

Logistic Regression

CS 530
Chapman
Spring 2021

Take message for rest of course:
Logistic Regression

- Logistic regression, lies at the intersection of regression and classification
- We will introduce the ubiquitous maximum-likelihood estimation
- We will introduce important notions in binary classification:
 - Confusion matrix
 - ROC curve
 - Pseudo- R^2 measures (McFadden's)
 - Does data trend according to logistic regression model?

1

2

TABLE OF CONTENTS

- Linear and Logistic Regression
 - Continuous and Categorical Data
- Maximum Likelihood: Fitting Logistic Regression
- Advantages & Disadvantages of Logistic Regression
- Multiple & Multinomial Logistic Regression
- Model Fit
 - Confusion Matrix: Accuracy, Sensitivity, Specificity
 - ROC Curve
 - McFadden's Pseudo-R²
 - Wald Test

Intro to Logistic Regression

After covering regression, we move on to classification with Logistic Regression, at the intersection of regression and classification.

Regression models relate a continuous (or almost continuous), dependent variable and other, similarly continuous, variables.

But how do we model relations between continuous variable & categorical variables (e.g., gender, species, alma mater, hometown)?

3

4

Intro to Logistic Regression

Example: Probability of defaulting on credit card payment given card balance

Probability of Default

Balance

[Image Source](#)

Intro to Logistic Regression

The x axis is like simple linear-regression. But the y axis is dichotomous (or, in this case, binary). So, we see that linear regression does not (and cannot) work well. The linear-regression model is

$$p(X) = \beta_0 + \beta_1 X$$

And, even if we insisted on fitting a linear regression, how can we interpret the resulting model? We cannot say that defaulting, as a binary condition, increases with the balance.

Logistic regression is a considerably better model for these kinds of data.

Probability of Default

Balance

[Image Source](#)

5

6

What is Logistic Regression?

To understand logistic regression, we need to understand the notion of *odds*.

The odds of an event, x , happening are defined as

$$\text{odds}(x) = \frac{p(x)}{1 - p(x)}.$$

Here $p(x)$ is the probability that the event occurs.

For one fair die, the odds of the result 3 is

$$\text{odds}(x = 3) = \frac{p(x = 3)}{1 - p(x = 3)} = \frac{1/6}{5/6} = \frac{1}{5}.$$

This is often reported as 1:5.



7

What is Logistic Regression?

Fitting a **linear regression** above, we use the model

$$p(X) = \beta_0 + \beta_1 X,$$

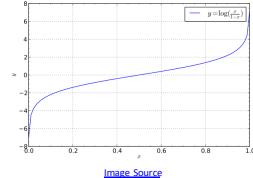
where $p(X)$ is the probability of an event occurring.

With **logistic regression**, we set the logarithm of the odds of X (also known as log-odds) equal to a linear function of X :

$$\text{logit}(p) = \log\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X.$$

The logit function is the inverse of the logistic function.

The logit function.



[Image Source](#)

8

What is Logistic Regression?

We can simplify this to a more explicit expression for $p(X)$:

$$\log\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X$$

$$\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$$

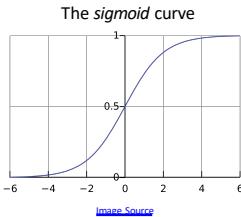
$$p(X) = [1 - p(X)]e^{\beta_0 + \beta_1 X}$$

$$p(X) = e^{\beta_0 + \beta_1 X} - p(X)e^{\beta_0 + \beta_1 X}$$

$$p(X)[1 + e^{\beta_0 + \beta_1 X}] = e^{\beta_0 + \beta_1 X}$$

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

This is the logistic function.

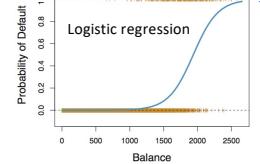
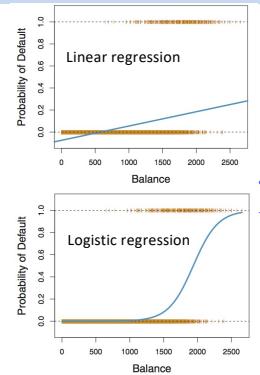


[Image Source](#)

What is Logistic Regression?

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

Applying logistic regression to the defaulting data above, we get the bottom right figure. The logistic regression plot, which fits a sigmoid function to the data, makes much more sense and is more easily interpretable.



[Image Source](#)

[Image Source](#)

9

What is Logistic Regression?

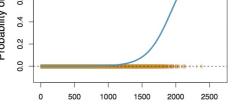
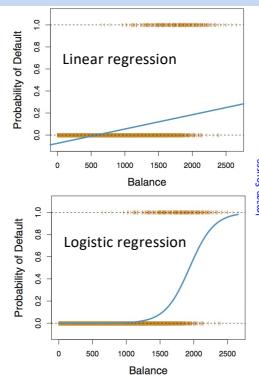
For linear regression, a change of 1 unit in x results in a change of β_1 units in y , regardless of x 's value:

$$p(X) = \beta_0 + \beta_1 X$$

For logistic regression, a change of 1 units in x results in a change of β_1 units in the log-odds of x , or a change of $e^{\beta_1 X}$ in the odds of x . So the rate of change of x changes with x :

$$\log\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X$$

$$\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$$



[Image Source](#)

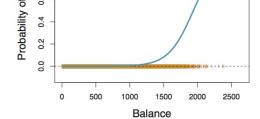
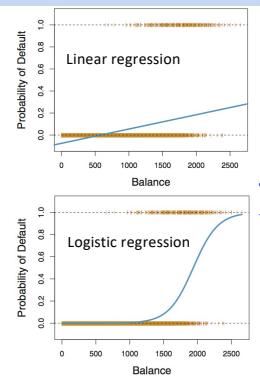
10

What is Logistic Regression?

How do we find optimal β_0 and β_1 for logistic regression?

For linear regression we used least squares on the RSS. We could use (non-linear) least squares here too.

But we will use a more general approach: *maximum (log) likelihood*



[Image Source](#)

11

12

Probability and likelihood in statistics

In statistics, Probability is used to describe future events, before the data become available.

$$P(\text{die}_1 = 6 \wedge \text{die}_2 = 6)$$



Likelihood is used after the data are available to fit the parameters of a given model to those data.

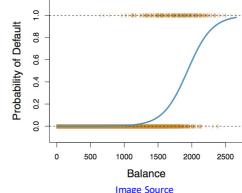
E.g.: Likelihood of fair dice given (6,6)×3



Maximum Likelihood: Fitting Logistic Regression

Given our available data, we estimate β_0 and β_1 such that $\hat{p}(x_i)$, for every trial i , will be as close as possible to the actual $p(x_i)$. In other words, we want $\hat{\beta}_0$ and $\hat{\beta}_1$ that

result in $\hat{p}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}}$ that is as close to 0 for no defaulting and as close to 1 for defaulting as possible.



13

Maximum Likelihood: Fitting Logistic Regression

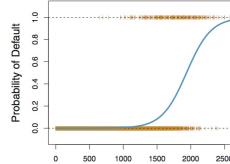
Given our available data, we estimate β_0 and β_1 such that $\hat{p}(x_i)$, for every trial i , will be as close as possible to the actual $p(x_i)$. In other words, we want $\hat{\beta}_0$ and $\hat{\beta}_1$ that

result in $\hat{p}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}}$ that is as close to 0 for no defaulting and as close to 1 for defaulting as possible.

This intuition can be formalized by the likelihood function:

$$\ell(\beta_0, \beta_1; \vec{x}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j)),$$

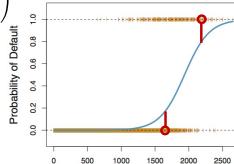
which we want to maximize.



14

Maximum Likelihood: Fitting Logistic Regression

$$\ell(\beta_0, \beta_1; \vec{x}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j)) \\ = \prod_{i:y_i=1} \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \prod_{j:y_j=0} \left(1 - \frac{e^{\beta_0 + \beta_1 x_j}}{1 + e^{\beta_0 + \beta_1 x_j}}\right)$$



15

Maximum Likelihood: Fitting Logistic Regression

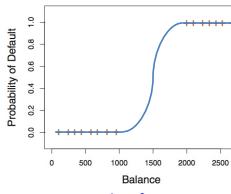
$$\ell(\beta_0, \beta_1; \vec{x}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j))$$

If the logistic regression curve fits the data perfectly, every term in every product of the likelihood function will be 1, giving us

$$\ell(\beta_0, \beta_1; \vec{x}) = 1.$$

This would maximize the likelihood function in our case.

Maximum likelihood is a very general method for parameter estimation for non-linear models. In fact, the least-squares approach for linear regression is a special case of maximum likelihood estimation.



16

Maximum Likelihood: Fitting Logistic Regression

$$\ell(\beta_0, \beta_1; \vec{x}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j))$$

Maximum likelihood is a very general method for parameter estimation for non-linear models. In fact, the least-squares approach for linear regression is a special case of maximum likelihood estimation.

17

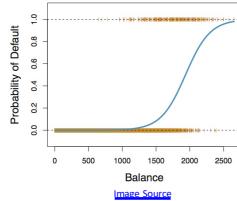
18

Maximum Likelihood: Fitting Logistic Regression

$$\ell(\beta_0, \beta_1; \vec{x}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j))$$

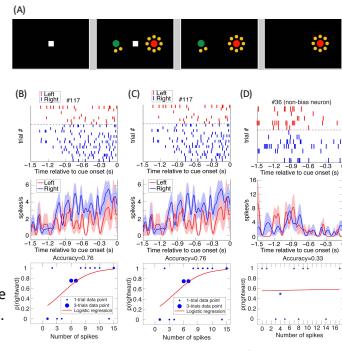
It is typically simpler to work with the log-likelihood function (numerical stability & analytic derivation), which turns products into sums:

$$\begin{aligned} L(\beta_0, \beta_1; \vec{x}) &= \ln[\ell(\beta_0, \beta_1; \vec{x})] \\ &= \ln \left[\prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j)) \right] \\ &= \ln \left[\prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i} \right] \\ &= \sum_{i=1}^n [y_i \ln(p(x_i)) + (1 - y_i) \ln(1 - p(x_i))] \end{aligned}$$



[Image Source](#)

Using Logistic Regression



Maoz et al., *Front Neurosci.*, 2013

19

Pros and Cons of Logistic Regression

What are the advantages of using logistic regression over linear regression for binary classification?

- Logistic function only gives values between 0 and 1
 - Logistic function “transitions quickly” from 0 to 1
 - Logistic regression does not require assumptions about normally distributed errors or homoscedasticity of errors
- Disadvantages?
- Logistic regression is an iterative process without a closed-form solution. So, it may not converge due to
 - completely separable classes
 - multicollinearity
 - sparse data
 - Logistic regression becomes cumbersome if we consider more than 2 categories in the dependent variable

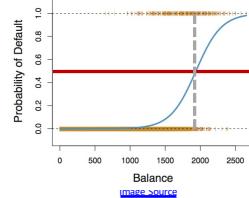
20

Prediction with Logistic Regression

Prediction based on a fitted logistic regression model—i.e., for given balance, how do we predict defaulting, a binary event?

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

The model returns the probability of default. We can then set a threshold (e.g., $p(X)=0.5$) and from that make a prediction (e.g., iff $\text{balance} > 1800$ predict default).



[Image Source](#)

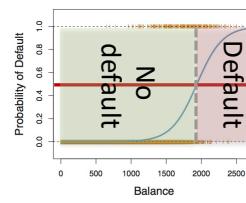
21

Prediction with Logistic Regression

Prediction based on a fitted logistic regression model—i.e., for given balance, how do we predict defaulting, a binary event?

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

The model returns the probability of default. We can then set a threshold (e.g., $p(X)=0.5$) and from that make a prediction (e.g., iff $\text{balance} > 1800$ predict default).



[Image Source](#)

22

Multiple Logistic Regression

What if we wanted to find the probability of default based on more variables: balance, credit history, & income? We could use multiple logistic regression.

With multiple logistic regression, we just add more terms to the linear expression (on the right-hand side):

$$\ln \left(\frac{p(\vec{X})}{1 - p(\vec{X})} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n = \beta_0 + \vec{\beta} \cdot \vec{X}$$

for $\vec{X} = (X_1, X_2, \dots, X_n)$ and $\vec{\beta} = (\beta_1, \beta_2, \dots, \beta_n)$. And similar algebra to before would give us:

$$p(\vec{X}) = \frac{e^{\beta_0 + \vec{\beta} \cdot \vec{X}}}{1 + e^{\beta_0 + \vec{\beta} \cdot \vec{X}}}.$$

And, again, finding good $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_n$ can be achieved using the maxim (log) likelihood method.

23

24

Multinomial Logistic Regression

What if we wanted to create a predictive model of mortgage payment from income, where our dependent variable, y , is categorical, but with more than two choices?

For example, say our subjects can pay (P), short sale (S), or foreclosure (F). And we had to predict which they will move.

This is termed multinomial logistic regression.



Multinomial Logistic Regression

For multinomial logistic regression, we pick a “pivot” outcome for y . In our example, we can pick $y = P$. Multinomial logistic regression then starts with equations for the non-pivot outcomes for y :

$$\ln \left[\frac{Pr(y = S)}{Pr(y = P)} \right] = \vec{\beta}_0^1 + \vec{\beta}^1 \cdot \vec{X},$$

$$\ln \left[\frac{Pr(y = F)}{Pr(y = P)} \right] = \vec{\beta}_0^2 + \vec{\beta}^2 \cdot \vec{X},$$

where the variables x_1, \dots, x_q represent incomes of different people, ...

25

Multinomial Logistic Regression

Modifying these equations, we get:

$$Pr(Y = S) = Pr(Y = P) e^{\vec{\beta}_0^1 + \vec{\beta}^1 \cdot \vec{X}},$$

$$Pr(Y = F) = Pr(Y = P) e^{\vec{\beta}_0^2 + \vec{\beta}^2 \cdot \vec{X}}$$

But,

$$Pr(Y = S) + Pr(Y = F) + Pr(Y = P) = 1.$$

So

$$Pr(Y = P) e^{\vec{\beta}_0^1 + \vec{\beta}^1 \cdot \vec{X}} + Pr(Y = P) e^{\vec{\beta}_0^2 + \vec{\beta}^2 \cdot \vec{X}} + Pr(Y = P) = 1.$$

And we get

$$Pr(Y = P) = \frac{1}{e^{\vec{\beta}_0^1 + \vec{\beta}^1 \cdot \vec{X}} + e^{\vec{\beta}_0^2 + \vec{\beta}^2 \cdot \vec{X}} + 1}.$$

26

Multinomial Logistic Regression

Multinomial logistic regression can similarly be extended beyond 3-class classification, to categorical variable y with any finite number of values.



However, a more popular approach for multi-class classification is LDA. That said, LDA's assumption of Gaussian distributions for all classes is stronger than logistic regression's assumption that the log-odds of \vec{X} (the variables) are a linear function of \vec{X} (those variables).

27

Questions from exit survey

- How are logistic and linear regression different?
- Multiple vs. multinomial logistic regression. How are they different? Why do we need both?

28

Model Fit: Accuracy

How do we test the goodness of fit for the logistic regression model?

The simplest and most intuitive way is to check the accuracy of the model. We set a threshold and check how often model predictions line up with our training data.

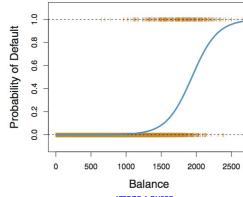


29

30

Model Fit: Accuracy

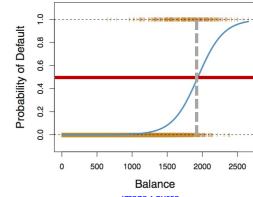
For the original default data and a decision threshold of 0.5



Model Fit: Accuracy

For the original default data and a decision threshold of

$$P(\text{Default}) = 0.5,$$



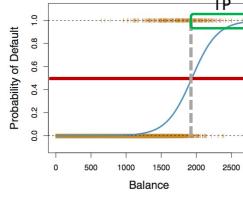
31

Model Fit: Accuracy

For the original default data and a decision threshold of

$$P(\text{Default}) = 0.5,$$

True positive (TP): Predicted default and default occurred.



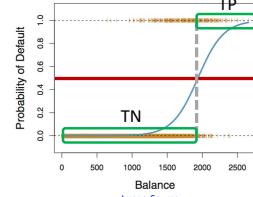
Model Fit: Accuracy

For the original default data and a decision threshold of

$$P(\text{Default}) = 0.5,$$

True positive (TP): Predicted default and default occurred.

True negative (TN): Predicted no default and no default occurred.



32

Model Fit: Accuracy

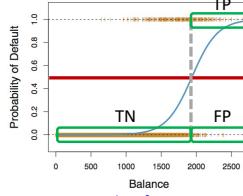
For the original default data and a decision threshold of

$$P(\text{Default}) = 0.5,$$

True positive (TP): Predicted default and default occurred.

True negative (TN): Predicted no default and no default occurred.

False positive (FP): Predicted default but no default occurred.



Model Fit: Accuracy

For the original default data and a decision threshold of

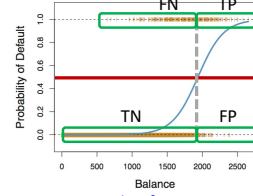
$$P(\text{Default}) = 0.5,$$

True positive (TP): Predicted default and default occurred.

True negative (TN): Predicted no default and no default occurred.

False positive (FP): Predicted default but no default occurred.

False negative (FN): Predicted no default but default occurred.



35

36

Model Fit: Confusion Matrix

More generally, we can describe accuracy for *binary classification* with a confusion matrix

False Positive = Type I Error

False Negative = Type II Error

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

	Predict True	Predict False
Actual True	TP	FN
Actual False	FP	TN

Model Fit: Confusion Matrix

Other measures stemming from the confusion matrix include:

True Positive Rate (hit rate, recall, sensitivity):

$$\text{TPR} = \frac{TP}{TP + FN}.$$

True Negative Rate (specificity, selectivity):

$$\text{TNR} = \frac{TN}{TN + FP}.$$

False Positive Rate (fall out, false alarm):

$$\text{FPR} = \frac{FP}{TN + FP} = 1 - \text{TNR}.$$

TPR and TNR: higher values mean better model fit.

37

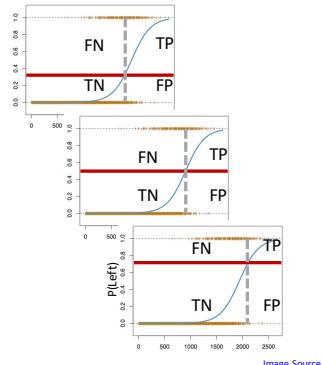
Model Fit: Confusion Matrix

TPR, TNR, & FPR all tradeoff and depend on the threshold we set.

$$\text{TPR} = \frac{TP}{TP+FN} \text{ (sensitivity)}$$

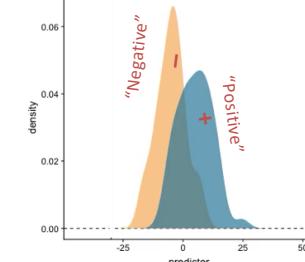
$$\text{TNR} = \frac{TN}{FP+TN} \text{ (specificity)}$$

$$\text{FPR} = \frac{FP}{TN+FP} = 1 - \text{TNR}$$



38

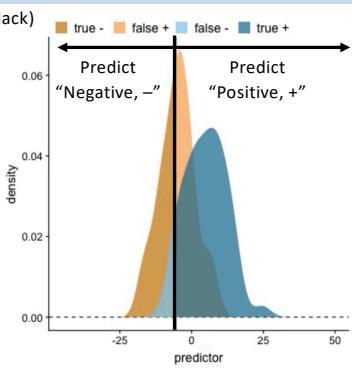
Model Fit: ROC Curve



39

Model Fit: ROC Curve

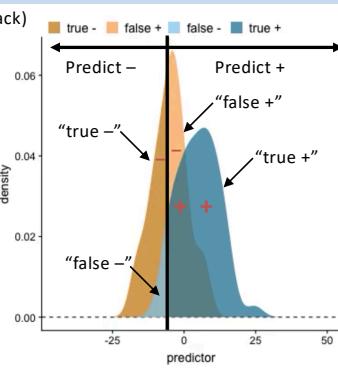
Predicted (in black)
Actual (in red)



40

Model Fit: ROC Curve

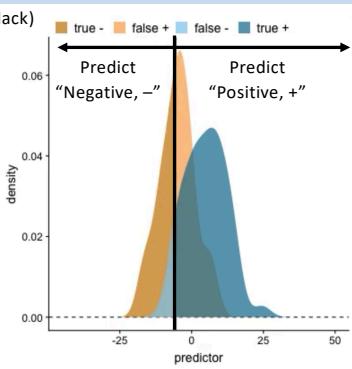
Predicted (in black)
Actual (in red)



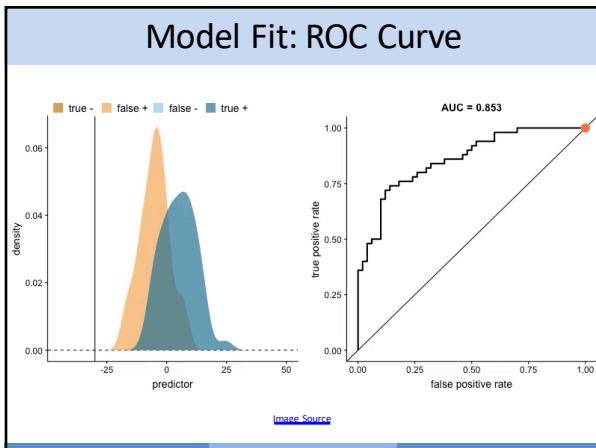
41

Model Fit: ROC Curve

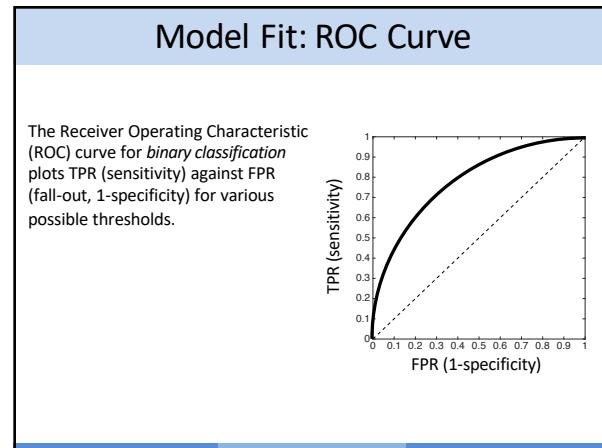
Predicted (in black)
Actual (in red)



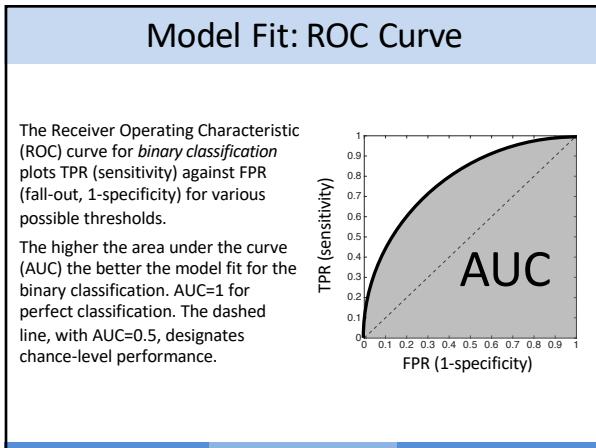
42



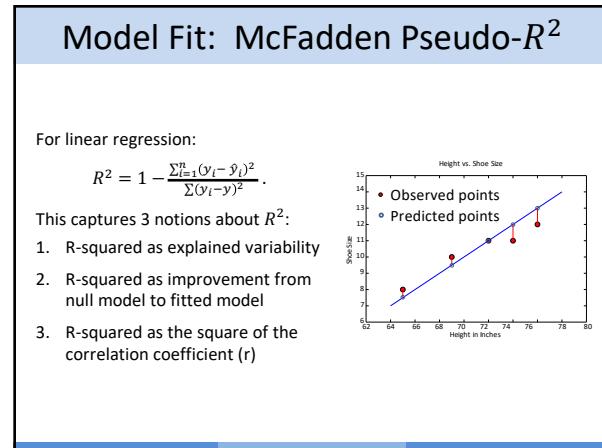
43



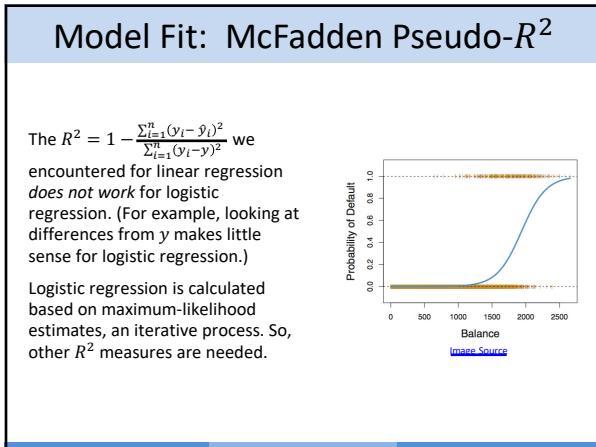
44



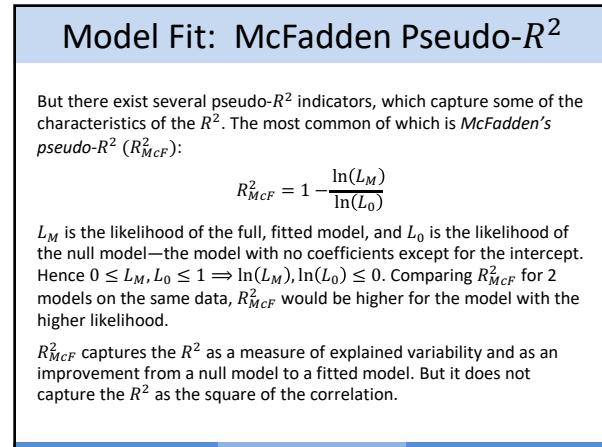
45



46



47



48

Model Fit: Wald Test

How statistically significant is each coefficient estimate $\hat{\beta}_j$ in our model? For linear regression we carried out t-tests on the variables. For logistic regression we can run a similar test based on the *z-statistic*, named the Wald test:

$$W_j = \frac{\hat{\beta}_j}{SE_{\hat{\beta}_j}},$$

where SE is the standard error of the mean. $|W_j| > 2$ (roughly) is significant at the 95% level.

49

Python code

```

1 from sklearn.datasets import load_iris
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LogisticRegression
4 X, y = load_iris(return_X_y=True)
5 clf = LogisticRegression(random_state=0, solver='lbfgs',
6                         multi_class='multinomial').fit(X, y)
7
8 # solver: algorithm to use in the optimization problem.
# For small datasets, 'liblinear' is a good choice, whereas 'sag' and 'saga'
# are faster for large ones.
9
10 # For multiclass problems, only 'newton-cg', 'sag', 'saga' and 'lbfgs'
# handle multinomial loss; 'liblinear' is limited to one-versus-rest schemes.
11
12
13
14
15
16
17 clf.predict(X[1:2, :])
array([0, 0])

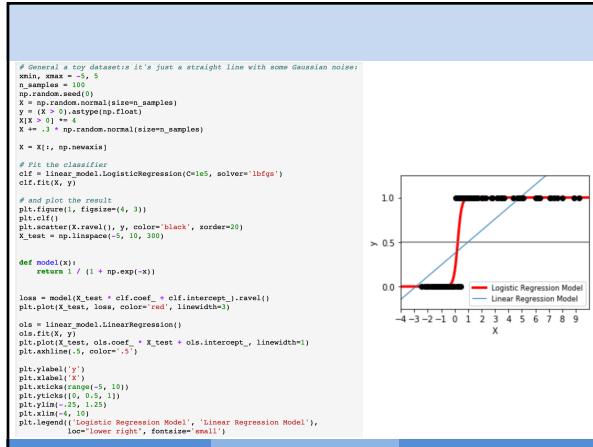
1 clf.predict_proba(X[1:2, :]) # The returned estimates for all classes are ordered by the label of classes.
2
3
4
5

array([[ 9.81805980e-01,  1.81940056e-02,  1.43354257e-08],
       [ 9.71803214e-01,  2.81967563e-02,  2.97462927e-08]])

1 clf.score(X, y) # Returns the mean accuracy on the given test data and labels.
0.9733333333333338

```

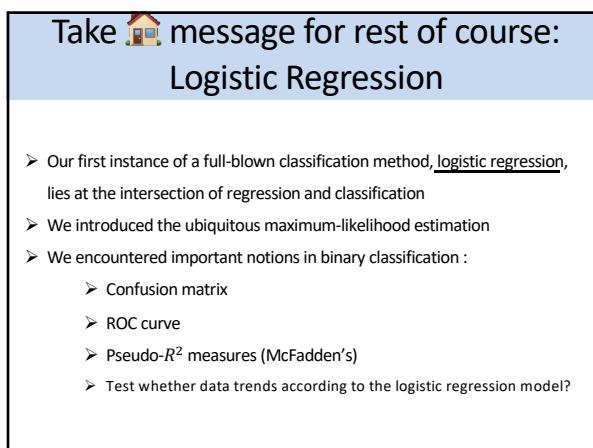
50



Summary

- For Logistic Regression $\log\left(\frac{p(X)}{1-p(X)}\right) = \beta_0 + \beta_1 X$ or $p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$. It is therefore a classification method rather than strictly a regression.
- Logistic regression is fit iteratively using the maximum likelihood method.
- It has extensions for multiple variables (multiple logistic regression) and multiple classes (multinomial logistic regression).
- Discussing model accuracy for binary classification, we encountered
 - Confusion Matrix: Accuracy, Sensitivity, Specificity
 - ROC Curve
 - McFadden Pseudo-R2
 - Wald Test

51



52