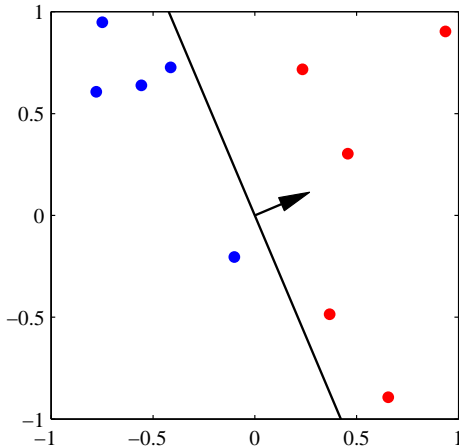


Perceptron

The whole idea behind linear classifiers is to directly represent the *decision* boundary (in binary classification problems). In decision trees and KNN, it is *implicit*.

We begin with the simplest kind of decision boundary: a hyperplane (a flat thing). We put points of each class on either side of the boundary. The mostly-vertical line in the figure below is a possible decision boundary between the blue (left) and red (right) points.



As a *representation* of a decision boundary, we store the vector that is orthogonal (aka “normal”) to the decision boundary. We usually call this \mathbf{w} (for “weights”... this will become clear soon). For instance, if the x-axis is the first coordinate and the y-axis is the second coordinate, the weight vector in the figure above (the arrow) might have value $\langle 0.4, 0.1 \rangle$, indicating that it is pointing right a lot and up a little.

The representation in terms of \mathbf{w} is overkill: if \mathbf{w} represents the hyperplane, then $\alpha\mathbf{w}$ also does, for any $\alpha > 0$ (it just makes the vector longer, but doesn’t change its direction). Thus, we’ll usually assume that $\|\mathbf{w}\| = 1$ (the 1 is for convenience; any positive constant would be fine).

This representation also assumes that the decision boundary passes through the origin. We will shift the hyperplane by adding an additional “bias” term b . In the above figure, if we wished to move the decision boundary in the direction of the arrow, we would use $b < 0$; to move it in the opposite direction, we use $b > 0$.

Given a representation of the decision boundary in terms of \mathbf{w} and b , we make our classification according to $\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$, where \mathbf{x} is an input and \hat{y} is the predicted class (note that we refer to our two classes as +1 and -1 instead of 0 and 1).

We can think about the hyperplane as being defined as the set $\{\mathbf{x} : \mathbf{w}^\top \mathbf{x} + b = 0\}$... in other words, it’s the set of *all* points for which our classifier is totally uncertain.

In general, if there exists a hyperplane that perfectly separates our data onto two sides, we refer to the data as *linearly separable*. Except in very peculiar circumstances, if data is linearly separable, then there are infinitely many separating hyperplanes (just jiggle the one in the figure a little bit any direction).

We usually want to pick out the “best” hyperplane and the criteria that is often chosen is the *maximal margin hyperplane*, which is essentially the one with the most “wiggle room.” Formally, the *margin* of a hyperplane (denoted γ is the distance to the closest point): $\gamma = \min_n y_n(\mathbf{w}^\top \mathbf{x}_n + b)$.

The next algorithm we discuss is an *error-correction* based algorithm. It maintains a current classifier h and uses this to make predictions. At each step, if h is right, it do nothing. If h is wrong, the update it so that it is “less wrong.”

Our representation is linear: we have a weight vector \mathbf{w} and a bias b and make predictions on \mathbf{x} of the form $\text{sign}[\mathbf{w}^\top \mathbf{x} + b]$.

The perceptron uses simple additive updates:

1. Initialize $\mathbf{w} = 0, b = 0$
2. For $n = 1 \dots$
 - (a) Obtain \mathbf{x}_n
 - (b) Predict $\hat{y}_n = \text{sign}[\mathbf{w}^\top \mathbf{x}_n + b]$
 - (c) Suffer loss ℓ_n (zero if $\hat{y}_n = y_n$, one otherwise)
 - (d) If $\ell_n = 0$, continue on as before. Otherwise, update:

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \hat{y}_n \mathbf{x}_n \\ b &\leftarrow b - \hat{y}_n\end{aligned}\tag{1}$$

The idea with the algorithm is that with binary classification, if \hat{y}_n was wrong, then $-\hat{y}_n$ is *right*. What we do in update (1) is to move (\mathbf{w}, b) closer to making the right prediction. Why? Let (\mathbf{w}', b') denote the updated classifier. Let's look at what we *would* predict, if we received the same \mathbf{x} again. Suppose that the truth is y (so $y = -\hat{y}$); we know we want $y[\mathbf{w}'^\top \mathbf{x} + b'] > 0$.

$$\begin{aligned}y[\mathbf{w}'^\top \mathbf{x} + b'] &= y(\mathbf{w} - \hat{y}\mathbf{x})^\top \mathbf{x} + y(b - \hat{y}) \\ &= y[\mathbf{w}^\top \mathbf{x} + b] - y(\hat{y}\mathbf{x}^\top \mathbf{x} + \hat{y}) \\ &= y[\text{old prediction}] - y\hat{y} \|\mathbf{x}\|^2 - y\hat{y} \\ &= y[\text{old prediction}] + \|\mathbf{x}\|^2 + 1\end{aligned}$$

In the last step, we used the fact that $y\hat{y} = -1$ if we made an error. So what's happening is we have our old prediction times y (which we know was negative—since we erred) plus $\|\mathbf{x}\|^2 + 1$, which is always positive. Which means that the whole thing is getting more positive. Which is exactly what we want.

This leads directly to the following theorem for *batch* applications:

Theorem (Block & Novikoff): if the training data is linearly separable with margin γ with weights \mathbf{u} , then Perceptron will terminate in $\|\mathbf{u}\|^2 \leq R^2/\gamma^2$ iterations (assuming $\|\mathbf{x}\| \leq R$ for all \mathbf{x}).

Proof sketch: (Suppose $R = 1$.) Let \mathbf{w}_k be the weights before the k th update, so $\mathbf{w}_1 = 0$ and if the k th error is on example n , then $y_n(\mathbf{w}_k^\top \mathbf{x}_n) \leq 0$. Then, $\mathbf{w}_{k+1}^\top \mathbf{u} = \mathbf{w}_k^\top \mathbf{u} + y_n(\mathbf{u}^\top \mathbf{x}_n) \geq \mathbf{w}_k^\top \mathbf{u} + \gamma$, so $\mathbf{w}_{k+1}^\top \mathbf{u} \geq k\gamma$. Also, $\|\mathbf{w}_{k+1}\|^2 = \|\mathbf{w}_k\|^2 + 2y_n(\mathbf{w}_k^\top \mathbf{x}_n) + \|\mathbf{x}_n\|^2 \leq \|\mathbf{w}_k\|^2 + 1$ so $\|\mathbf{w}_{k+1}\|^2 \leq k$. Thus $\sqrt{k} \geq \|\mathbf{w}_{k+1}\| \geq \mathbf{w}_{k+1}^\top \mathbf{u} \geq k\gamma$; then algebra.

For the inseparable case, just force it to be separable (add a feature).

Voting: keep a copy of each value of \mathbf{w} and make predictions by voting ($\text{sign}[\sum_i \text{sign}[\mathbf{w}_i^\top \mathbf{x}]]$).

Averaging: use averaged weights to make prediction ($\text{sign}[(\sum_i \mathbf{w}_i)^\top \mathbf{x}]$).