



COMP7/8118 M50

Data Mining

Logistic Regression & Supervised Learning

Xiaofei Zhang

Slides compiled from Jiawei Han and Raymond C.W. Wong's work

THE UNIVERSITY OF
MEMPHIS

Classification via regression

- Instead of predicting the **class** of an record we want to **predict the probability of the class** given the record
- The problem of **predicting continuous values** is called **regression** problem
- General approach: find a continuous function that models the continuous points.

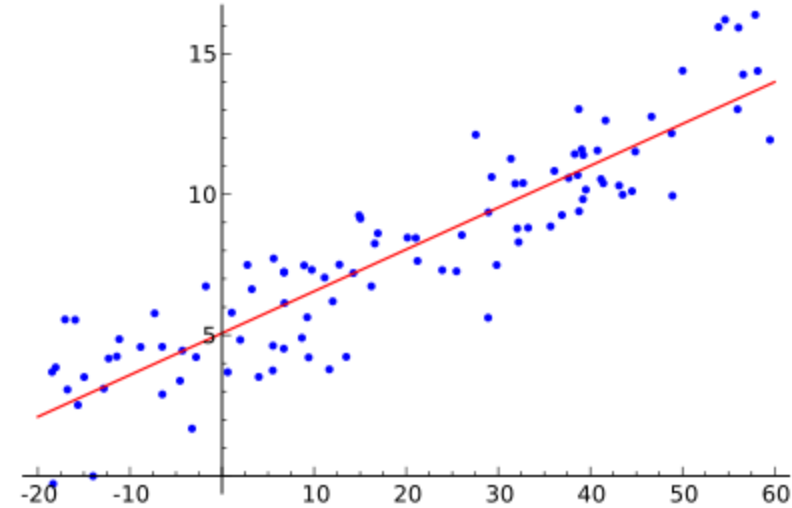
Example: Linear regression

- Given a dataset of the form $\{(x_1, y_1), \dots, (x_n, y_n)\}$ find a linear function that given the vector x_i predicts the y_i value as $y'_i = w^T x_i$

- Find a vector of weights w that minimizes the sum of square errors

$$\sum_i (y'_i - y_i)^2$$

- Several techniques for solving the problem.



Classification via regression

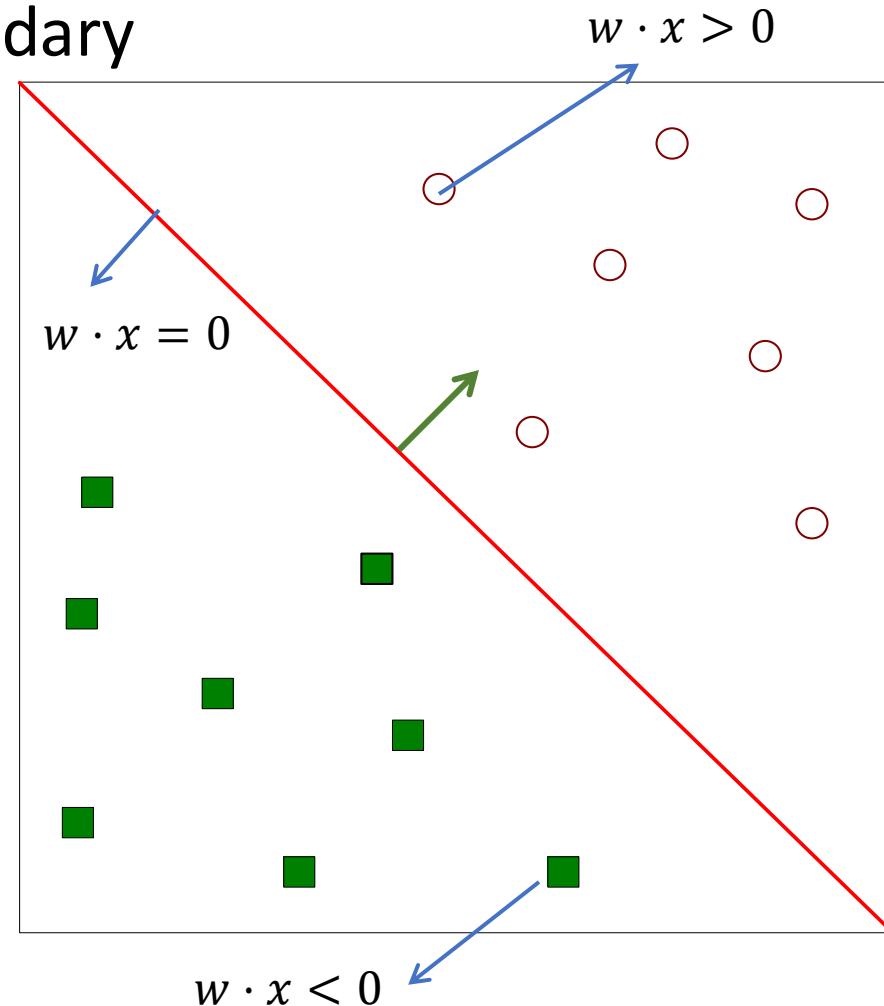
- Assume a linear classification boundary

For the positive class the **bigger** the **value of $w \cdot x$** , the further the point is from the classification boundary, the higher our **certainty** for the membership to the **positive class**

- Define $P(C_+ | x)$ as an **increasing** function of $w \cdot x$

For the negative class the **smaller** the **value of $w \cdot x$** , the further the point is from the classification boundary, the higher our **certainty** for the membership to the **negative class**

- Define $P(C_- | x)$ as a **decreasing** function of $w \cdot x$



Logistic Regression

The **logistic function**

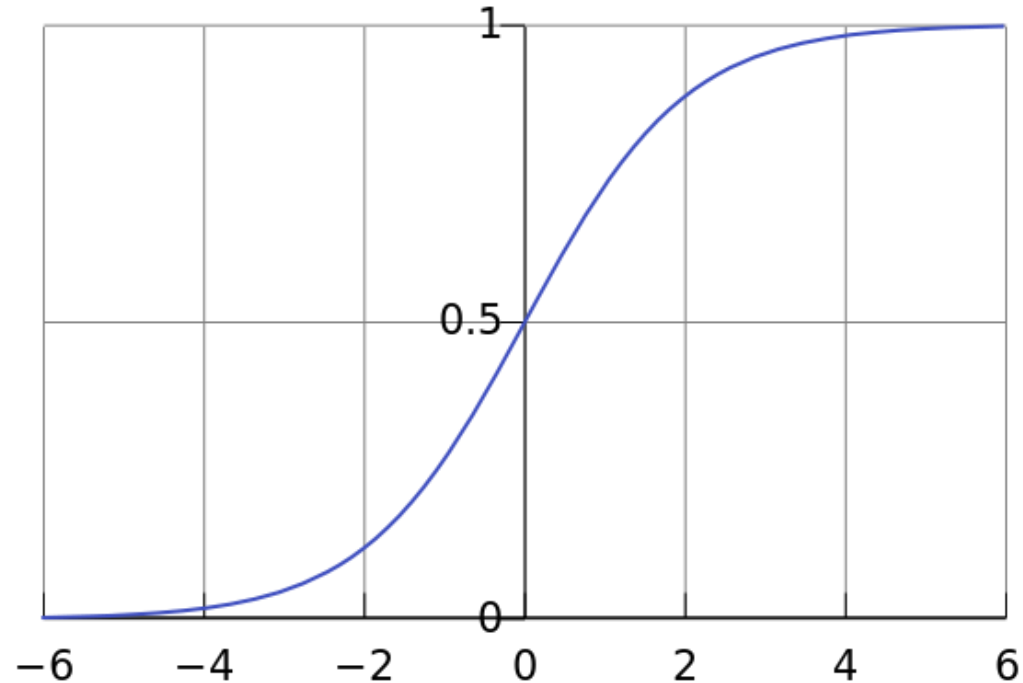
$$f(t) = \frac{1}{1 + e^{-t}}$$

$$P(C_+|x) = \frac{1}{1 + e^{-w \cdot x}}$$

$$P(C_-|x) = \frac{e^{-w \cdot x}}{1 + e^{-w \cdot x}}$$

$$\log \frac{P(C_+|x)}{P(C_-|x)} = w \cdot x$$

Linear regression on the **log-odds ratio**



Logistic Regression: Find the vector **w** that **maximizes the probability** of the observed data

Logistic Regression in one dimension

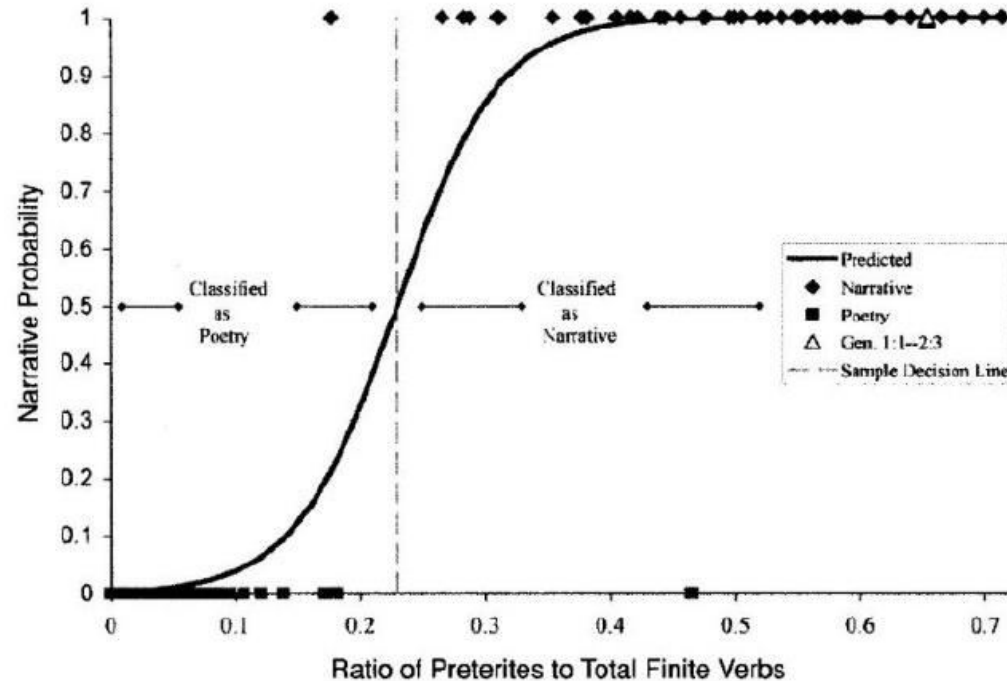


Figure 10-3. The solid curved line is called a logistic regression curve. The vertical axis measures the probability that an Old Testament passage is narrative, based on the use of preterite verbs. The probability is zero for poetry and unity or one for narrative. Passages with high preterite verb counts, falling to the right of the vertical dotted line, are likely narrative. The triangle on the upper right represents Genesis 1:1-2:3, which is clearly literal, narrative history.

Logistic Regression

- Produces a **probability estimate** for the **class membership** which is often very useful.
- The **weights** can be useful for understanding the **feature importance**.
- Works for relatively large datasets
- Fast to apply.

Supervised Learning

Learning

- **Supervised Learning**: learn a model from the data using **labeled data**.
 - **Classification** and **Regression** are the prototypical examples of supervised learning tasks. Other are possible (e.g., ranking)
- **Unsupervised Learning**: learn a model – extract structure from **unlabeled data**.
 - **Clustering** and **Association Rules** are prototypical examples of unsupervised learning tasks.
- **Semi-supervised Learning**: learn a model for the data using both **labeled and unlabeled** data.

Supervised Learning Steps

- **Model** the problem
 - What is you are trying to predict? What kind of optimization function do you need? Do you need classes or probabilities?
- Extract **Features**
 - How do you find the right features that help to discriminate between the classes?
- Obtain **training data**
 - Obtain a collection of labeled data. Make sure it is large enough, accurate and representative. Ensure that classes are well represented.
- Decide on the **technique**
 - What is the right technique for your problem?
- **Apply** in practice
 - Can the model be trained for very large data? How do you test how you do in practice? How do you improve?

Modeling the problem

- Sometimes it is not obvious. Consider the following three problems
 - Detecting if an email is spam
 - Categorizing the queries in a search engine
 - Ranking the results of a web search

Feature extraction

- Feature extraction, or **feature engineering** is the most tedious but also the most important step
 - How do you separate the players of the Greek national team from those of the Swedish national team?
- One line of thought: throw features to the classifier and the classifier will figure out which ones are important
 - **More features**, means that you need **more training data**
- Another line of thought: **Feature Selection**: Select carefully the features using various functions and techniques
 - Computationally intensive

Training data

- An overlooked problem: How do you get **labeled data** for training your model?
 - E.g., how do you get training data for ranking?
 - Chicken and egg problem
- Usually requires a lot of manual effort and domain expertise and carefully planned labeling
 - Results are not always of high quality (lack of expertise)
 - And they are not sufficient (low coverage of the space)
- Recent trends:
 - Find a **source** that generates the labeled data for you.
 - **Crowd-sourcing** techniques

Dealing with small amount of labeled data

- **Semi-supervised learning** techniques have been developed for this purpose.
- **Self-training**: Train a classifier on the data, and then feed back the high-confidence output of the classifier as input
- **Co-training**: train two “independent” classifiers and feed the output of one classifier as input to the other.
- **Regularization**: Treat learning as an optimization problem where you define relationships between the objects you want to classify, and you exploit these relationships
 - Example: Image restoration

Technique

- The choice of technique depends on the problem requirements (do we need a probability estimate?) and the problem specifics (does independence assumption hold? do we think classes are linearly separable?)
- For many cases finding the right technique may be trial and error
- For many cases the exact technique does not matter.

Big Data Trumps Better Algorithms

- If you have enough data then the algorithms are not so important
- The web has made this possible.
 - Especially for text-related tasks
 - Search engine uses the **collective human intelligence**

Google lecture: [Theorizing from the Data](https://www.youtube.com/watch?v=nU8DcBF-qo4)
(<https://www.youtube.com/watch?v=nU8DcBF-qo4>)

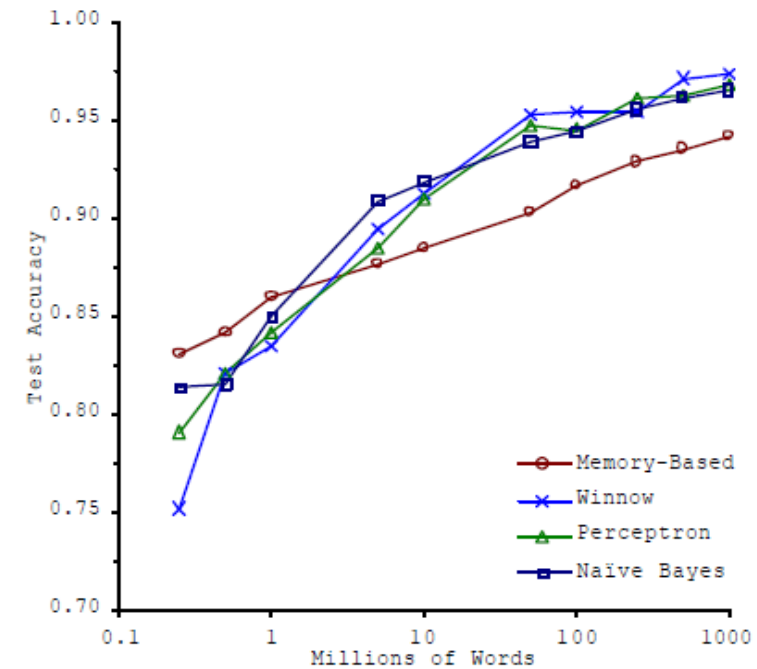


Figure 1. Learning Curves for Confusion Set Disambiguation

Apply-Test

- How do you **scale** to very large datasets?
 - Distributed computing – **map-reduce** implementations of machine learning algorithms (Mahout, over Hadoop)
- How do you test something that is running online?
 - You cannot get labeled data in this case
 - **A/B testing**
- How do you deal with changes in data?
 - **Active learning**