# DMML, 3 March 2020

Classification — Supervised learning

Sentiment Analysis

Naively — each word has a score (+ve or -ve)
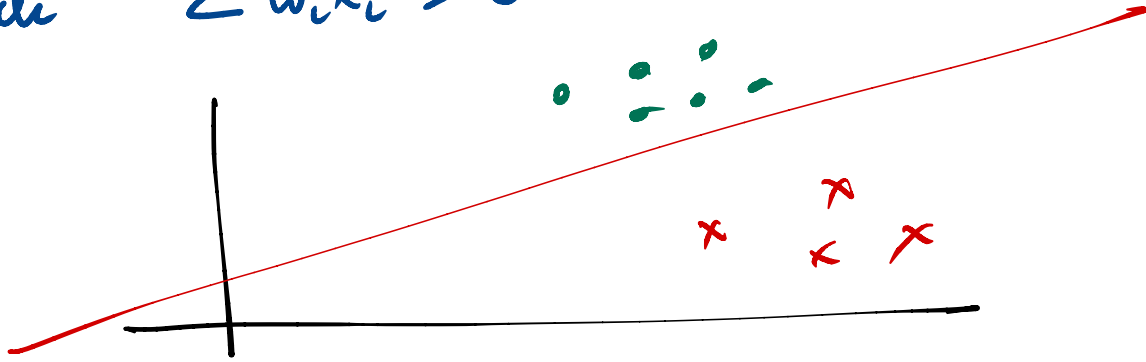
Compute weighted sum of scores of new review

If above threshold — positive

score of word 2

Words are features/attributes : $x_1, x_2, \ldots, x_N$

Compute $\sum w_i x_i$ for words in the document

Check $\quad \sum w_i x_i > t$



Linear separator

Geometric interpretation of data

Separate categorie by a hyperplane

Data items are of the form $\bar{x} = \langle x_1, x_2 \ldots, x_m \rangle$

$\quad x_i \in \mathbb{R}$

Assume classification allows linear separability

$\quad \exists \, \bar{w} = \langle w_1, w_2, \ldots, w_m \rangle$ such that

$\qquad$ For each positive $\bar{x}$, $\quad \bar{w} \cdot \bar{x} > t$

$\qquad\qquad$ negative $\bar{x}$, $\quad \bar{w} \cdot x < t$

How do we find $\bar{w}$ ?

# Instead

$$\overline{W} \cdot \overline{x} - t > 0 \qquad \text{Positive } \overline{x}$$

$$\overline{W} \cdot \overline{x} - t < 0 \qquad \text{Negative } \overline{x}$$

Expand $\overline{W}$ to $\langle \overline{W}, -t \rangle \quad \sim \hat{W}$ ] → write as
$\overline{x}$ to $\langle \overline{x}, 1 \rangle \quad \sim \hat{x}$ $\hat{W}, \overline{x}$
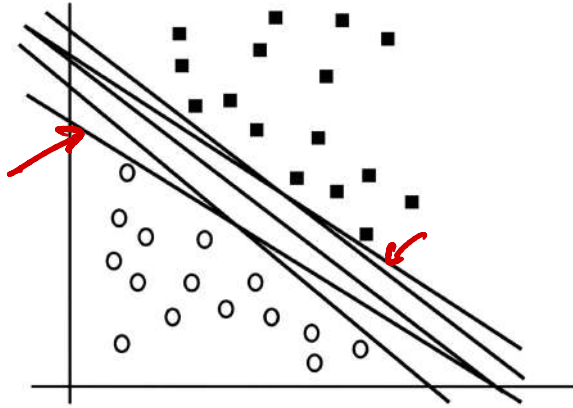
$$\hat{W} \hat{x} > 0$$

$$\hat{W} \hat{x} < 0$$

Each input $\bar{x}_i$ has a label — Yes/No

$$l_i = +1 \quad -1$$

Now $\forall i. \quad \bar{w} \cdot \bar{x}_i \cdot l_i > 0$

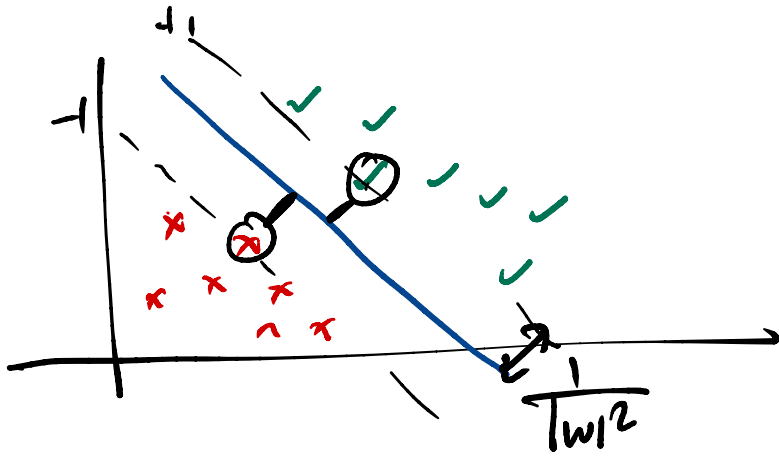Direct way is to use linear programming

But, there is a simpler way

Many possible $\bar{w}$ may work

Right now, we are happy to find any one

Want $w^*$ s.t. $w^* \cdot x_i \cdot l_i > 0$ $\forall i$

Scale $w^*$ s.t. $w^* \cdot x_i \cdot l_i > 1$ $\forall i$

Distance of nearest point is $\dfrac{1}{|w_i^*|^2}$

$$\frac{1}{|w|^2}$$

# Algorithm $\qquad$ <span style="color:red">[Perceptron]</span>

$$W \leftarrow 0$$

while there exists $x_i$ s.t. $W \cdot x_i \cdot \ell_i \leq 0$

$$W \leftarrow W + x_i \ell_i = \begin{array}{l} W + x_i \text{ if } x_i \text{ positive} \\ W - x_i \text{ if } x_i \text{ negative} \end{array}$$

Why does this converge?

## Theorem

Suppose $\exists w^*$ s.t $w^* \cdot x_i \cdot l_i > 1$ $\forall i$

Perceptron algorithm finds $w$ s.t. $w \cdot x_i \cdot l_i > 0$ $\forall i$

in at most $r^2 |w^*|^2$ updates, where $r = \max |x_i|$

No guarantee on "quality" of $w$

## Proof

Keep track of $w^T w^*$ , $|w|^2$

right-side note
$w$ current estimate

$w^*$ is assumed

1. Each update to $w$ increases $w^T w^*$ by at least 1

$$(w + x_i l_i)^T w^* = w^T w^* + x_i^T l_i w^*$$

at least +1

Suppose $l_i = +1$    $w^* x_i^T > 1$    $> 1$

$l_i = -1$    $w^* x_i^T < -1$    $> 1$

2. Each update to $w$ increases $|w|^2$ by at most $r^2$

$$(w + x_i l_i)^T (w + x_i l_i) = |w|^2 + 2 x_i^T l_i w + |x_i l_i|^2$$

$\leq 0$

$$\leq |w|^2 + |x_i|^2$$

Suppose we update $w$ $m$ times

$$w^T w^* \geq m \qquad \text{(grows by at least 1 each update)}$$

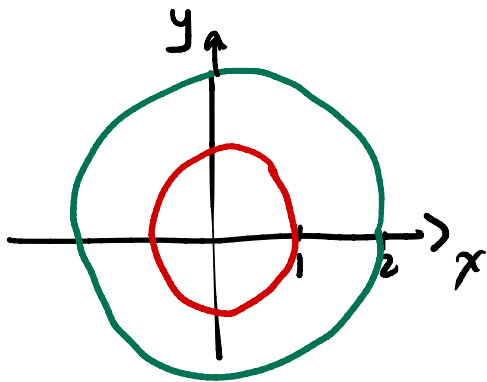$$|w|^2 \leq mr^2 \qquad \text{(each update increases by at most } r^2\text{)}$$

$$m \leq |w||w^*|$$

$$\frac{m}{|w^*|} \leq |w| \leq \sqrt{m}\, r$$

$$\sqrt{m} \leq r \cdot |w^*|$$

$$m \leq r^2 |w^*|^2$$

But what about the assumption of linear separability?



Good points — circle of radius 2

Bad points — circle of radius 1

Not linearly separable.

Geometric transformation.

$$\langle x, y \rangle \longrightarrow \langle x, y, x^2 + y^2 \rangle$$

$$z = 1.5 \quad \text{plane separates}$$

Given a transformation $\bar{x} \longmapsto \varphi(\bar{x})$

## Perceptron

$$W = 0 \pm x_i \pm x_j \pm x_k \pm \cdots$$

$$W = \underbrace{\langle u_1, u_2, \ldots, u_m \rangle \langle x_1, x_2, \ldots, x_m \rangle}_{x_{tram}}$$

Classify a new $\bar{z}$

$$\bar{w} \cdot \bar{z} > 0 \ ? \qquad \left( \bar{u} \cdot \bar{x}_{tram} \right) \cdot \bar{z} > 0$$

$$\left( \bar{u} \cdot \overline{\varphi(\bar{x})} \right) \cdot \varphi(\bar{z}) > 0$$

All we need to know about $\varphi$ is how to compute dot products $\varphi(\bar{x}_i) \cdot \varphi(\bar{x}_j)$

Suppose we have a function $k(\bar{x}_i, \bar{x}_j)$

s.t. $\forall \bar{x}_i, \bar{x}_j \quad k(\bar{x}_i, \bar{x}_j) = \varphi(\bar{x}_i) \cdot \varphi(\bar{x}_j)$

Such a k is called a kernel function

$\exists \varphi$ s.t. $k(\bar{x}_i, \bar{x}_j) = \varphi(\bar{x}_i) \cdot \varphi(\bar{x}_j) \quad \forall \bar{x}_i, \bar{x}_j$

**Strategy** Try out various kernels (don't worry about what $\varphi$ looks like)

When is $k$ a kernel?

Minimum requirement — symmetric

Non constructive characterization — Mercer's Theorem

Constructive definition in terms of positive semidefinite matrices