1) perceptron dates back to early 50s. biologically inspired (and can be viewed as a "single neuron" network)

L7

Mar. 31st (Midterm)
HW 4, due Apr. 2nd

2) support vector machines: funny name (and not intuitive). as we discussed the reason for "support vectors" comes by examining the dual. in some sense, the optimal model (w) is "supported" on a subset of the data points
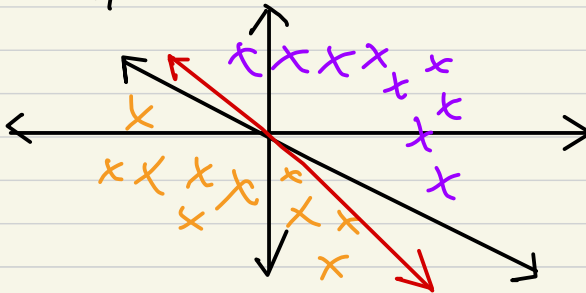
Recap:

$\begin{cases} K\text{-}NN \\ \text{The classification Algorithm} \end{cases}$

Support vector machines (SVMs)

kernel methods

## Support Vector Machines (SVMs)



there could be multiple lines (separators)
Infinite!

Q: What is the "best" separator b/w the classes?

Revisiting the Perceptron

Dataset $\{(x_1, y_1), (x_2, y_2), (x_3, y_3) \cdots \}$

Find $W$ s.t. $\text{sign} \langle W, x_i \rangle = y_i$

such that

Rewriting this :

$$\ell_i(w) = \begin{cases} 0 & \text{, if } \mathrm{sign}\langle w, x_i\rangle = y_i \\ 1 & \text{, otherwise} \end{cases}$$

$\Rightarrow y_i \cdot \mathrm{sign}\langle w, x_i\rangle \geq 1$   why not $\geq 0$

$\Rightarrow y_i \langle w, x_i\rangle \geq 1$

$\Rightarrow 1 - y_i \langle w, x_i\rangle < 1$

Intuition: 1) if $y_i$ & $\langle w, x_i\rangle$ are of same sign

$$\begin{cases} 0 \quad \text{AND } \langle w, x_i\rangle \text{ is large} \\ \qquad\quad \text{then no loss} \\ 1 \quad \text{2) if } y_i \text{ & } \langle w, x_i\rangle \text{ are of opposite signs, then loss} \geq 1 \end{cases}$$

3) everything else in between 0 & 1

$$L(w) = \sum_{i=1}^{n} \ell_i(w)$$

$$= \sum_{i=1}^{n} \max(0, 1 - y_i\langle w, x_i\rangle)$$

Hinge Loss

$f(z) = \max(0, 1-z)$

$$\ell_i(w) = \begin{cases} 0 & \text{, } y\langle w, x_i\rangle \geq 1 \\ 1 - y_i\langle w, x_i\rangle & \text{otherwise} \end{cases}$$

$$\partial \ell_i(w) = \begin{cases} 0 & \text{, } y\langle w, x_i\rangle \geq 1 \\ -y_i x_i & \text{, otherwise} \end{cases}$$

(Sub) gradient descent :

不是严格意义
的梯度 / GD ...

$$w_{k+1} \longleftarrow w_k - \eta_k \, \partial L(w_k)$$

$$= \begin{cases} w_k & \text{, } y_i\langle w, x_i\rangle \geq 1 \\ w_k + \eta \, y_i x_i & \text{, otherwise} \end{cases}$$

## SVM

$$L_{SVM}(w) = L_{hinge}(w) + \frac{\lambda}{2}\|w\|_2^2 \quad \text{support vector}$$

$$= \sum_{i=1}^{n} \max(0, 1-y_i\langle w, x_i\rangle) + \frac{\lambda}{2}\|w\|_2^2$$

"primal SVM"



(Primal)     $\min\limits_{w} L_{SVM}(w)$

Solution: $W^*$

(Dual)     $\max\limits_{\alpha} -\sum_{i=1}^{n} \alpha_i$

$$-\sum y_i y_j \alpha_i \alpha_j \langle x_i, x_j\rangle$$

such that     $0 \leq \alpha_i \leq \frac{1}{\lambda}$

{ Hinge loss
{ regularizer
SVM

Solution:     $\alpha^*$

$$\boxed{W^* = \sum_{i=1}^{n} \alpha_i^* y_i x_i}$$

Intuition: SVM model is a linear combination of
data points

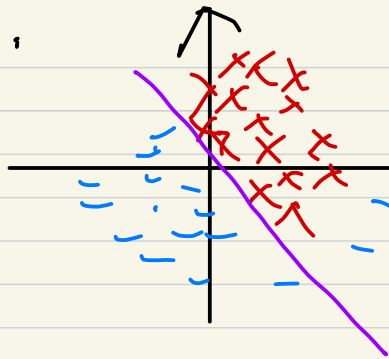wherever $\alpha_i^* = 0$, solution $W^*$ doesn't depend on $x_i$

∴ The data points $x_i$ for

are called "support vectors"

Separability.

Separable

(could be small mistake)

( could be no-100% $\checkmark$ )

this is not linear separable.

Ex 1

Ex 2

EX] :  $(X_1, X_2) \rightarrow (X_1, X_2, \sqrt{X_1^2 + X_2^2})$

blue is outer and higher
with third dimension

EX2 : do we need 3rd dimension ?  No!

look at graph.   red in $I \cdot III$ phase
blue in $II \ IV$ phase

$(X_1, X_2) \longrightarrow X_1 X_2$

Embed the data into a different feature space in which

kernel methods :

$$x \longrightarrow \phi(x)$$

$\phi(\cdot)$ :   a transition of the data called
the kernel Mapping

Typical feature transformations :

* original feature
* =RBS Quadratic features
  $$(x_1, x_2) \longrightarrow (x_1^2, x_2^2, x_1 x_2)$$
* Higher order polynomials
  $$(x_1 x_2) \longrightarrow (x_1^2 \cdot x_2^2, x_1 x_2, x_1^3, x_2^3,$$
  $$x_1^2 x_2, x_1 x_2^2, x_1 x_2),$$
  $$etc.$$

$$(x_1, x_2, \cdots, x_d) \xrightarrow{\quad Quadratic \quad}$$
$$(x_1^2, x_2^2, \cdots, x_d^2, x_1 x_2, x_1 x_3, \cdots, x_1 x_d, \cdots$$
$$\cdots \quad x_{d-1} x_d)$$

approximately $d^2$ seems unpractical

Going back to circle example
$$(x_1, x_2) \longrightarrow (x_1, x_2, x_1^2, x_2^2, x_1 x_2)$$
$$w \qquad (0, 0, 1, 1, 0)$$

Circle

$$\langle w, \phi(x) \rangle = x_1^2 + x_2^2$$

$$f(x) = \text{sign}(\langle w, \phi(x)\rangle - \text{radius})$$

XOR $\quad w \in 0,0,0,0,1)$

$$f(x) = \text{sgn}(\langle w, \phi(x)\rangle)$$

What is important is <u>not</u> the <u>feature mapping</u>, but the ability to take dot products with the feature mapping

$\longrightarrow$ "kernel trick"

implicitly define feature mapping via the <u>dot product</u>

$\longrightarrow$ kernel dot product

or kernel inner product

Example of kernel dot products

☆ Regular dot product

$$(x_1, x_2) \rightarrow (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2)$$

☆ Quadratic dot product **lift dimension, but make dot product more simple**

$$\langle \phi(x), \phi(y)\rangle = k(x,y) \rightarrow (1 + \langle x, y\rangle)^2$$

$$(y_1, y_2) \rightarrow (1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2, y_1^2, y_2^2)$$

**important**

Dot product $= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 y_1 x_2 y_2$
$$+ x_1^2 y_1^2 + x_2^2 y_2^2$$

$$\overset{?}{=} (1 + x_1 y_1 + x_2 y_2)^2 = (1 + \langle [x_1, x_2], [y_1, y_2]\rangle)^2$$

☆ Cubic dot product
$$k(x,y) \rightarrow (1 + \langle x, y\rangle)^3$$

☆ Exponential
$$k(x,y) = \exp(-\|x-y\|^2) \quad \text{"Gaussian kernel"} \quad \text{"Radial basis function"} / RBF$$