

Classification with Logistic Regression

Chad Wakamiya Spring 2020

Berkeley SCET

Agenda

Classification

Introduction to types of classification and set up.

Logistic Regression

The logistic regression formula and intuition.

Multiclass Classification

Extending logistic regression for datasets with multiple features.



Classification

Classification

CET

Berkeley :

Classification is the problem of assigning observations to one or more categories.



Multiclass Classification

Involves more than 2 classes



Examples

Binary

- Spam detection
- Churn/no churn customer retention
- Develop Diabetes/don't
- Default/repay loan

Multiclass

- Image recognition
- Natural language • processing

Classification

We have **features** and **labels** data (X, Y):

 $(x_1, y_1) \\ (x_2, y_2) \\ \dots \\ (x_n, y_n)$ Features Actu

Actual Labels

- Features: \boldsymbol{x}_i is a vector (or even matrix) for each data element

• For a picture: $\mathcal{X}_i = [32 \times 32 \times 3]$: array of numbers

• Actual Labels: $y_i \in \{0, 1\}$ \circ If picture i is a dog, $y_i = 1$ \circ If picture i is a cat, $y_i = 0$

Berkeley SCET



Logistic Regression

Logistic Function

The **logistic function** $\sigma(t)$ can be used to classify binary observations.

- \circ When t is large, $\sigma(t)
 ightarrow 1$
- \circ When t is small, $\sigma(t)
 ightarrow 0$

$$\sigma(t)=rac{1}{1+e^{-t}}$$



https://en.wikipedia.org/wiki/Sigmoid_function



Logistic Regression vs. Linear Regression



Logistic Regression vs. Linear Regression



Logistic Regression vs. Linear Regression



Logistic Regression



- Notice that $0 < \sigma(t) < 1$ for all real numbers t, so we can use the logistic function to model the *probability* that an observation belongs to a certain class.
- If $t = w_0 + w_1 x$ we can use the logistic function to write:

$$P(Y=1|x)=\sigma(t)=rac{1}{1+e^{-(w_0+w_1x)}}$$

Probability the image is a dog given features x

Example: $P(Y=1|x) = \frac{1}{2} \text{ when } w_0 + w_1 x = 0$

Berkeley SCET



Logistic Regression Threshold



Logistic Regression Example

1. **Problem:** Will student *i* pass given *i* studies x_i hours?

 $\hat{y}_i \!=\! \left\{ \begin{matrix} 1, & \text{if pass} \\ 0, & \text{if fail} \end{matrix} \right.$

- 2. **Model:** We use this curve to predict the probability that the student would pass given x hours of study
- 3. **Classify:** If Prob > 0.5, we predict the student will pass the exam



Data : $x =$ hours studied:	Hours	0.50	0.75	1.00	1.25	1.50	1.75	1.75	2.
y = pass/not pass:	Pass	0	0	0	0	0	0	1	0



Multi-feature Logistic Regression



The logistic function can be extended for multiple features (d dimensions) if we write:

$$t = x_i^T W = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

$$P(Y=1|x_i)=\sigma(x_i^TW)=rac{1}{1+e^{-(w_0+w_1x_{i,1}+w_2x_{i,2}+\cdots+w_dx_{i,d})}}$$

Probability the image is a dog given features x_i

Berkeley SCET

Linear Model

Training a Logistic Regression Classifier (Find *W*)

Steps to Train a Classifier Model

1. Choose our **model** to estimate y_i

 $f_{\scriptscriptstyle W}(x_{i}) = \hat{y}_{i} = \begin{cases} 1, & \text{if Prob} > \text{threshold} \\ 0, & \text{otherwise} \end{cases}$

- 2. Define a loss function (*L*) \rightarrow Allows for scoring each sample point (picture) x_i
- Optimize across the parameter space
 (W) to minimize the loss function to some small threshold



Goal: Find the best values for the model parameters W.



1. Choose a Model (Review)

We select a logistic regression model and set the threshold to 0.5:

$$f_{\scriptscriptstyle W}(x_{\scriptscriptstyle i})=\hat{y}_{\scriptscriptstyle i}= \begin{cases} 1, & \text{if} \ p_{\scriptscriptstyle i}>0.5\\ 0, \ \text{otherwise} \end{cases}$$

where:
$$p_i = P(Y=1|x_i) = \sigma(x_i^T W)$$

Probability the image is in class 1 given features x_i

CET

Berkeley

Next: Find the parameters $W = [w_0, w_1, w_2, \dots, w_d]$

2. Define a Loss Function

 $p_i = P(Y=1|x_i) = \sigma(x_i^T W)$

Logistic regression commonly uses the **Cross Entropy** loss function to score predictions:

$$L_i = -y_i \ln(p_i) - (1-y_i) \ln(1-p_i)$$

- If the predicted label is **wrong the loss is large** and if the predicted label is **right the the loss is small**.
- Since y_i is binary there are 2 cases:

I.
$$y_i = 0 \Rightarrow L_i = -\ln(1-p_i)$$

II. $y_i = 1 \Rightarrow L_i = -\ln(p_i)$
II. $y_i = 1 \Rightarrow p_i$ is near $0 \Rightarrow$ Loss is large
If $\hat{y}_i = 1 \Rightarrow p_i$ is near $1 \Rightarrow$ Loss is -0
 $L_i = \frac{1}{p_i} \int_{1}^{2} f(p_i) = -\ln(1-p_i)$
 $L_i = \frac{1}{p_i} \int_{1}^{1.5} f(p_i) = -\ln(p_i)$
 $L_i = \frac{1}{p_i} \int_{1}^{1.5} f(p_i) = -\ln(p_i)$
 $L_i = \frac{1}{p_i} \int_{1}^{1.5} f(p_i) = -\ln(p_i)$

3. Optimize Across the Parameter Space

• We want the *W* with the **lowest average loss** across all data points in our training or test set.

Average Loss
$$=rac{1}{n}\sum_{i=1}^n L_i$$
 .

• Formally this can be written as:

Berkeley SCET

$$W^* = \operatorname*{argmin}_W \underbrace{\frac{1}{n} \sum_{i=1}^n -y_i \ln(\sigma(x_i^T W)) - (1 - y_i) \ln(1 - \sigma(x_i^T W))}_{\mathsf{Average Loss}}$$

• With regularization, the average loss function is convex \Rightarrow solve for W^*

Logistic Regression Classifier Summary

1. Choose our **model** to estimate y_i

$$f_{W}(\mathbf{x}_{i}) = \begin{cases} 1, & \text{if } p_{i} > 0.5 \\ 0, & \text{otherwise} \end{cases}$$

$$p_i = \sigma(x_i^T W)$$

2. Define a loss function (*L*)

 \rightarrow Allows for scoring each sample point (picture) x_i

$$L_i=-y_i\ln(p_i)-(1-y_i)\ln(1-p_i)$$

3. Optimize across the parameter space (*W*) to minimize the loss function to some small threshold

$$W^* = \operatorname*{arg\,min}_W \frac{1}{n} \sum_{i=1}^n -y_i \ln(\sigma(x_i^T W)) - (1 - y_i) \ln(1 - \sigma(x_i^T W))$$

Berkeley SCET

Multiclass Classification

Multiclass Classification

Logistic regression can be applied to solve multiclass problems.

Common Approaches

Berkeley

- 1. One-vs-Rest (One-vs-All)
- 2. Softmax Regression (Multinomial Logistic Regression)



Involves more than 2 classes



One-vs-Rest (One-vs-All)

For each class, build a logistic regression to find the probability the observation belongs to that class. For each data point, predict the class with the highest probability.



Softmax Regression (Multinomial Logistic Regression)

In softmax regression the probability that a data point belongs to each class is calculated by:



If K=2, softmax regression reduces to the same binary logistic regression formulas we saw earlier. Check out this <u>overview of softmax regression</u> for the proof.

Berkeley

References

- Data^x (IEOR 135/290) Ikhlaq Sidhu and Arash Nourian)
 - The content presented in this lecture draws on materials by the IEOR 135 course instructors.
- Logistic Regression
 - <u>https://en.wikipedia.org/wiki/Logistic_regression</u>
- Data Science Principles and Techniques (DS 100 at UC Berkeley) Ani Adhikari, Joseph E. Gonzalez
 - o <u>http://www.ds100.org/sp20/syllabus/</u>
- Softmax Regression
 - <u>http://deeplearning.stanford.edu/tutorial/supervised/SoftmaxRegression/</u>
- One-vs-all
 - <u>https://www.coursera.org/lecture/machine-learning/multiclass-classification-one-vs-all-68Pol</u>



Images for Notebook

Splitting the Dataset



Confusion Matrix

	Predicted Positive (1)	Predicted Negative (0)
Actually	True	False
Positive (1)	Positive	Negative
Actually	False	True
Negative (0)	Positive	Negative