# The perceptron algorithm

## Frank Rosenblatt (1928 - 1971)

Two-class model; input $\vec{x}$ => feature $\vec{\phi}(\vec{x})$ vector

typically, (bias) $\phi_0(\vec{x}) = 1$ is included

Consider $y(\vec{x}) = f(\vec{w}^T \cdot \vec{\phi}(\vec{x}))$, where

$$f(a) = \begin{cases} +1 & , a \geq 0 \\ -1 & , a < 0 \end{cases}$$

Assume $t = \begin{cases} +1 & c_1 \\ -1 & c_2 \end{cases}$

Error function : perceptron criterion

Note that we seek $\vec{w}$ s.t. all patterns in $c_1$ have $\vec{w}^T \cdot \vec{\phi}(\vec{x}_n) > 0$, and all patterns in $c_2$ have $\vec{w}^T \cdot \vec{\phi}(\vec{x}_n) < 0$.

Then $\vec{w}^T \cdot \vec{\phi}(\vec{x}_n) t_n > 0$ in both classes

Let's assign $\emptyset$ error to all correctly classified points (patterns), and $-\vec{w}^T \cdot \vec{\phi}(\vec{x}_n) t_n$ to all missclassified points. Then

$$E(\vec{w}) = - \sum_{n \in N} \vec{w}^T \cdot \vec{\phi}(\vec{x}_n) t_n$$

↑ set of all misclassified points

needs to be minimized.

Note that $E(\vec{w}) > 0$ since $\sum_{new}$ sums over misclassified patterns only.

-1-

$E(\vec{w})$ is a piecewise linear ⏜contribution of nth t⏜ $\underbrace{\sim \vec{w}}$ in $\underbrace{}_{\text{is}}$ in $E(\vec{w})$
regions of $\vec{w}$ space where $\vec{x}_n$ is
misclassified, and $\emptyset$ otherwise.

Now optimise $E(\vec{w})$ by gradient descent:
cycle through n and at each step,
update

$$\vec{w}^{\,(\tau+1)} = \vec{w}^{\,(\tau)} - \eta \, \vec{\nabla}_n E(\vec{w}) =$$

with annotations: step number ↓, learning rate ↓

$$= \vec{w}^{\,(\tau)} + \eta \, \vec{y}(\vec{x}_n) t_n$$

Set $\eta = 1 \iff y(\vec{x})$ is insensitive to the choice of $\eta$.

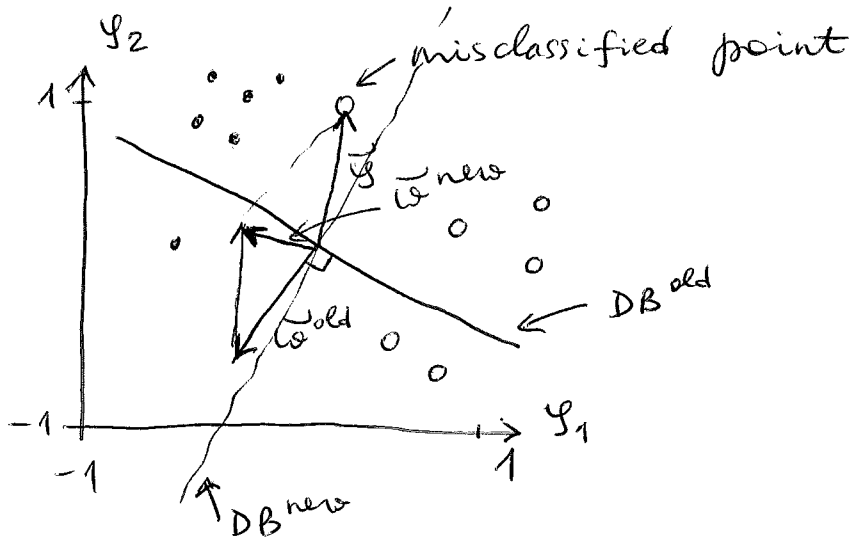Note that $\overset{\vec{y}(\vec{x}_n)}{\overbrace{\phantom{xx}}}$

$$-\vec{w}^{\,(\tau+1)T} \vec{y}_n t_n = -\vec{w}^{\,(\tau)T} \vec{y}_n t_n - \underbrace{(\vec{y}_n t_n)^T (\vec{y}_n t_n)}_{> 0} < -\vec{w}^{\,(\tau)T} \vec{y}_n t_n.$$

Thus the error always decreases for
a given term, but some other patterns
may become misclassified $\Rightarrow E(\vec{w})$ is
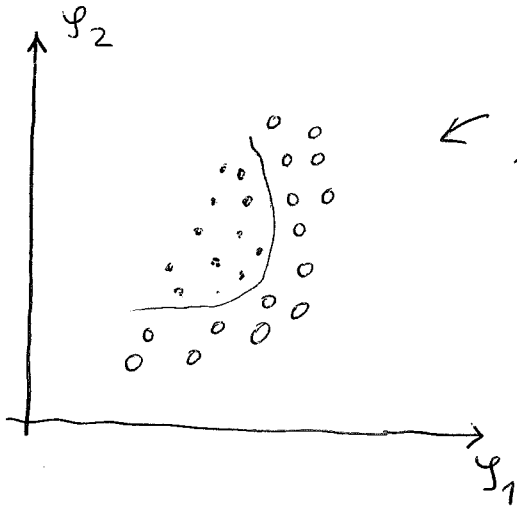not guaranteed to be reduced at each
step.

However, one can argue that the
perceptron finds a solution (but may
always
need a lot of steps to do so) for
linearly separable problems.
[The solution may not be unique]

−2−

Ex.



Difficulties: not probabilistic, does not generalize to $k > 2$ classes easily.



← lack of convergence on linearly non-separable datasets

# Probabilistic models

Class-conditional density: $p(\vec{x} \mid C_k)$
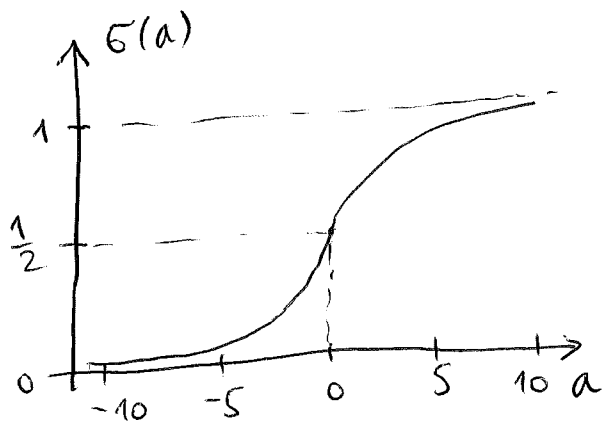$$k = 1, \ldots, K$$

Class prior: $p(C_k)$

Class posterior: $p(C_k \mid \vec{x})$

### $K = 2$ case:

$$p(C_1 \mid \vec{x}) = \frac{p(\vec{x} \mid C_1)\, p(C_1)}{p(\vec{x} \mid C_1)\, p(C_1) + p(\vec{x} \mid C_2)\, p(C_2)} =$$

$$= \frac{1}{1 + \underbrace{\dfrac{p(\vec{x} \mid C_2)\, p(C_2)}{p(\vec{x} \mid C_1)\, p(C_1)}}_{\text{"} e^{-a}}} = \sigma(a)$$

$\uparrow$ sigmoid f'n

$$a = \log \frac{p(\vec{x} \mid C_1)\, p(C_1)}{p(\vec{x} \mid C_2)\, p(C_2)}$$



Note that

$$\sigma(-a) = \frac{1}{1 + e^{a}}$$

"
$$1 - \sigma(a) = 1 - \frac{1}{1 + e^{-a}} =$$

$$= \frac{e^{-a}}{1 + e^{-a}} = \frac{1}{e^{a} + 1}$$

also, $\quad a = \log \left( \dfrac{\sigma}{1 - \sigma} \right)$

$\underbrace{\phantom{a = \log \left( \dfrac{\sigma}{1-\sigma} \right)}}_{\text{logit f'n}}$

## $K > 2$ case:

$$p(C_k | \vec{x}) = \frac{p(\vec{x}|C_k)\, p(C_k)}{\sum_{j=1}^{k} p(x|C_j)\, p(C_j)} =$$

$$= \frac{e^{a_k}}{\sum_j e^{a_j}} \quad , \text{ if } \quad a_k = \log[\, p(\vec{x}|C_k)\, p(C_k)\,]$$

$\underbrace{\phantom{= \frac{e^{a_k}}{\sum_j e^{a_j}}}}$ softmax f'n:

$$\text{if } a_k \gg a_j \;,\; \forall_j \neq k$$

$$\frac{e^{a_k}}{\sum_j e^{a_j}} \simeq 1 = p(C_k|\vec{x}) \implies p(C_j|\vec{x}) \simeq 0, \quad \forall_j \neq k$$

Now assume that

$$p(\vec{x}|C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x}-\vec{\mu}_k)^T \overset{\downarrow \text{same for all classes for now}}{\Sigma^{-1}}(\vec{x}-\vec{\mu}_k)}$$

$$D = \#\text{ dim's in } \vec{x}$$

(k=2 case:)

Then $\quad p(C_1|\vec{x}) = \sigma\left( \log \frac{p(\vec{x}|C_1)}{p(\vec{x}|C_2)} + \log \frac{p(C_1)}{p(C_2)} \right)$

$\underbrace{\phantom{\log \frac{p(\vec{x}|C_1)}{p(\vec{x}|C_2)} + \log \frac{p(C_1)}{p(C_2)}}}$

$$\log \frac{p(C_1)}{p(C_2)} - \frac{1}{2}(\vec{x}-\vec{\mu}_1)^T \underset{\Sigma^{-1}}{\Sigma^{-1}}(\vec{x}-\vec{\mu}_1) + \frac{1}{2}(\vec{x}-\vec{\mu}_2)^T \Sigma^{-1}(\vec{x}-\vec{\mu}_2) =$$

$$= \frac{1}{2}\vec{x}^T \underset{\Sigma^{-1}}{\vec{\mu}_1} + \frac{1}{2}\vec{\mu}_1^T \underset{\Sigma^{-1}}{\vec{x}} - \frac{1}{2}\vec{x}^T \underset{\Sigma^{-1}}{\vec{\mu}_2} + \frac{1}{2}\vec{\mu}_2^T \underset{\Sigma^{-1}}{\vec{x}} +$$

$$+ \log \frac{p(C_1)}{p(C_2)} - \frac{1}{2}\vec{\mu}_1^T \Sigma^{-1}\vec{\mu}_1 + \frac{1}{2}\vec{\mu}_2^T \Sigma^{-1}\vec{\mu}_2 =$$

$$= \frac{1}{2}(\Sigma^{-1}\vec{\mu}_1)^T \vec{x} + \frac{1}{2}\underbrace{(\Sigma^{-1}\vec{\mu}_1)^T}_{\vec{\mu}_1^T(\Sigma^{-1})^T}\vec{x} - \frac{1}{2}(\Sigma^{-1}\vec{\mu}_2)^T \vec{x} - \frac{1}{2}(\Sigma^{-1}\vec{\mu}_2)^T \vec{x} + \ldots$$

$\overset{?}{(\Sigma^{-1})^T} = \Sigma^{-1}$

Defining $\begin{cases} \vec{\omega} = \Sigma^{-1}(\vec{\mu}_1 - \vec{\mu}_2), \\ \omega_0 = -\frac{1}{2}\vec{\mu}_1^T \Sigma^{-1}\vec{\mu}_1 + \frac{1}{2}\vec{\mu}_2^T \Sigma^{-1}\vec{\mu}_2 + \log\frac{p(C_1)}{p(C_2)}, \end{cases}$
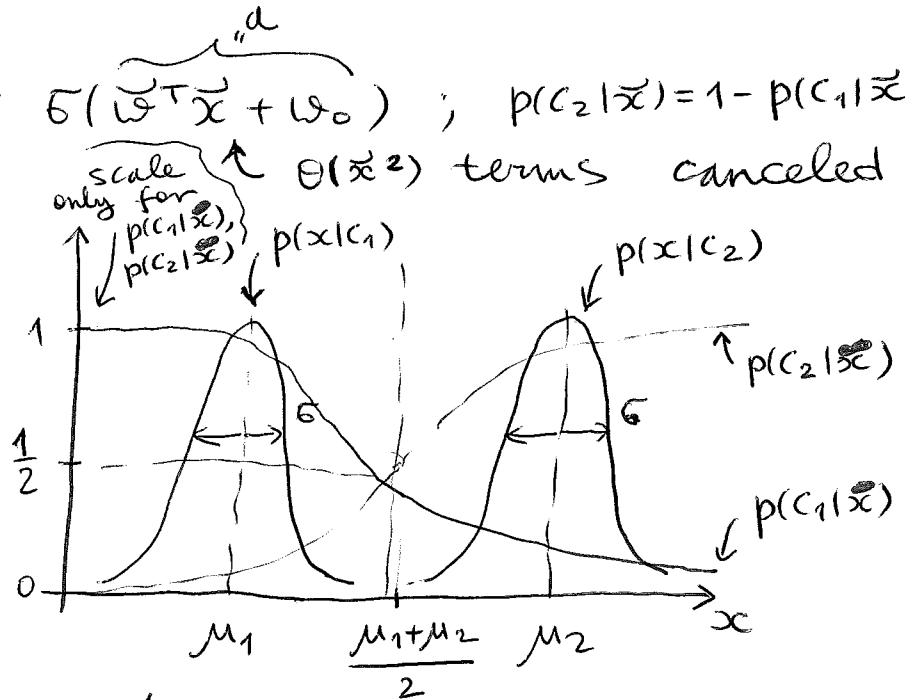
we obtain:

$$p(C_1|\vec{x}) = \sigma(\overbrace{\vec{\omega}^T\vec{x} + \omega_0}^{\text{"}a}) \quad ; \quad p(C_2|\vec{x}) = 1 - p(C_1|\vec{x})$$

$\uparrow \mathcal{O}(\vec{x}^2)$ terms canceled

<u>$D=1$ exc.:</u>

assume

$p(C_1) = p(C_2) = \frac{1}{2}$



Here, $\Sigma^{-1} \Rightarrow \frac{1}{\sigma^2}$ ;

$\begin{cases} \omega = \dfrac{\mu_1 - \mu_2}{\sigma^2}, \\ \omega_0 = \dfrac{1}{2\sigma^2}(\mu_2^2 - \mu_1^2) \end{cases} \Rightarrow a = \dfrac{(\mu_1 - \mu_2)x}{\sigma^2} + \dfrac{1}{2\sigma^2}(\mu_2^2 - \mu_1^2)$

$a = 0$ when $x = \dfrac{\mu_1 + \mu_2}{2}$,

s.t. $\sigma(a) = \frac{1}{2}$.

DB: surfaces at which

$$p(C_1|\vec{x}) = p(C_2|\vec{x})$$

$\qquad \| \qquad\qquad \| \qquad \swarrow 1 - \sigma(a) = \sigma(-a)$

$\sigma(\vec{\omega}^T\vec{x} + \omega_0) \qquad \sigma(-\vec{\omega}^T\vec{x} - \omega_0)$

Then $\vec{\omega}^T\vec{x} + \omega_0 = -\vec{\omega}^T\vec{x} - \omega_0$,

$\vec{\omega}^T\vec{x} + \omega_0 = 0 \Leftarrow$ linear eq'n for DB

In our 1D example,

$$\underbrace{\omega \cdot x + \omega_0}_{a} = 0 \quad \text{is satisfied by}$$

$$x = \frac{\mu_1 + \mu_2}{2} \Leftarrow DB$$

In fact, all contours of $p(C_1|\vec{x})$ & $p(C_2|\vec{x})$ are given by functions linear in $\underline{\underline{x}}$.

(K > 2 case:)

$$d_k(\vec{x}) = -\frac{D}{2}\log(2\pi) - \frac{1}{2}\log|\Sigma| - \frac{1}{2}(\vec{x}-\vec{\mu}_k)^T \Sigma^{-1}(\vec{x}-\vec{\mu}_k) +$$

$$+ \log p(C_k) = \quad \overbrace{\phantom{xxx}}^{\vec{\omega}_k}$$

$$= -\frac{1}{2}\vec{x}^T\Sigma^{-1}\vec{x} + (\overbrace{\Sigma^{-1}\vec{\mu}_k})^T\vec{x} \overbrace{- \frac{1}{2}\vec{\mu}_k^T\Sigma^{-1}\vec{\mu}_k}^{\omega_{0,k}} -$$

$$- \frac{D}{2}\log(2\pi) - \frac{1}{2}\log|\Sigma| + \log p(C_k) =$$

$$= -\frac{1}{2}\vec{x}^T\Sigma^{-1}\vec{x} + \vec{\omega}_k\vec{x} + \omega_{0,k} .$$

DB's are still linear however:

$$p(C_k|\vec{x}) = p(C_j|\vec{x}) \qquad [j \neq k]$$

$$\Downarrow$$

$$e^{a_k} = e^{a_j}, \quad \text{or} \quad a_k = a_j:$$

$$\vec{\omega}_k^T\vec{x} + \omega_{0,k} = \vec{\omega}_j^T\vec{x} + \omega_{0,j} .$$

$$\uparrow \qquad O(\vec{x}^2) \text{ terms cancel}$$

linear eq'n in $\vec{x}$, gives $D-1$ dimensional decision surface (DB)

Finally, if $p(\vec{x}|C_k) = \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$,
the DB's at which the 2 largest
posterior probs are equal, are given
by quadratic f's of $\vec{x}$:

(2D)

$x_2$

$p(\vec{x}|C_1) = \mathcal{N}(\vec{x}|\vec{\mu}_1, \Sigma)$                    DB

$p(\vec{x}|C_3) = \mathcal{N}(\vec{x}|\vec{\mu}_3, \Sigma')$

$\#\Sigma$

$p(\vec{x}|C_2) = \mathcal{N}(\vec{x}|\vec{\mu}_2, \Sigma)$

$x_1$