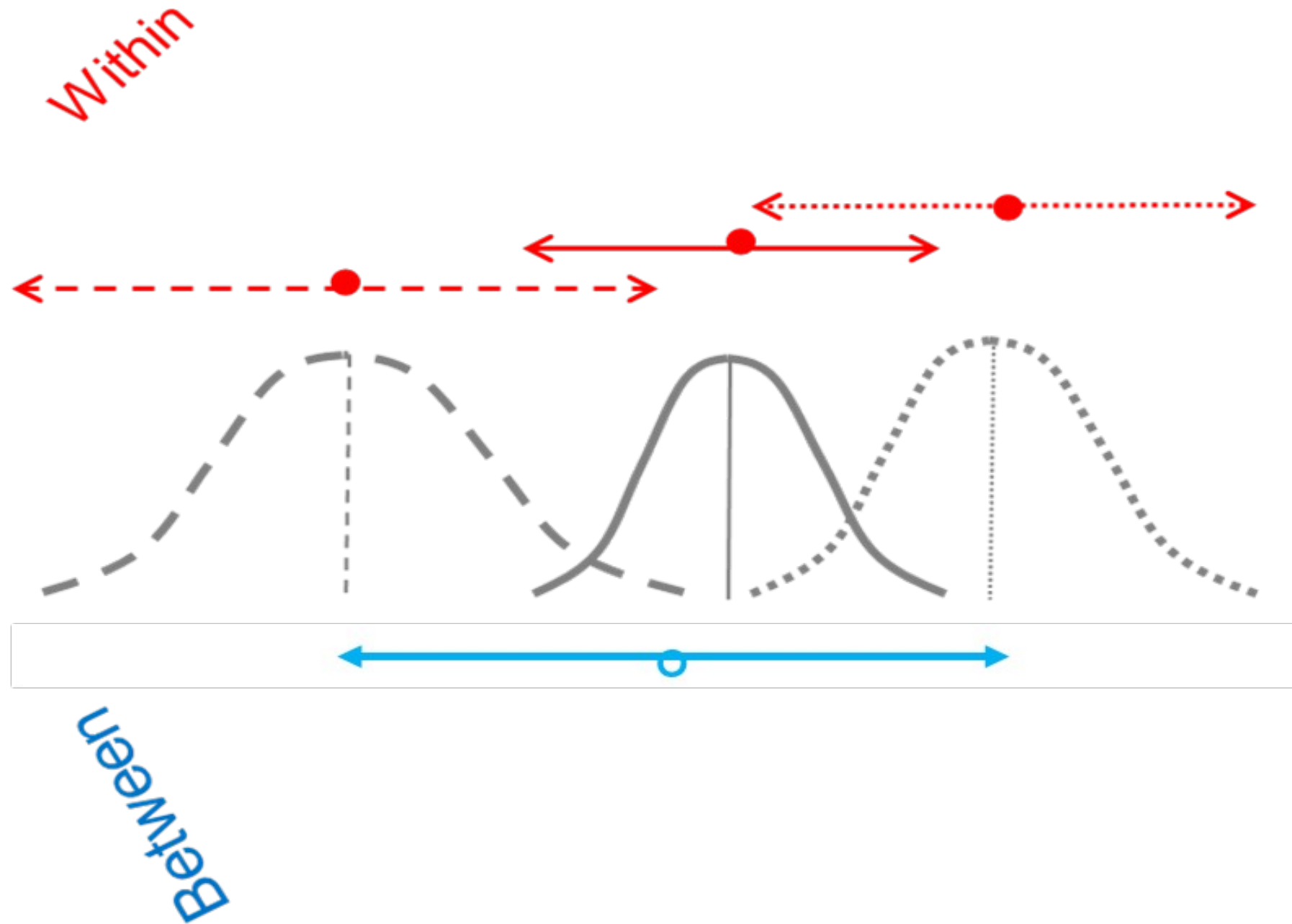
2024 Cohort

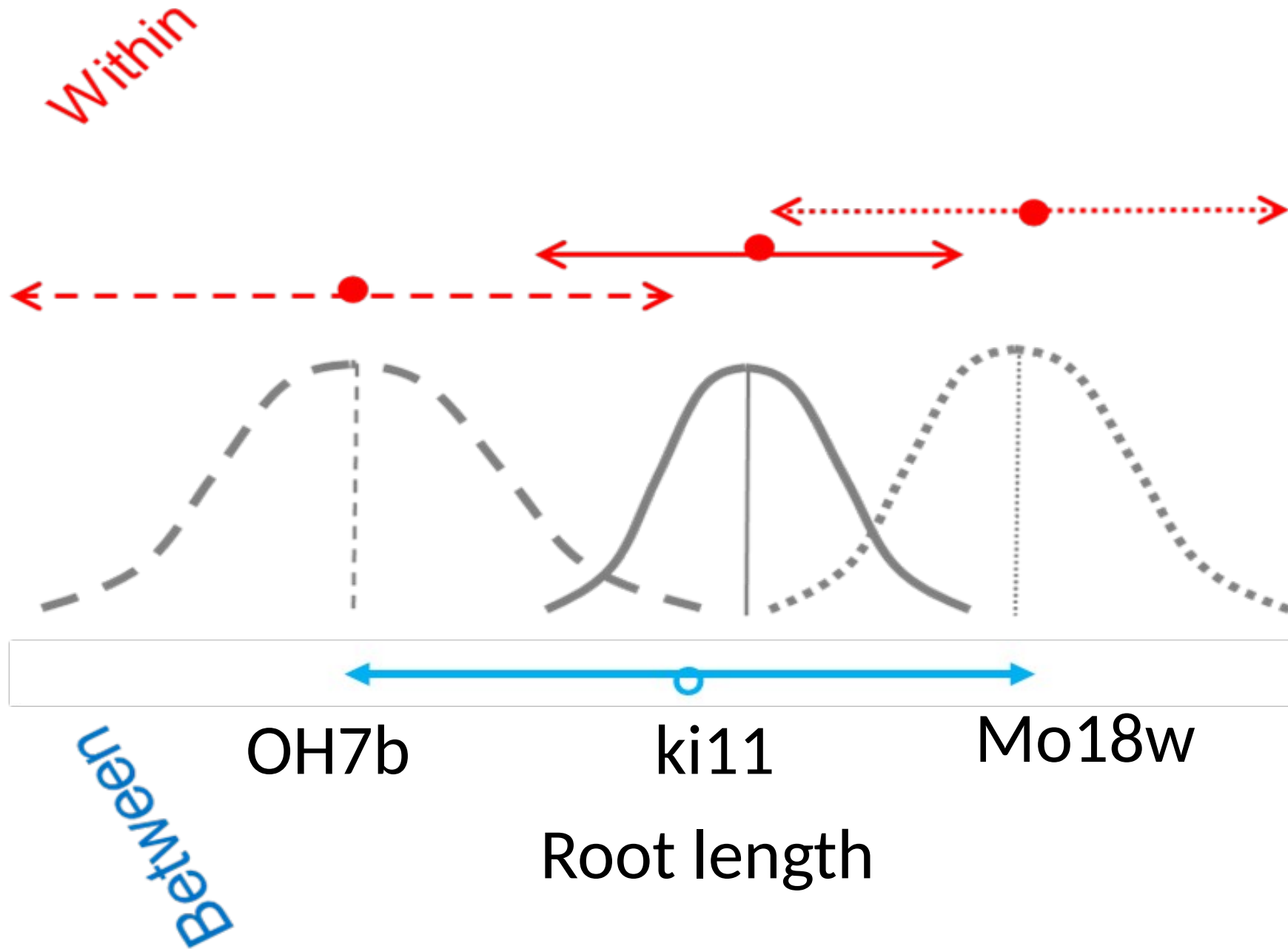
Discussion points for today

- What does ANOVA stand for?
- Why do we use ANOVA?
- How do we form a hypothesis?
- How do we test a hypothesis?
- How do we structure an ANOVA?
- What are the assumptions of an ANOVA?
- Almost there! One more step: How do we compare means?
- How do we interpret the results?

-What does ANOVA stand for?

ANalysis
Of
VAriance

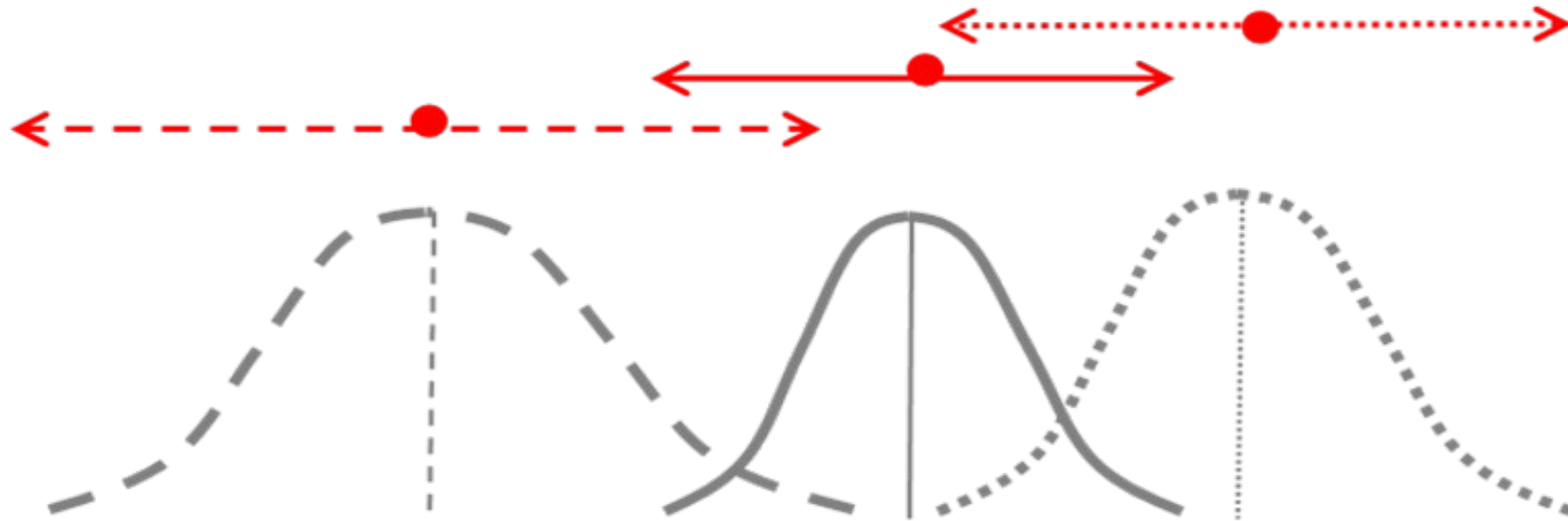




-Why do we conduct ANOVAs?

Within

Answer the question: Do the means of more than **two groups** differ *significantly*?



Between

CML322

B73

Mow18

Root length

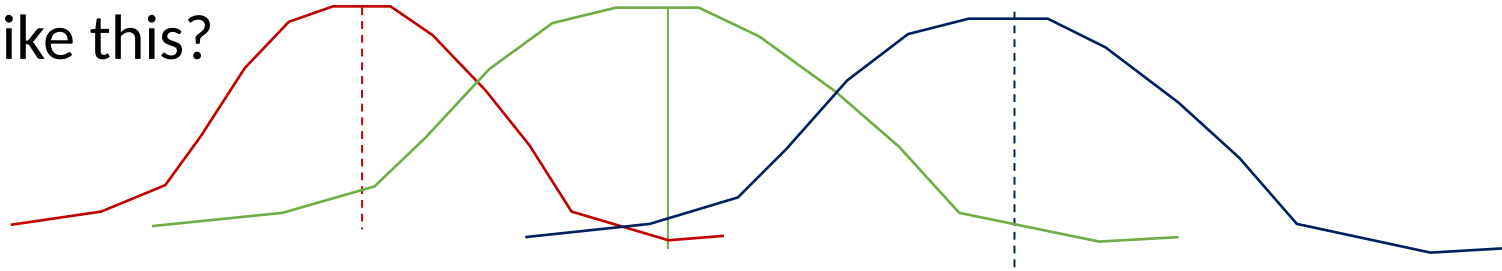
Which groups are more likely different?

CML322

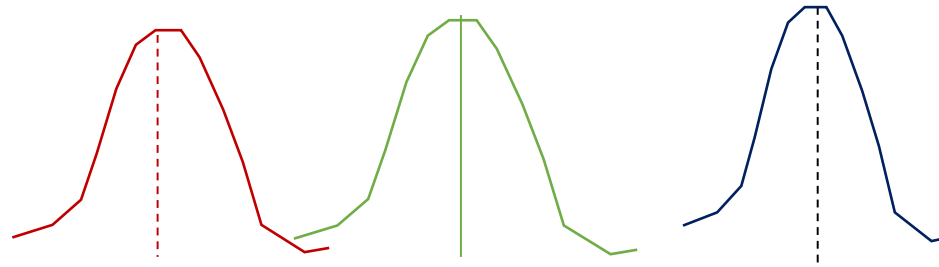
B73

MO18w

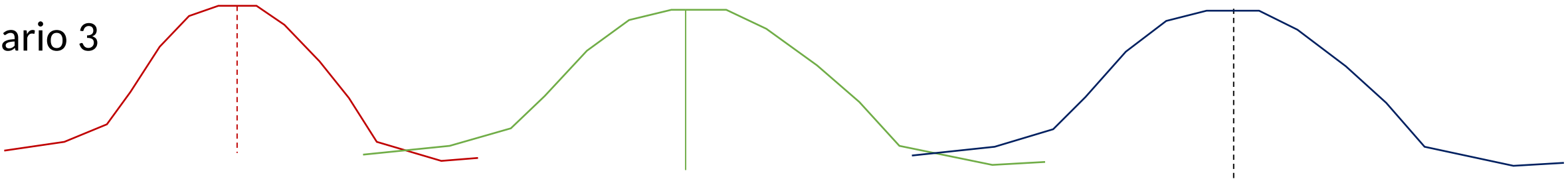
If this data looked like this?
Scenario 1



Scenario 2



Scenario 3



-How do we form a hypothesis?

$H_0: \mu_1 = \mu_2 = \mu_3$ $H_1: \mu_i \neq \mu_j$ for at least one pair of i and j .

H means *hypothesis*

$_0$ means *null (or no difference)*

μ means *average*

-How do we form a hypothesis?

$H_0: \mu_1 = \mu_2 = \mu_3$ $H_1: \mu_i \neq \mu_j$ for at least one pair of i and j .

H means *hypothesis*

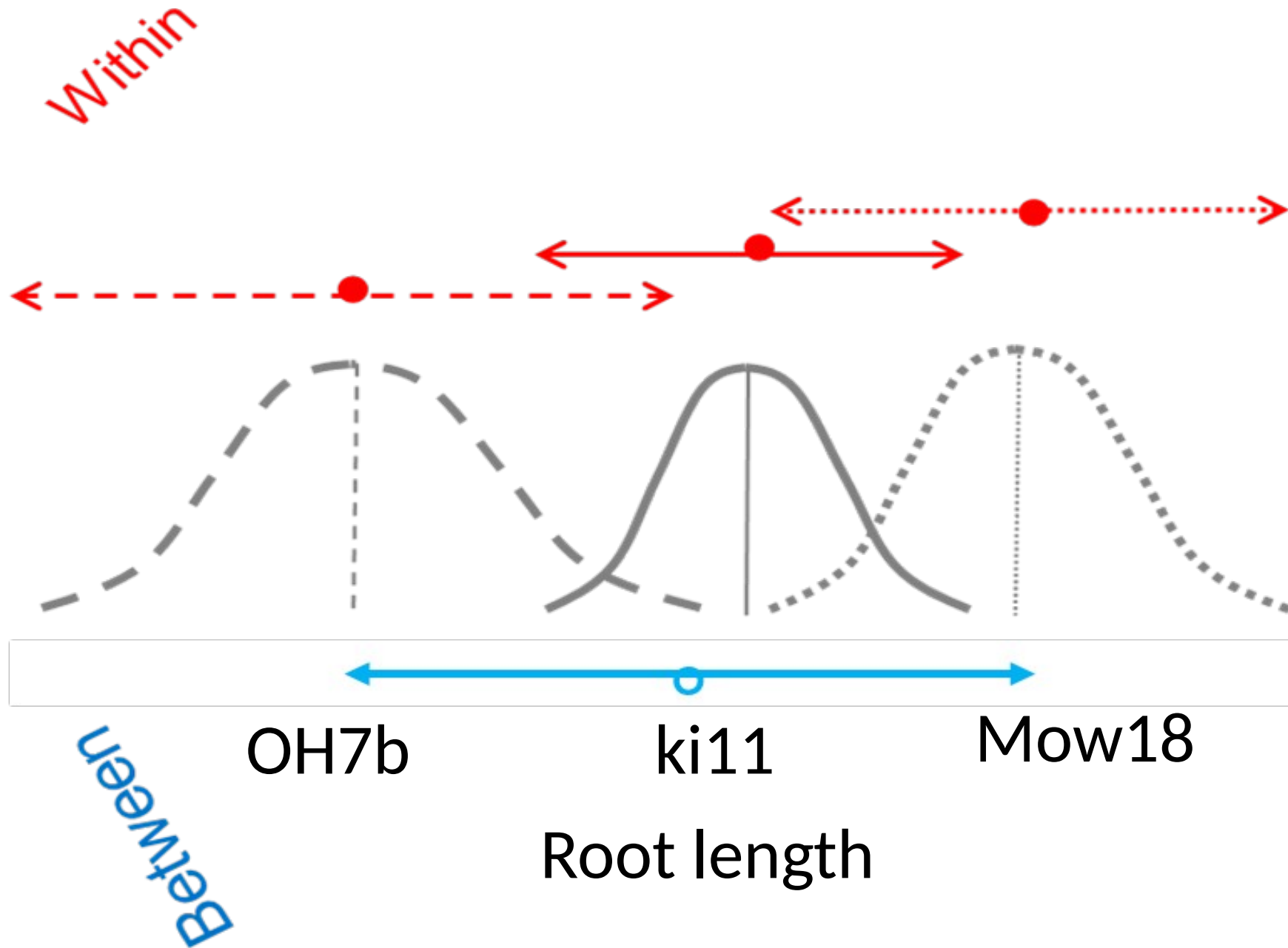
$_0$ means *null (or no difference)*

μ means *average*

Questions: How many groups in this example? How many groups in your research project?

-How do we test a hypothesis?

Think back to its name...



-How do we test a hypothesis?

The answer is in the name...

$$\text{Statistic} = \frac{\text{Variance between treatments}}{\text{Variance within treatments}}$$

-How do we test a hypothesis?

The answer is in the name...

$$\text{Statistic} = \frac{\text{Variance between treatments}}{\text{Variance within treatments}}$$

How does our calculated statistic compare to the critical value for our given sample size and number of groups with a particular confidence level?

The larger the statistic, the more evidence for differences between groups

-How do we structure an ANOVA?

Clue, we have seen this function before. It is the basis for most of our statistics.

-How do we structure an ANOVA?

Correct, like a linear regression! *But with categories rather than continuous variables as the independent variable*

$$y_{ik} = \mu + \alpha_k + \varepsilon_{ik}$$

μ = grand mean

α_k = an effect of treatment for group k

ε_{ik} = a person i 's residual within group k

Calculation of Sum of Squares: Now we do ANOVA by hand

Group (X)	Score (Y)							
1	15							
1	16							
1	14							
1	13							
1	12							
2	26							
2	25							
2	23							
2	20							
2	21							
3	10							
3	9							
3	9							
3	6							
3	6							

Between group Within group

Nah, just joking!

Calculation of Sum of Squares: But if we did ANOVA by hand

Group (X)	Score (Y)	\bar{Y}	$\bar{Y}_k - \bar{Y}$	$Y - \bar{Y}_k$	Y^2	\bar{Y}^2	$(\bar{Y}_k - \bar{Y})^2$	$(Y - \bar{Y}_k)^2$
1	15	15	14-15	'15-14	225	225	1	1
1	16	15	14-15	'16-14	256	225	1	4
1	14	15	14-15	'14-14	196	225	1	0
1	13	15	14-15	'13-14	169	225	1	1
1	12	15	14-15	'12-14	144	225	1	4
2	26	15	23-15	26-23	676	225	64	9
2	25	15	23-15	25-23	625	225	64	4
2	23	15	23-15	23-23	529	225	64	0
2	20	15	23-15	20-23	400	225	64	9
2	21	15	23-15	21-23	441	225	64	4
3	10	15	'8-15	'10-8	100	225	49	4
3	9	15	'8-15	'9-8	81	225	49	1
3	9	15	'8-15	'9-8	81	225	49	1
3	6	15	'8-15	'6-8	36	225	49	4
3	6	15	'8-15	'6-8	36	225	49	4

Between group = 570

Within group = 50

R output

```
> anova$Group<-as.factor(anova$Group)
> anovares<-lm(Score~Group, data=anova)
> anova(anovares)
```

Analysis of Variance Table

Response: Score

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Group	2	570	285.000	68.4	2.751e-07 ***
Residuals	12	50	4.167		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R code

```
lm.APA_primary <- lm(APA~ Genotype_ID, data = subset(roots_proc,  
Root_type = "Primary"))
```

```
anova(lm.APA_primary)
```

-What are the assumptions of an ANOVA?

Is it a good model? Or do we need to transform data?

-What are the assumptions of an ANOVA?

Is it a good model? Or do we need to transform data?

Normality: The population(s) from which the samples are drawn is normally distributed.

Test with residual histogram and PP or QQ plot.

In R: `hist(model)`

-What are the assumptions of an ANOVA?

Is it a good model? Or do we need to transform data?

Normality: The population(s) from which the samples are drawn is normally distributed.

Test with residual histogram and PP or QQ plot.

In R: `hist(residuals(model))`

Homogeneity: The variance of the groups is assumed to be equal

Examine with residual vs. the predicted values plot

In R: `plot(model)`

-What are the assumptions of an ANOVA?

Is it a good model? Or do we need to transform data?

Normality: The population(s) from which the samples are drawn is normally distributed.

Test with residual histogram and PP or QQ plot.

In R: `hist(residuals(model))`

Homogeneity: The variance of the groups is assumed to be equal

Examine with residual vs. the predicted values plot

In R: `plot(model)`

Independence of observations: The measurements are independent and random

Growth series: We would need to use repeated measures ANOVA

-Let us evaluate the model

Is it a good model? Or do we need to transform data?

R code

```
lm.APA_primary <- lm(APA~ Genotype_ID, data = subset(roots_proc,  
Root_type == "Primary"))
```

```
plot(lm.APA_primary)
```

```
hist(residuals(lm.APA_primary))
```

R code with transformation

```
lm.APA_primary <- lm(log(APA)~ Genotype_ID, data =  
subset(roots_proc, Root_type == "Primary"))
```

```
plot(lm.APA_primary)
```

```
hist(residuals(lm.APA_primary))
```

```
anova(lm.APA_primary)
```

Interpretation: What is a “significant” difference?

First, look at p-value

- ... p-value tells us about probability based on your data [means & deviation]
- ... of getting a value more extreme than dataset if null hypothesis were true

Second, select a significance level

- ... typically 0.05 (by convention)
- ... indicating that an extreme outcome is unlikely under null hypothesis

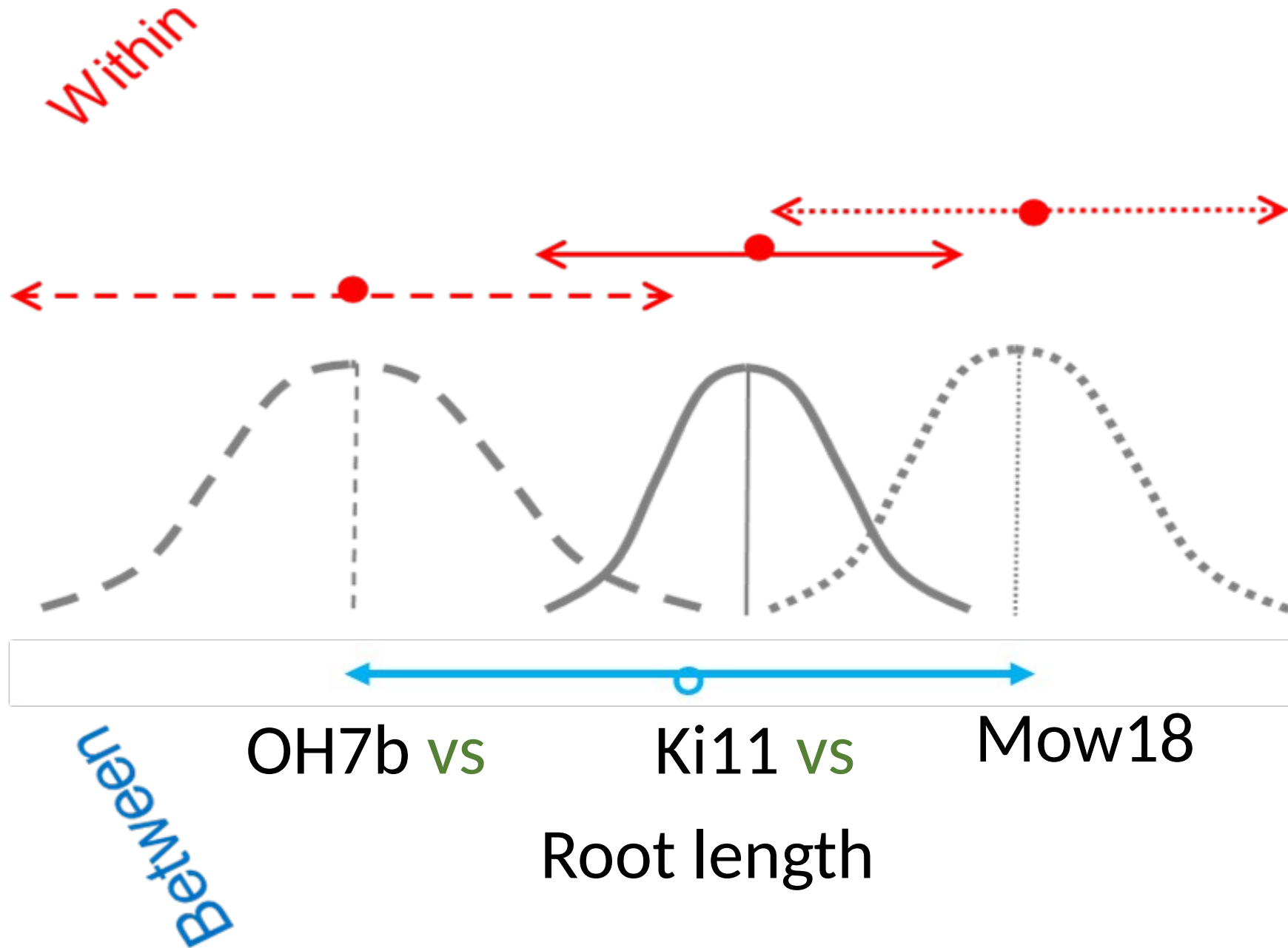
Third, accept or reject null hypothesis

- ... example “p-value < 0.05 , and so reject null hypothesis”
- ... also confident that the means will differ 95% of the time

-Almost there! One more step:

How do we compare means?

- ANOVA results only tells us whether there is *a* significant mean difference
- But the test does not tell us where the difference is
- And so, we need to explore all pair-wise comparisons of means



R code with transformation

```
lm.APA_primary <- lm(log(APA)~ Genotype_ID, data =  
subset(roots_proc, Root_type == "Primary"))
```

```
plot(lm.APA_primary)
```

```
hist(lm.APA_primary)
```

```
anova(lm.APA_primary)
```

```
cld(summary(glht(lm.APA_primary, mcp(Genotype_ID = "Tukey")),  
test=adjusted("bonferroni")), level=0.05, decreasing = TRUE)
```