

Perceptron

Jack David Carson - February 3, 2025

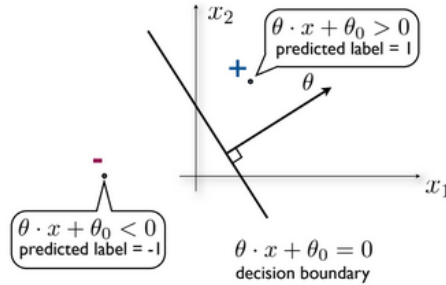
If we have some linear classifier that we wish to use to predict the value y corresponding to a feature vector $\mathbf{x}^{(i)} \in \mathbb{R}^d$ for $y \in \{-1, 1\}$, we can first formalize the problem as trying to select an ideal classifier $\hat{h} \in \mathcal{H}$, the set of all classifiers. In functional form we can say $h : \mathbb{R}^d \rightarrow \{1, -1\}$ and that classifier \hat{h} has some parameters θ and evaluate its training error over training data $S_n = \{\mathbf{x}^{(i)}, y^{(i)}, i = 1, \dots, n\}$

$$E_n(h) = \frac{1}{n} \sum_{i=1}^n [[h(\mathbf{x}^{(i)}) \neq y^{(i)}]]$$

where $[[\text{true}]] = 1$ is an error and $[[\text{false}]] = 0$ is a correct classification. We define the actual classification $h(\mathbf{x}^{(i)})$ in terms of its parameters and a “bias” term θ_0 as

$$h(\mathbf{x}, \theta) = \text{sgn}(\theta^\top \mathbf{x} + \theta_0) = \begin{cases} +1 & \text{for } \theta^\top \mathbf{x} + \theta_0 > 0 \\ -1 & \text{for } \theta^\top \mathbf{x} + \theta_0 \leq 0 \end{cases}$$

In this case the classifier represents a **hyperplane** through \mathbb{R}^d . If $\theta_0 = 0$ then the hyperplane intersects the origin.



We say that the training data S^n is **linearly separable** if we can design a linear classifier that is correct for all $\mathbf{x}^{(i)}, y^{(i)}$.

0.1 Perceptron Algorithm

The **Perceptron Algorithm** is the simplest algorithm we can use to optimize the parameters of our linear classifier. Simply, we compute the product of the parameters with the observation as in $h(\mathbf{x})$ and, if the sign of $y^{(i)}$ does not equal the sign of $\theta^\top \mathbf{x}^{(i)} + \theta_0$, then we know the observation is incorrect, and can update our parameters by $\theta = \theta + y^{(i)} \mathbf{x}^{(i)}$. In algorithmic notation:

```

1: function PERCEPTRON(D, T)
2:   ▷ Initialize parameters
3:    $\theta \leftarrow 0$ 
4:    $\theta_0 \leftarrow 0$ 
5:
6:   for  $t = 1, \dots, T$  do
7:     for  $i = 1, \dots, n$  do
8:       if  $y^{(i)}(\theta^T \mathbf{x}^{(i)} + \theta_0) \leq 0$  then
9:         ▷ Update weight and bias
10:         $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$ 
11:         $\theta_0 \leftarrow \theta_0 + y^{(i)}$ 
12:
13:   return  $(\theta, \theta_0)$ 

```

0.2 Perceptron Convergence

In order to understand convergence we can first specify a linearly seperable dataset formally as **Definition:** $S_n = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, n\}$ is **linearly seperable** with margin γ if there exists some margin

$$y^{(i)}(\theta^T \mathbf{x}^{(i)} + \theta_0) \geq \gamma \|\theta_0\|$$

for *all* training points. We can further this understanding by considering the projection of a point $\mathbf{x}^{(i)}$. The projection must of course exist on the hyperplane governed by $\theta^T \mathbf{x}_0^{(i)} + \theta_0 = 0$ for projected point \mathbf{x}_0 . Since this exists on our hyperplane the orthogonal complement (the distance from $\mathbf{x}_0^{(i)}$ to $\mathbf{x}^{(i)}$ is simply

$$d = \|\mathbf{x}^{(i)} - \mathbf{x}_0^{(i)}\|$$

So, given that there exists some

- a) θ^* such that $y^{(i)}(\theta^{*T} \mathbf{x}^{(i)} + \theta_0^*) \geq \gamma \|\theta_0^*\|$
- b) All examples are bounded by $\|\mathbf{x}^{(i)}\| \leq R$ to ensure finite vectors

The **perceptron convergence theorem** states that the perceptron algorithm will make at most R^2/γ^2 mistakes on the way to find the correct classifier. Remarkably, this does **not** depend on the length of the feature vector nor the number of examples!