

DS 4400

Machine Learning and Data Mining I

Alina Oprea
Associate Professor
Khoury College of Computer Science
Northeastern University

October 13 2020

Outline

- Logistic regression
 - Cross-entropy objective
 - Gradient descent for logistic regression
- Project discussion
- Evaluation of classifiers
 - Metrics
 - ROC curves
- Linear Discriminant Analysis (LDA)

Logistic Regression

- Setup

- Training data: $\{x_i, y_i\}$, for $i = 1, \dots, N$
- Labels: $y_i \in \{0, 1\}$

- Goals

- Learn $P(Y = 1|X = x)$

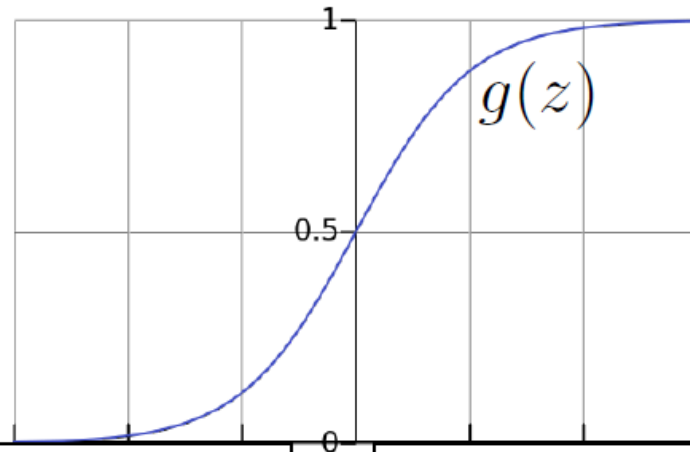
- Highlights

- Probabilistic output
- At the basis of more complex models (e.g., neural networks)
- Supports regularization (Ridge, Lasso)
- Can be trained with Gradient Descent

Logistic Regression

$$h_{\theta}(x) = g(\theta^T x)$$

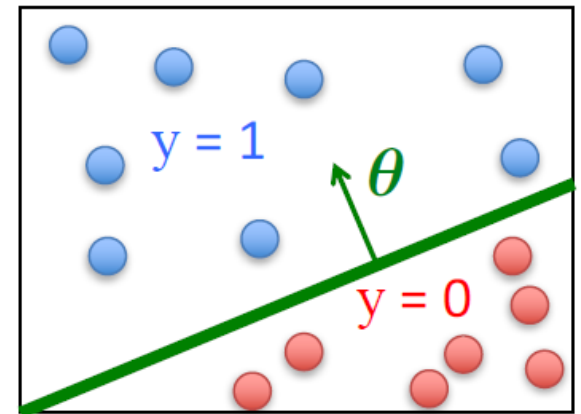
$$g(z) = \frac{1}{1 + e^{-z}}$$



$\theta^T x$ should be large negative values for negative instances

$\theta^T x$ should be large positive values for positive instances

- Assume a threshold and...
 - Predict $Y = 1$ if $h_{\theta}(x) \geq 0.5$
 - Predict $Y = 0$ if $h_{\theta}(x) < 0.5$



Logistic Regression is a linear classifier!

Cross-Entropy Objective

$$P(Y = y_i | X = x_i; \theta) = h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))\end{aligned}$$

Logistic regression objective

$$\min_{\theta} J(\theta)$$

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

Gradient Descent for Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

Want $\min_{\theta} J(\theta)$

- Initialize θ
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update
for $j = 0 \dots d$

Computing Gradients

- Derivative of sigmoid

- $g(z) = \frac{1}{1+e^{-z}}; g'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = g(z)(1 - g(z))$

- Derivative of hypothesis

- $h_{\theta}(x) = g(\theta^T x) = g(\theta_j x_j + \sum_{k \neq j} \theta_k x_k)$

- $\frac{\partial h_{\theta}(x)}{\partial \theta_j} = \frac{\partial g(\theta^T x)}{\partial \theta_j} x_j = g(\theta^T x)(1 - g(\theta^T x))x_j$

- Derivation of C_i

- $\frac{\partial C_i}{\partial \theta_j} = y_i \frac{1}{h_{\theta}(x_i)} g(\theta^T x_i)(1 - g(\theta^T x_i))x_{ij} -$
 $(1 - y_i) \frac{1}{1 - h_{\theta}(x_i)} g(\theta^T x_i)(1 - g(\theta^T x_i))x_{ij}$
 $= (y_i - h_{\theta}(x_i))x_{ij}$

Gradient Descent for Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

Want $\min_{\theta} J(\theta)$

- Initialize θ
- Repeat until convergence (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

Gradient Descent for Logistic Regression

Want $\min_{\theta} J(\theta)$

- Initialize θ
- Repeat until convergence (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

This looks IDENTICAL to Linear Regression!

- However, the form of the model is very different:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Regularized Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- We can regularize logistic regression exactly as before:

$$\begin{aligned} J_{\text{regularized}}(\boldsymbol{\theta}) &= J(\boldsymbol{\theta}) + \lambda \sum_{j=1}^d \theta_j^2 \\ &= J(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2 \end{aligned}$$

L2 regularization

Also supports L1 regularization

Outline

- Logistic regression
 - Cross-entropy objective
 - Gradient descent for logistic regression
- Project discussion
- Evaluation of classifiers
 - Metrics
 - ROC curves
- Linear Discriminant Analysis (LDA)

Classifier Evaluation

- Classification is a supervised learning problem
 - Prediction is binary or multi-class
- Classification techniques
 - **Linear classifiers**
 - Perceptron (online or batch mode)
 - Logistic regression (probabilistic interpretation)
 - **Instance learners**
 - kNN: need to store entire training data
- Cross-validation should be used for parameter selection and estimation of model error

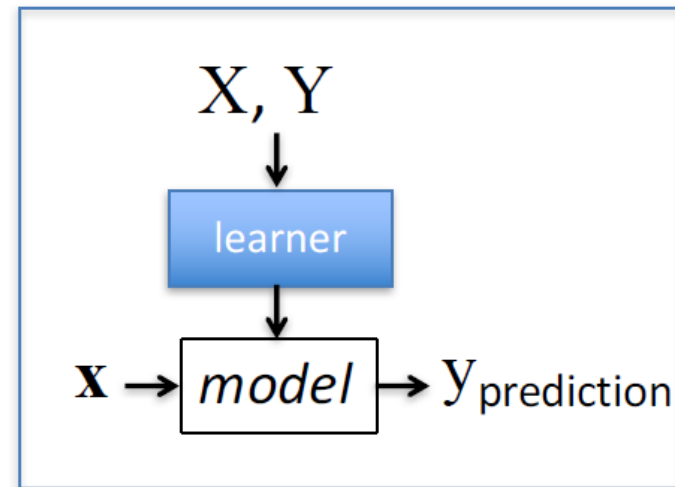
Evaluation of classifiers

Given: labeled training data $X, Y = \{\langle \mathbf{x}_i, y_i \rangle\}_{i=1}^n$

- Assumes each $\mathbf{x}_i \sim \mathcal{D}(\mathcal{X})$

Train the model:

$model \leftarrow classifier.train(X, Y)$



Apply the model to new data:

- Given: new unlabeled instance $\mathbf{x} \sim \mathcal{D}(\mathcal{X})$

$y_{\text{prediction}} \leftarrow model.predict(\mathbf{x})$

Classification Metrics

$$\text{accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ test instances}}$$

$$\text{error} = 1 - \text{accuracy} = \frac{\# \text{ incorrect predictions}}{\# \text{ test instances}}$$

- Training set accuracy and error
- Testing set accuracy and error

Confusion Matrix

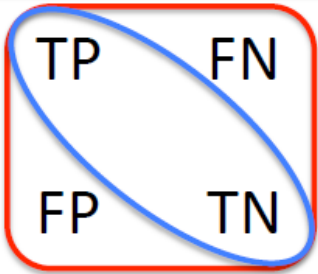
Given a dataset of P positive instances and N negative instances:

| | | Predicted Class | |
|--------------|-----|-----------------|----|
| | | Yes | No |
| Actual Class | Yes | TP | FN |
| | No | FP | TN |

Accuracy and Error

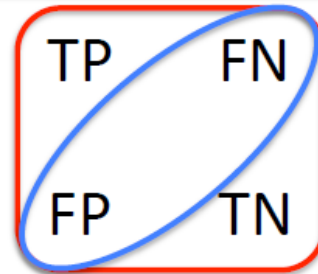
Given a dataset of P positive instances and N negative instances:

| | | Predicted Class | |
|--------------|-----|-----------------|----|
| | | Yes | No |
| Actual Class | Yes | TP | FN |
| | No | FP | TN |



$$\text{accuracy} = \frac{TP + TN}{P + N}$$

| | | Predicted Class | |
|--------------|-----|-----------------|----|
| | | Yes | No |
| Actual Class | Yes | TP | FN |
| | No | FP | TN |



$$\begin{aligned}\text{error} &= 1 - \frac{TP + TN}{P + N} \\ &= \frac{FP + FN}{P + N}\end{aligned}$$

Confusion Matrix

- Given a dataset of P positive instances and N negative instances:

| | | Predicted Class | |
|--------------|-----|-----------------|----|
| | | Yes | No |
| Actual Class | Yes | TP | FN |
| | No | FP | TN |

$$\text{accuracy} = \frac{TP + TN}{P + N}$$

- Imagine using classifier to identify positive cases (i.e., for information retrieval)

$$\text{precision} = \frac{TP}{TP + FP}$$

Probability that classifier predicts positive correctly

$$\text{recall} = \frac{TP}{TP + FN}$$

Probability that actual class is predicted correctly

Why One Metric is Not Enough

Assume that in your training data, Spam email is 1% of data, and Ham email is 99% of data

- Scenario 1
 - Have classifier always output HAM!
 - What is the accuracy? 99%
- Scenario 2
 - Predict one SPAM email as SPAM, all other emails as legitimate
 - What is the precision? 100%
- Scenario 3
 - Output always SPAM!
 - What is the recall? 100%

Precision & Recall

Precision

- the fraction of positive predictions that are correct
- $P(\text{is pos} | \text{predicted pos})$

$$\text{precision} = \frac{TP}{TP + FP}$$

Recall

- fraction of positive instances that are identified
- $P(\text{predicted pos} | \text{is pos})$

$$\text{recall} = \frac{TP}{TP + FN}$$

-
- You can get high recall (but low precision) by only predicting positive
 - Recall is a non-decreasing function of the # positive predictions
 - Typically, precision decreases as either the number of positive predictions or recall increases
 - Precision & recall are widely used in information retrieval

F-Score

- Combined measure of precision/recall tradeoff

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- This is the harmonic mean of precision and recall
 - In the F_1 measure, precision and recall are weighted evenly
 - Can also have biased weightings that emphasize either precision or recall more ($F_2 = 2 \times \text{recall}$; $F_{0.5} = 2 \times \text{precision}$)
- Limitations:
 - F-measure can exaggerate performance if balance between precision and recall is incorrect for application
 - Don't typically know balance ahead of time

A Word of Caution

- Consider binary classifiers A, B, C:

| | | A | | B | | C | |
|-------------|---|-----|-----|-----|-----|------|-----|
| | | 1 | 0 | 1 | 0 | 1 | 0 |
| Predictions | 1 | 0.9 | 0.1 | 0.8 | 0 | 0.78 | 0 |
| | 0 | 0 | 0 | 0.1 | 0.1 | 0.12 | 0.1 |

- Clearly A is useless, since it always predicts 1
- B is slightly better than C
 - less probability mass wasted on the off-diagonals
- But, here are the performance metrics:

| Metric | A | B | C |
|-----------|-------|-------|--------|
| Accuracy | 0.9 | 0.9 | 0.88 |
| Precision | 0.9 | 1.0 | 1.0 |
| Recall | 1.0 | 0.888 | 0.8667 |
| F-score | 0.947 | 0.941 | 0.9286 |

Acknowledgements

- Slides made using resources from:
 - Andrew Ng
 - Eric Eaton
 - David Sontag
- Thanks!