

Intro to Neural Networks

Odelia Schwartz
Luis G Sanchez Giraldo

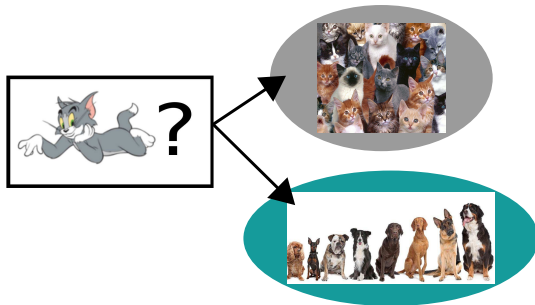
University of Miami

March 24th, 2016

Object Features and the Classification Problem

Classification Problem:

Given a set of objects, where each object belongs to one of two classes, we want to find a rule or function that can predict to which class each object from the set belongs to.



Object Features and the Classification Problem (cont...)

Features for Object Representation:

Each object is represented as by a set of numbers called features. Mathematically, an object can be represented as a d -dimensional vector $\mathbf{x} = \{x_1, x_2, \dots, x_d\}$, where d is the number of features.

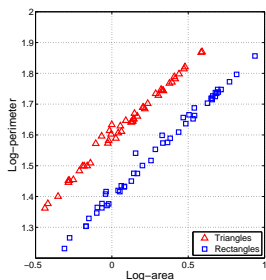
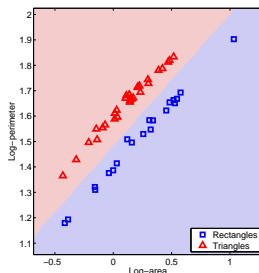


Figure : Representation of rectangles and triangles by their area and perimeter. The *log*-scale is used for mathematical convenience.

Object Features and the Classification Problem (cont...)

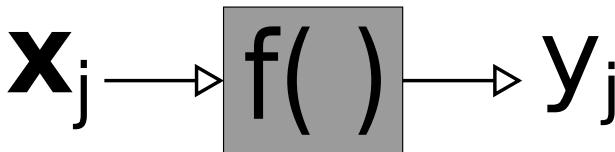
The assumption is that objects that belong together in the same class will have similar values for their features. Therefore, if we want to predict the class of an object, we can hope to do so by looking at the values of its features. However, it is important to notice that features do interact and that these **feature interactions (structure)** are important in identifying the object's class. For example, a simple rule to separate triangles from rectangles would be to divide the log-perimeter, log-volume plane into two regions.



Rosenblatt's Perceptron

Rosenblatt's Perceptron is one of the earliest models for learning with supervision, or what is commonly known as **supervised learning**.

In **supervised learning**, we are given a set $\{(\mathbf{x}_j, y_j)\}_{j=1}^N$ of feature-class pairs. The goal is to express the relation between \mathbf{x}_j and y_j , for all $j = 1, \dots, N$, as a function f that receives \mathbf{x}_j as input, and returns y_j as output. Learning consist of finding such an f within a given collection of functions.



Rosenblatt's Perceptron (cont...)

In particular, the perceptron implements the following function to solve the classification problem.

$$f(\mathbf{x}_j) = \text{sign} \left(\sum_{i=1}^d w_i x_{ij} + b \right). \quad (1)$$

As we can see, changing the values of w_i and b gives us different functions and thus, it defines a collection of functions. Learning amounts to finding \mathbf{w} and b that “best capture” the input-output relation.

Perceptron Algorithm for Learning f

consist of the following steps:

1. Let $j = 1$, and initialize $\{w_i\}_{i=1}^d$ and b (\mathbf{w} and b can be initialized with zeros).

2. Compute

$$f(\mathbf{x}_j) = \text{sign} \left(\sum_{i=1}^d w_i x_{ij} + b \right). \quad (2)$$

3. If $f(\mathbf{x}_j)$ is **NOT** equal to y_j , update \mathbf{w} and b as follows:

$$w_i \leftarrow w_i + y_j x_{ij} \quad (3)$$

$$b \leftarrow w_i + y_j \quad (4)$$

4. Repeat 2 for next value of j , that is, $j = j + 1$. When $j = N$, restart the counter to $j = 1$. If for all $j = 1, \dots, N$, $f(\mathbf{x}_j)$ is equal to y_j , stop iterating.

Perceptron Algorithm for Learning f (cont...)

Why is this supposed to work? (Basic Idea)

Suppose we have updated \mathbf{w} and b , based on the rule given above, using the j th input-output pair (\mathbf{x}_j, y_j) . The linear part of the updated function evaluated at \mathbf{x}_j is given by:

$$\sum_{i=1}^d w_i x_{ij} + b = \sum_{i=1}^d \hat{w}_i x_{ij} + \hat{b} + y_j \left(\sum_{i=1}^d x_{ij} x_{ij} + 1 \right), \quad (5)$$

where $\hat{\mathbf{w}}$ and \hat{b} are the parameters before the update.

As we can see, the update rule applies a positive or negative correction to the current function to enforce agreement with the actual class label y_j .

Going Beyond the Perceptron

The perceptron works at separating objects given a set of features (numerical descriptors). However, in this framework a lot of human expertise is required for crafting such features.

- ▶ What if features can be learned from exemplars?
- ▶ How can features be learned from exemplars?

Going Beyond the Perceptron (cont...)

Consider the following generalizations of the function in the perceptron

1.

$$f(\mathbf{x}_j) = h\left(\sum_{i=1}^d w_i x_{ij} + b\right), \quad (6)$$

where h is a nonlinear function. In this case the values of the outputs of f do not need to be restricted to -1 and $+1$ as with sign function.

2. To compare the output of $f(\mathbf{x}_j)$ with the actual label values y_j , we define a loss function:

$$\mathcal{L}(f(\mathbf{x}_j), y_j) \quad (7)$$

The loss function can be used to guide the learning since it penalizes errors made by f .

Going Beyond the Perceptron (cont...)

Common examples of loss functions are:

- ▶ Zero-One loss

$$\begin{cases} 1 & \text{if } \text{sign}[f(\mathbf{x}_j)] \neq y_j, \\ 0 & \text{if } \text{sign}[f(\mathbf{x}_j)] = y_j. \end{cases} \quad (8)$$

- ▶ Squared loss

$$(y_j - f(\mathbf{x}_j))^2. \quad (9)$$

- ▶ Hinge loss

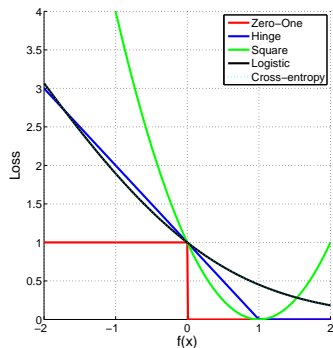
$$\max\{0, 1 - y_j f(\mathbf{x}_j)\}. \quad (10)$$

- ▶ Cross-entropy loss. In this case $f(\mathbf{x}_j)$ is transformed to a function $h(f(\mathbf{x}_j))$ with outputs in the range $[0, 1]$

$$-\frac{y_j + 1}{2} \log(f(\mathbf{x}_j)) - \frac{1 - y_j}{2} \log(1 - f(\mathbf{x}_j)). \quad (11)$$

Going Beyond the Perceptron (cont...)

Values of the loss functions for $y_j = 1$ vs $f(\mathbf{x}_j)$



Going Beyond the Perceptron (cont...)

Given a set $\{(\mathbf{x}_j, y_j)\}_{j=1}^N$ of feature-class pairs learning can be accomplished by minimizing the average loss on this set. Having differentiable loss and activation functions provide mathematically a simple framework to use gradient-based minimization techniques to the process of learning. For instance, using the logistic function for h in (6) and the cross-entropy loss (11) yields an iterative procedure.

Going Beyond the Perceptron (cont...)

Gradient descent for cross-entropy loss with logistic sigmoid activation.

1. Initialize $\{w_i\}_{i=1}^d$ and b (\mathbf{w} and b can be initialized with zeros).
2. Let

$$f(\mathbf{x}_j) = \frac{1}{1 + \exp(-z_j)}, \text{ where } z_j = \left(\sum_{i=1}^d w_i x_{ij} + b \right), \quad (12)$$

3. and

$$\Delta w_i = \frac{\partial}{\partial w_i} \mathcal{L}(f(\mathbf{x}_j), y_j) = \left(f(\mathbf{x}_j) - \frac{y_j + 1}{2} \right) x_{ij} \quad (13)$$

$$\Delta b = \frac{\partial}{\partial b} \mathcal{L}(f(\mathbf{x}_j), y_j) = \left(f(\mathbf{x}_j) - \frac{y_j + 1}{2} \right) \quad (14)$$

4. Update parameters. Similarly to the perceptron

$$w_i \leftarrow w_i - \mu \Delta w_i \text{ and } b \leftarrow b - \mu \Delta b \quad (15)$$