

Chapter 9: Multiple and logistic regression

Wen-Han Hwang

(Slides primarily developed by Mine Çetinkaya-Rundel from OpenIntro.)



Institute of Statistics
National Tsing Hua University
Taiwan

Outline

- 1 Introduction to multiple regression
- 2 Model selection
- 3 Checking model conditions using graphs

Introduction to multiple regression

Multiple regression

- Simple linear regression: Bivariate - two variables: Y and X
- Multiple linear regression: Multiple variables: Y and X_1, X_2, X_3, \dots

Poverty vs Region

- Poverty: % poverty in each state
- Region: West vs East
- Region4: northeast, midwest, west, or south

Poverty Analysis Across US States by Region

State	Poverty Rate (%)	Region (West/East)	Region4
California	12.5	West	Region4West
New York	14.2	East	Region4East
Texas	15.6	West	Region4South
Florida	13.3	East	Region4South
...

- This table categorizes states by poverty rates, broader regional division (West vs East), and the specific 'Region4' classification.
- Each state is further analyzed by how it fits into the 'Region4' context, aiding in detailed comparative analysis.

Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: **region**
 - **Reference level:** East (east=0)
 - West=1 if the state belongs to the West region; otherwise, West=0
- **Intercept:** 11.17%
 - Represents the estimated average poverty percentage in Eastern states.
 - This is the value obtained when West=0.
- **Slope:** 0.38%
 - Indicates that the estimated average poverty percentage in Western states is 0.38% higher than in Eastern states.
 - Thus, the estimated average poverty percentage in Western states is $11.17 + 0.38 = 11.55\%$.
 - This is the value obtained when West=1.

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) is the reference level?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

- (a) northeast
- (b) midwest
- (c) west
- (d) south
- (e) cannot tell

- **Midwest (D1):** Equals 1 if the state is in the Midwest, 0 otherwise. Northeast is the reference.
- **West (D2):** Equals 1 if the state is in the West, 0 otherwise. Northeast is the reference.
- **South (D3):** Equals 1 if the state is in the South, 0 otherwise. Northeast is the reference.

Using these dummy variables, any state from the northeast will have the values 0 for all three dummies (D1, D2, D3), serving as the reference category. This approach allows you to compare the effects of being in the midwest, west, or south relative to being in the northeast in your regression model.

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) has the lowest poverty percentage?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

- (a) northeast
- (b) midwest
- (c) west
- (d) south
- (e) cannot tell

Weights of books

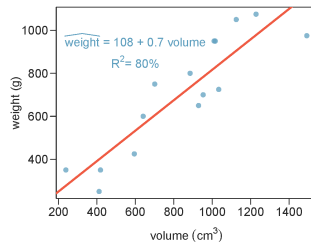
	weight (g)	volume (cm ³)	cover
1	800	885	hc
2	950	1016	hc
3	1050	1125	hc
4	350	239	hc
5	750	701	hc
6	600	641	hc
7	1075	1228	hc
8	250	412	pb
9	700	953	pb
10	650	929	pb
11	975	1492	pb
12	350	419	pb
13	950	1010	pb
14	425	595	pb
15	725	1034	pb



From: Maindonald, J.H. and Braun, W.J. (2nd ed., 2007) "Data Analysis and Graphics Using R"

Weights of books (cont.)

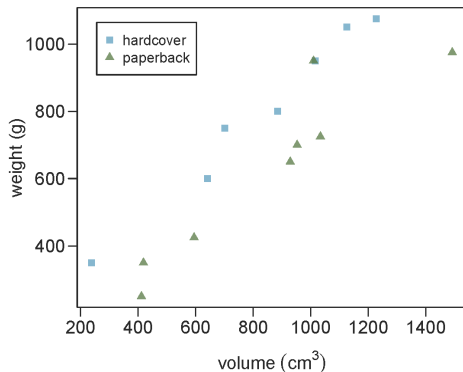
The scatterplot shows the relationship between weights and volumes of books as well as the regression output. Which of the below is correct?



- (a) Weights of 80% of the books can be predicted accurately using this model.
- (b) Books that are 10 cm³ over average are expected to weigh 7 g over average.
- (c) The correlation between weight and volume is $R = 0.80^2 = 0.64$.
- (d) The model underestimates the weight of the book with the highest volume.

Weights of hardcover and paperback books

Can you identify a trend in the relationship between volume and weight of hardcover and paperback books?



Linear model

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	197.96	59.19	3.34	0.01
volume	0.72	0.06	11.67	0.00
cover:pb	-184.05	40.49	-4.55	0.00

$$\widehat{weight} = 197.96 + 0.72 \text{ volume} - 184.05 \text{ cover} : pb$$

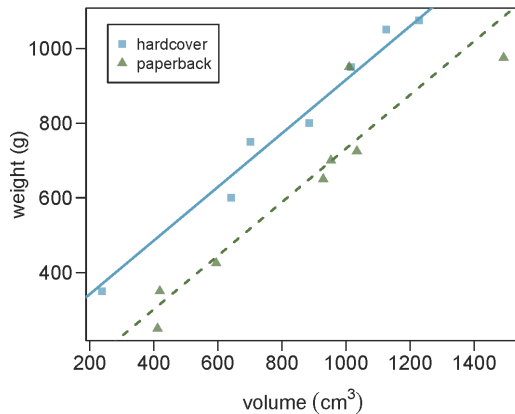
- For **hardcover** books: plug in **0** for cover

$$\begin{aligned}\widehat{weight} &= 197.96 + 0.72 \text{ volume} - 184.05 \times 0 \\ &= 197.96 + 0.72 \text{ volume}\end{aligned}$$

- For **paperback** books: plug in **1** for cover

$$\begin{aligned}\widehat{weight} &= 197.96 + 0.72 \text{ volume} - 184.05 \times 1 \\ &= 13.91 + 0.72 \text{ volume}\end{aligned}$$

Visualising the linear model



Modeling weights of books using volume and cover type

```
lm(weight ~ volume + cover)
```

```
model <- lm(weight ~ volume + cover)
summary(model)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	197.96284	59.19274	3.344	0.005841	**
volume	0.71795	0.06153	11.669	6.6e-08	***
cover:pb	-184.04727	40.49420	-4.545	0.000672	***

Residual standard error: 78.2 on 12 degrees of freedom
Multiple R-squared: 0.9275, Adjusted R-squared: 0.9154
F-statistic: 76.73 on 2 and 12 DF, p-value: 1.455e-07

Interpretation of the regression coefficients

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	197.96	59.19	3.34	0.01
volume	0.72	0.06	11.67	0.00
cover:pb	-184.05	40.49	-4.55	0.00

- **Slope of volume:** All else held constant, books that are 1 more cubic centimeter in volume tend to weigh about 0.72 grams more.
- **Slope of cover:** All else held constant, the model predicts that paperback books weigh 184 grams lower than hardcover books.
- **Intercept:** Hardcover books with no volume are expected on average to weigh 198 grams.
 - Obviously, the intercept does not make sense in context. It only serves to adjust the height of the line.

Prediction

Which of the following is the correct calculation for the predicted weight of a paperback book that is 600 cm³?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	197.96	59.19	3.34	0.01
volume	0.72	0.06	11.67	0.00
cover:pb	-184.05	40.49	-4.55	0.00

- (a) $197.96 + 0.72 * 600 - 184.05 * 1$
- (b) $184.05 + 0.72 * 600 - 197.96 * 1$
- (c) $197.96 + 0.72 * 600 - 184.05 * 0$
- (d) $197.96 + 0.72 * 1 - 184.05 * 600$

Another example: Modeling kid's test scores

Predicting cognitive test scores of three- and four-year-old children using characteristics of their mothers. Data are from a survey of adult American women and their children - a subsample from the National Longitudinal Survey of Youth.

	kid_score	mom_hs	mom_iq	mom_work	mom_age
1	65	yes	121.12	yes	27
⋮					
5	115	yes	92.75	yes	27
6	98	no	107.90	no	18
⋮					
434	70	yes	91.25	yes	25

Gelman, Hill. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. (2007) Cambridge University Press.

• $lm(score \sim hs + iq + work + age)$

Interpreting the slope

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	19.59	9.22	2.13	0.03
mom_hs:yes	5.09	2.31	2.20	0.03
mom_iq	0.56	0.06	9.26	0.00
mom_work:yes	2.54	2.35	1.08	0.28
mom_age	0.22	0.33	0.66	0.51

All else held constant, kids with mothers whose IQs are one point higher tend to score on average 0.56 points higher.

Kids whose moms haven't gone to HS, did not work during the first three years of the kid's life, have an IQ of 0 and are 0 yrs old are expected on average to score 19.59. Obviously, the intercept does not make any sense in context.

Interpreting the slope

What is the correct interpretation of the slope for mom_work?

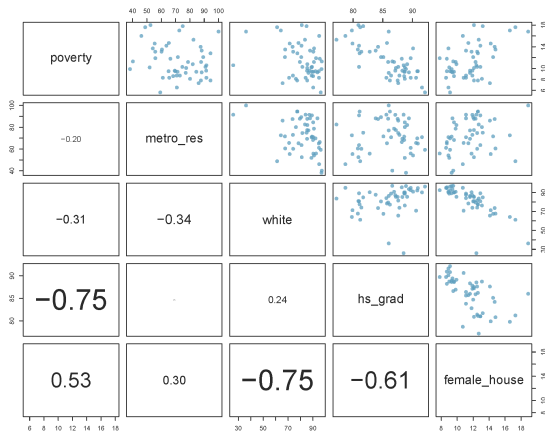
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	19.59	9.22	2.13	0.03
mom_hs:yes	5.09	2.31	2.20	0.03
mom_iq	0.56	0.06	9.26	0.00
mom_work:yes	2.54	2.35	1.08	0.28
mom_age	0.22	0.33	0.66	0.51

All else being equal, kids whose moms worked during the first three years of the kid's life

- Ⓐ are estimated to score 2.54 points lower
- Ⓑ are estimated to score 2.54 points higher

than those whose moms did not work.

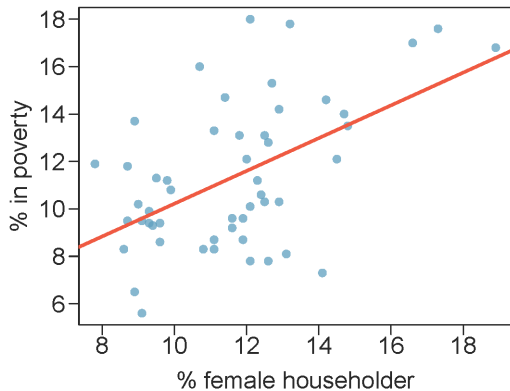
Modeling poverty



- **metro_res (%)**: The percentage of people living in urban areas. Higher values suggest more urban populations.
- **white (%)**: The percentage of people who are white in the population. Higher percentages indicate a larger white demographic.
- **hs_grad (%)**: The percentage of people who have graduated from high school. Higher values show a more educated population.
- **female_house (%)**: The percentage of households led by women. Higher percentages reflect more female-headed households.

Predicting poverty using % female householder

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	3.31	1.90	1.74	0.09
female_house	0.69	0.16	4.32	0.00



$$R = 0.53$$

$$R^2 = 0.53^2 = 0.28$$

Another look at R^2

R^2 can be calculated in three ways:

- 1 square the correlation coefficient of x and y (how we have been calculating it)
- 2 square the correlation coefficient of y and \hat{y}
- 3 based on definition:

$$R^2 = \frac{\text{explained variability in } y}{\text{total variability in } y}$$

Using [ANOVA](#) we can calculate the explained variability and total variability in y .

Sum of squares

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
female_house	1	132.57	132.57	18.68	0.00
Residuals	49	347.68	7.10		
Total	50	480.25			

Sum of squares of y : $SS_{Total} = \sum (y - \bar{y})^2 = 480.25 \rightarrow$ total variability

Sum of squares of residuals: $SS_{Error} = \sum e_i^2 = 347.68 \rightarrow$ unexplained variability

Sum of squares of x : $SS_{Model} = SS_{Total} - SS_{Error} \rightarrow$ explained variability
 $= 480.25 - 347.68 = 132.57$

$$R^2 = \frac{\text{explained variability}}{\text{total variability}} = \frac{132.57}{480.25} = 0.28 \checkmark$$

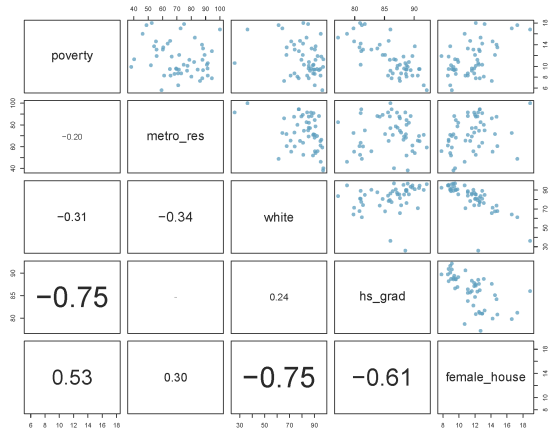
Predicting poverty using % female hh + % white

Linear model:	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.58	5.78	-0.45	0.66
female_house	0.89	0.24	3.67	0.00
white	0.04	0.04	1.08	0.29

ANOVA:	Df	Sum Sq	Mean Sq	F value	Pr(>F)
female_house	1	132.57	132.57	18.74	0.00
white	1	8.21	8.21	1.16	0.29
Residuals	48	339.47	7.07		
Total	50	480.25			

$$R^2 = \frac{\text{explained variability}}{\text{total variability}} = \frac{132.57 + 8.21}{480.25} = 0.29$$

Does adding the variable white to the model add valuable information that wasn't provided by female_house?



Collinearity between explanatory variables

poverty vs. % female head of household

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	3.31	1.90	1.74	0.09
female_house	0.69	0.16	4.32	0.00

poverty vs. % female head of household and % female hh

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.58	5.78	-0.45	0.66
female_house	0.89	0.24	3.67	0.00
white	0.04	0.04	1.08	0.29

Collinearity between explanatory variables (cont.)

- Two predictor variables exhibit **collinearity** when they are correlated, which can complicate the estimation of models.
Remember: Predictors are also referred to as explanatory or *independent* variables, but this does not imply independence among these variables.
- Adding predictors that are highly correlated with each other to a model is generally avoided because it often contributes little additional value. Instead, a **parsimonious** model, which is the simplest effective model, is preferred.
- It is challenging to completely avoid collinearity in observational data. In contrast, experimental data can be designed to minimize correlation among predictors.

R^2 vs. adjusted R^2

	R^2	Adjusted R^2
Model 1 (Single-predictor)	0.28	0.26
Model 2 (Multiple)	0.29	0.26

- When any variable is added to the model R^2 increases.
- But if the added variable doesn't really provide any new information, or is completely unrelated, adjusted R^2 does not increase.

Adjusted R^2

Adjusted R^2

$$R_{adj}^2 = 1 - \left(\frac{SS_{Error}}{SS_{Total}} \times \frac{n-1}{n-p-1} \right)$$

where n is the number of cases and p is the number of predictors (explanatory variables) in the model.

- Because p is never negative, R_{adj}^2 will always be smaller than R^2 .
- R_{adj}^2 applies a penalty for the number of predictors included in the model.
- Therefore, we choose models with higher R_{adj}^2 over others.

Calculate adjusted R^2

ANOVA:	Df	Sum Sq	Mean Sq	F value	Pr(>F)
female_house	1	132.57	132.57	18.74	0.0001
white	1	8.21	8.21	1.16	0.2868
Residuals	48	339.47	7.07		
Total	50	480.25			

$$\begin{aligned} R_{adj}^2 &= 1 - \left(\frac{SS_{Error}}{SS_{Total}} \times \frac{n-1}{n-p-1} \right) \\ &= 1 - \left(\frac{MS_{Error}}{MS_{Total}} \right) \\ &= 1 - \left(\frac{7.07}{480.25/50} \right) \\ &= 0.26 \end{aligned}$$

Model selection

Beauty in the classroom

- Data: Student evaluations of instructors' beauty and teaching quality for 463 courses at the University of Texas.
- Evaluations conducted at the end of semester, and the beauty judgements were made later, by six students who had not attended the classes and were not aware of the course evaluations (2 upper level females, 2 upper level males, one lower level female, one lower level male).

Hamermesh & Parker. (2004) Beauty in the classroom: instructors' pulchritude and putative pedagogical productivity? Economics Education Review.

Variables in Instructor Evaluation Model

- **Response Variable:**

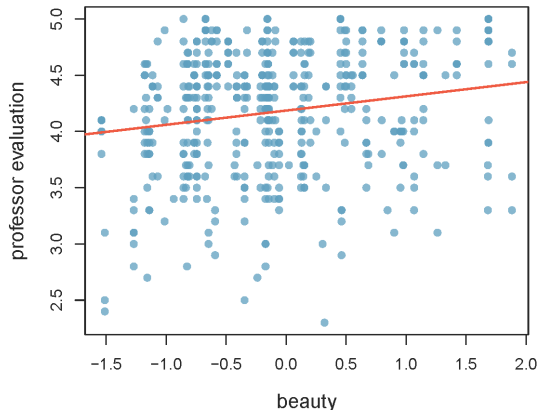
- **Score** - Measures the teaching quality of the instructors as evaluated by students.

- **Predictor Variables:**

- **Beauty** - Attractiveness rating of the instructor.
- **Gender.Male** - Binary indicator (1 if male, 0 otherwise).
- **Age** - Age of the instructor.
- **Formal.Yes** - Binary indicator of whether the instructor's photo shows them in formal attire (1 for yes, 0 for no).
- **Lower.Yes** - Binary indicator of whether the course is a lower division course (1 for yes, 0 for no).
- **Native.Non English** - Binary indicator (1 if the instructor is a non-native English speaker, 0 otherwise).
- **Minority.Yes** - Binary indicator (1 if the instructor belongs to a minority group, 0 otherwise).
- **Students** - Number of students enrolled in the course.
- **Tenure** - Categorical variable with levels: Non-tenure track, Tenure track, Tenured.

Professor rating vs. beauty

Professor evaluation score (higher score means better) vs. beauty score (a score of 0 means average, negative score means below average, and a positive score above average):



Which of the below is correct based on the model output?

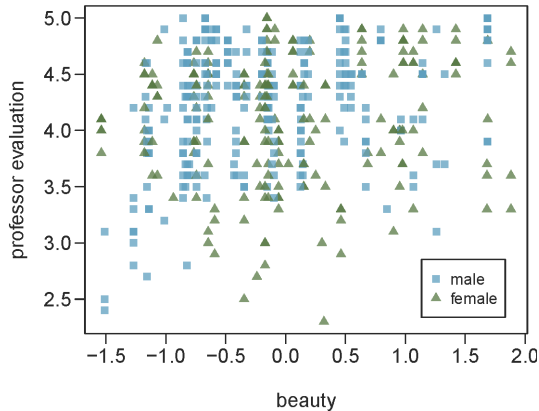
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.19	0.03	167.24	0.00
beauty	0.13	0.03	4.00	0.00

$R^2 = 0.0336$

- Ⓐ Model predicts 3.36% of professor ratings correctly.
- Ⓑ Beauty is not a significant predictor of professor evaluation.
- Ⓒ Professors who score 1 point above average in their beauty score are tend to also score 0.13 points higher in their evaluation.
- Ⓓ 3.36% of variability in beauty scores can be explained by professor evaluation.
- Ⓔ The correlation coefficient could be $\sqrt{0.0336} = 0.18$ or -0.18 , where the sign cannot be sure.

Exploratory analysis

For a given beauty score, are male professors evaluated higher, lower, or about the same as female professors?



Professor rating vs. beauty + gender

For a given beauty score, are male professors evaluated higher, lower, or about the same as female professors?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.09	0.04	107.85	0.00
beauty	0.14	0.03	4.44	0.00
gender.male	0.17	0.05	3.38	0.00

$$R_{adj}^2 = 0.057$$

- (a) higher
- (b) lower
- (c) about the same

Full model

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.6282	0.1720	26.90	0.00
beauty	0.1080	0.0329	3.28	0.00
gender.male	0.2040	0.0528	3.87	0.00
age	-0.0089	0.0032	-2.75	0.01
formal.yes ¹	0.1511	0.0749	2.02	0.04
lower.yes ²	0.0582	0.0553	1.05	0.29
native.non english	-0.2158	0.1147	-1.88	0.06
minority.yes	-0.0707	0.0763	-0.93	0.35
students ³	-0.0004	0.0004	-1.03	0.30
tenure.tenure track ⁴	-0.1933	0.0847	-2.28	0.02
tenure.tenured	-0.1574	0.0656	-2.40	0.02

- tenure is a categorical variable: non-tenure, tenure track, tenured.

Hypotheses

Just as the interpretation of the slope parameters take into account all other variables in the model, the hypotheses for testing for significance of a predictor also takes into account all other variables.

$H_0: B_i = 0$ when other explanatory variables are included in the model.

$H_A: B_i \neq 0$ when other explanatory variables are included in the model.

Assessing significance: numerical variables

The p-value for age is 0.01. What does this indicate?

	Estimate	Std. Error	t value	Pr(> t)
...				
age	-0.0089	0.0032	-2.75	0.01
...				

- (a) Since p-value is positive, higher the professor's age, the higher we would expect them to be rated.
- (b) If we keep all other variables in the model, there is strong evidence that professor's age is associated with their rating.
- (c) Probability that the true slope parameter for age is 0 is 0.01.
- (d) There is about 1% chance that the true slope parameter for age is -0.0089.

Assessing significance

Which predictors do not seem to meaningfully contribute to the model, i.e. may not be significant predictors of professor's rating score?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.6282	0.1720	26.90	0.00
beauty	0.1080	0.0329	3.28	0.00
gender.male	0.2040	0.0528	3.87	0.00
age	-0.0089	0.0032	-2.75	0.01
formal.yes	0.1511	0.0749	2.02	0.04
lower.yes	0.0582	0.0553	1.05	0.29
native.non english	-0.2158	0.1147	-1.88	0.06
minority.yes	-0.0707	0.0763	-0.93	0.35
students	-0.0004	0.0004	-1.03	0.30
tenure.tenure track	-0.1933	0.0847	-2.28	0.02
tenure.tenured	-0.1574	0.0656	-2.40	0.02

Model selection strategies

Based on what we've learned so far, what are some ways you can think of that can be used to determine which variables to keep in the model and which to leave out?

Backward-elimination

- 1 Start with the full model
- 2 Drop one variable at a time and record R_{adj}^2 of each smaller model
- 3 Pick the model with the highest increase in R_{adj}^2
- 4 Repeat until none of the models yield an increase in R_{adj}^2

Backward-elimination

Full	beauty + gender + age + formal + lower + native + minority + students + tenure	0.0839
Step 1	gender + age + formal + lower + native + minority + students + tenure	0.0642
	beauty + age + formal + lower + native + minority + students + tenure	0.0557
	beauty + gender + formal + lower + native + minority + students + tenure	0.0706
	beauty + gender + age + lower + native + minority + students + tenure	0.0777
	beauty + gender + age + formal + native + minority + students + tenure	0.0837
	beauty + gender + age + formal + lower + minority + students + tenure	0.0788
	beauty + gender + age + formal + lower + native + students + tenure	0.0842
	beauty + gender + age + formal + lower + native + minority + tenure	0.0838
Step 2	beauty + gender + age + formal + lower + native + minority + students	0.0733
	gender + age + formal + lower + native + students + tenure	0.0647
	beauty + age + formal + lower + native + students + tenure	0.0543
	beauty + gender + formal + lower + native + students + tenure	0.0708
	beauty + gender + age + lower + native + students + tenure	0.0776
	beauty + gender + age + formal + native + students + tenure	0.0846
	beauty + gender + age + formal + lower + native + tenure	0.0844
	beauty + gender + age + formal + lower + native + students	0.0725
Step 3	gender + age + formal + native + students + tenure	0.0653
	beauty + age + formal + native + students + tenure	0.0534
	beauty + gender + formal + native + students + tenure	0.0707
	beauty + gender + age + native + students + tenure	0.0786
	beauty + gender + age + formal + students + tenure	0.0756
	beauty + gender + age + formal + native + tenure	0.0855
	beauty + gender + age + formal + native + students	0.0713
Step 4	gender + age + formal + native + tenure	0.0667
	beauty + age + formal + native + tenure	0.0553
	beauty + gender + formal + native + tenure	0.0723
	beauty + gender + age + native + tenure	0.0806
	beauty + gender + age + formal + tenure	0.0773
	beauty + gender + age + formal + native	0.0713

step function in R

Best model: beauty + gender + age + formal + native + tenure

Call:

```
lm(formula = profevaluation ~ beauty + gender + age + formal +  
    native + tenure, data = d)
```

Coefficients:

(Intercept)	beauty	gendermale
4.628435	0.105546	0.208079
age	formalyes	nativenon english
-0.008844	0.132422	-0.243003
tenuretenure track	tenuretenured	
-0.206784	-0.175967	

Forward selection

- ① Start with regressions of response vs. each explanatory variable
- ② Pick the model with the highest R^2_{adj}
- ③ Add the remaining variables one at a time to the existing model, and once again pick the model with the highest R^2_{adj}
- ④ Repeat until the addition of any of the remaining variables does not result in a higher R^2_{adj}

- Backward elimination with the p-value approach:
 - ① Start with the full model
 - ② Drop the variable with the highest p-value and refit a smaller model
 - ③ Repeat until all variables left in the model are significant
- Forward selection with the p-value approach:
 - ① Start with regressions of response vs. each explanatory variable
 - ② Pick the variable with the lowest significant p-value
 - ③ Add the remaining variables one at a time to the existing model, and pick the variable with the lowest significant p-value
 - ④ Repeat until any of the remaining variables does not have a significant p-value

Backward-elimination: p – value approach

Step	Variables included & p-value									
Full	beauty 0.00	gender male 0.00	age 0.01	formal yes 0.04	lower yes 0.29	native nonenglish 0.06	minority yes 0.35	students 0.30	tenure tenure track 0.02	tenure tenured 0.02
Step 1	beauty 0.00	gender male 0.00	age 0.01	formal yes 0.04	lower yes 0.38	native nonenglish 0.03		students 0.34	tenure tenure track 0.02	tenure tenured 0.01
Step 2	beauty 0.00	gender male 0.00	age 0.01	formal yes 0.05		native nonenglish 0.02		students 0.44	tenure tenure track 0.01	tenure tenured 0.01
Step 3	beauty 0.00	gender male 0.00	age 0.01	formal yes 0.06		native nonenglish 0.02			tenure tenure track 0.01	tenure tenured 0.01
Step 4	beauty 0.00	gender male 0.00	age 0.01			native nonenglish 0.06			tenure tenure track 0.01	tenure tenured 0.01
Step 5	beauty 0.00	gender male 0.00	age 0.01						tenure tenure track 0.01	tenure tenured 0.01

Best model: beauty + gender + age + tenure

Adjusted R^2 vs. p-value approaches

- The two approaches are similar, but they sometimes lead to different models, with the adjusted R^2 approach tending to include more predictors in the final model.
- When we care about understanding which variables are statistically significant predictors of the response, or if there is interest in producing a simpler model at the potential cost of a little prediction accuracy, then the p-value approach is preferred.
- Regardless of the approach we use, our job is not done after variable selection – we must still verify the model conditions are reasonable.

step function in R

The step function in R does similar backward elimination process, however it uses a different metric called AIC (Akaike Information Criterion) instead of adjusted R^2 to do the model selection.

- $\text{AIC} = n \log\left(\frac{\text{RSS}}{n}\right) + 2 |\beta|$

- $\frac{\text{RSS}}{n} = \hat{\sigma}^2$

Checking model conditions using graphs

Modeling conditions

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$$

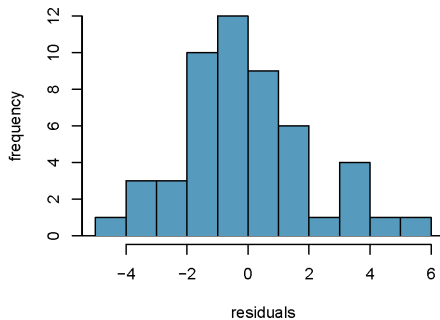
The model depends on the following conditions

- ① residuals are nearly normal (less important for larger data sets)
- ② residuals have constant variability
- ③ residuals are independent
- ④ each variable is linearly related to the outcome

We often use graphical methods to check the validity of these conditions, which we will go through in detail in the following slides.

(1) nearly normal residuals

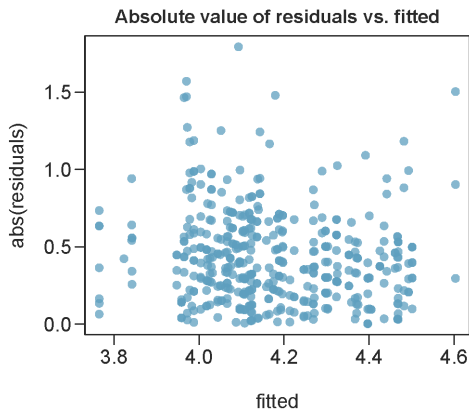
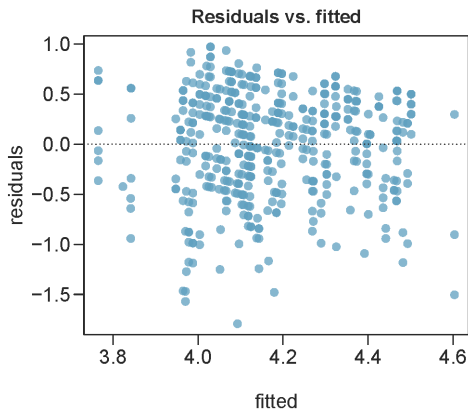
normal probability plot and/or histogram of residuals:



Does this condition appear to be satisfied?

(2) constant variability in residuals

scatterplot of residuals and/or absolute value of residuals vs. fitted (predicted):



Does this condition appear to be satisfied?

Checking constant variance - recap

- When we did simple linear regression (one explanatory variable) we checked the constant variance condition using a plot of **residuals vs. x** .
- With multiple linear regression (2+ explanatory variables) we checked the constant variance condition using a plot of **residuals vs. fitted**.

Why are we using different plots?

In multiple linear regression there are many explanatory variables, so a plot of residuals vs. one of them wouldn't give us the complete picture.

(3) independent residuals

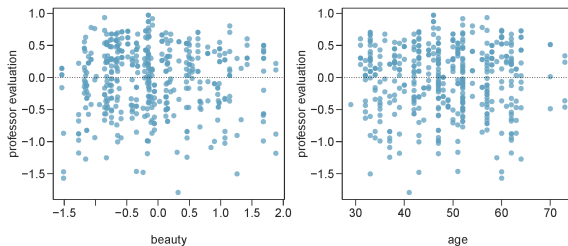
scatterplot of residuals vs. order of data collection:



Does this condition appear to be satisfied?

(4) linear relationships

scatterplot of residuals vs. each (numerical) explanatory variable:



Does this condition appear to be satisfied?

Note: We use residuals instead of the predictors on the y-axis so that we can still check for linearity without worrying about other possible violations like collinearity between the predictors.

Several options for improving a model

- Transforming variables
- Seeking out additional variables to fill model gaps
- Using more advanced methods that would account for challenges around inconsistent variability or nonlinear relationships between predictors and the outcome

Transformations

If the concern with the model is non-linear relationships between the explanatory variable(s) and the response variable, transforming the response variable can be helpful.

- Log transformation ($\log y$)
- Square root transformation (\sqrt{y})
- Inverse transformation ($1/y$)
- Truncation (cap the max value possible)

It is also possible to apply transformations to the explanatory variable(s), however such transformations tend to make the model coefficients even harder to interpret.

Models can be wrong, but useful

All models are wrong, but some are useful. - George Box

- No model is perfect, but even imperfect models can be useful, as long as we are clear and report the model's shortcomings.
- If conditions are grossly violated, we should not report the model results, but instead consider a new model, even if it means learning more statistical methods or hiring someone who can help.