

Perceptrons

Virgil Pavlu October 3, 2014

1 The perceptron

Lets suppose we are (as with regression regression) with $(\mathbf{x}_i, y_i); i = 1, \dots, m$ the data points and labels. This is a classification problem with two classes $y \in \{-1, 1\}$

Like with regression we are looking for a linear predictor (classifier)

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}\mathbf{w} = \sum_{d=0}^D x^d w^d$$

(we added the $x^0 = 1$ component so we can get the free term w^0) such that $h_{\mathbf{w}}(\mathbf{x}) \geq 0$ when $y = 1$ and $h_{\mathbf{w}}(\mathbf{x}) \leq 0$ when $y = -1$.

On $y = -1$ data points: given that all \mathbf{x} and y are numerical, we will make the following transformation: when $y = -1$, we will reverse the sign of the input; that is replace \mathbf{x} with $-\mathbf{x}$ and $y = -y$. Then the condition $h_{\mathbf{w}}(\mathbf{x}) \leq 0$ becomes $h_{\mathbf{w}}(\mathbf{x}) \geq 0$ for all data points.

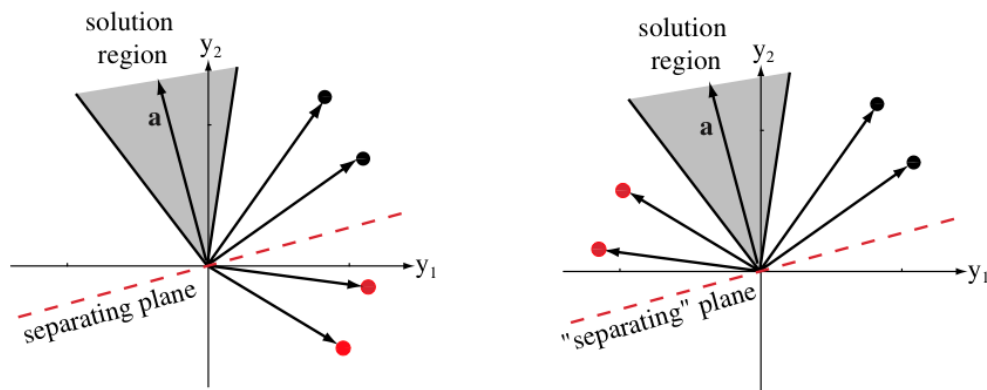


Figure 1: transforming $y = -1$ datapoints into $y = 1$ datapoints

The perceptron objective function is a combination of the number of miss-classification points and how bad the miss-classification is

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in M} -h_{\mathbf{w}}(\mathbf{x}) = \sum_{\mathbf{x} \in M} -\mathbf{x}\mathbf{w}$$

where M is the set of miss-classified data points. Note that each term of the sum is positive, since miss-classified implies $\mathbf{w}\mathbf{x} < 0$. Using gradient descent, we first differentiate J

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{\mathbf{x} \in M} -\mathbf{x}^T$$

then we write down the gradient descent update rule

$$\mathbf{w} := \mathbf{w} + \lambda \sum_{\mathbf{x} \in M} \mathbf{x}^T$$

(λ is the learning rate). The batch version looks like

1. init \mathbf{w}
2. LOOP
3. get M = set of missclassified data points
4. $\mathbf{w} = \mathbf{w} + \lambda \sum_{\mathbf{x} \in M} \mathbf{x}^T$
5. UNTIL $|\lambda \sum_{\mathbf{x} \in M} \mathbf{x}| < \epsilon$

Assume the instances are linearly separable. Then we can modify the algorithm

1. init \mathbf{w}
2. LOOP
3. get M = set of missclassified data points
4. for each $\mathbf{x} \in M$ do $\mathbf{w} = \mathbf{w} + \lambda \mathbf{x}^T$
5. UNTIL M is empty

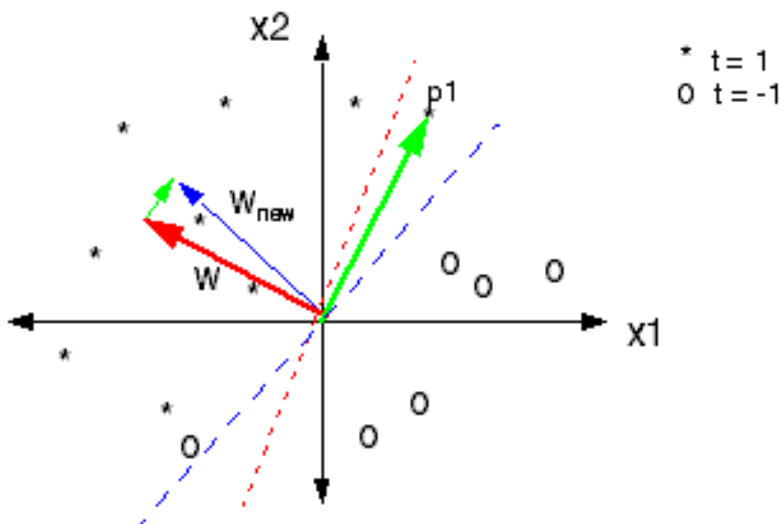


Figure 2: perceptron update: the plane normal w moves in the direction of misclassified x until the x is on the correct side.

Intuitively, the update $w^{new} = w^{old} + x$ for misclassified points x is the following: if x is on the wrong side of the plane $\langle w, x \rangle = 0$, it means that the normal vector to the plane, w , is on the opposite side to

x . The update essentially moves w in the direction of x ; as long as x remains on the wrong side, w moves towards it until it w and x are on the same side of the plane (thus x is correctly classified).

Proof of perceptron convergence Assuming data is linearly separable, or there is a solution $\bar{\mathbf{w}}$ such that $\mathbf{x}\bar{\mathbf{w}} > 0$ for all \mathbf{x} .

Lets call \mathbf{w}_k the \mathbf{w} obtained at the k -th iteration (update). Fix an $\alpha > 0$. Then

$$\mathbf{w}_{k+1} - \alpha\bar{\mathbf{w}} = (\mathbf{w}_k - \alpha\bar{\mathbf{w}}) + \mathbf{x}_k^T$$

where \mathbf{x}_k is the datapoint that updated \mathbf{w} at iteration k . Then

$$\|\mathbf{w}_{k+1} - \alpha\bar{\mathbf{w}}\|^2 = \|\mathbf{w}_k - \alpha\bar{\mathbf{w}}\|^2 + 2\mathbf{x}_k(\mathbf{w}_k - \alpha\bar{\mathbf{w}}) + \|\mathbf{x}_k\|^2 \leq \|\mathbf{w}_k - \alpha\bar{\mathbf{w}}\|^2 - 2\mathbf{x}_k\alpha\bar{\mathbf{w}} + \|\mathbf{x}_k\|^2$$

Since $\mathbf{x}_k\bar{\mathbf{w}} > 0$ all we need is an α sufficiently large to show that this update process cannot go on forever. When it stops, all datapoints must be classified correctly.

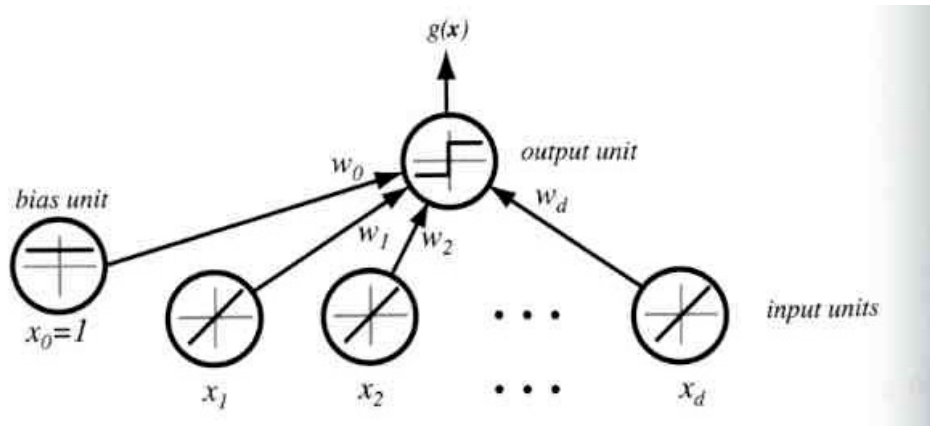


Figure 3: bias unit