

Unit 14: Analysis of Variance (ANOVA)

Case Studies:

 Is there an <u>association</u> between **political affiliation** (republican/democrat/independent/"no preference"/"other party") and **age**?

Purpose of this Lectures



- 1. Analyses for Associations
- 2. Association Analyses Summary: One Categorical Explanatory Variable (with >2 Levels)-> Numerical Response Variable (One Way)
- 3. Association Analyses Summary: One Categorical Explanatory Variable (with >2 Levels)-> Numerical Response Variable (Another Way)
- Modeling the Association Between a Categorical Explanatory Variable (with >2 levels) and a Numerical Response Variable (in the <u>Sample</u>).
 - **4.1.** Relationship between the sample intercept and slopes and the sample means.
- Modeling the Association Between a Categorical Explanatory Variable (with >2 levels) and a Numerical Response Variable (in the <u>Population</u>).
 - 5.1. ANOVA: Relationship between the population intercept and slopes and the population means

Additional Resources

Section 5.5 in Diez, Barr, and Cetinkaya-Rundel, (2015), *OpenIntro Statistics* https://www.openintro.org/download.php?file=os3&redirect=/stat/textbook/os3.php

1. ANALYSES FOR ASSOCIATIONS

Questions to consider, when selecting an analysis to test an association.



2. <u>Association Analysis Summary:</u> (ONE WAY TO PERFORM THIS ANALYSIS)

<u>Response</u>: Numerical

EXPLANATORY: ONE CATEGORICAL WITH >2 LEVELS

	Type of Variables Involved in the Association Test	Explanatory Variables: One Categorical Variable with >2 levels Response Variable: Numerical Variable				
Research Questions about Associations	Example	Is there an <u>association</u> between political affiliation (republican/democrat/independent/"no preference"/"other party") and age ?				
	Type of Association (Way to Quantify Association)	Multiple Linear Regression Model (<u>linear relationship</u> between categorical explanatory variable (x) and response variable (y))				
Descriptive Analytics	How to <u>Describe</u> an Association in a <u>Sample</u> ?	1. Multiple Linear Regression Model: • $\hat{y} = \widehat{\beta_0} + \widehat{\beta_1}x_1 + \widehat{\beta_2}x_2 + \dots + \widehat{\beta_p}x_p$ • R^2 of the model				
	When is this analysis <i>(for the sample)</i> appropriate to use?	Linearity condition is met				
Inferential	How to <u>Infer</u> an Association for a <u>Population?</u>	 Conduct inference on: A single population parameter β_i <u>All</u> population parameters β₁,, β_p A subset of population parameters β_{p+1},, β_{p+q} 				
Statistics	When is this analysis <i>(for the population)</i> appropriate to use?	 Linearity condition is met Constant variance of residuals condition is met. Residuals are normal (and centered at 0). Residuals are independent. No-Multicollinearity condition is met. 				
Predictive	Making Predictions	Use your multiple linear regression line to make predictions $\hat{y} = \hat{\beta_0} + \hat{\beta_1}x_1 + \hat{\beta_2}x_2 + \dots + \hat{\beta_p}x_p$				
Analytics	How to quantify the performance of your prediction(s)?	 <u>Individual Data Point</u>: residual <u>All Data</u>: root mean square error (RMSE) 				

3. <u>Association Analysis Summary:</u> (<u>ANOTHER WAY</u> TO PERFORM THIS ANALYSIS)

<u>Response</u>: Numerical

EXPLANATORY: ONE CATEGORICAL WITH >2 LEVELS

	Type of Variables Involved in the Association Test	Explanatory Variables: One Categorical Variable with >2 levels Response Variable: Numerical Variable				
Research Questions about Associations	Example	Is there an <u>association</u> between political affiliation (republican/democrat/independent/"no preference"/"other party") and age ?				
	Type of Association (Way to Quantify Association)	Does <u>at least one pair of political affiliations (</u> out of the five given) have <u>mean</u> <u>ages</u> that are <u>different</u> ?				
Descriptive	How to <u>Describe</u> an Association in a <u>Sample</u> ?	 \$\overline{x_1}\$, \$\overline{x_2}\$, \$\dots\$, \$\overline{x_p}\$ at least one pair of these sample means are different. \$\overline{p}\$ is the number of levels of the categorical variable \$\overline{x_i}\$ is the sample mean (the numerical variable) for the <i>ith</i> level of the categorical variable. 				
Analytics	When is this analysis <i>(for the sample)</i> appropriate to use?	Sample distribution <i>i</i> symmetric and unimodal (for <i>i=1,2,,p</i>).				
Inferential	How to <u>Infer</u> an Association for a <u>Population?</u>	$\begin{array}{l} \hline \textbf{Conduct an Analysis of Variance (ANOVA)} \\ H_0: \mu_1 = \mu_2 = \cdots = \mu_p \\ H_A: \mu_i \neq \mu_j \ for \ at \ least \ one \ pair \ of \ levels \\ \bullet p \ is \ the \ number \ of \ levels \ of \ the \ categorical \ variable \\ \bullet \bar{x}_i \ is \ the \ sample \ mean \ (the \ numerical \ variable) \ for \ the \ ith \ level \ of \ the \ categorical \ variable. \end{array}$				
Statistics	When is this analysis <i>(for the population)</i> appropriate to use?	Conditions for ANOVA Inference1. Each of the p sample distributions are random.2. $n_i < 10\%$ of the population distribution.3. Each sample distribution is approximately normal.4. $s_1 \approx s_2 \approx \cdots \approx s_p$ 5. Each of the p samples are independent of each other.				
Predictive Analytics	Making Predictions How to quantify the performance of your prediction(s)?	N/A N/A				

4. <u>MODELING</u> THE ASSOCIATION BETWEEN A CATEGORICAL EXPLANATORY VARIABLE (WITH >2 LEVELS) AND A NUMERICAL RESPONSE VARIABLE <u>IN THE SAMPLE</u>

Ex: We would first like to model the association between political association and age, first using a <u>multiple linear regression model</u> (with age as the response variable and party as the explanatory variable.) In this analysis, we will again examine our Pew Survey dataset (from 2017) which has randomly sampled n=1465 adults living in the U.S., which has five possible values in the "party" variable: Democrat, Independent, No preference, Other party, and Republican.

If we were to set up a multiple linear regression model for this analysis, how many slopes would the model have?

	party	age
0	Democrat	18.0
1	Independent	18.0
2	No preference (VOL.)	19.0
3	Other party (VOL.)	32.0
4	Republican	18.0

Using Python we find that the model yields the following output.

OLS Regression Results

Dep. Variable	:	aç	ge	R-squ	ared:		0.052
Model	:	OL	.S Ad	j. R-squ	ared:		0.049
Method	: Lea	ast Square	es	F-sta	tistic:		19.82
Date	: Mon, 3	0 Mar 202	20 Prob	(F-stat	istic):		6.66e - 16
Time	:	10:27:8	55 Lo	g-Likeli	hood:		-6261.1
No. Observations	:	146	65		AIC:	1.	253e+04
Df Residuals	:	146	60		BIC:	1.	256e+04
Df Model	:		4				
Covariance Type	:	nonrobu	st				
	coef	std err	t	P> t	[0.0	25	0.975]
Intercept	50.4991	0.758	66.618	0.000	49.0	12	51.986
party[T.Ind]	-3.6914	1.073	-3.440	0.001	- 5.7	96	-1.587
party[T.No_Pref]	-7.3527	2.821	-2.606	0.009	- 12.8	87	-1.818
party[T.Other]	-5.8991	7.819	-0.754	0.451	- 21.2	37	9.439
party[T.Rep]	6.2775	1.183	5.306	0.000	3.9	57	8.598
Omnibus:	130.613	Durbi	n-Watso	n: 1	1.725		
Prob(Omnibus):	0.000	Jarque-	Bera (JB): 40).798		
Skew:	-0.017		Prob(JB): 1.38	8e-09		
Kurtosis:	2.183		Cond. N	0.	19.0		

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

With the information in the summary output table, we can find the intercept and slopes for this multiple linear regression model *for the sample data*.

age = 50.4991 - 3.6914*party*[*T*.*Ind*] - 7.3527*party*[*T*.*No_Pref*] - 5.8991*party*[*T*.*Other*] + 6.2775*party*[*T*.*Rep*] By interpreting our intercept and slopes, we can use this multiple linear regression model to say the following things about the sample.

- $\widehat{\beta_0}$ Interpretation: We expect the mean age of **Democrats** in the sample to be **50.4991.**
- <u>\$\heta_1\$ Interpretation:</u> We expect the mean age of **Independents** in the sample to be 3.6914 years lower than the mean age of Democrats in the sample (ie. **46.808**=50.4991-3.6914).
- <u>β₂ Interpretation</u>: We expect the mean age of people with **no party preference** in the sample to be 7.3527 years lower than the mean age of Democrats in the sample (ie. **43.1464**=50.4991-7.3527).
- <u>\$\heta_3\$ Interpretation:</u> We expect the mean age of people from "other parties" in the sample to be 5.8991 years lower than the mean age of Democrats in the sample (ie. 44.6=50.4991-5.8991).
- <u>\$\heta_4\$ Interpretation:</u> We expect the mean age of **Republicans** in the sample to be 6.2775
 years higher than the mean age of Democrats in the sample (ie.

 56.7766=50.4991+6.2775).

4.1. RELATIONSHIP BETWEEN THE <u>SAMPLE INTERCEPT AND SLOPES</u> OF THIS MULTIPLE LINEAR REGRESSION AND THE <u>SAMPLE MEANS</u>

Let's find the actual sample mean age of each of the five political affiliations in this Pew dataset and compare them to what we said when we interpreted our slopes and intercept.

Sample Mean Ages

	age
party	
Dem	50.499051
Ind	46.807619
No_Pref	43.146341
Other	44.600000
Rep	56.776567

Each of these sample mean ages matches up to what we said when we interpreted the intercepts and slopes of the multiple linear regression line!

In general, this relationship that we observed between the sample slopes and intercepts and the sample means will hold whenever you have a multiple linear regression model with a single explanatory variable that is categorical.

5. <u>MODELING</u> THE ASSOCIATION BETWEEN A CATEGORICAL EXPLANATORY VARIABLE (WITH >2 LEVELS) AND A NUMERICAL RESPONSE VARIABLE <u>IN THE POPULATION</u>

Now let's think about what the multiple linear regression equation would look like if we modelled the relationship between age (response variable) and political affiliation (explanatory variable) in the *population* of all adults living in the U.S.

 $\widehat{age} = \beta_0 + \beta_1 party[T.Ind] + \beta_2 party[T.Ind] + \beta_3 party[T.No \ pref] + \beta_4 party[T.Rep]$

Do we have sufficient evidence to suggest that at least one of the four population slopes is non-zero?

1. First let's formulate the hypotheses. $H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$

 H_A : At least one $\beta_i \neq 0$ (for i=1,2,3,4)

2. Next, let's check the conditions for conducting inference.



b. Constant Variability of Residuals Condition?



c. Normality of Residuals Condition (with Mean 0)?



d. Independence of residuals condition?

e. No multicollinearity condition?

OLS Regression Res	sults						
Dep. Variable:		age	R	squared:	0	.052	
Model:		OLS	Adj. R·	squared:	0	.049	
Method:	Least	Squares	F	-statistic:	1	9.82	
Date:	Mon, 05 A	Apr 2021	Prob (F-	statistic):	6.66	e-16	
Time:		16:03:02	Log-Li	kelihood:	-62	61.1	
No. Observations:		1465		AIC:	1.253	e+04	
Df Residuals:		1460		BIC:	1.256	e+04	
Df Model:		4					
Covariance Type:	no	onrobust					
		coe	f std eri	r t	P> t	[0.025	0.975]
	Intercept	coe 50.4991	f stderi 1 0.758	t 66.618	P> t 	[0.025 49.012	0.975] 51.986
party[T.In	Intercept dependent]	coe 50.4991 -3.6914	f std er 1 0.758 1 1.073	t 66.618 -3.440	P> t 0.000 0.001	[0.025 49.012 -5.796	0.975] 51.986 -1.587
party[T.In	Intercept dependent] nce (VOL.)]	coe 50.4991 -3.6914 -7.3527	f std err 1 0.758 4 1.073 7 2.821	t 66.618 6 -3.440 -2.606	P> t 0.000 0.001 0.009	[0.025 49.012 -5.796 -12.887	0.975] 51.986 -1.587 -1.818
party[T.In party[T.No prefere party[T.Other pa	Intercept dependent] nce (VOL.)] arty (VOL.)]	coe 50.4991 -3.6914 -7.3527 -5.8991	f std err 0.758 1 0.758 1 1.073 7 2.821 1 7.819	t 666.618 -3.440 -2.606 -0.754	P>[t] 0.000 0.001 0.009 0.451	[0.025 49.012 -5.796 -12.887 -21.237	0.975] 51.986 -1.587 -1.818 9.439
party[T.In party[T.No prefere party[T.Other pa party[T.F	Intercept dependent] nce (VOL.)] arty (VOL.)] Republican]	coe 50.4991 -3.6914 -7.3527 -5.8991 6.2775	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183	t 66.618 -3.440 -2.606 -0.754 5.306	P>[t] 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus:	Intercept dependent] nce (VOL.)] arty (VOL.)] Republican] 130.613	coe 50.4991 -3.6914 -7.3527 -5.8991 6.2775 Durbin-V	f std err 1 0.758 1 1.073 7 2.821 1 7.819 5 1.183 Vatson:	 t 66.618 -3.440 -2.606 -0.754 5.306 1.725 	P>[t] 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus: Prob(Omnibus):	Intercept dependent] nce (VOL.)] arty (VOL.)] Republican] 130.613 0.000 J	coe 50.4991 -3.6914 -7.3527 -5.8991 6.2775 Durbin-V arque-Be	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183 Vatson: ra (JB):	t 8 66.618 9 -3.440 -2.606 0 -0.754 3 5.306 1.725 40.798	P>[t] 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus: Prob(Omnibus): Skew:	Intercept dependent] nce (VOL.)] arty (VOL.)] Republican] 130.613 0.000 J -0.017	coe 50.4991 -3.6914 -7.3527 -5.8991 6.2775 Durbin-V arque-Be Pr	f std err 0.758 1.073 2.821 1.7.819 5.1.183 Vatson: ra (JB): ob(JB):	 t 66.618 -3.440 -2.606 -0.754 5.306 1.725 40.798 1.38e-09 	P>[t] 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598

4. Finally, use the p-value (and a significance level of $\alpha = 0.05$) to make a conclusion about your hypotheses.

5.1. ANOVA: RELATIONSHIP BETWEEN THE <u>POPULATION SLOPES</u> OF THIS MULTIPLE LINEAR REGRESSION AND THE <u>POPULATION MEANS</u>

In section 4, we saw that there was a relationship between the <u>sample intercept and the 4</u> sample slopes of the multiple linear regression model (*ie*. $\hat{\beta}_0$, $\hat{\beta}_1$, $\hat{\beta}_2$, $\hat{\beta}_3$, $\hat{\beta}_4$) and the <u>5 sample means ages</u> corresponding to the 5 political affiliations (ie.

 $\overline{x}_{dem}, \overline{x}_{ind}, \overline{x}_{no \ pref}, \overline{x}_{no \ party}, \overline{x}_{rep}).$

Is there also a relationship between the **population** intercept and the 4 sample slopes of the multiple linear regression model (β_0 , β_1 , β_2 , β_3 , β_4) and the <u>5 population means ages</u> corresponding to the 5 political affiliations (ie. μ_{dem} , μ_{ind} , $\mu_{no \ pref}$, $\mu_{no \ party}$, μ_{rep})?

Yes!

<u>Response Variable</u>: Numerical <u>Explanatory Variable</u>: Categorical (*w* levels)

 $\widehat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{w-1} x_{w-1}$

 $\mu_1, \mu_2, \ldots, \mu_w$

Hypotheses

 $H_0: \beta_1 = \beta_2 = \dots = \beta_{w-1} = 0$ $H_A: \text{at least one } \beta_i \neq 0 \text{ for i=1,2,...w-1}$

Test Statistic

 $test \ statistic = \frac{SSR/(w-1)}{SSE/(n-(w-1)-1)}$

 $p - value = P(F_{w-1,n-(w-1)-1} \ge test \ statistic)$

P-value

Hypotheses $H_0: \mu_1 = \mu_2 = \dots = \mu_w$

 $H_A: \mu_i \neq \mu_j$ for at least one pair (i, j) between 1 and w

Test Statistic

P-value

Do we have sufficient evidence to suggest that at least one pair of political affiliations (Democrat, Republican, Independent, No preference, Other party) have average ages *out of all adults living in the U.S.* that are different?

- 1. First let's formulate the hypotheses. $H_0: \mu_{dem} = \mu_{rep} = \mu_{ind} = \mu_{no \ pref} = \mu_{no \ party}$ $H_A: \mu_{group} \neq \mu_{another \ group}$ for at least one pair of political affiliations groups
- 2. **Next, let's check the conditions for conducting inference** (on the intercept and slopes $\beta_0, \beta_1, \beta_2, \beta_3, \beta_4$ of the population model $age = \beta_0 + \beta_1 party[T.Ind] + \beta_2 party[T.Ind] + \beta_3 party[T.No pref] + \beta_4 party[T.Rep])$
 - a. Linearity Condition?



b. Constant Variability of Residuals Condition?



c. Normality of Residuals Condition (with Mean 0)?



d. Independence of residuals condition?

e. No multicollinearity condition?

OLS Regression Res	sults						
Dep. Variable:		age	R-	squared:	0	.052	
Model:		OLS	Adj. R-	squared:	0	.049	
Method:	Leas	t Squares	F	-statistic:	1	9.82	
Date:	Mon, 05	6 Apr 2021	Prob (F-	statistic):	6.66	e-16	
Time:		16:03:02	Log-Lil	kelihood:	-62	61.1	
No. Observations:		1465		AIC:	1.253€	e+04	
Df Residuals:		1460		BIC:	1.256	e+04	
Df Model:		4					
Covariance Type:		nonrobust					
		coe	f std err	t t	P> t	[0.025	0.975]
	Interce	coe pt 50.499	f std err 1 0.758	t 66.618	P> t 0.000	[0.025 49.012	0.975] 51.986
party[T.In	Interce dependen	coe pt 50.499 t] -3.6914	f std err 1 0.758 4 1.073	t 66.618 -3.440	P> t 0.000 0.001	[0.025 49.012 -5.796	0.975] 51.986 -1.587
party[T.In	Interce dependen nce (VOL.	coe pt 50.499 (t] -3.691 (.)] -7.352	f std err 1 0.758 4 1.073 7 2.821	t 66.618 -3.440 -2.606	P>[t] 0.000 0.001 0.009	[0.025 49.012 -5.796 -12.887	0.975] 51.986 -1.587 -1.818
party[T.In party[T.No prefere party[T.Other pa	Interce dependen nce (VOL. arty (VOL.	coe pt 50.499 (t] -3.691 (.)] -7.352 (.)] -5.899	f stderr 1 0.758 4 1.073 7 2.821 1 7.819	t 666.618 -3.440 -2.606 -0.754	P> t 0.000 0.001 0.009 0.451	[0.025 49.012 -5.796 -12.887 -21.237	0.975] 51.986 -1.587 -1.818 9.439
party[T.In party[T.No prefere party[T.Other pa party[T.F	Interce dependen nce (VOL. arty (VOL. Republica	coe pt 50.499 .t] -3.6914 .)] -7.352 .)] -5.899 n] 6.277	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183	t 66.618 -3.440 -2.606 -0.754 5.306	P>[t] 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus:	Interce dependen nce (VOL. arty (VOL. Republican 130.613	coe pt 50.499 t] -3.691 .)] -7.352 .)] -5.899 n] 6.277 Durbin-V	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183 Watson:	 66.618 -3.440 -2.606 -0.754 5.306 1.725 	P> t 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus: Prob(Omnibus):	Interce dependen nce (VOL. arty (VOL. Republican 130.613 0.000	coe pt 50.499 tt] -3.691 .)] -7.352 .)] -5.899 n] 6.277 Durbin-V Jarque-Be	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183 Watson: ara (JB):	 t 66.618 -3.440 -2.606 -0.754 5.306 1.725 40.798 	P> t 0.000 0.001 0.009 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598
party[T.In party[T.No prefere party[T.Other pa party[T.F Omnibus: Prob(Omnibus): Skew:	Interce dependen nce (VOL. arty (VOL. Republican 130.613 0.000 -0.017	coe pt 50.499 t] -3.691)] -7.352)] -5.899 n] 6.277 Durbin-V Jarque-Be Pr	f std err 1 0.758 4 1.073 7 2.821 1 7.819 5 1.183 Watson: tra (JB): rob(JB):	 66.618 -3.440 -2.606 -0.754 5.306 1.725 40.798 1.38e-09 	P> t 0.000 0.001 0.451 0.000	[0.025 49.012 -5.796 -12.887 -21.237 3.957	0.975] 51.986 -1.587 -1.818 9.439 8.598

4. Finally, use the p-value (and a significance level of $\alpha = 0.05$) to make a conclusion about your hypotheses.