

- 3.4. Relationship between bias, variance, overfitting, underfitting, and mean squared error of a model
- 3.5. Goal of selecting a model that will make good predictions on new data
- 4. Goal: Find a Parsimonious Model
- 5. More about Fitting a Logistic Regression Model
 - 5.1. How are the optimal values of $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ determined in a logistic regression model?
 - 5.2. Where do we find the optimal log-likelihood function value for a given logistic regression model?
- 6. Model Selection with Log Likelihood Ratio Test
- 7. Model Selection with AIC and BIC

Additional Resources

Sections 5.1, 6.1, and 6.2 Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. An Introduction to Statistical Learning: with Applications in R. New York :Springer, 2013. https://www.ime.unicamp.br/~dias/Intoduction%20to%20Statistical%20Learning.pdf

1. <u>PROBLEM</u>: WHAT ARE SOME ALGORITHMS THAT CAN HELP US <u>EFFICIENTLY</u> IDENTIFY MODELS THAT HAVE "HIGH PARSIMONY"?

1. <u>Review</u>: What methods have learned so far that help us determine which set of explanatory variables should be included in the "most parsimonious" logistic regression model?

2. <u>Review</u>: In general, why would we want a model that is parsimonious?

3. <u>Total Number of Possible Models</u>: If we have *P* available explanatory variables to model in a logistic regression model, how many possible logistic regression models can we make and evaluate whether they are the "most parsimonious"?

4. <u>Unit 20 Problem</u>: What we will explore in this unit is some algorithms that can help us *efficiently* identify which explanatory variables to include in a regression model such that the resulting model will have **high parsimony**?

2. <u>One Idea</u>: Use Backwards Elimination or Forward Selection Algorithms

We can use what we call a **backwards elimination algorithm** to help us select a good combination of explanatory variables to include in *a model* that achieve some goal we are trying to reach. For instance, what we fill in the **black** blanks below is how we can use a backwards elimination algorithm that finds a model with

_____, however you can use this same structure select a model that optimizes some other metric as well.

Backwards Elimination Algorithm

Goal: Find the model with the ______.

Steps:

1. Fit a "current model" and find the _____ of this model.

In the beginning, your "current model" should include ______ possible explanatory variables you are

considering.

2. For each explanatory variable in your "current model" do the following:

a. Fit a "test model".

Your "test model" should be every explanatory variable in the "current model" ______ the

explanatory variable you are considering.

- b. Find the ______ of this "test model".
- 3. If NONE of the "test models" from step (2) had a ______ that was ______ than the

______ for the "current model", then STOP THE ALGORITHM, and return the "current

model" as your "final model".

4. Otherwise, choose the "test model" from step (2) that had the ______

and set your new "current model" to be this "test model". Then go back to step (2).

We can use what we call a **forward selection algorithm** to help us select a good combination of explanatory variables to include in *a model* that achieve some goal we are trying to reach. For instance, what we fill in the **black** blanks below is how we can use a backwards elimination algorithm that finds a model with

_____, however you can use this same structure select a model that optimizes some other

metric as well.

Forward Selection Algorithm

Goal: Find the model with the ______

Steps:

1. Fit a "current model" and find the _____ of this model.

In the beginning, your "current model" should include ______ possible explanatory variables you are considering.

2. For each explanatory variable in your "current model" do the following:

a. Fit a "test model".

Your "test model" should be every explanatory variable in the "current model" ______ the

explanatory variable you are considering.

- b. Find the _____ of this "test model".
- 3. If NONE of the "test models" from step (2) had a ______ that was ______ than the

______ for the "current model", then STOP THE ALGORITHM, and return the "current

model" as your "final model".

4. Otherwise, choose the "test model" from step (2) that had the ______,

and set your new "current model" to be this "test model". Then go back to step (2).

Caveat: Forward selection algorithms and backwards elimination algorithms that try to find a model with a
______, are not guaranteed to find the model with the ______
____out of all possible models!

Pros of Using Backwards Elimination (or Forward Selection) to Help find a Parsimonious Model

Cons of Using Backwards Elimination (or Forward Selection) to Help find a Parsimonious Model

See the exercise in the Jupyter notebook section 2.

3. ANOTHER IDEA: USE A REGULARIZATION TERM IN YOUR REGRESSION MODEL

We can also use a *single* model with a **regularization term** to help us select a good combination of explanatory variables to include in *a model* that achieve some goal we are trying to reach.

3.1 <u>Recap</u> of Non-Regularized Logistic Regression Models and AIC/BIC

Non-Regularized Regression Objective Function

1. <u>Main Problem</u>: Recall from previous units, when we fit a logistic regression model (ie. a **non-regularized logistic regression model)** we seek to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that ______ the log likelihood function LLF().

$$\begin{split} \max_{\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}} LLF(\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}) &= \max_{\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}} \{y_{1} \log(p_{1}) + (1 - y_{1}) \log(1 - p_{1}) + \dots + y_{n} \log(p_{n}) + (1 - y_{n}) \log(1 - p_{n}) \}, \\ \end{split}$$

$$\begin{aligned} & \text{Where, } p_{i} = \frac{e^{\hat{\beta}_{0} + \hat{\beta}_{1} x_{i1} + \dots + \hat{\beta}_{p} x_{ip}}{1 + e^{\hat{\beta}_{0} + \hat{\beta}_{1} x_{i1} + \dots + \hat{\beta}_{p} x_{ip}} \end{split}$$

We call the LLF the **objective function** as it is the function we are trying to optimize (ie. either maximize or minimize) in a given problem.

2. **Equivalent Problem**: This problem is equivalent to trying to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that _______ "-2" times the log likelihood function, ie. -2*LLF().

$$\begin{split} \min_{\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}} - 2 \cdot LLF(\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}) &= \min_{\hat{\beta}_{0}, \hat{\beta}_{1}, \dots, \hat{\beta}_{p}} - 2 \cdot [y_{1} \log(p_{1}) + (1 - y_{1}) \log(1 - p_{1}) + \dots + y_{n} \log(p_{n}) + (1 - y_{n}) \log(1 - p_{n})], \\ \end{split}$$

$$\begin{aligned} & \text{Where, } p_{i} = \frac{e^{\hat{\beta}_{0} + \hat{\beta}_{1} X_{i1} + \dots + \hat{\beta}_{p} X_{ip}}{1 + e^{\hat{\beta}_{0} + \hat{\beta}_{1} X_{i1} + \dots + \hat{\beta}_{p} X_{ip}} \end{split}$$

Model Selection

Also, recall from previous units that the AIC for a given model (that has already been _____) is calculated as:

$$AIC = -2 \cdot LLF + 2 \cdot p.$$

Parsimonious Models

Ideally, we want the AIC of an ______ to be _____, which we can get when the:

- LLF is ______, and the
- p (ie. number of slopes) is _____.

Penalty Terms

We call ______ the **penalty term** in the AIC equation, as the ______ slopes we have in the model, the more the AIC function for the model is penalized.

Main Idea Behind Regularization Terms

Use an **objective function** that tries to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that:

- 1. _____"-2" times the log likelihood function, (ie. -2*LLF()) and
- 2. also uses a **penalty term** (ie. a **regularization term)** that penalizes models with a large number of ______

slopes for the explanatory variables.

 $\min_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p} -2 \cdot LLF(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) + penalty term$

Main Types of Regularization Terms

There are three common types of regularization terms that we will discuss in this unit.

- 1. Ridge Regression regularization term (L2 Penalty)
- 2. LASSO regularization term (L1 Penalty)
- 3. Elastic net regularization term (Combination of L2 and L1 Penalty)

3.2 <u>REGULARIZATION TERM 1:</u> RIDGE REGRESSION (L2 PENALTY)

In a **Logistic Ridge Regression model**, we seek to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that minimize the following objective function (for a given $\lambda > 0$).

$$\underset{\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}}{\text{MIN}} - 2 \cdot LLF(\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}) + \lambda(\hat{\beta}_{1}^{2} + \hat{\beta}_{2}^{2} + \cdots + \hat{\beta}_{p}^{2})$$

The best solutions select values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$, in which

- $LLF(\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p)$ is _____, which means the model will be encouraged to have
- $\hat{\beta}_1^2 + \hat{\beta}_2^2 + \dots + \hat{\beta}_p^2$ is _____, which means the model will be encouraged to have _____.

A higher value of λ means

• The model will be ______ likely to select values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ in which $\hat{\beta}_1^2 + \hat{\beta}_2^2 + \cdots +$

 $\hat{\beta}_p^2$ is low, at the expense of creating models with high value for $LLF(\hat{\beta}_{0,\hat{\beta}_1}, ..., \hat{\beta}_p)$.

3.3 <u>REGULARIZATION TERM 2:</u> LASSO (L1 PENALTY)

In a **LASSO logistic regression model**, we seek to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that minimize the following objective function (for a given $\lambda > 0$).

$$\underset{\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}}{\text{MIN}} - 2 \cdot LLF(\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}) + \lambda(|\hat{\beta}_{1}| + |\hat{\beta}_{2}| + \cdots + |\hat{\beta}_{p}|)$$

The best solutions select values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$, in which

- $LLF(\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p)$ is _____, which means the model will be encouraged to have
- $|\hat{\beta}_1| + |\hat{\beta}_2| + \dots + |\hat{\beta}_p|$ is ______, which means the model will be encouraged to have ______.

A higher value of λ means

• The model will be ______ likely to select values of $\hat{\beta}_{0,}\hat{\beta}_{1},...,\hat{\beta}_{p}$ in which $|\hat{\beta}_{1}| + |\hat{\beta}_{2}| + \cdots + |\hat{\beta}_{p}|$ is low, at the expense of creating models with high value for $LLF(\hat{\beta}_{0,}\hat{\beta}_{1},...,\hat{\beta}_{p})$.

3.4 COMPARING LASSO (L1) PENALTY TO THE RIDGE REGRESSION (L2) PENALTY

Why would you want to use ridge regression penalty over a LASSO penalty, and vice versa?

Ridge Regression (L2) Penalty $\widehat{m{eta}}_1^2 + \widehat{m{eta}}_2^2 + \cdots + \widehat{m{eta}}_p^2$

LASSO (L1) Penalty $|\widehat{\beta}_1| + |\widehat{\beta}_2| + \dots + |\widehat{\beta}_p|$

3.5 <u>REGULARIZATION TERM 3:</u> ELASTIC NET (L2 AND L1 PENALTY COMBINATION)

In an **Elastic Net Logistic regression model**, we seek to find the optimal values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$ that minimize the following objective function (for a given $\lambda > 0$ and $0 \le \alpha \le 1$)

$$\underset{\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}}{\text{MIN}} - 2 \cdot LLF\left(\hat{\beta}_{0},\hat{\beta}_{1},\ldots,\hat{\beta}_{p}\right) + \lambda \left[\alpha\left(\left|\hat{\beta}_{1}\right| + \left|\hat{\beta}_{2}\right| + \cdots + \left|\hat{\beta}_{p}\right|\right) + \frac{1-\alpha}{2}\left(\hat{\beta}_{1}^{2} + \hat{\beta}_{2}^{2} + \cdots + \hat{\beta}_{p}^{2}\right)\right]$$

The best solutions select values of $\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p$, in which

- $LLF(\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p)$ is _____, which means the model will be encouraged to have
- $|\hat{\beta}_1| + |\hat{\beta}_2| + \dots + |\hat{\beta}_p|$ AND $(\hat{\beta}_1^2 + \hat{\beta}_2^2 + \dots + \hat{\beta}_p^2)$ are _____, which means the model will be encouraged to have _____.

A α closer to 1 means

• The model will be ______ likely to select values of $\hat{\beta}_{0,}\hat{\beta}_{1},...,\hat{\beta}_{p}$ that resemble solutions to the

A α closer to 0 means

• The model will be ______ likely to select values of $\hat{\beta}_{0,}\hat{\beta}_{1},...,\hat{\beta}_{p}$ that resemble solutions to the

A higher value of λ means

• The model will be ______ likely to select values of $\hat{\beta}_{0,}\hat{\beta}_{1}, \dots, \hat{\beta}_{p}$ in which $|\hat{\beta}_{1}| + |\hat{\beta}_{2}| + \dots + |\hat{\beta}_{p}|$ AND $(\hat{\beta}_{1}^{2} + \hat{\beta}_{2}^{2} + \dots + \hat{\beta}_{p}^{2})$ are low, at the expense of creating models with high value for $LLF(\hat{\beta}_{0,}\hat{\beta}_{1}, \dots, \hat{\beta}_{p})$.

3.5 How to Use Regularized Regression to select a Parsimonious Model

- 1. Fit a *single* regression model with a **regularization term.**
- 2. [OPTIONAL STEP]: Examine the resulting slopes of this this fitted model with **regularization term**, and use this to determine generate some ideas of possible "reduced models" to test.
- 3. [OPTIONAL STEP]: Fit each of these "reduced models" as well as the "full model" (ie. the model with all possible explanatory variables you are considering) and use the AIC, BIC, and/or log-likelihood ratio test to determine which model is best.

Go to the Unit 20 notebook, section 3 for an exercise in using regularized logistic regression models.