# Contents

## 0.1   Linear Separator

Assuming that red and blue datasets represents points $X_1$ and $X_2$, then the two sets $X_1$ and $X_2$ are linearly separable if there exists $(n + 1)$ real numbers $w_1, w_2, \ldots, w_n, k$

- such that every point in $X_1$ satisfies $\sum_{i=1}^{n} w_i x_i < k$

- such that every point in $X_2$ satisfies $\sum_{i=1}^{n} w_i x_i > k$

Binary classification $y_i \in \{-1, 1\}$ can be viewed as the task of separating classes in feature space.

- Hypothesis class of linear decision surfaces is $f(x_i) = \text{sign}(\mathbf{w}^T \mathbf{x_i} + b)$.

- Without loss of generality, we assume that $b = 0$. Thus, we get the simplified $f(x_i) = \text{sign}(\mathbf{w}^T \mathbf{x_i})$.

- $(y_i)(\mathbf{w}^T \mathbf{x_i}) > 0 \Leftrightarrow$ data point $x_i$ is correctly classified.

  - Remember, $y_i$ is counting as 1 or -1.

## 0.2   Perceptron Algorithm

- Set time $t = 1$, start with vector $\mathbf{w_1} = \vec{0}$.

- Given example $\mathbf{x}$, predict positive iff (if and only if) $\mathbf{w_1} \cdot \mathbf{x} \geq \mathbf{0}$.

- On a mistake, update as follows:

  - Mistake on positive, then update $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \mathbf{x}$.
  - Mistake on negative, then update $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \mathbf{x}$.

## 0.3   Geometric Margin

The margin of example  . . .

The margin $\gamma$ of a set of examples $S$ w.r.t (with respect to) a linear separator $\mathbf{w}$ is the largest margin over points $\mathbf{x} \in S$. Theorem: If the data has a margin $\gamma$ and all points lie inside a ball of radius $R$, then the Perceptron algorithm makes $\leq \frac{R}{\gamma^2}$ mistakes.

## 0.4   Support Vector Machine

Support vector machines (SVMs) are supervised max-margin models with associated learning algorithms.

- Good generalization in theory.

- Good generalization in practice.

- Work well with few training instances.

- Find globally best model.

- Efficient algorithms.

- Amenable to the kernel trick.

## 0.5   Optimal Linear Separator

Which of the linear separators is optimal?

## 0.6   Classification Margin

Examples closest to the hyperplane are support vectors. Margin $\rho$ of the separator is the distance between support vectors.

## 0.7   Maximizing the Margin

- Better Generalization – A larger margin allows the SVM to better generalize to new, unseen data, leading to higher predictive accuracy.

- Improved Robustness – A larger margin can lead to improved robustness against noise and outliers in the training data, as it allows for greater tolerance of misclassified examples.

- Reducing Overfitting – A larger...

## 0.8   Linear SVM

Let training set $\left\{(\mathbf{x}_i, y_i)_{i=1\ldots n}, \mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, 1\}\right\}$ be separated by a hyperplane with margin $\rho$. Then for each training example $(\mathbf{x}_i, y_i)$

$$\mathbf{w}^T\mathbf{x}_i + b \geq 1 \qquad \text{if } y_i = 1$$
$$\Leftrightarrow y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1$$
$$\mathbf{w}^T\mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

Geometrically, the distance between the 2 hyperplanes can be expressed as:

$$\rho = \frac{2}{||w||} \tag{1}$$

Then we can formulate the quadratic optimization problem:

Find $\mathbf{w}$ and $b$ such that

$$\rho = \frac{2}{||\mathbf{w}||}$$

is maximized and for all $(\mathbf{x}_i, y_i), i = 1\ldots n : y_i\left(\mathbf{w}^T\mathbf{x}_i + b\right) \geq 1$

Which can be reformulated as:

$\mathbf{x}_i, y_i$, find $\mathbf{w}$ and $b$ such that

Minimize $Q(w) = \frac{1}{2}||\mathbf{w}||^2 = \frac{1}{2}\mathbf{w}^T\mathbf{w}$

subject to   ...

## 0.9   Lagrangian Duality

- Need to optimize a quadratic function subject to linear constraints.

- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.

- Solution involves constructing dual problem where Lagrange multipliers $a_i$ is associated with all inequality constraint in primal (original) problem:

$\forall i$, find $a_1, \ldots, a_n$ such that ... subject to $a_i \geq 0$