# Classification

## Perceptron — First proposed in 1962

$(x_i, y_i), i = 1, 2, \ldots N$



$w^T x + w_0 = 0$

$w/\|w\|_2$

$C_1$

$x - x_0$

$y_i = -1$

$C_2$

$y_i = +1$

$w^T x_0 + w_0 = 0$

$\Rightarrow w^T x_0 = -w_0$

Distance of $x$ to hyperplane

Signed Distance $\overset{\text{of } x}{\to}$ to the hyperplane

$$\left(\frac{w}{\|w\|}\right)^T (x - x_0)$$

$$= \frac{w^T x - w^T x_0}{\|w\|} = \frac{w^T x + w_0}{\|w\|}$$

For class $C_1$,

$$w^T x_i + w_0 > 0 \quad \text{for all points } x_i \text{ in } C_1 \text{ that are correctly classified}$$

For class $C_2$, $w^T x_i + w_0 < 0$ for all points in $C_2$ that are correctly classified

For points that are correctly classified $(x_i$ in $C_1$ or $C_2)$

$$y_i (w^T x_i + w_0) > 0$$

For misclassified points,

$$y_i (w^T x_i + w_0) < 0$$

Perceptron Criterion:

$$\min_{w, w_0} D(w, w_0) = \boxed{-\sum_{i \in M} y_i (w^T x_i + w_0)} \to M \text{ indexes the misclassified points}$$

$$\nabla_w D(w, w_0) = -\sum_{x \in M} y_i x_i = -\left(\sum_{i \in M \cap C_1} x_i - \sum_{j \in M \in C_2} x_j\right)$$

$$\nabla_{w_0} D(w, w_0) = -\sum_{i \in M} y_i = -(N_1 - N_2)$$

$N_i$ be the number of misclassified points in $C_i$

## Perceptron Update Rule

$$\begin{bmatrix} w \\ w_0 \end{bmatrix} \leftarrow \begin{bmatrix} w \\ w_0 \end{bmatrix} - \eta \begin{bmatrix} -y_i x_i \\ -y_i \end{bmatrix} \quad \text{if } x_i \text{ is misclassified}$$

$$\hookrightarrow \begin{bmatrix} w \\ w_o \end{bmatrix} \leftarrow \begin{bmatrix} w \\ w_o \end{bmatrix} + \eta \begin{bmatrix} y_i x_i \\ y_i \end{bmatrix} \quad \text{...Not gradient descent}$$
but stochastic gradient descent

Gradient Descent:

$$\begin{bmatrix} w \\ w_o \end{bmatrix} = \begin{bmatrix} w \\ w_o \end{bmatrix} + \eta \begin{bmatrix} \sum_i y_i x_i \\ \sum_i y_i \end{bmatrix}$$

<u>Stochastic Gradient Descent</u> : Look at $x_i$ and if $x_i$
is misclassified then go in direction of the negative gradient
(but only contribution to gradient by $x_i$)

<u>Perceptron Algorithm</u>    $(x_i, y_i)$, $i = 1, 2, \dots n$  is the training data

Start with some $w, w_o$

Repeat

   for $i=1$ to $n$

      if $y_i(w^T x_i + w_o) < 0$  then

         $w \leftarrow w + \eta y_i x_i$

         $w_o \leftarrow w_o + \eta y_i$

      end if

   end for

until there are no misclassifications (mistakes)
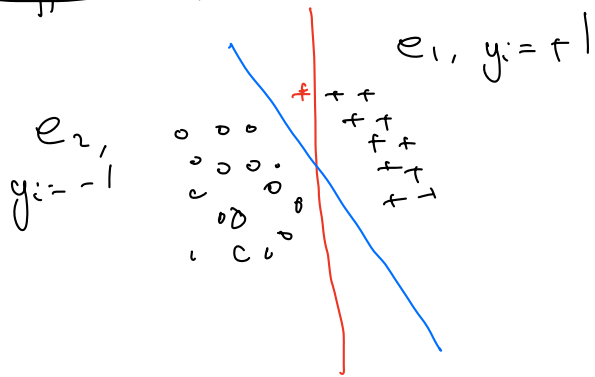within the for loop

Perceptron algorithm is guaranteed to find <u>a</u> separating
hyperplane <u>if</u> the data is linearly separable

<u>Drawbacks of Perceptron</u>

1. If the data is linearly separable, the hyperplane that is
   output by the perceptron depends on the order in which
   points (data) is presented to the algorithm.

2. Number of iterations might be large

3. If the classes are not linearly separable, then the

algorithm might not converge — cycles can develop that are not easy to detect

## Support Vector Machines (SVM)



$e_1, y_i = +1$

$e_2, y_i = -1$

Signed distance of $x$ to hyperplane

$$\frac{w^T x - w^T x_0}{\|w\|} = \frac{w^T x + w_0}{\|w\|}$$

$$y_i \left( \frac{w^T x + w_0}{\|w\|} \right) - \text{Distance of training point } x_i \text{ to the hyperplane}$$

Suppose we put the requirement that each of these distances is greater than $C$

$$\frac{y_i (w^T x + w_0)}{\|w\|} \geq C$$

$$\Rightarrow y_i (w^T x + w_0) \geq C \cdot \|w\|$$

In linear support vector machines, the goal is to maximize $C$:

$$\underset{w, w_0}{\text{maximize}} \quad C$$

such that $y_i (w^T x_i + w_0) \geq C \|w\|$, $i = 1, 2, \dots n$

Fix $C \|w\| = 1$ (since we can arbitrarily scale $w$ and $w_0$)

$$\underset{w, w_0}{\text{maximize}} \quad 1/\|w\|_2$$

such that $y_i (w^T x_i + w_0) \geq 1$, $i = 1, 2, \dots n$

Primal SVM Problem
(Constrained
Optimization
Problem)

$$\underset{w, w_0}{\text{Minimize}} \quad \|w\|_2^2$$

such that $y_i (w^T x_i + w_0) \geq 1$, $i = 1, 2, \dots n$

$1/\|w\|$ $\quad 2/\|w\|$ $\rightarrow$ margin

Hyperplane:
$w^T x + w_0 = 0$

SVM finds hyperplane with maximum margin

## General Constrained Optimization Problem

$$\min_x f_0(x)$$
$$\text{st} \quad f_i(x) \le 0 \quad, \quad i = 1, 2, \ldots m$$
$$\quad\quad h_i(x) = 0 \quad, \quad i = 1, 2, \ldots, p$$

① Primal Problem

Here $x \in \mathbb{R}^n$ (i.e., $f_i : \mathbb{R}^n \to \mathbb{R}$, $h_i : \mathbb{R}^n \to \mathbb{R}$)

Lagrangian $\quad L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \to \mathbb{R}$

$$L(\underline{x}, \underline{\lambda}, \underline{v}) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} v_i h_i(x)$$

where $\lambda_i$ is Lagrange multiplier for i-th inequality constraint
& $v_i$ is Lagrange multiplier for i-th equality constraint

$\underline{\lambda}$ and $\underline{v}$ are also called dual variables

## Lagrange Dual Function

$$g(\lambda, v) = \inf_x L(x, \lambda, v) = \inf_\lambda \left( f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} v_i h_i(x) \right)$$

$\lambda$ & $v$ are dual feasible if $\lambda \ge 0$ & $g(\lambda, v) > -\infty$

Fact 1 : $g(\lambda, v) \le p^*$ for any dual feasible $\lambda, v$
(where $p^*$ is optimal value of the primal problem ①, i.e., $p^* = f_0(x^*)$)

Fact 2 : If $\exists$ dual feasible $\lambda^*, v^*$ $(\lambda^* \ge 0)$ and primal feasible $x^*$ such that
$$g(\lambda^*, v^*) = p^* = f_0(x^*)$$ then strong duality is said to hold