# Introduction to Machine Learning
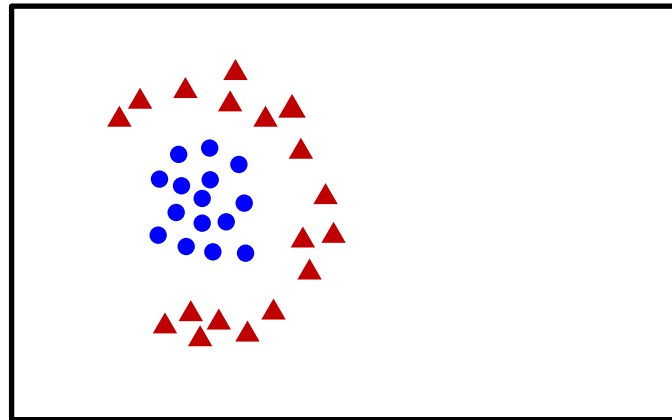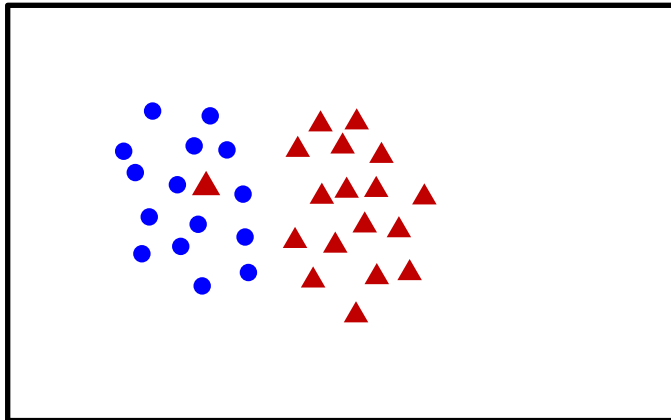## Eran Halperin, Yishay Mansour, Lior Wolf
## 2013-2014
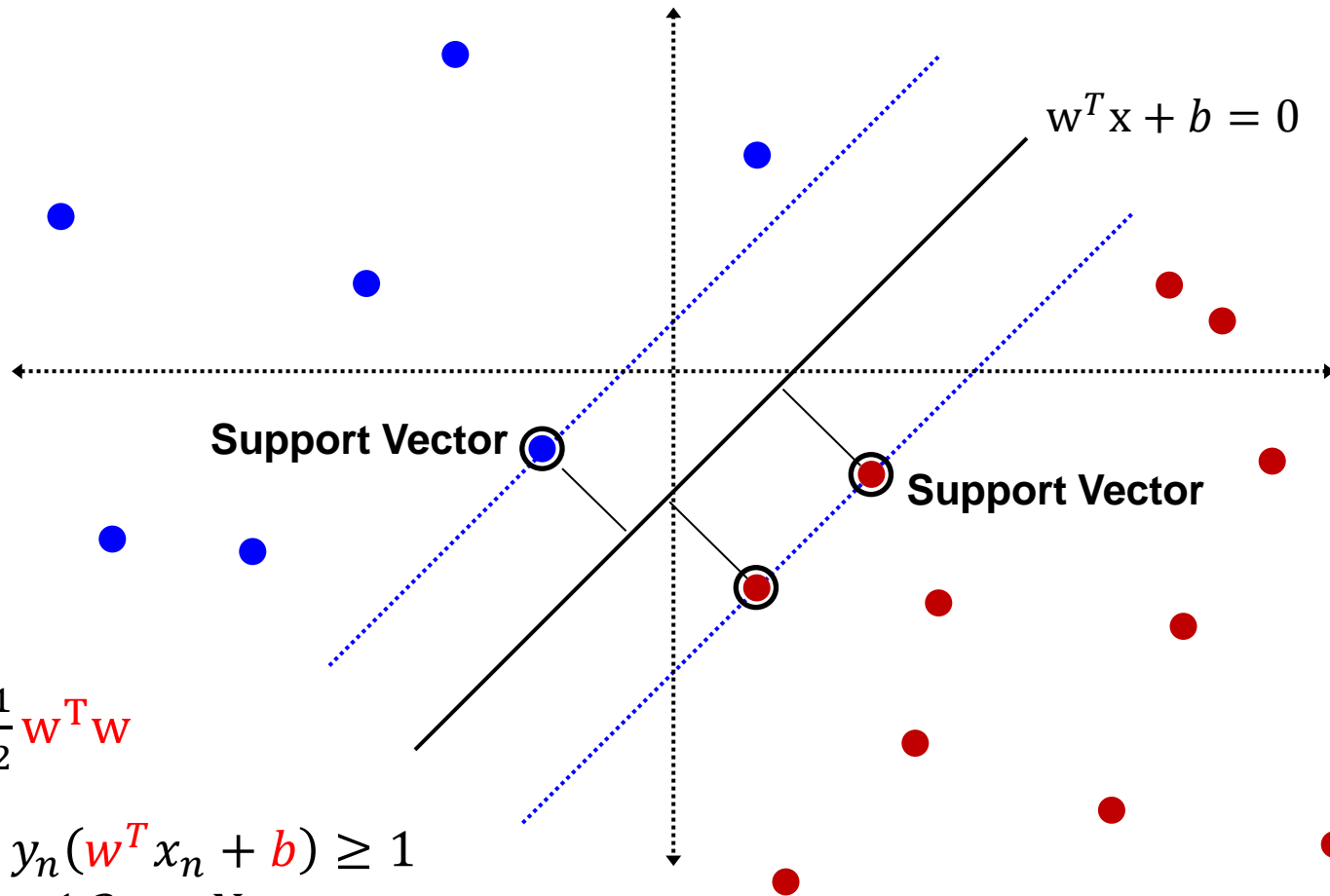
# Lecture 7: Kernels

# Outline

- Following discussions from last class
  - How many support vectors are there anyhow?
  - Positive definite matrices

- Support Vector machine (SVM) classifier
  - The kernel trick
  - Which kernels to use
  - Contructing kernels
  - SVD and kernel SVD

# Number of support vectors

linearly separable data*:  #SV $\leq$ d+1 = VC-dim

$$w^T x + b = 0$$

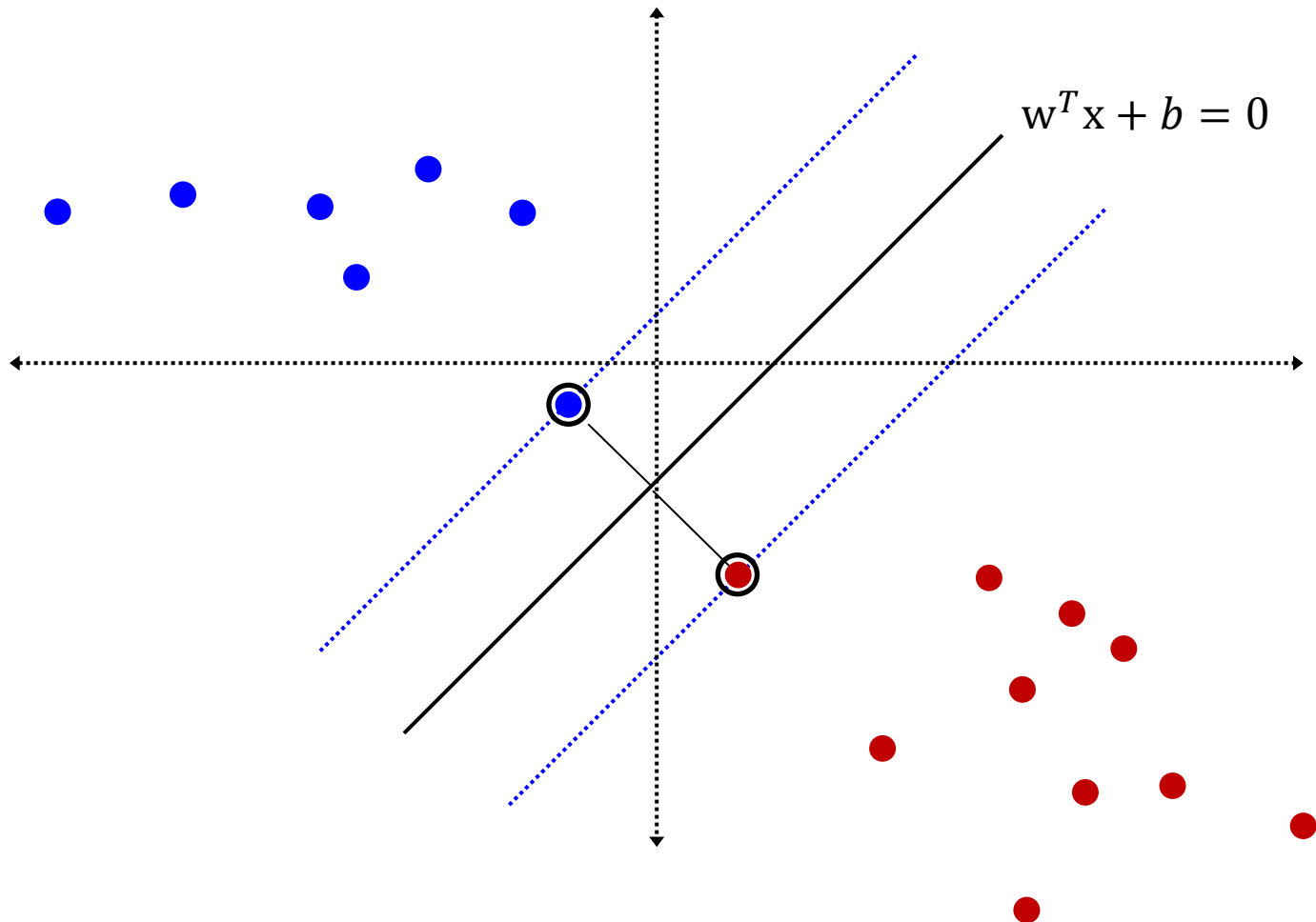**Support Vector**

**Support Vector**

Minimize $\frac{1}{2} w^T w$

subject to $y_n(w^T x_n + b) \geq 1$
$\qquad n = 1, 2, \ldots, N$

$w \in \mathrm{R}^d, b \in \mathrm{R}$

* Could be much more in ``degenerate cases''

# #SV=2 is sometimes enough



$\mathrm{w}^T\mathrm{x} + b = 0$

# Positive (Semi) Definite (PD/PSD) Matrices

(1) The *n×n* matrix **A** is positive definite if and only if:

$$\boldsymbol{Y}^T \boldsymbol{A} \boldsymbol{Y} > 0, \quad \forall \boldsymbol{Y} \neq \boldsymbol{0}$$

‣ The *n×n* matrix **A** is positive semi-definite if and only if:

$$\boldsymbol{Y}^T \boldsymbol{A} \boldsymbol{Y} \geq 0, \quad \forall \boldsymbol{Y} \neq \boldsymbol{0}$$

semi

(2) $\boldsymbol{A}$ is positive definite $\iff \exists \boldsymbol{P}$ s.t. $\boldsymbol{A} = \boldsymbol{P}\boldsymbol{P}^T, |\boldsymbol{P}| \neq \boldsymbol{0}$

$\boldsymbol{A}$ is positive definite $\Rightarrow$ it is symmetric

# Eigenvalues of PD Matrices

- Given the *n×n* matrix **A**, there are *n* eigenvalues λ and vectors **X≠0** where

$$AX = \lambda X$$

$$\begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix}^T A \begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix} = \text{diag}(\lambda_1, \lambda_2, \ldots, \lambda_n)$$

▸ If $\lambda_i > 0 \quad \forall i \in [1, n] \Longleftrightarrow A$ is positive definite

▸ If $\lambda_i \geq 0 \quad \forall i \in [1, n] \Longleftrightarrow A$ is positive semi-definite

$A$ is positive definite $\Rightarrow |A| > 0$

# Positive (Semi) Definite (PD/PSD) Matrices

The *n×n* matrix **A** is positive definite if and only if:

1. $Y^T A Y > 0, \quad \forall Y \neq 0$

2. $\exists P$ s.t. $A = P P^T, |P| \neq 0$

3. $\lambda_i \geq 0 \quad \forall i \in [1, n]$

# The dual formulation

$$\min_{\propto} \frac{1}{2} \propto^{\mathrm{T}} \begin{bmatrix} y_1 y_1 x_1^{\mathrm{T}} x_1 & y_1 y_2 x_1^{\mathrm{T}} x_2 & ... & y_1 y_N x_1^{\mathrm{T}} x_N \\ y_2 y_1 x_2^{\mathrm{T}} x_1 & y_2 y_2 x_2^{\mathrm{T}} x_2 & ... & y_2 y_N x_2^{\mathrm{T}} x_N \\ ... & ... & ... & ... \\ y_N y_1 x_N^{\mathrm{T}} x_1 & y_N y_2 x_N^{\mathrm{T}} x_2 & ... & y_N y_N x_N^{\mathrm{T}} x_N \end{bmatrix} \propto + \underbrace{\left(-1^{\mathrm{T}}\right)}_{\text{linear}} \propto$$

$$\underbrace{\phantom{M}}_{M}$$

subject to    $y^{\mathrm{T}} \underbrace{\propto}_{\text{linear constraint}} = 0$

$$\underbrace{0}_{\text{lower bounds}} \quad \leq \quad \propto \quad \leq \quad \underbrace{C}_{\text{upper bounds}}$$

# The dual formulation

$$\min_{\propto} \frac{1}{2} \propto^{\mathrm{T}} \underbrace{\begin{bmatrix} y_1 y_1 \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_1 & y_1 y_2 \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_2 & ... & y_1 y_N \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_N \\ y_2 y_1 \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_1 & y_2 y_2 \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_2 & ... & y_2 y_N \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_N \\ ... & ... & ... & ... \\ y_N y_1 \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_1 & y_N y_2 \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_2 & ... & y_N y_N \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_N \end{bmatrix}}_{\mathrm{M=A^tA}} \propto + \underbrace{\left(-1^{\mathrm{T}}\right)}_{\text{linear}} \propto$$

$$\begin{bmatrix} y_1 y_1 \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_1 & y_1 y_2 \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_2 & ... & y_1 y_N \mathrm{x}_1^{\mathrm{T}} \mathrm{x}_N \\ y_2 y_1 \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_1 & y_2 y_2 \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_2 & ... & y_2 y_N \mathrm{x}_2^{\mathrm{T}} \mathrm{x}_N \\ ... & ... & ... & ... \\ y_N y_1 \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_1 & y_N y_2 \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_2 & ... & y_N y_N \mathrm{x}_N^{\mathrm{T}} \mathrm{x}_N \end{bmatrix} = [y_1 \mathrm{x}_1, y_2\, \mathrm{x}_2, ..., y_N\, \mathrm{x}_N]^t [y_1 \mathrm{x}_1, y_2\, \mathrm{x}_2, ..., y_N\, \mathrm{x}_N]$$

# Support Vector Machine



$$\text{w}^T\text{x} + b = 0$$

**Support Vector**

**Support Vector**

**w**

$$f(x) = \sum_i \propto_i y_i \left( \text{x}_i^\text{T} \text{x} \right) + b$$
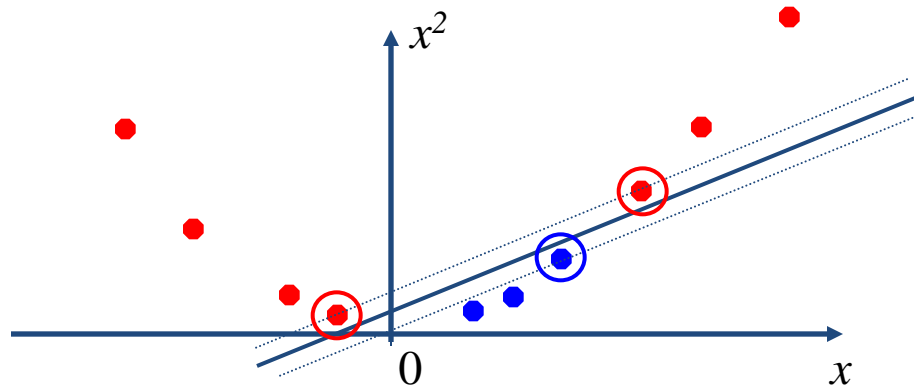
support vectors

# Nonlinear SVMs

- Datasets that are linearly separable work out great:
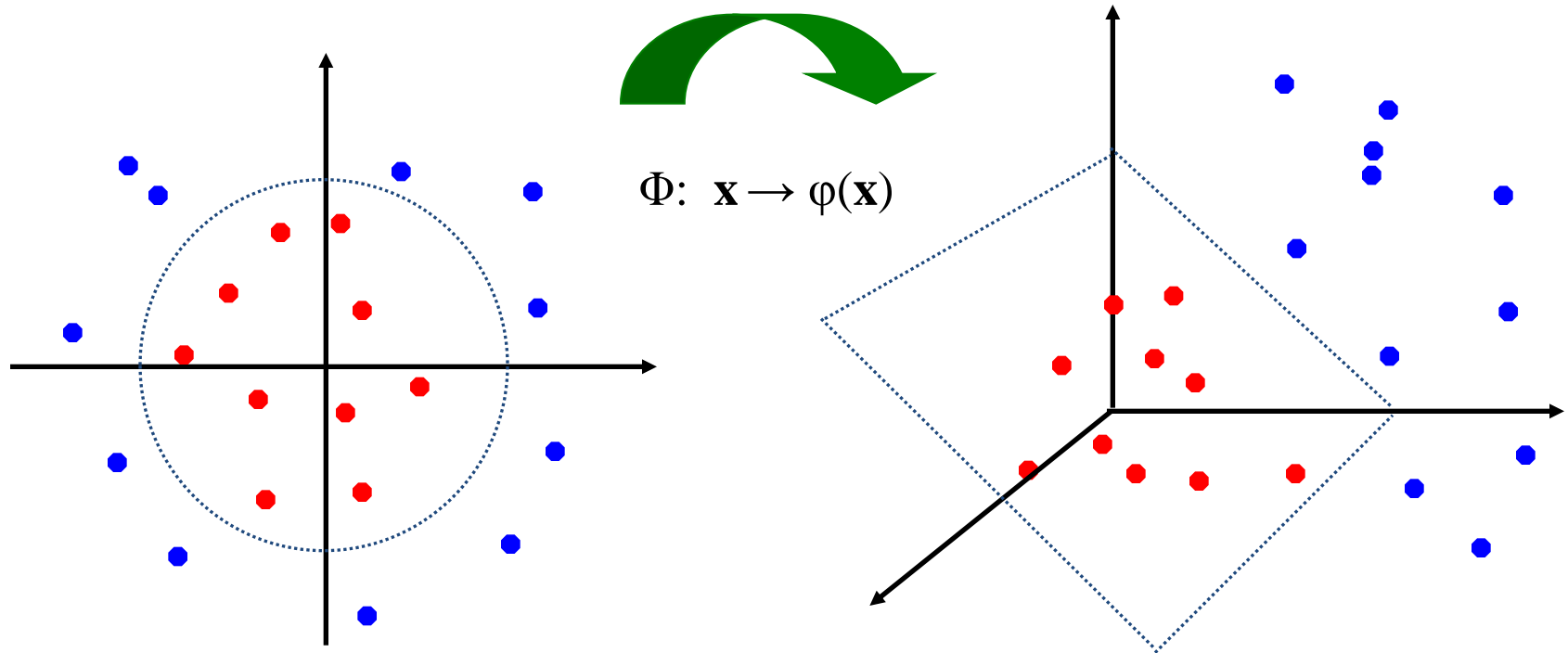


- But what if the dataset is just too hard?



- We can map it to a higher-dimensional space:



Slide credit: Andrew Moore

# Another example (2D)

# Nonlinear SVMs

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$

# A potential problem

- If we map the input vectors into a <span style="color:red">very</span> high-dimensional feature space, optimizing the SVM and even classification might become computationally intractable
  - The mathematics is the same
  - The vectors have a huge number of components
  - Taking the dot product of two vectors is very expensive
  - What would happen to the primal QP?

$$\text{Minimize } \frac{1}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}$$

$$\text{subject to } y_n(w^T\phi(x_n) + b) \geq 1$$
$$n = 1, 2, \ldots, N$$

$$\mathbf{w} \in \mathrm{R}^?, b \in \mathrm{R}$$

# A potential problem

- If we map the input vectors into a <span style="color:red">very</span> high-dimensional feature space, optimizing the SVM and even classification might become computationally intractable
  - The mathematics is the same
  - The vectors have a huge number of components
  - Taking the dot product of two vectors is very expensive
  - What would happen to the primal QP?
  - What would happen to the dual QP?

$$\min_{\propto} \frac{1}{2} \propto^{\mathrm{T}} \begin{bmatrix} y_1 y_1 \phi(\mathrm{x}_1)^T \phi(\mathrm{x}_1) & y_1 y_2 \phi(\mathrm{x}_1)^T \phi(\mathrm{x}_2) & \dots & y_1 y_N \phi(\mathrm{x}_1)^T \phi(\mathrm{x}_N) \\ y_2 y_1 \phi(\mathrm{x}_2)^T \phi(\mathrm{x}_1) & y_2 y_2 \phi(\mathrm{x}_2)^T \phi(\mathrm{x}_2) & \dots & y_2 y_N \phi(\mathrm{x}_2)^T \phi(\mathrm{x}_N) \\ \dots & \dots & \dots & \dots \\ y_N y_1 \phi(\mathrm{x}_N)^T \phi(\mathrm{x}_1) & y_N y_2 \phi(\mathrm{x}_N)^T \phi(\mathrm{x}_2) & \dots & y_N y_N \phi(\mathrm{x}_N)^T \phi(\mathrm{x}_N) \end{bmatrix} \propto + (-1^{\mathrm{T}}) \propto$$

# A potential problem

- If we map the input vectors into a <span style="color:red">very</span> high-dimensional feature space, optimizing the SVM and even classification might become computationally intractable
  - The mathematics is the same
  - The vectors have a huge number of components
  - Taking the dot product of two vectors is very expensive
  - What would happen to the primal QP?
  - What would happen to the dual QP?
  - And during classification?

$$f(x) = \sum_i \propto_i y_i (\phi(x_i)^T x) + b \qquad\qquad f(x) = w^T x + b$$

Dual decision rule                    Primal decision rule

# Where is the $\phi$ "feature" space?

Dual optimization:

$$\min_{\propto} \frac{1}{2} \propto^T \begin{bmatrix} y_1 y_1 \phi(x_1)^T \phi(x_1) & y_1 y_2 \phi(x_1)^T \phi(x_2) & \dots & y_1 y_N \phi(x_1)^T \phi(x_N) \\ y_2 y_1 \phi(x_2)^T \phi(x_1) & y_2 y_2 \phi(x_2)^T \phi(x_2) & \dots & y_2 y_N \phi(x_2)^T \phi(x_N) \\ \dots & \dots & \dots & \dots \\ y_N y_1 \phi(x_N)^T \phi(x_1) & y_N y_2 \phi(x_N)^T \phi(x_2) & \dots & y_N y_N \phi(x_N)^T \phi(x_N) \end{bmatrix} \propto + (-1^T) \propto$$

subject to $\quad y^T \propto = 0 \qquad 0 \quad \leq \quad \propto \quad \leq \quad C$
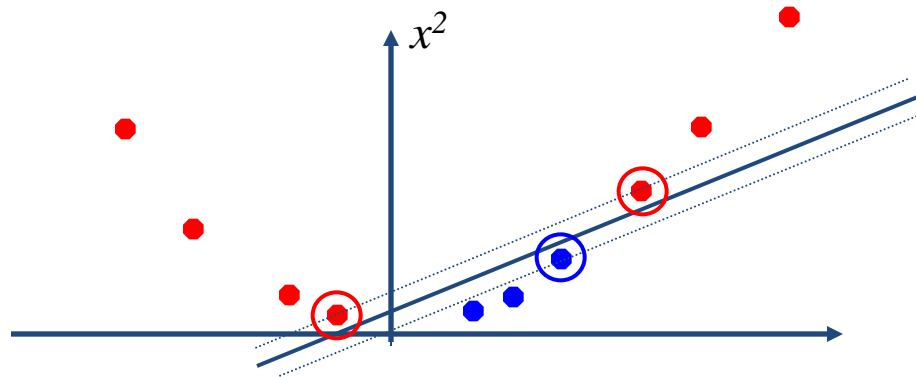
w?

$$w = \sum_{i \in SV} \propto_i y_i \phi(x_i) \qquad\qquad f(x) = \sum_i \propto_i y_i (\phi(x_i)^T \phi(x)) + b$$

b?

$$f(x) = \sum_i \propto_i y_i (\phi(x_i)^T \phi(x)) + b = y$$

# The "Kernel trick"

Linear SVM

$$f(x) = \sum_i \propto_i y_i (\mathrm{x}_i^\mathrm{T} \mathrm{x}) + b$$

Non-linear SVM

$$f(x) = \sum_i \propto_i y_i (\phi(\mathrm{x}_i)^T \phi(\mathrm{x})) + b$$

Define the "kernel function" K

$$K(x', x'') = \phi(\mathrm{x}')^T \phi(\mathrm{x}'')$$

then

$$f(x) = \sum_i \propto_i y_i K(\mathrm{x}_i,^T \mathrm{x}) + b$$

# Where is the $\phi$ "feature" space?

Dual optimization:

$$\min_{\propto} \frac{1}{2} \propto^T \begin{bmatrix} y_1 y_1 \phi(x_1)^T \phi(x_1) & y_1 y_2 \phi(x_1)^T \phi(x_2) & \ldots & y_1 y_N \phi(x_1)^T \phi(x_N) \\ y_2 y_1 \phi(x_2)^T \phi(x_1) & y_2 y_2 \phi(x_2)^T \phi(x_2) & \ldots & y_2 y_N \phi(x_2)^T \phi(x_N) \\ \ldots & \ldots & \ldots & \ldots \\ y_N y_1 \phi(x_N)^T \phi(x_1) & y_N y_2 \phi(x_N)^T \phi(x_2) & \ldots & y_N y_N \phi(x_N)^T \phi(x_N) \end{bmatrix} \propto + (-1^T) \propto$$

K

subject to $\quad y^T \propto = 0 \qquad 0 \quad \leq \quad \propto \quad \leq \quad C$

w?

$$w = \sum_{i \in SV} \propto_i y_i \phi(x_i) \qquad\qquad f(x) = \sum_i \propto_i y_i \left( \phi(x_i)^T \phi(x) \right) + b$$

K

b?

$$f(x) = \sum_i \propto_i y_i \left( \phi(x_i)^T \phi(x) \right) + b = y$$

K

# Nonlinear kernel: Example

- Consider the mapping $\varphi(x) = (x, x^2)$



$$\varphi(x) \cdot \varphi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2 y^2$$

$$K(x, y) = xy + x^2 y^2$$

# Computing K(x,x') without explicitly computing $\phi(x)$

For example: 2nd order polynomial kernel in 2d

$$K(x, x') = (1 + x^t x')^2 = (1 + x_1 x_1' + x_2 x_2')^2 =$$
$$= 1 + x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_1' + 2x_2 x_2' + 2x_1 x_1' x_2 x_2'$$

$$\text{K(x, x')} = \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1 x_2 \end{bmatrix}^T \begin{bmatrix} 1 \\ x_1'^2 \\ x_2'^2 \\ \sqrt{2}x_1' \\ \sqrt{2}x_2' \\ \sqrt{2}x_1' x_2' \end{bmatrix} = \phi(x)^T \phi(x)$$

# Popular kernels

$$K(x', x'') = \phi(x')^T \phi(x'')$$

| Name | params | Kernel eqation $K(x',x'')$ | Non-linear mapping $\phi(x)$ |
|------|--------|----------------------------|------------------------------|
| Linear | | $(x')^t x''$ | x |
| Polinomial | D | $(1+(x')^t x'')^D$ | All polynomials up to degree D in the elements of the vector x |
| Gaussian==RBF | $\sigma$ | $\exp(-||x'-x''||^2/(2\sigma^2))$ | Infinite dimensional vector |

# Popular kernels

$$K(x', x'') = \phi(\mathrm{x}')^T \phi(\mathrm{x}'')$$

| Name | params | Kernel eqation K(x',x'') | Non-linear mapping $\phi(x)$ |
|---|---|---|---|
| Linear | | $(x')^t x''$ | x |
| Polinomial | D | $(1+(x')^t x'')^D$ | All polynomials up to degree D in the elements of the vector x |
| Gaussian==RBF | $\sigma$ | exp(-\|\|x'-x''\|\|$^2$/$(2\sigma^2)$) | Infinite dimensional vector |

$$K(x,y) = \left(\sum_{i=1}^{n} x_i y_i + 1\right)^2 = \sum_{i=1}^{n} x_i^2 y_i^2 + \sum_{i=2}^{n}\sum_{j=1}^{i-1} \sqrt{2}x_i y_i \sqrt{2}x_j y_j + \sum_{i=1}^{n} \sqrt{2}x_i \sqrt{2}y_i + 1$$

$$\varphi(x) = \langle x_n^2, \ldots, x_1^2, \sqrt{2}x_n x_{n-1}, \ldots, \sqrt{2}x_n x_1, \sqrt{2}x_{n-1}x_{n-2}, \ldots, \sqrt{2}x_{n-1}x_1, \ldots, \sqrt{2}x_2 x_1, \sqrt{2}x_n, \ldots, \sqrt{1}x_1, 1\rangle$$

Complexity does not depend on D! (take log multiply and exponent)

# Popular kernels

$$K(x', x'') = \phi(x')^T \phi(x'')$$

| Name | params | Kernel eqation K(x',x'') | Non-linear mapping $\phi(x)$ |
|------|--------|--------------------------|------------------------------|
| Linear | | $(x')^t x''$ | x |
| Polinomial | D | $(1+(x')^t x'')^D$ | All polynomials up to degree D in the elements of the vector x |
| Gaussian==RBF | $\sigma$ | $\exp(-||x'-x''||^2/(2\sigma^2))$ | Infinite dimensional vector |

$$\overline{\overline{K}}_{ij} = K(x_i, x_j)$$

$$\overline{\overline{K}} = [\phi(x_1)\ \phi(x_2)\ \dots \phi(x_N)]^t\ [\phi(x_1)\ \phi(x_2)\ \dots \phi(x_N)]$$

$$rank(\overline{\overline{K}}) = \text{rank}([\phi(x_1)\ \phi(x_2)\ \dots \phi(x_N)])$$

$$take\ \overline{\overline{K}} = \begin{bmatrix} 1 & \epsilon & \epsilon \\ \epsilon & 1 & \epsilon \\ \epsilon & \epsilon & \ddots \end{bmatrix}$$

# Proper kernels

- Symmetric

  $K(x_i, x_j) = K(x_j, x_i)$

- Positive definite kernel

  $\forall N, \forall x_1, \ldots, x_N, \forall c \in R^N, \qquad c^t \overline{\overline{K}} c > 0$

- But in practice, we don't necessarily need PSD kernels..

# Constructing proper kernels

- $K_3(x',x'') = K_1(x',x'') + K_2(x',x'')$

$$\phi_3(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \end{bmatrix}$$

- $K_3(x',x'') = K_1(x',x'') * K_2(x',x'')$

$(\phi_1(x')^t \phi_1(x''))(\phi_2(x'')^t \phi_2(x')) = \text{tr}(\phi_1(x')^t \phi_1(x'')\phi_2(x'')^t \phi_2(x'))=$

$\text{tr}((\phi_2(x')\phi_1(x')^t)( \phi_1(x'')\phi_2(x'')^t)) = <\text{VEC}(\phi_2(x')\phi_1(x')^t), \text{VEC}(\phi_2(x'')\phi_1(x'')^t)>$

# Distances and kernels

- Suppose we want to apply kNN in kernel space

  In input space:

  $$\left\|a - b\right\|^2 = (a - b)^t(a - b) = a^t a - 2a^t b + b^t b$$

  Similarly, in "feature space"
  $$\left\|\phi(a) - \phi(b)\right\|^2 = (\phi(a) - \phi(b))^t(\phi(a) - \phi(b)) =$$
  $$= \phi(a)^t\phi(a) - 2\phi(a)^t\phi(b) + \phi(b)^t\phi(b) =$$
  $$= K(a,a) - 2K(a,b) + K(b,b)$$

Generalized Gaussian kernel for histograms:

$$K(h_1, h_2) = \exp\left(-\frac{1}{A}D(h_1, h_2)^2\right)$$

- L1 distance: $D(h_1, h_2) = \sum_{i=1}^{N} |h_1(i) - h_2(i)|$

- L2 distance: $D^2(h_1, h_2) = \sum_{i=1}^{N} \left(h_1(i) - h_2(i)\right)^2$

- L-inf distance: $D(h_1, h_2) = \max_{1 \leq i \leq N} |h_1(i) - h_2(i)|$

- $\chi^2$ distance: $D(h_1, h_2) = \sum_{i=1}^{N} \frac{\left(h_1(i) - h_2(i)\right)^2}{h_1(i) + h_2(i)}$

- Hellinger distance: $D^2(h_1, h_2) = \sum_{i=1}^{N} \left(\sqrt{h_1(i)} - \sqrt{h_2(i)}\right)^2$

- Mahalanobis distance: $D^2(h_1, h_2) = \left(h_1 - h_2\right)^T S^{-1} \left(h_1 - h_2\right)$

# The Intersection Kernel

Histogram Intersection kernel between histograms *a, b*

$$K(a, b) = \sum_{i=1}^{n} min(a_i, b_i)$$

$$a_i \geq 0$$
$$b_i \geq 0$$

*K* small -> *a, b* are different
*K* large -> *a, b* are similar

Intro. by Swain and Ballard 1991 to compare color histograms.
Odone et al 2005 proved positive definiteness.

# Demonstration of Positive Definiteness

Histogram Intersection kernel  between histograms *a, b*

$$K(a,b) = \sum_{i=1}^{n} min(a_i, b_i) \qquad \begin{array}{c} a_i \geq 0 \\ b_i \geq 0 \end{array}$$

To see that $min(a_i, b_i)$ is positive definite,

represent *a, b* in "Unary", *n* is written as *n* ones in a row:

$$min(a_i, b_i) = \left\langle a_{i\,\text{unary}}, b_{i\,\text{unary}} \right\rangle$$

$$min(3, 5) = \left\langle (1, 1, 1, 0, 0), (1, 1, 1, 1, 1) \right\rangle = 3$$

# The Trick

log( #support vectors ) x #dimensions

Decision function is $\text{sign}\left(h(x)\right)$ where:

$$h(x) = \sum_{j=1}^{\#\text{sv}} \alpha^j \left( \sum_{i=1}^{\#\text{dim}} \min(x_i, x_i^j) \right) + b$$

$$= \sum_{i=1}^{\#\text{dim}} \left( \sum_{j=1}^{\#\text{sv}} \alpha^j \min(x_i, x_i^j) \right) + b$$

$$= \sum_{i=1}^{\#\text{dim}} h_i(x_i)$$

$$h_i(x_i) = \sum_{j=1}^{\#\text{sv}} \alpha^j \min(x_i, x_i^j) + b$$

$$= \sum_{x_i^j < x_i} \alpha^j x_i^j + \left( \sum_{x_i^j \geq x_i} \alpha^j \right) x_i$$

Just sort the support vector values in each coordinate, and pre-compute

To evaluate, find position of $x_i$ in the sorted support vector values $x_i^j$ (cost: log #sv) look up values, multiply & add

Slide credit: Jitendra Malik

# Singular Value Decomposition (SVD)

- Handy mathematical technique that has application to many problems

- Given any $m \times n$ matrix **A**, it decomposes it to three matrices **U**, **V**, and **W** such that

$$\mathbf{A} = \mathbf{U}\,\mathbf{W}\,\mathbf{V}^{\mathsf{T}}$$

**U** is $m \times n$ and orthonormal

**W** is $n \times n$ and diagonal

**V** is $n \times n$ and orthonormal

# SVD

Matlab: [U,W,V]=svd(A,0)

$$
\mathbf{A} = \mathbf{U} \begin{pmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{pmatrix} \mathbf{V}^{\mathrm{T}}
$$

- The $w_i>0$ are called the singular values of **A** and are sorted
- If **A** is singular, some of the $w_i$ will be 0
- *rank*(**A**) = number of nonzero $w_i$
- SVD is unique (unless some $w_i$ are equal)

# SVD and Inverses

- **$A^{-1} = (V^T)^{-1} W^{-1} U^{-1} = V W^{-1} U^T$**
  - Using fact that inverse = transpose for orthogonal matrices
  - Note: $W^{-1}$ is also diagonal with elements one over those of **W**

- Pseudoinverse: if $w_i = 0$, set $1/w_i$ to 0 (!)
  - Defined for all (even non-square, singular, etc.) matrices
  - Equal to $(A^T A)^{-1} A^T$ if $A^T A$ invertible

- Solving **Ax**=**b** by least squares
  **x**=pinv(**A**)*b

# SVD and Eigenvectors

- Let $\mathbf{A}=\mathbf{U}\mathbf{W}\mathbf{V}^\mathrm{T}$, and let $x_i$ be $i^\mathrm{th}$ column of $\mathbf{V}$
- Consider $\mathbf{A}^\mathrm{T}\mathbf{A}x_i$:

$$\mathbf{A}^\mathrm{T}\mathbf{A}x_i = \mathbf{V}\mathbf{W}^\mathrm{T}\mathbf{U}^\mathrm{T}\mathbf{U}\mathbf{W}\mathbf{V}^\mathrm{T}x_i = \mathbf{V}\mathbf{W}^2\mathbf{V}^\mathrm{T}x_i = \mathbf{V}\mathbf{W}^2 \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{V} \begin{pmatrix} 0 \\ \vdots \\ w_i^2 \\ \vdots \\ 0 \end{pmatrix} = w_i^2 x_i$$

- So elements of $\mathbf{W}$ are sqrt(eigenvalues) and columns of $\mathbf{V}$ are eigenvectors of $\mathbf{A}^\mathrm{T}\mathbf{A}$
- Similarly, the columns of U are the eigenvectors of $\mathbf{A}\mathbf{A}^\mathrm{T}$, and $\mathbf{diag(W)}$ are sqrt(eigenvalues $\mathbf{A}\mathbf{A}^\mathrm{T}$)

# Kernel SVD

Let A=$[\phi(x_1), \phi(x_2), \ldots, \phi(x_N)]$

How do we compute the SVD decomposition U,W, V$^T$?

Mental framework – A is of size $\infty \times N$

AA$^t$ is of size $\infty \times \infty$ and U is also $\infty \times N$

but A$^t$A is of size $N \times N$ and we can compute V and W

$A = UWV^t \rightarrow U = AVW^{-1}$

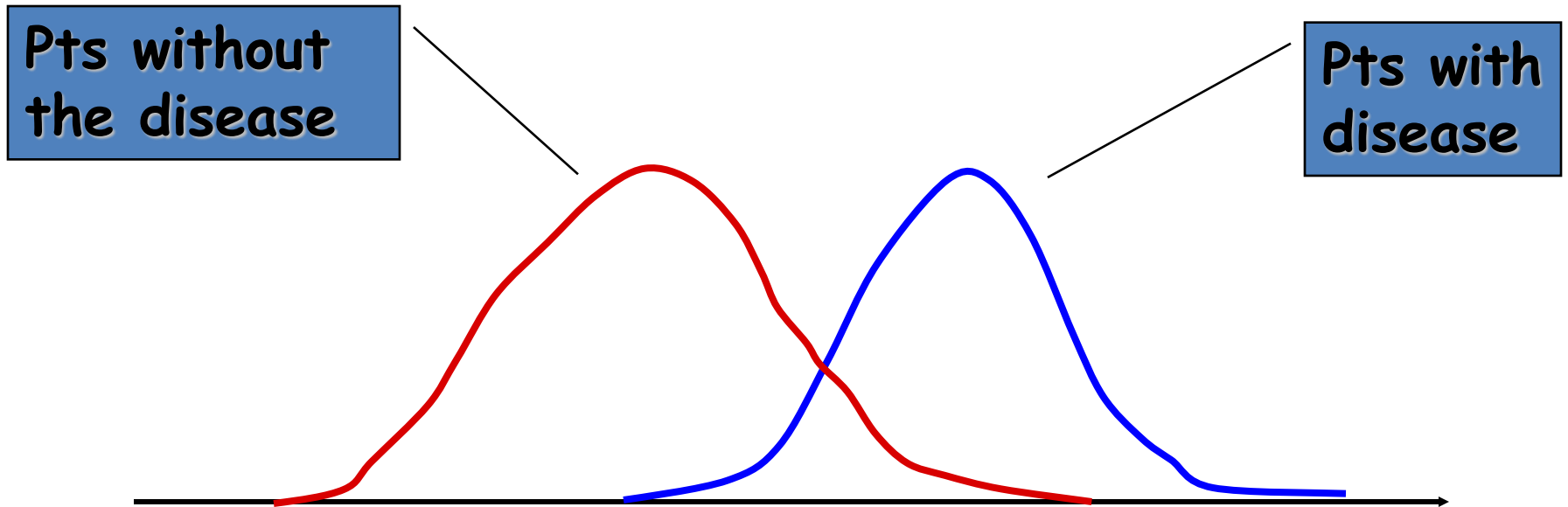and we can compute, e.g., $U^t \phi(x) = W^{-1} V^t A^t \phi(x)$

$= W^{-1} V^t \begin{bmatrix} k(x_1, x) \\ k(x_2, x) \\ \vdots \end{bmatrix}$

Applications: (1) **A**$^{-1}$ = **V W**$^{-1}$ **U**$^T$
(2) kernel PCA

# True label vs. classifier result

|  | Prediction=-1 | Prediction=1 |
|---|---|---|
| classifier / Real label | | |
| No disease (D = -1) | ☺ True negative | ✗ False positive |
| Disease (D = +1) | ✗ Miss | ☺ True positive |

Original slide credit: Darlene Goldstein

# Specific Example



Pts without the disease

Pts with disease

Test Result

$$f(\mathrm{x}) = \mathrm{w}^{\mathrm{t}}\mathrm{x} + b$$

# Threshold

Call these patients "negative"  |  Call these patients "positive"



Test Result

# Some definitions ...



Call these patients "negative"

Call these patients "positive"

**True Positives**

Test Result

**without the disease**
**with the disease**

Call these patients "negative"     Call these patients "positive"

Test Result

False Positives

**without the disease**
**with the disease**

Original slide credit: Darlene Goldstein

Call these patients "negative" | Call these patients "positive"

True negatives

Test Result

**without the disease**
**with the disease**

Original slide credit: Darlene Goldstein

# Moving the Threshold: right



"−"

"+"

Test Result

**without the disease**
**with the disease**

Original slide credit: Darlene Goldstein

# Moving the Threshold: left



"−"

"+"

Test Result

without the disease
with the disease

Original slide credit: Darlene Goldstein

# ROC curve



True Positive Rate

100%

0%

0%    False Positive Rate    100%

# ROC curve comparison
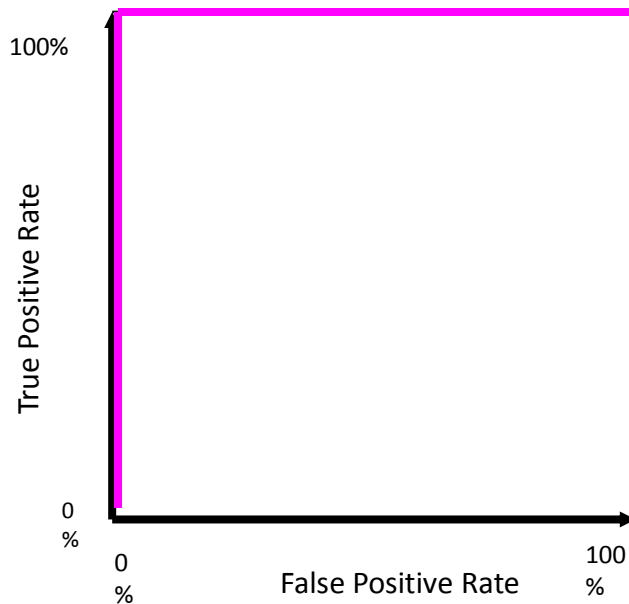
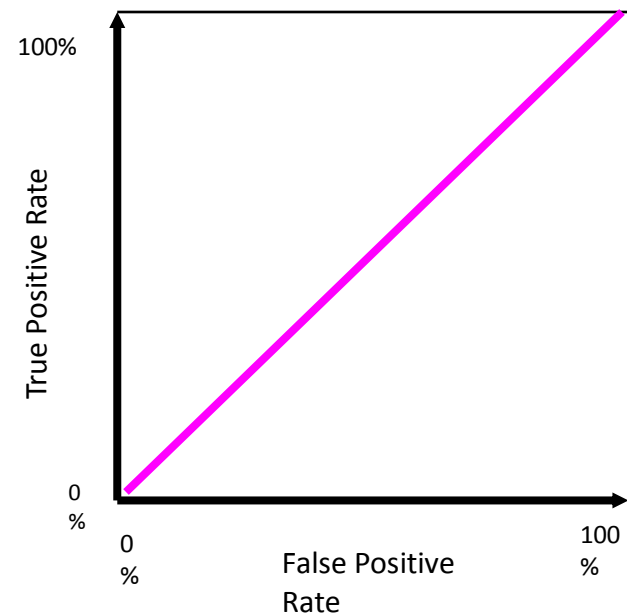## A good classifier:



## A poor classifier:

# ROC curve extremes

Best Classifier:



The distributions don't overlap at all

Worst Classifier:



The distributions overlap completely