

# **Unit 14**: Logistic Regression Variable Selection

### **Case Studies:**

- To introduce the concept of <u>using training data to build a</u>
   model and <u>using test data to test a model for it's predictive</u>
   <u>capabilities</u> we will, again, examine the relationship between
   a:
  - Categorical response variable: support for a certain opinion (favor/not in favor) and an
  - Explanatory variables:
    - <u>Sex</u>
    - Party, and
    - Age

### **Summary of Concepts:**

- Definitions/Properties
  - Maximum Likelihood Estimation
  - O AIC
  - O BIC
- Modeling
  - O How can we select the "best" explanatory variables to include in our logistic regression model?
    - Using log-likelihood ratio test to determine this.
    - Using AIC and BIC to determine this.

### **Unit 14: Logistic Regression Variable** Selection

Previously we have been building models manually either by having specific variables in mind, or by making targeted comparisons between models with different variables. In this section we begin to explore more automated methods for modeling. The key concept is to embed a a model in a larger class of potential models and tune the models within this class. This tuning process is called \*\*learning\*\* the model, and starts us on the road to machine learning.

A big issue in model selection is the temptation to fit bigger and bigger models in order to improve the fit to the training data. This tendancy is called **overfitting**. By overfitting the data at hand, we risk losing the ability to generalize the results to future data or larger populations, because the model is too fine tuned to the data at hand.

Learning methods are designed to counteract the tendancy to overfit the data. A simple approach introduced in the previous section is to split the data randomly into training and testing subsets of the data. We do all the model building on the training data, and then assess the model using the test data.

### **Topic 0: Review of Methods to Deal with Overfitting that We've Already Learned**

### For Linear Regression Models

For linear regression models (such as the one below), what is a method that we talked about in the past for testing whether multiple explanatory variables were "needed" in the model?

**Ex**: Considering whether 
$$x_2$$
 and  $x_5$  are needed in the linear regression model below.  $\hat{y} = \hat{\beta_0} + \hat{\beta_1} x_1 + \hat{\beta_2} x_2 + \hat{\beta_3} x_3 + \hat{\beta_4} x_4 + \hat{\beta_5} x_5$ 

### For Logistic Regression Models

# Topic 1 (Definitions/Theory): How to Describe the Tradeoff Between Overfitting a Model and Underfitting a Model

### **Bias**

We define the bias of a model  $\hat{Y} = \hat{\beta_0} + \hat{\beta_1} x_1 + \hat{\beta_2} x_2 + \ldots + \hat{\beta_p} x_p$  as:

$$Bias(\hat{Y}) = E[\hat{Y}] - \mu$$

Idea: Bias is likely to be higher in models that are: \_\_\_\_

### **Variance**

We define the variance of a model  $\hat{Y} = \hat{\beta_0} + \hat{\beta_1} x_1 + \hat{\beta_2} x_2 + \ldots + \hat{\beta_p} x_p$  as:

$$Var(\hat{Y}) = E[(\hat{Y} - E(\hat{Y}))^2]$$

Idea: Variance is likely to be higher in models that are: \_\_\_\_

The idea is that simpler models might be biased due to some missing variables or transformations, so  $E[\hat{Y}] \neq \mu$ , but if the bias is not too large compared to the variance reduction they provide, the mean square error can be improved over larger, less biased models with larger variance. If we go too far in this direction, however, the bias will overtake the variance. So we expect there will be some optimal model between the two extremes.

### Mean Square Error (MSE)

Finally, we define the mean square error of a model as

$$E[(\hat{Y}-\mu)^2]$$

A key modeling aim is to find an effective compromise between bias reduction and variance reduction, for example, by searching for models with small **mean square error** for prediction, such a compromise might be found.

Fundamental bias-variance decomposition for model prediction  $\hat{Y}$ :

### **Bias-Variance Tradeoff Relationship**

$$MSE(\hat{Y}) = E[(\hat{Y} - \mu)^2] =$$

$$=$$

$$=$$

$$= E[(\hat{Y} - E(\hat{Y}))^2] + [E(\hat{Y}) - \mu]^2$$

$$= Var(\hat{Y}) + Bias^2(\hat{Y}).$$

### Relationship

This section explores several methodologies useful in model selection, aimed at addressing the overfit/underfit challenge:

- Log-Likelihood-Ratio Tests for comparing nested logistic regression models; analogous to F-tests in ANOVA
- · Information criteria such as AIC and BIC that trade off model fit with model complexity
- Train/Test data splitting to evaluate model based classifiers for sensitivy, specificity and accuracy

### Python libraries and functions:

```
statsmodels.api
statsmodels.formula.api
logit
scipy.stats
bernoulli
chi2
norm
sklearn.model_selection
train_test_split
sklearn.metrics
accuracy_score
confusion_matrix
roc_curve
roc_auc_score
```

### Topic 2: Maximum Likelihood Estimation for Determining Logistic Regression Model Parameters

Recall that in linear regression modeling it can be useful to test between two models using an analysis of variance F test, which compares the residual sums of squares for two, nested models. It allows us to test multiple parameters within one hypothesis test.

In logistic regression modeling, the F test is no longer applicable. However, the same general testing idea is possible by comparing log-likelihoods between two nested models. The change in log-likelihood is used as a large sample chi-square test of the null hypothesis that the simpler model is adequate.

# How are the optimal values of $\hat{\beta_0}$ , $\hat{\beta_1}$ , ..., $\hat{\beta_p}$ determined in a logistic regression model?

### 1. Assumption Behind Logistic Regression

```
• Each y_1, y_2,...,y_n are independent.

• y_i \sim Bern(p_i), where

• \log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots \beta_p X_{ip}, for i=1,2,\ldots,n.

• Put another way: p_i = \frac{e^{\beta_0+\beta_1 X_{i1}+\cdots\beta_p X_{ip}}}{1+e^{\beta_0+\beta_1 X_{i1}+\cdots\beta_p X_{ip}}}
```

### 2. Representing the Response Variable Probabilities

How can we represent the probability mass function of  $y_i$ , given the explanatory variable values and a given logistic regression model

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}, \quad \text{for} \quad i = 1, 2, \dots, n. \text{ (ie. } p_i = \frac{e^{\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}}})$$
?

Answer:  $P(y_i|\beta_0, \beta_1, \dots, \beta_n, X_i) = p_i^{y_i} (1 - p_i)^{1 - y_i}$ 

How can we represent the JOINT probability mass function of  $y_1$ ,  $y_2$ ,..., $y_n$  given the explanatory variable values and a given logistic regression model

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}, \quad \text{for} \quad i = 1, 2, \dots, n. \text{ (ie. } p_i = \frac{e^{\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}}})?$$

**Answer:** In binary response models such as logistic regression the **likelihood function (LF)** is the joint probability mass function of the responses viewed as a function of the parameters. For a logit model with independent Bernoulli responses, the likelihood function has the form:

$$LF(\beta_0, \beta_1, \dots, \beta_p) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1 - y_i}$$

where

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}, \quad \text{for} \quad i = 1, 2, \dots, n.$$

## 3. Goal: Determine the best values of $\hat{\beta_0}$ , $\hat{\beta_1}$ , ..., $\hat{\beta_p}$ that maximize the likelihood function.

Maximize with respect to  $\beta_0, \beta_1, \dots, \beta_p$ 

$$LF(\beta_0, \beta_1, \dots, \beta_p) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1 - y_i}$$

where

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}, \quad \text{for} \quad i = 1, 2, \dots, n.$$

Can we transform the likelihood function into another function that gives us the same optimal values of  $\overset{\wedge}{\beta_0}$ ,  $\overset{\wedge}{\beta_1}$ , ...,  $\overset{\wedge}{\beta_p}$ , but is easier to take the derivative of?

**Answer:** The logarithmic tranformation converts the product to a sum of log values, the log-likelihood function (LLF):  $LLF(\beta_0,\beta_1,\ldots,\beta_p) = \sum_{i=1}^n \{y_i \log(p_i) + (1-y_i) \log(1-p_i)\}.$ 

## 4. Where can we find this optimal value to the log-likelihood function?

The result reported in the model summary (listed as 'Log-Likelihood') is the optimized value computed by **maximum likelihood estimation**:

$$llf = LLF(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) = \sum_{i=1}^n \{ y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i) \}$$

## Topic 3: Log-likelihood Ratio Test: For Comparing two Logistic Regression Models

### **Step 1: Two Nested Logistic Regression Models**

### Model 1:

$$\operatorname{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p.$$

Model 0: ...

$$logit(p) = log(\frac{p}{1-p}) =$$

In order to test between two logit models, Model 0 and Model 1, where Model 0 is a special case of Model 1 obtained by setting some regression coefficients equal to zero.

### Step 2: Set up two hypotheses

Consider one is a special case of the other we can compare their log-likelihood ratios. Consider testing:

 $H_0$ : Model 0 is correct,

 $H_A$ : Model 0 is incorrect because at least one missing coefficient from Model 1 is not ze

### **Step 3: Test Statistic**

A general result from large sample theory is if  $H_0$  is true, then twice the difference in negative log-likelihoods

$$llr = -2\left(llf_0 - llf_1\right)$$

# Step 4: What distribution is this test statistic an observation from (and what degrees of freedom should it have)?

This test statistic has an approximate Chi-square distribution with degrees of freedom equal to the difference in the numbers of parameters for the two models. Like the central limit theorem, this approximation works better for larger sample size n.

### Some Properties of Chi-Squared Distribution:

- 1. Positive distribution.
- 2. Parameter that Determines Shape of Distribution: degrees of freedom value

Step 5: Use this test statistic and Chi-Squared distribution to find the p-value that corresponds to these hypotheses. Use the p-value to make a conclusion about the hypotheses.

Applying this test in our example lets us evaluate multiple coefficients at the same time to determine whether we can reduce to the simpler model.

### Some Properties of Calculating the p-value with Chi-Squared Distribution for this Test

1. Just use a right-tail.

Here's how it works.

## Topic 4: Pew Research Survey Example (Model Selection with Log Likelihood Ratio Test)

Using Log-likelihood Ratio Test For Comparing two Logistic Regression Models

In an earlier section we considered two models for predicting a favorable opinion of border wall construction in the Pew Research Survey of February 2017. Let's load the data and the two models and first see how we can test between the two models. The idea is analogous to the ANOVA method for comparing two linear regression models.

### Preprocessing and data validation

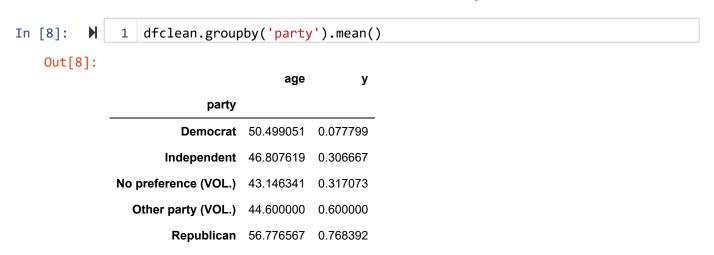
```
In [1]:
         H
                 import numpy as np
              2
                 import pandas as pd
                 import zipfile as zp
              3
                 import statsmodels.api as sm
                 import statsmodels.formula.api as smf
In [2]:
         M
                 zf = zp.ZipFile('../data/Feb17-public.zip')
              1
                 missing_values = ["NaN", "nan", "Don't know/Refused (VOL.)"]
              2
                 df = pd.read_csv(zf.open('Feb17public.csv'),
                                  na_values=missing_values)[['age', 'sex', 'q52', 'party'
                 # reduce q52 responses to two categories
In [3]:
         H
              1
                 # and create binary reponse variable
                 df['q52'][df['q52']!='Favor'] = 'Not_favor'
                df['y'] = df['q52'].map({'Not_favor':0,'Favor':1})
                # use cleaned data without records that have missing values
              6 dfclean = df.dropna()
In [4]:
              1 dfclean.head()
   Out[4]:
                age
                       sex
                                q52
                                         party y
             0 80.0 Female
                            Not favor Independent 0
             1 70.0 Female Not_favor
                                      Democrat 0
             2 69.0 Female Not favor Independent 0
             3 50.0
                                     Republican 1
                      Male
                               Favor
             4 70.0 Female Not favor
                                      Democrat 0
In [5]:
         H
              1 dfclean['party'].value_counts()
   Out[5]: Democrat
                                     527
            Independent
                                     525
            Republican
                                     367
            No preference (VOL.)
                                      41
            Other party (VOL.)
                                       5
            Name: party, dtype: int64
                dfclean['sex'].value_counts()
In [6]:
   Out[6]: Male
                       760
            Female
                       705
            Name: sex, dtype: int64
```

```
In [7]: ► 1 dfclean.describe()
```

Out[7]:

	age	У
count	1465.000000	1465.000000
mean	50.522867	0.341297
std	17.843611	0.474307
min	18.000000	0.000000
25%	35.000000	0.000000
50%	52.000000	0.000000
75%	65.000000	1.000000
max	96.000000	1.000000

# <u>Descriptive Analytics Question:</u> Is the proportion of people that support the border wall different for at least one pair of political parties *in the sample*?



We can see that the proportion of 'favor' responses varies quite a bit between party affiliations, by looking at the mean values for 'y'. In each subgroup, the sample mean of y equals the proportion who favored building the wall.

Inference Question: Is the proportion of people that support the border wall different for at least one pair of political parties in the population of all adults that live in the U.S.?

### Use a Full model and reduced model for log-likelihood-ratio test

Recall that 'party' is a categorical variable with 5 categories. If we wish to test the null hypothesis of no party effects, we need a 4 degree of freedom test. For this we can use the log-likelihood-ratio test.

### Step 1: Set up a full model and a null model.

- Null Model (Model 0):
  - Response = Support for Border Wall
  - Explanatory Variables:
    - age
    - sex
- Full Model (Model 1):
  - Response = Support for Border Wall
  - Explanatory Variables:
    - age
    - sex
    - party

First we fit the null and full model:

### Step 2: Set up the null and alternative hypotheses.

 $H_0$ : Model 0 is correct,

 $H_A$ : Model 0 is incorrect because the missing 'party' coefficient in model 0 is not zero.

### **Step 3: Calculate the test statistic.**

We don't need to display the summaries to perform the test, but it is informative to review the model summaries to understand the variables. The maximized log-likelihood is shown in the model summary as 'Log-Likelihood'.

### Step 3a: Extract the log-likelihoods for the two models:

In [10]: model0.summary() Out[10]: Logit Regression Results Dep. Variable: No. Observations: 1465 Model: Logit **Df Residuals:** 1462 2 Method: MLE Df Model: Thu, 19 Nov 2020 Pseudo R-squ.: 0.03557 Date: Time: 11:05:08 Log-Likelihood: -906.92 converged: LL-Null: -940.37 True **Covariance Type:** LLR p-value: 2.960e-15 nonrobust [0.025 0.975] coef std err Z P>|z| 0.196 -10.637 Intercept -2.0818 0.000 -2.465 -1.698 0.5415 0.000 0.318 0.765 sex[T.Male] 0.114 4.750 age 0.0220 0.003 6.770 0.000 0.016 0.028 In [11]: model1.summary() Out[11]: Logit Regression Results Dep. Variable: No. Observations: 1465 **Df Residuals:** 1458 Model: Logit Method: MLE **Df Model:** 6 Thu, 19 Nov 2020 Pseudo R-squ.: 0.2738 Date: Time: 11:05:08 Log-Likelihood: -682.88 converged: True LL-Null: -940.37 **Covariance Type:** LLR p-value: 4.971e-108 nonrobust coef std err P>|z| [0.025 0.975] Z Intercept -3.5261 0.281 -12.536 0.000 -4.077 -2.975 party[T.Independent] 1.6843 0.191 8.796 0.000 1.309 2.060 4.807 0.000 2.566 party[T.No preference (VOL.)] 1.8226 0.379 1.079 party[T.Other party (VOL.)] 0.938 0.002 4.732 2.8930 3.083 1.054 party[T.Republican] 3.5862 0.206 17.435 0.000 3.183 3.989 sex[T.Male] 0.3721 0.137 2.712 0.007 0.103 0.641 0.0168 0.004 4.305 0.000 0.009 0.024 age

In [12]: ▶ 1 model0.llf, model1.llf

## Step 3b: Use these log-likelihoods to calculate the likelihood ratio test statistic.

Just be careful to get the multiplier (-2) right so the chi-sqaure approximation works correctly.

```
In [14]:
          H
                  # Extract log-likelihood function values
               2
                 # and model degrees of freedom from each model
               3
                  11f0, df0 = model0.llf, model0.df_model
                  llf1, df1 = model1.llf, model1.df_model
                  # take differences
                  llr, dfdiff = -2*(llf0 - llf1), df1 - df0
               7
                  # display results
               8
                  pd.DataFrame({'-2*11f': [-2*11f0, -2*11f1, llr],
               9
                                 'df_model': [df0, df1, dfdiff]},
              10
                               index=['model0','model1', 'diff'])
    Out[14]:
                           -2*Ilf df_model
```

# -2\*Ilf df\_model model0 1813.836471 2.0 model1 1365.759089 6.0 diff 448.077382 4.0

## Step 4: Calculate the degrees of freedom for the chi-squared distribution that this test statistic is an observation from.

Why was df = 4 in this analysis?

### **Step 5: Calculate the p-value and make a conclusion.**

Summarize the test with calculated p-value using chi-square distribution

**Conclusion:** We definitely reject the null hypothesis and favor Model 1 over Model 0. Party affiliation is a significant factor associated with the response to question 52 in the survey.

### **Topic 5: Model Selection with AIC and BIC**

## Two metrics that measure the balance between having a good model fit and a small number of variables.

AIC and BIC are criteria for evaluating a model that combine the likelihood assessment of fit with a penalty for complex models. Historically they were derived from different perspectives.

### **Akaike Information Criterion (AIC)**

The Akaike Information Criterion (AIC), has the form

$$AIC = -2 * ll f + 2 * p$$

### How to use it:

where p is the same as the model degrees of freedom. Small values are considered better than large values, so minimizing AIC favors larger likelihoods and simpler models, while trying to balance these two goals.

### **Bayesian Information Criterion (BIC)**

The Bayesian Information Criterion (BIC) is related but uses different relative weighting of likelihood and complexity:

$$BIC = -2 * llf + p * log(n).$$

### How to use it:

AIC vs. BIC
Again, models with smaller values are better than models with larger values. Both methods enforce favoring simpler models among those with simimlar fit overall, and they help prevent overfitting the model because of the complexity penalty.

BIC tends to favor simplicity more heavily than does AIC due to its heavier penalty for large p.

### What AIC and BIC can be used for:

**Use cases:** AIC and/or BIC are often used to guide variable selection when multiple exogensous variables are considered for inclusion in the model. This enables us to compare a whole series of models and try to find a reasonable tradeoff between bias and variance, i.e., goodness of fit and model complexity.

### What AIC and BIC cannot be used for:

**Evaluation of predictive accuracy:** Although model selection criteria like AIC and BIC can help avoid overfitting and underfitting the data, they do no provide us with assessment of classification performance. In order to evaluate the model selected by these criteria or related strategies, it is still necessary to use some version of the train/test method, where the training data are used for the model building process, and the test data are reserved for predictive evaluation only.

# **Topic 6: Pew Research Survey Example (Model Selection with AIC and BIC)**

In the current implementation of the statsmodels logit api, both of these criteria are available from the model fitting results. Here's a summary for our two models of the Pew survey data for predicting favorable or unfavorable opinions of the border wall:

In [17]: ▶ 1 model1.summary()

Out[17]:

Logit Regression Results

Dep. Variable:	у	No. Observations:	1465
Model:	Logit	Df Residuals:	1458
Method:	MLE	Df Model:	6
Date:	Thu, 19 Nov 2020	Pseudo R-squ.:	0.2738
Time:	11:05:08	Log-Likelihood:	-682.88
converged:	True	LL-Null:	-940.37
Covariance Type:	nonrobust	LLR p-value:	4.971e-108

	coef	std err	z	P> z	[0.025	0.975]
Intercept	-3.5261	0.281	-12.536	0.000	-4.077	-2.975
party[T.Independent]	1.6843	0.191	8.796	0.000	1.309	2.060
party[T.No preference (VOL.)]	1.8226	0.379	4.807	0.000	1.079	2.566
party[T.Other party (VOL.)]	2.8930	0.938	3.083	0.002	1.054	4.732
party[T.Republican]	3.5862	0.206	17.435	0.000	3.183	3.989
sex[T.Male]	0.3721	0.137	2.712	0.007	0.103	0.641
age	0.0168	0 004	4 305	0.000	0.009	0.024

```
In [18]:
                    model0.summary()
    Out[18]:
               Logit Regression Results
                   Dep. Variable:
                                              y No. Observations:
                                                                       1465
                         Model:
                                            Logit
                                                      Df Residuals:
                                                                       1462
                        Method:
                                                         Df Model:
                                                                          2
                                            MLE
                           Date: Thu, 19 Nov 2020
                                                    Pseudo R-squ.:
                                                                     0.03557
                          Time:
                                         11:05:08
                                                   Log-Likelihood:
                                                                     -906.92
                     converged:
                                            True
                                                          LL-Null:
                                                                     -940.37
                Covariance Type:
                                       nonrobust
                                                      LLR p-value: 2.960e-15
                              coef std err
                                                z P>|z|
                                                          [0.025 0.975]
                  Intercept -2.0818
                                     0.196 -10.637 0.000 -2.465 -1.698
                sex[T.Male]
                            0.5415
                                             4.750 0.000
                                                          0.318
                                                                 0.765
                                     0.114
                       age
                            0.0220
                                     0.003
                                             6.770 0.000
                                                          0.016
                                                                 0.028
In [19]:
                    model1.aic, model0.aic, model1.bic, model0.bic
    Out[19]:
               (1379.7590888950426,
                1819.8364712252783,
                1416.7863625452007,
                1835.7053027896318)
In [20]:
            H
                    pd.DataFrame({'-2*11f': [-2*11f0, -2*11f1],
                 2
                                     'df_model': [df0, df1],
                 3
                                     'AIC': [model0.aic, model1.aic],
                                     'BIC': [model0.bic, model1.bic]},
                 4
                 5
                                   index=[0,1])
    Out[20]:
                          -2*IIf df_model
                                                              BIC
                                                 AIC
                   1813.836471
                                         1819.836471
                                                      1835.705303
                                     2.0
                   1365.759089
                                     6.0
                                         1379.759089 1416.786363
```

### **Conclusion:**

Both AIC and BIC favor Model 1. This suggests that Model 0 is too simple, so the bias due to omitted variables is too large for this model compared to Model 1.

# Topic 7: Determining What our Null Model Should Be (Given a Full Model): Backwards Elimination Algorithm

### Generate the features matrix and response data from a logit model

For illustration in an example with many explanatory variables we generate binary response data with 20 explanatory variables. First we set up the coefficient vector for the simulation model.

### **Generating Simulated Data**

Next we generate a random features matrix, using numpy matrix operations to form the matrix.

Use numpy matrix multiplication to form the log-odds model, and exponentiate to get the vector of n odds for the responses.

```
In [25]: | # compute the odds of 1 for n observations
2 # use numpy matrix multiplication to make this easier
3 odds = np.exp(b0 + np.matmul(X, bvec))
4 odds.shape
Out[25]: (200,)
```

Convert the odds vector to the population probability vector for the n 0/1 responses. Then use the bernoulli.rvs function to generate the responses from the model.

To set up for formula based modeling, we assign names to the columns of X.

```
1 # Add y to the data frame
In [29]: ▶
              2 df['y'] = y
              3 display(df.shape, df.iloc[0:5,:7], \
                        df.iloc[0:5, 7:14], df.iloc[0:5, 15:])
            (200, 21)
```

	X1	X2	Х3	X4	X5	Х6	X7
0	1.624345	-0.611756	-0.528172	-1.072969	0.865408	-2.301539	1.744812
1	-1.100619	1.144724	0.901591	0.502494	0.900856	-0.683728	-0.122890
2	-0.191836	-0.887629	-0.747158	1.692455	0.050808	-0.636996	0.190915
3	-0.754398	1.252868	0.512930	-0.298093	0.488518	-0.075572	1.131629
4	-0.222328	-0.200758	0.186561	0.410052	0.198300	0.119009	-0.670662

	X8	Х9	X10	X11	X12	X13	X14
0	-0.761207	0.319039	-0.249370	1.462108	-2.060141	-0.322417	-0.384054
1	-0.935769	-0.267888	0.530355	-0.691661	-0.396754	-0.687173	-0.845206
2	2.100255	0.120159	0.617203	0.300170	-0.352250	-1.142518	-0.349343
3	1.519817	2.185575	-1.396496	-1.444114	-0.504466	0.160037	0.876169
4	0.377564	0.121821	1.129484	1.198918	0.185156	-0.375285	-0.638730

	X16	X17	X18	X19	X20	у
0	-1.099891	-0.172428	-0.877858	0.042214	0.582815	0
1	-0.012665	-1.117310	0.234416	1.659802	0.742044	1
2	0.586623	0.838983	0.931102	0.285587	0.885141	0
3	-2.022201	-0.306204	0.827975	0.230095	0.762011	0
4	0.077340	-0.343854	0.043597	-0.620001	0.698032	0

### 1. FOR ASSESSING PREDICTIVE POWER OF NEW OBSERVATIONS: Split the data into training data and test data

## 2. FOR CREATING A MODEL THAT IS "PARSIMONIOUS": Model the training data: are 20 variables necessary?

### **Full Model**

```
In [32]:
                    mod0 = smf.logit(
                 1
                 2
                         'v \sim X1+X2+X3+X4+X5+X6+X7+X8+X9+X10
                 3
                         +X11+X12+X13+X14+X15+X16+X17+X18+X19+X20',
                 4
                         data=df_train).fit()
               Optimization terminated successfully.
                         Current function value: 0.176600
                         Iterations 10
In [33]:
                 1 # model information
                 2 mod0.summary().tables[0]
    Out[33]:
               Logit Regression Results
                   Dep. Variable:
                                              y No. Observations:
                                                                       160
                         Model:
                                           Logit
                                                     Df Residuals:
                                                                       139
                        Method:
                                           MLE
                                                        Df Model:
                                                                        20
                          Date: Thu, 19 Nov 2020
                                                   Pseudo R-squ.:
                                                                    0.7376
                                        11:05:09
                                                   Log-Likelihood:
                          Time:
                                                                    -28.256
                                                         LL-Null:
                     converged:
                                           True
                                                                    -107.68
                Covariance Type:
                                       nonrobust
                                                     LLR p-value: 1.251e-23
```

```
In [34]:
                     # model coefficient summary table
                     mod0.summary().tables[1]
    Out[34]:
                              coef std err
                                                z P>|z|
                                                           [0.025 0.975]
                 Intercept -1.3640
                                     0.504
                                           -2.709
                                                   0.007
                                                           -2.351
                                                                  -0.377
                                     0.700
                       X1
                            2.4668
                                            3.524
                                                   0.000
                                                           1.095
                                                                   3.839
                       X2
                            2.5231
                                     0.752
                                            3.355 0.001
                                                           1.049
                                                                   3.997
                       X3
                            2.6866
                                     0.737
                                            3.645 0.000
                                                           1.242
                                                                   4.131
                       X4
                            2.1278
                                     0.767
                                             2.773
                                                   0.006
                                                           0.624
                                                                   3.632
                                            3.541
                       X5
                            2.1866
                                     0.618
                                                   0.000
                                                           0.976
                                                                   3.397
                       X6
                           -1.3311
                                     0.503
                                            -2.645 0.008
                                                          -2.317
                                                                  -0.345
                                                   0.000
                           -2.4016
                                     0.653
                                            -3.677
                                                           -3.682
                                                                  -1.121
                       X8
                           -1.6166
                                     0.576
                                           -2.806
                                                   0.005
                                                          -2.746
                                                                  -0.487
                       X9
                           -2.0160
                                     0.710
                                            -2.838
                                                   0.005
                                                          -3.409
                                                                  -0.624
                      X10
                           -2.2436
                                     0.721
                                            -3.112 0.002
                                                           -3.657
                                                                  -0.831
                      X11
                            1.4669
                                     0.551
                                             2.660
                                                   0.008
                                                           0.386
                                                                   2.548
                      X12
                            0.7205
                                     0.432
                                             1.670 0.095 -0.125
                                                                   1.566
                      X13
                            0.8124
                                     0.451
                                             1.802 0.072 -0.071
                                                                   1.696
                      X14
                            0.1165
                                     0.412
                                             0.283
                                                   0.777
                                                          -0.690
                                                                   0.923
                      X15
                            0.2603
                                     0.459
                                             0.567
                                                   0.571
                                                          -0.640
                                                                   1.160
                      X16
                            0.6365
                                     0.489
                                             1.302 0.193
                                                          -0.321
                                                                   1.594
                      X17
                                     0.392
                                             0.960
                                                   0.337
                                                           -0.392
                            0.3760
                                                                   1.144
                      X18
                           -0.1898
                                     0.545
                                            -0.348 0.728 -1.258
                                                                   0.878
                      X19
                                             2.248 0.025
                            1.0728
                                     0.477
                                                           0.137
                                                                   2.008
                      X20
                            0.8356
                                     0.489
                                             1.710 0.087 -0.122
                                                                   1.793
```

Here are AIC and BIC for the model:

```
In [35]: ▶ 1 (mod0.aic, mod0.bic)
Out[35]: (98.51189376685147, 163.09054388676185)
```

### 3. Backwards Elimination Algorithm Ideas

Which Null Model Should we Select to Compare to the Full Model? Let's compare a simpler model. There are many possible models ( $2^{20}$ ), so how can we process them? An old idea is to use the coefficient tests to help filter variables.

```
In [36]:
                   mod0.pvalues.sort_values()
    Out[36]: X7
                            0.000236
              Х3
                            0.000267
              X5
                            0.000399
              Х1
                            0.000425
              Χ2
                            0.000794
              X10
                            0.001857
              Х9
                            0.004545
              X8
                            0.005020
              X4
                            0.005551
              Intercept
                            0.006752
              X11
                            0.007813
              Х6
                            0.008164
              X19
                            0.024573
              X13
                            0.071623
              X20
                            0.087216
              X12
                            0.095010
              X16
                            0.192819
              X17
                            0.337042
              X15
                            0.570756
              X18
                            0.727601
              X14
                            0.777220
              dtype: float64
                  mod0.pvalues[mod0.pvalues < 0.05]</pre>
In [37]:
    Out[37]: Intercept
                            0.006752
              X1
                            0.000425
              Χ2
                            0.000794
              Х3
                            0.000267
              Х4
                            0.005551
              X5
                            0.000399
              X6
                            0.008164
              Χ7
                            0.000236
              X8
                            0.005020
              Х9
                            0.004545
              X10
                            0.001857
              X11
                            0.007813
              X19
                            0.024573
              dtype: float64
```

## One Idea: Null Model = Only the Variables that have Statistically Significant p-values in the Full Model

Let's compare the model that only keeps these "significant" variables.

BIC is reduced. AIC is about the same. Let's try to go further.

```
mod1.pvalues[mod1.pvalues < 0.05]</pre>
In [39]:
    Out[39]: Intercept
                           0.010066
              Х1
                            0.000045
              Χ2
                           0.000043
              Х3
                           0.000013
              Χ4
                           0.000229
              X5
                           0.000192
              X6
                           0.008512
              Χ7
                           0.000085
              X8
                           0.000584
              Х9
                           0.000647
              X10
                           0.000220
              X11
                           0.002939
              X19
                           0.023514
              dtype: float64
```

## Another Idea: Null Model = Drop Even More Explanatory Variables, Starting with the Ones that Have the Highest p-values

Try dropping this least significant variable to see what happens.

```
In [41]:
                  mod2 = smf.logit('y \sim X1+X2+X3+X4+X5+X6+X7+X8+X9+X10+X11',
                                    data=df_train).fit()
               3
                 (mod2.aic, mod2.bic)
             Optimization terminated successfully.
                       Current function value: 0.245458
                       Iterations 9
   Out[41]: (102.5466613953034, 139.44874717810933)
         AIC and BIC increase a bit. Try dropping one more...
In [42]:
               1 mod2.pvalues[mod2.pvalues==max(mod2.pvalues)]
    Out[42]: Intercept
                           0.013388
              dtype: float64
In [43]:
                 mod2.pvalues.sort_values()
   Out[43]: X3
                           0.000015
             X2
                           0.000018
             Х1
                           0.000032
             Χ7
                           0.000081
             X5
                           0.000210
             X10
                           0.000236
             Χ4
                           0.000363
             X8
                           0.000818
             Х9
                           0.001048
             X11
                           0.005432
             X6
                           0.008977
              Intercept
                           0.013388
              dtype: float64
         What happens if we drop X6?
                  mod3 = smf.logit('y \sim X1+X2+X3+X4+X5+X7+X8+X9+X10')
In [44]:
               2
                                    data=df train).fit()
               3
                 (mod3.aic, mod3.bic)
             Optimization terminated successfully.
                       Current function value: 0.297416
                       Iterations 8
```

Even more increase. Looks like we can't reduce the model beyond mod1, based on these criteria.

Out[44]: (115.17324407941844, 145.9249822317567)

	aic	bic
0	98.511894	163.090544
1	98.738294	138.715554
2	102.546661	139.448747
3	115.173244	145.924982

### Conclusion

According to BIC, model 1 is the best. According to AIC it's very close between mod0 and mod1. Here's the model summary for mod1:

In [46]: ► 1 mod1.summary()

Out[46]:

Logit Regression Results

Dep. V	ariable:		у	No. C	Observat	ions:	160
	Model:		Logit		Df Resid	uals:	147
ı	Method:		MLE		Df M	odel:	12
	Date:	Thu, 19 N	Nov 2020	Ps	eudo R-	squ.:	0.6623
	Time:		11:05:09	Lo	g-Likelih	ood:	-36.369
con	verged:		True		LL-	Null:	-107.68
Covariano	e Type:	n	onrobust		LLR p-v	alue:	1.766e-24
	coef	std err	z	P> z	[0.025	0.975	5]
Intercept	-0.9374	0.364	-2.574	0.010	-1.651	-0.22	3
X1	2.1566	0.529	4.080	0.000	1.120	3.19	3
X2	2.0498	0.501	4.090	0.000	1.068	3.03	2
Х3	1.9472	0.447	4.358	0.000	1.071	2.82	3
X4	1.7873	0.485	3.685	0.000	0.837	2.73	8
X5	1.4869	0.399	3.729	0.000	0.705	2.26	8
X6	-0.9173	0.349	-2.631	0.009	-1.601	-0.23	4
X7	-1.9646	0.500	-3.930	0.000	-2.944	-0.98	5
X8	-1.4560	0.423	-3.439	0.001	-2.286	-0.62	6
Х9	-1.7072	0.500	-3.411	0.001	-2.688	-0.72	6
X10	-1.7636	0.477	-3.695	0.000	-2.699	-0.82	8
X11	1.1117	0.374	2.974	0.003	0.379	1.84	4
X19	0.7941	0.351	2.265	0.024	0.107	1.48	1

Compared to the simulation model that generated the data we see that the best fitted model is missing X8 and includes X14, which we know to have a zero coefficient from the simulation model. This is an example of the effects of sample variation in model building.

### 4. Evaluate selected model as a classifier on test data

Let's compute the accuracies of the models as classifiers. We'll use the predictive probability as the classification score use the classification rule:

$$\hat{y} = \begin{cases} 1, & \text{if } \hat{p} \ge 0.5\\ 0, & \text{if } \hat{p} < 0.5 \end{cases}$$

We compute the predictive probabilities for Model 1:

```
In [47]:
              1 phat1 = mod1.predict(exog=df_test)
              2 phat1[0:10]
   Out[47]: 95
                    0.347583
             15
                    0.050317
             30
                    0.001393
             158
                    0.999849
             128
                    0.016493
             115
                    0.001555
                    0.000541
             69
             170
                    0.000013
             174
                    0.020079
             45
                    0.079501
             dtype: float64
```

## **New Statistic Assessing Predictive Power: Classification Accuracy**

### **Classification Accuracy for Model 1**

Here's a function to compute the classification accuracy, which is the overall fraction correctly classified.

The classification accuracy for Model 1 is estimated to be 92.5%. This is the combination of the true positives rate and the true negatives rate, if we view 1's as positive and 0's as negative.

#### **Comparing Classification accuracy for all Models?**

```
# bind together all the predictive probabilities into a matrix
In [50]:
           H
               1
               2
                  phat_matrix = np.array([mod0.predict(exog=df_test),
               3
                                          mod1.predict(exog=df test),
               4
                                          mod2.predict(exog=df test),
                                          mod3.predict(exog=df_test)])
               5
               6
                  phat matrix.shape
   Out[50]: (4, 40)
In [51]:
               1 | phat matrix[2][0:10] # compare with values above
   Out[51]: array([3.09116764e-01, 1.02689627e-01, 1.21609512e-03, 9.99795672e-01,
                     1.93365814e-02, 6.22526164e-03, 6.48536897e-04, 1.61225639e-05,
                     2.16677944e-02, 8.21426220e-02])
                  accuracy_list = []
In [52]:
               1
           H
               2
                  for i in range(0,4):
               3
                      accuracy_list.append(
               4
                          accuracy_score(y_true=df_test['y'],
               5
                                          y_pred=1*(phat_matrix[i] >= pthresh),
                                          normalize=True)
               6
               7
                      )
               1 accuracy list
In [53]:
    Out[53]: [0.875, 0.925, 0.875, 0.875]
In [54]:
                  pd.DataFrame({'aic': [mod0.aic, mod1.aic, mod2.aic, mod3.aic],
           H
               2
                                'bic': [mod0.bic, mod1.bic, mod2.bic, mod3.bic],
               3
                                'accuracy': accuracy list},
               4
                                index=[0,1,2,3]
   Out[54]:
                       aic
                                 bic accuracy
                  98.511894
                           163.090544
              0
                                         0.875
                  98.738294
                           138.715554
                                         0.925
              1
              2 102.546661
                           139.448747
                                         0.875
              3 115.173244 145.924982
                                         0.875
```

**Conclusion using Classification Accuracy** 

For this test set there is no much difference between these models in tersm of classification accuracy, though the model with smallest AIC had the highest accuracy.

### Other Predictive Power Statistics: Sensitivity, specificity

Accuracy is a blunt measure that depends on the overall fraction of each category as well as the sensitivity and specificity. We can break out the component sensitivity and specificity as illustrated in the previous section.

Here's a function used in the previous section for that purpose, modified to include accuracy, and to return a single row data frame.

```
In [55]:
                  from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score
                  def senspec(y, score, thresh, index=0):
In [56]:
           H
               1
                      yhat = 1*(score >= thresh)
               2
               3
                      tn, fp, fn, tp = confusion matrix(y true=y, y pred=yhat).ravel()
               4
                      sens = tp / (fn + tp)
               5
                      spec = tn / (fp + tn)
               6
                      accuracy = (tn+tp)/(tn+fp+fn+tp)
               7
                      return pd.DataFrame({'tn':[tn],
               8
                                             'fp':[fp],
               9
                                             'fn':[fn],
              10
                                             'tp':[tp],
                                             'sens':[sens],
              11
              12
                                             'spec':[spec],
                                             'accuracy':[accuracy]})
              13
                  # sensitivy and specificity for Model 0
In [57]:
               1
                  senspec(df_test['y'], phat_matrix[0], 0.5)
    Out[57]:
                 tn fp fn tp
                                  sens
                                          spec accuracy
                    1 4 13 0.764706 0.956522
              0 22
                                                   0.875
```

#### Out[58]:

_		tn	fp	fn	tp	sens	spec	accuracy
-	0	22	1	4	13	0.764706	0.956522	0.875
	1	23	0	3	14	0.823529	1.000000	0.925
	2	22	1	4	13	0.764706	0.956522	0.875
	3	22	1	4	13	0.764706	0.956522	0.875

## Model Comparison Using Sensitivity, Specificity, and Accuracy for all 4 Models.

The accuracy, sensitivity and specificity are all better for Model 1 (minimum BIC model) versus the others.

#### Remark on the selection of variables

From the simulation model we know that all 20 variables had some nonzero population coefficients, so why are some not significant? And why are they removed by the AIC/BIC criteria?

- First, note that the effect sizes for variables X11-X20 are small compared to the effects of X1-X10. With the sample size of 200, small effects often are not statistically significant due to the large standard errors compared to the estimates. We don't have enough power to detect those small effects.
- Second, the predictive performance of the model can be sometimes be improved by removing seemingly significant variables due to the reduced burden of estimation. Having fewer coefficients to estimate can decrease variance and improve mean square error for prediction as long as we retain enough highly informative variables.

STAT 207, Victoria Ellison and Douglas Simpson, University of Illinois at Urbana-Champaign

n [ ]: <b>N</b> 1			
n [ ]: ▶ 1			
In [ ]: ▶ 1		s m	-
	In I I •	N	1