

Applied Statistical Analysis

EDUC 6050

Week 13

Finding clarity using data

Today

Categorical Outcomes

Categorical Outcomes

For simple research questions
Not controlling for other factors
Doesn't provide a lot of information
(ie., only tells us difference or not)



General Requirements

1. One or more categorical variables

Goodness of Fit

Test of Independence

ID	X	Y
1	0	0
2	2	1
3	1	0
4	2	1
5	0	1
6	0	1
7	2	0
8	1	0

Hypothesis Testing with Chi Square (Independence)

The same 6 step approach!

1. Examine Variables to Assess Statistical Assumptions
2. State the Null and Research Hypotheses (symbolically and verbally)
3. Define Critical Regions
4. Compute the Test Statistic
5. Compute an Effect Size and Describe it
6. Interpreting the results

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
2. Appropriate measurement of variables for the analysis
3. Expected frequency 5+

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data

2. Appropriate
for the a

3. Expected

Individuals are independent of each other (one person's scores does not affect another's)

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
- 2. Appropriate measurement of variables for the analysis**
3. Expected frequency 5+



Here we need interval/ratio outcome

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence
2. Appropriate test for the analysis



Variance around the line should be roughly equal across the whole line

3. Expected frequency 5+

1 Examine Variables to Assess Statistical Assumptions

Examining the Basic Assumptions

1. **Independence:** random sample
2. **Appropriate measurement:** know what your variables are
3. **Expected frequency 5+:** Check expected frequencies

2

State the Null and Research Hypotheses (symbolically and verbally)

Hypothesis Type	Symbolic	Verbal	Difference between means created by:
Research Hypothesis	$OF \neq EF$	Observed frequency is not equal to expected frequency	True relationship
Null Hypothesis	$OF = EF$	Observed frequency is the same as the expected frequency	Random chance (sampling error)

3 Define Critical Regions

How much evidence is enough to believe the null is not true?

generally based on an $\alpha = .05$

Use software's p-value to judge if it is below .05

4

Compute the Test Statistic

Jamovi Tutorial

5

Compute an Effect Size and Describe it

"Phi" → $\phi = \sqrt{\frac{\chi^2}{n}}$ *Cramer's* $\phi = \sqrt{\frac{\chi^2}{n(df)}}$

ϕ	<i>Cramer's</i> ϕ	Estimated Size of the Effect
Close to .1	Depends	Small
Close to .3	on df	Moderate
Close to .5	(pg 557)	Large

6

Interpreting the results

“The voters’ opinions of the president’s policies were associated with the voters’ political affiliations, $\chi^2(2, N = 58) = 16.40$, $p = .02$, $\phi = .53$. More democrats and fewer republicans approved of the president’s policies than would be expected by chance.” – pg 577.


Logistic Regression

Intro to Logistic Regression

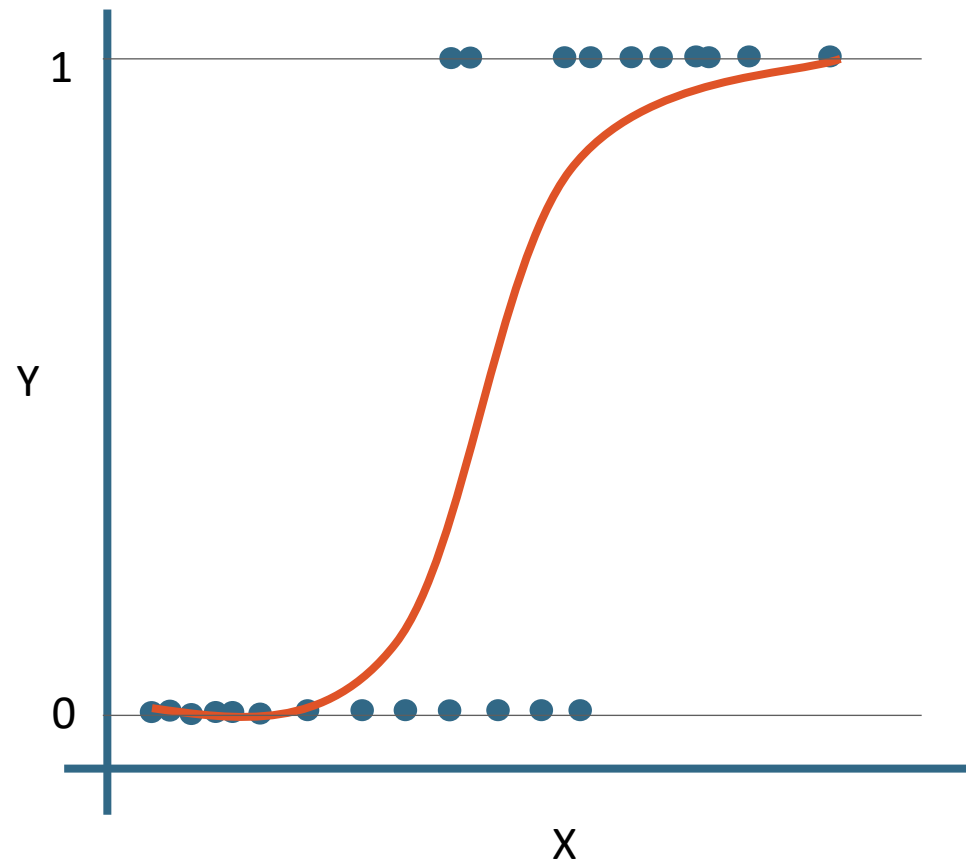
So far, we have always wanted continuous outcome variables

But what if our outcome is a categorical variable??

Logistic Regression is just like linear regression but works with binary (dichotomous) outcomes

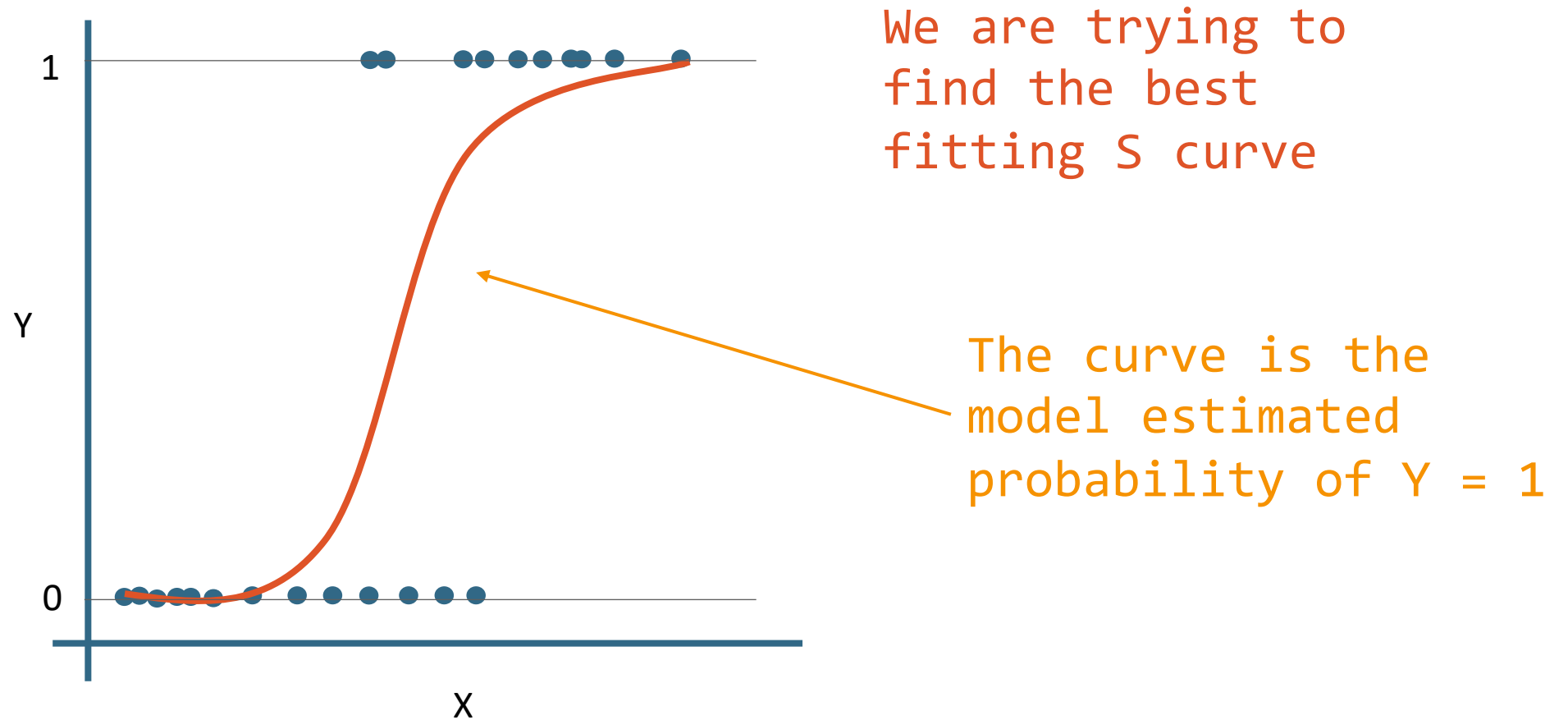
- Substance Use or Not
 - Cancer or Not
 - Buy it or Not
- 

Logic of Logistic Regression



We are trying to find the best fitting S curve

Logic of Logistic Regression



Logistic Regression

Simple

- Only one predictor in the model
- Tells you if that one predictor is associated with the odds of $Y = 1$

Multiple

- More than one variable in the model
- Tells you if, while holding the other variables constant, if that predictor is associated with the odds of $Y = 1$

Logistic Regression

- Logistic does what regression does but with a little bit of **mathematical magic**



$$\boxed{\mathit{logit}(Y)} = \beta_0 + \beta_1 X + \epsilon$$

Logistic Regression

- Logistic does what regression does but with a little bit of **mathematical magic**



$$\boxed{\mathit{logit}(Y)} = \beta_0 + \beta_1 X + \epsilon$$

intercept

slope

Logistic Regression

$$\mathit{logit}(Y) = \beta_0 + \beta_1 X + \epsilon$$

Example

We have two variables, X and Y . X is continuous, Y is binary. We want to know if increases/decreases in X are associated (or predict) changes in the chance of Y equaling 1.

Logistic Regression

- It is trying to predict the outcome accurately using the information from the predictor
- Better prediction tells us that the predictor(s) is/are more strongly related to the outcome

General Requirements

1. Two or more variables,
2. Outcome needs to be binary
3. Others can be continuous or categorical

ID	X	Y
1	8	0
2	6	1
3	9	1
4	7	1
5	7	0
6	8	0
7	5	1
8	5	0

Hypothesis Testing with Logistic Regression

The same 6 step approach!

1. Examine Variables to Assess Statistical Assumptions
2. State the Null and Research Hypotheses (symbolically and verbally)
3. Define Critical Regions
4. Compute the Test Statistic
5. Compute an Effect Size and Describe it
6. Interpreting the results

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
2. Appropriate measurement of variables for the analysis
3. Normality of distributions
4. Homoscedastic

1

Examine Variables to Assess Statistical Assumptions


Basic Assumptions

1. Independence of data

2. Appropriate for the analysis

3. Normality

4. Homoscedastic



Individuals are independent of each other (one person's scores does not affect another's)

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
- 2. Appropriate measurement of variables for the analysis**
3. Normality of distributions
4. Homoscedasticity



Here we need nominal outcome

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence
 2. Appropriateness for the analysis
 3. Normality of distributions
 4. Homoscedastic
- Residuals should be normally distributed

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data

2. Appropriateness for the test

3. Normality

4. Homoscedastic

Variance around the line should be roughly equal across the whole line

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
2. Appropriate measurement of variables for the analysis
3. Normality of distributions
4. Homoscedastic
- 5. Logistic Relationship**
- 6. No omitted variables**

1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data

2. Appropriateness of the model for the data

3. Normality of distributions

4. Homoscedastic

5. Logistic Relationships

6. No omitted variables

The “S-shaped” curve should fit to the data



1

Examine Variables to Assess Statistical Assumptions

Basic Assumptions

1. Independence of data
2. Appropriate measurement of variables for the
3. Normality
4. Homoscedasticity
5. Logistic

Any variable that is related to both the predictor and the outcome should be included in the regression model

6. No omitted variables

1

Examine Variables to Assess Statistical Assumptions

Examining the Basic Assumptions

1. **Independence:** random sample
2. **Appropriate measurement:** know what your variables are
3. **Normality:** Histograms, Q-Q, skew and kurtosis
4. **Homoscedastic:** Scatterplots
5. **Logistic:** Scatterplots
6. **No Omitted:** check correlations, know the theory

2

State the Null and Research Hypotheses (symbolically and verbally)

Hypothesis Type	Symbolic	Verbal	Difference between means created by:
Research Hypothesis	$\beta \neq 0$	X predicts Y	True relationship
Null Hypothesis	$\beta = 0$	There is no <i>real</i> relationship.	Random chance (sampling error)

3 Define Critical Regions

How much evidence is enough to believe the null is not true?

generally based on an $\alpha = .05$

Use software's p-value to judge if it is below .05

4

Compute the Test Statistic

The screenshot shows a software interface with a menu bar containing 'Data' and 'Analyses'. Below the menu bar are icons for 'Exploration', 'T-Tests', 'ANOVA', 'Regression', 'Frequencies', and 'Factor'. A data table is visible with columns 'nam', 'prod1', 's', 'marr', and 'gend'. A dropdown menu is open over the table, listing options: 'Correlation Matrix', 'Linear Regression', 'Logistic Regression', '2 Outcomes' (with 'Binomial' as a sub-option), and 'N Outcomes' (with 'Multinomial' as a sub-option). An orange arrow points from a text box to the '2 Outcomes Binomial' option.

	nam	prod1	s	marr	gend
1	Michael	2	8	0	
2	Pam	3	7	1	
3	Jim	3	8	1	
4	Dwight	5	8	0	
5	Stanley	4	4	1	
6	Phyllis	4	4	1	
7	Creed	1	4	0	
8	Meredith	3	5	4	0
9	Oscar	5	7	7	0
10	Angela	4	5	7	0
11	Kevin	2	6	2	0
12	Stella	3	5	5	0
13	Ryan	2	2	5	0
14	Toby	4	1	6	0
15	Andy	3	5	7	0
16	Jan	4	6	6	1
17	April	1	6	4	1
18	Andy	1	2	2	1
19	Leslie	5	8	7	0
20	Ron	3	8	7	0
21	Tom	2	5	5	0
22	Donna	2	7	6	0
23	Ben	5	8	5	0
24	Chris	4	6	8	0
25	Gary (Larry, ...)	3	5	3	1
26	Jean Ralphio	1	1	2	0
27	Mona Lisa	1	1	1	0
28	Ann	5	8	8	0
29	Kyle	3	5	2	1

Click on
"2 Outcomes
Binomial"

4

Compute the Test Statistic

The screenshot shows the SPSS interface for a Binomial Logistic Regression analysis. The dialog box on the left shows the dependent variable 'subs' and the covariate 'inco'. The results window on the right displays model fit measures and coefficients.

Binomial Logistic Regression

Dependent Variable: subs

Covariates: inco

Binomial Logistic Regression Results

Model Fit Measures

Model	Deviance	AIC	R ² _{McF}
1	29.8	33.8	0.238

Model Coefficients

Predictor	Estimate	SE	Z	p
Intercept	2.1381	1.3809	1.55	0.122
inco	-0.0805	0.0333	-2.42	0.016

Note. Estimates represent the log odds of "subs = 1" vs. "subs = 0"

Outcome goes here

Continuous predictors go here

Categorical predictors go here

Other model options

4

Continuous Predictor

Model Coefficients

Predictor	Estimate	SE	Z	p	Odds ratio	95% Confidence Interval	
						Lower	Upper
Intercept	2.1381	1.3809	1.55	0.122	8.483	0.566	127.060
Income	-0.0805	0.0333	-2.42	0.016	0.923	0.864	0.985

Note. Estimates represent the log odds of "subs = 1" vs. "subs = 0"

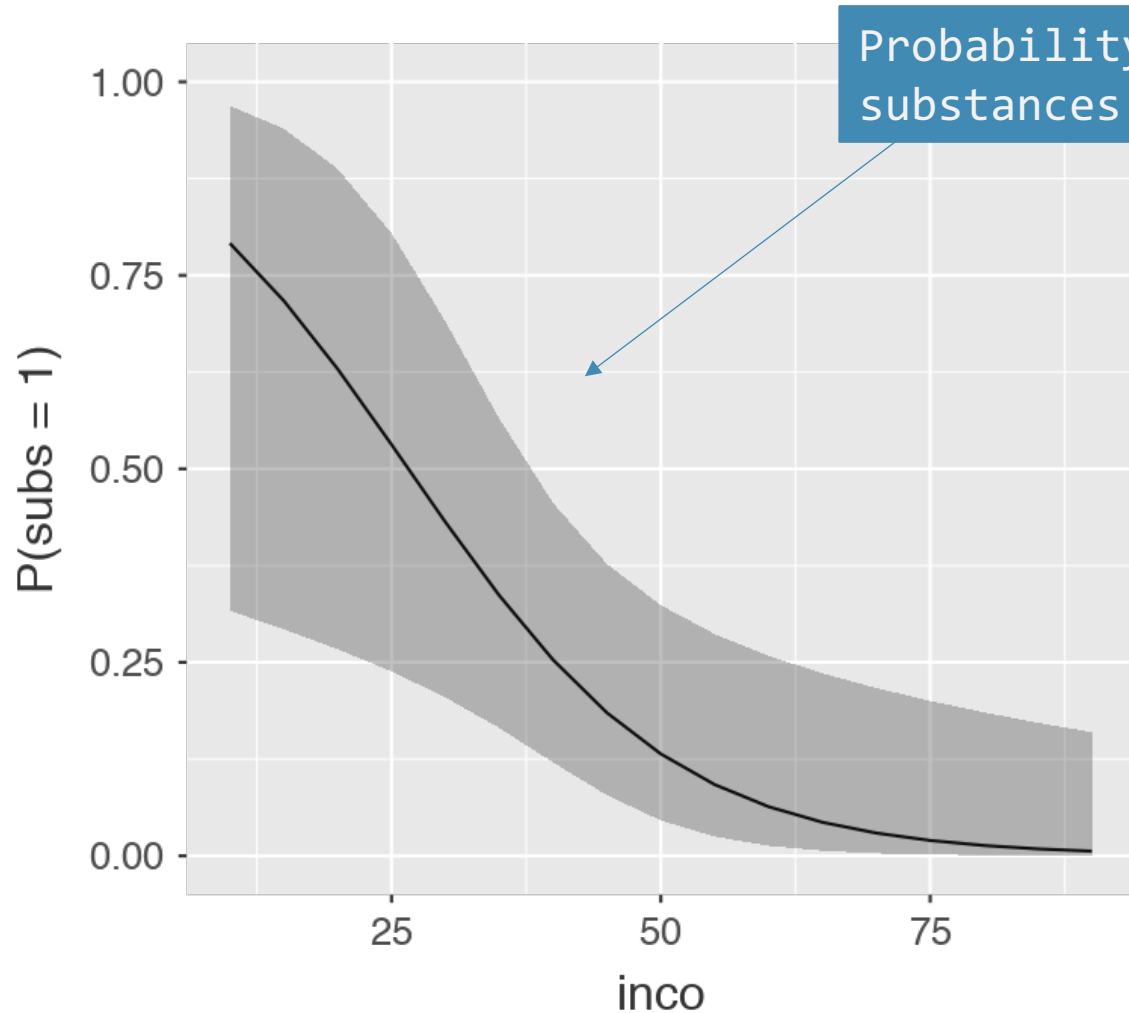
Estimate in "log-odds" units

Significant

The odds ratio is below 1 so as income increases, the odds of using substances decreases by $\sim 1 - .923 = .077$ (7.7% decrease)

4

Continuous Predictor



How well can we predict substance use with just income?

Classification Table – subs

Observed	Predicted		% Correct
	0	1	
0	29	1	96.7
1	5	3	37.5

Note. The cut-off value is set to 0.5

4

Categorical Predictor

Model Coefficients

Predictor	Estimate	SE	Z	p	Odds ratio	95% Confidence Interval	
						Lower	Upper
Intercept	-1.504	0.553	-2.721	0.007	0.222	0.0752	0.657
The Office – Parks and Rec	0.405	0.799	0.507	0.612	1.500	0.3131	7.186

Show:

The Office –
Parks and Rec

0.405

0.799

0.507

0.612

1.500

0.3131

7.186

Note: Estimates represent the log odds of "subs = 1" vs. "subs = 0"

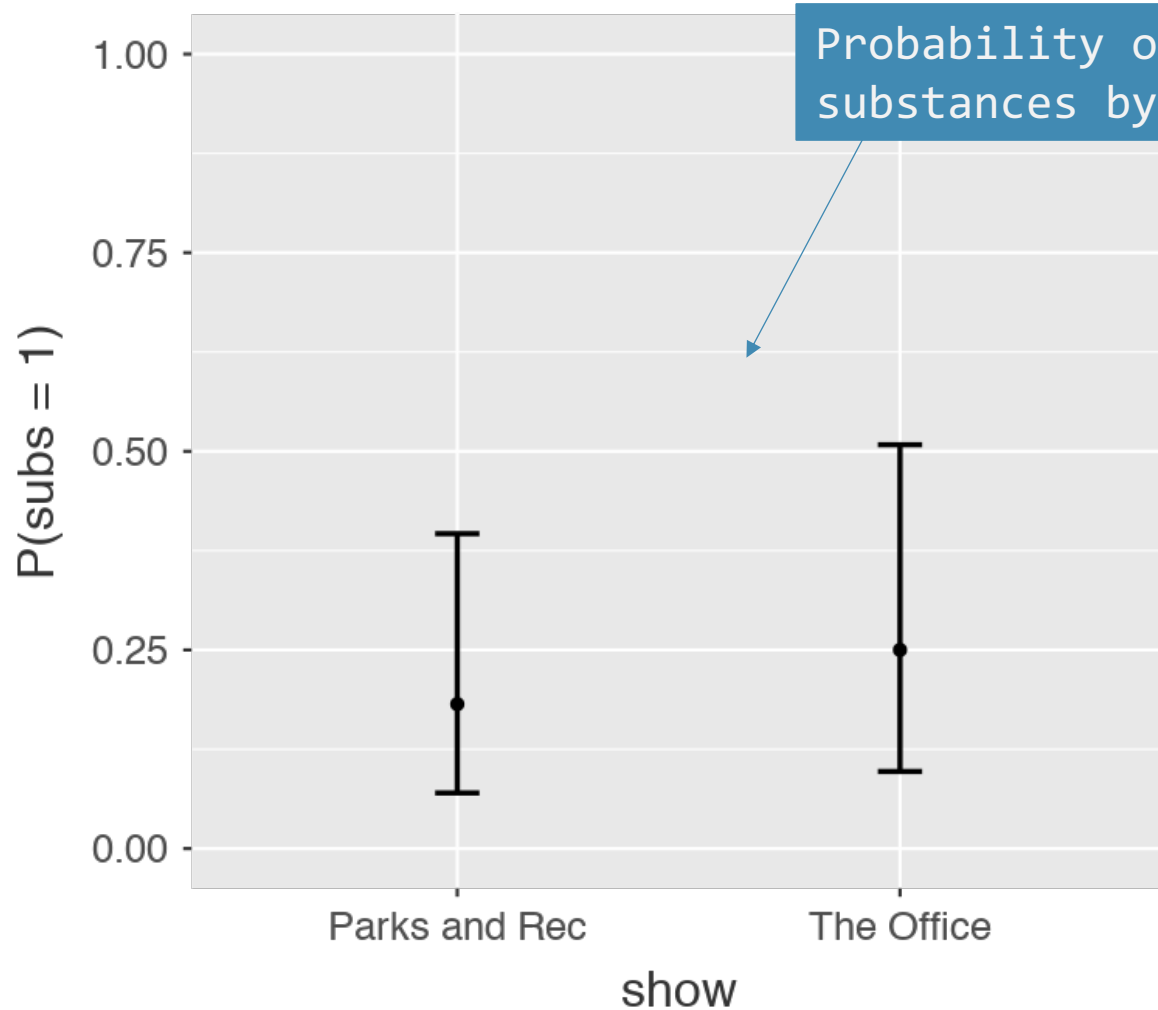
Estimate in "log-odds" units

Not Significant

The odds ratio is above 1 so individuals on The Office have an odds of using substances 50% ($1.5 - 1 = .5 = 50\%$) higher than PR

4

Categorical Predictor



How well can we predict substance use with just income?

Classification Table – subs

Observed	Predicted		% Correct
	0	1	
0	30	0	100
1	8	0	0.00

Note. The cut-off value is set to 0.5

5

Compute an Effect Size and Describe it

One of the main effect sizes for regression is R^2

$$\text{Odds Ratio} = \frac{\text{Odds of } Y \text{ when } X \text{ is one unit higher}}{\text{Odds of } Y \text{ when } X \text{ is not one unit higher}}$$

6

Interpreting the results

The logistic regression analysis showed that income significantly predicted the odds of substance use (OR = .923, $p = .016$). As income increased by \$1000, the odds of using substances decreased by 7.7%.

Multiple Logistic Regression

Multiple Logistic Regression

More than one predictor in the same model

This change the interpretation just a little:

Slope is now the change in the odds of $Y = 1$ for a one unit change in X , while holding the other predictors constant.

Multiple Regression

Provides us with a few more things to think about

1. Variable Selection
2. Assumption Checks (much more difficult in logistic regression)
3. Multi-collinearity
4. Interactions

Variable Selection

Several Approaches

1. Forward
2. Backward
3. Lasso
4. Covariates then predictor of interest

Variable Selection when theory isn't clear

Several Approaches

1. Forward
2. Backward
3. Lasso
4. Covariates then predictor of interest

I'd recommend these two

Assumption Checks

Difficult (we won't cover it in this class)

Jamovi doesn't provide many checks (only collinearity)

Multi-Collinearity

When two or more predictors are very related to each other or are linear combinations of each other

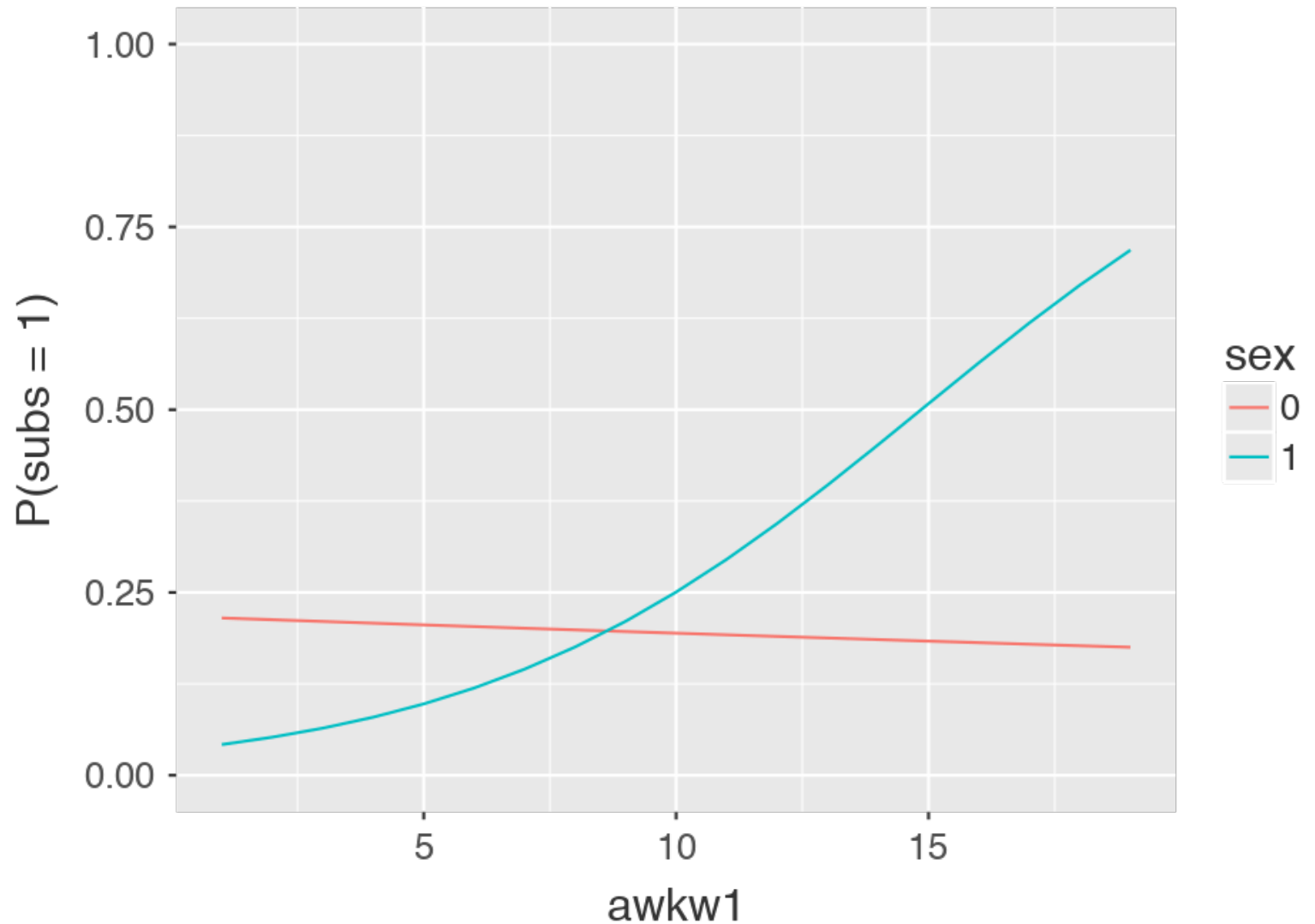
Check correlations

Dummy codes are correct (Jamovi does this automatically)

Interactions

Just as we do
in linear
models

Can have 2+
variables in
the interaction



Interactions

The screenshot shows the Jamovi software interface. The top navigation bar includes 'Data' and 'Analyses' tabs. Below the navigation bar are icons for various statistical tests: Exploration, T-Tests, ANOVA, Regression, Frequencies, and Factor. The main window is titled 'Binomial Logistic Regression'. On the left, a list of variables includes 'sex', 'race', 'inco', 'chil', 'alco', 'spor', 'depr1', 'awkw1', and 'prod2'. The 'Factors' section contains 'show' and 'Married'. The 'Predictors' section contains 'show' and 'Married'. The 'Blocks' section shows 'Block 1' with 'show' and 'Married' and 'Block 2' with 'variables here'. A context menu is open over the 'Interaction' option, showing 'Main Effects', 'All 2 way', 'All 3 way', 'All 4 way', and 'All 5 way'. An orange arrow points from a text box to the 'Interaction' option. On the right, the 'Estimated Marginal Means' plot shows the probability of 'subs = 1' for two shows: 'Parks and Rec' and 'The Office'. The y-axis ranges from 0.00 to 1.00. The plot shows that 'The Office' has a higher probability (around 0.2) compared to 'Parks and Rec' (around 0.14). Below the plot is a table of 'Estimated Marginal Means - show'.

show	Probability	SE	95% Confidence Interval	
			Lower	Upper
Parks and Rec	0.138	0.0834	0.0387	0.387
The Office	0.191	0.1085	0.0564	0.483

Can tell Jamovi to do an interaction

Questions?

Please post them to the
discussion board before
class starts

End of Pre-Recorded Lecture Slides

In-class discussion slides



Application

Example Using
The Office/Parks and Rec Data Set

Hypothesis Test with
Logistic Regression