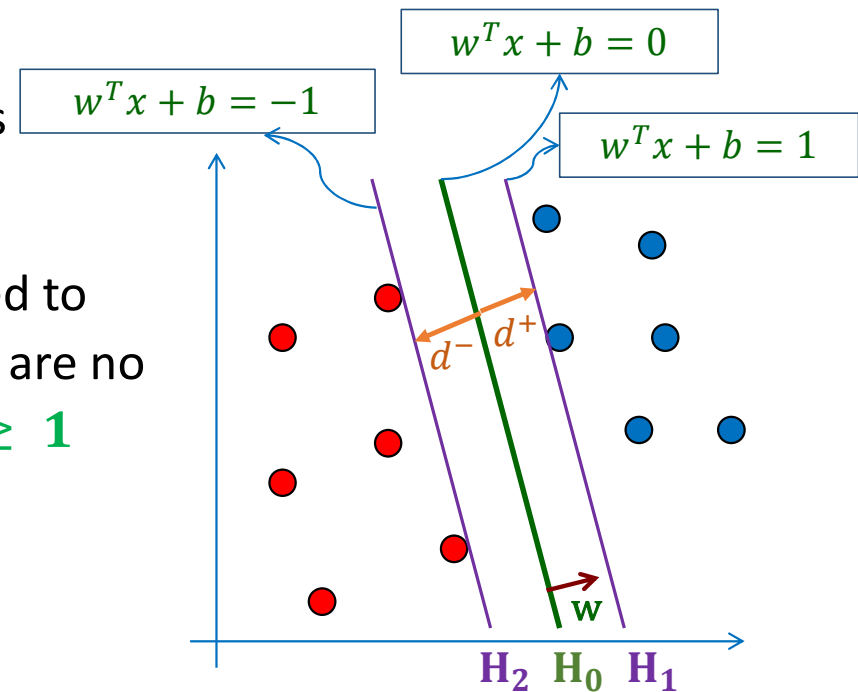# Non Linear SVM

Kernel Trick

# Maximizing the Margin

- We want a classifier (linear separator) with as big a margin as possible.

- In order to maximize the margin, we thus need to minimize $||w||$. With the condition that there are no datapoints between $H_1$ and $H_2$ : $y_i\,(w^T x_i) \geq 1$

- Minimize $J(w) = \frac{1}{2} w^T w$, subject to:
$$\forall_i, \; y_i(w^T x_i + b) \geq 1$$

- Constrained quadratic optimization problem solved by the **Lagrangian multipler method**



$w^T x + b = 0$
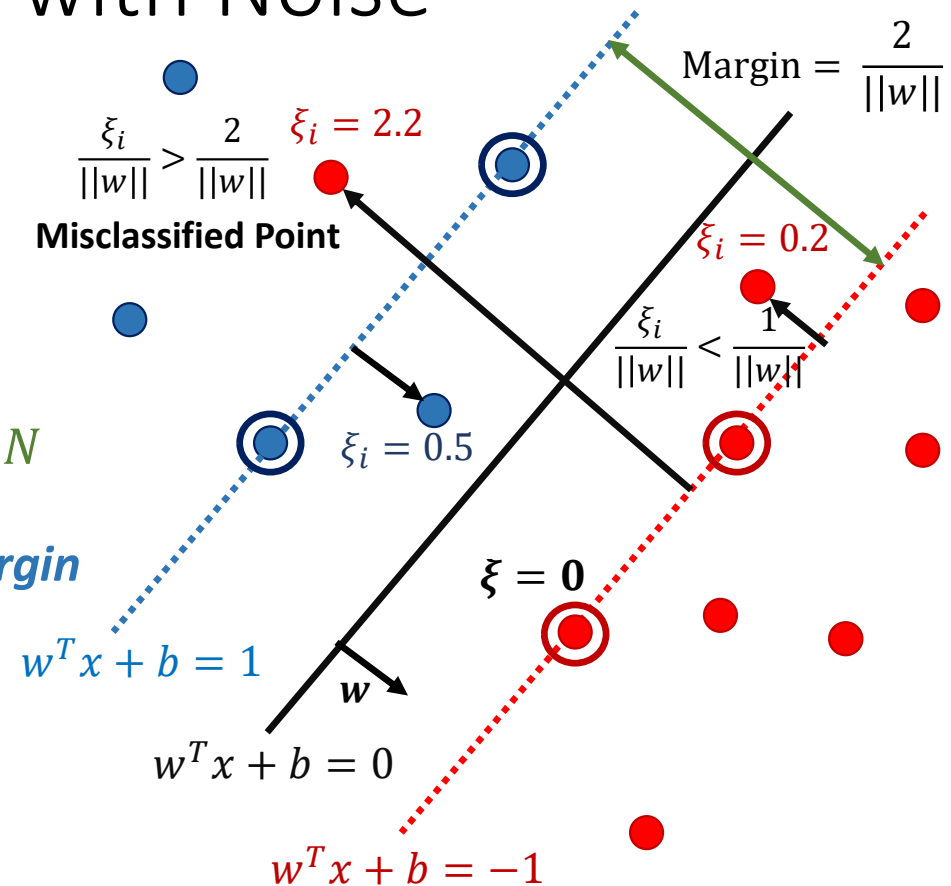
$w^T x + b = -1$

$w^T x + b = 1$

$d^-$ $d^+$

$w$

$H_2$ $H_0$ $H_1$

# Learning Maximum Margin with Noise

- The error terms $\xi_N^{'s}$ are incorporated into our optimization problem by:

$$\min_{w,b} ||w||^2 + C \sum_{i=1}^{N} \xi_i$$

such that $y_i(w^T x_i + b) \geq 1 - \xi_i, \ i = 1, \ldots, N$

- The solution to this problem is called *soft margin support vector classification*

$\text{Margin} = \dfrac{2}{||w||}$

$\dfrac{\xi_i}{||w||} > \dfrac{2}{||w||}$

$\xi_i = 2.2$

**Misclassified Point**

$\xi_i = 0.2$

$\dfrac{\xi_i}{||w||} < \dfrac{1}{||w||}$

$\xi_i = 0.5$

$\xi = 0$

$w^T x + b = 1$

$w$

$w^T x + b = 0$

$w^T x + b = -1$

# SVM Solution

Involves computing the ***inner products*** $x_i^T x_j$ between all training points

1. Maximize: $Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \boxed{x_i^T x_j}$, subject to $0 \leq \alpha_i \leq C \ \forall i$, $\sum_{i=1}^{N} \alpha_i y_i = 0$

- $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i;$    (very few $\alpha_i$ s are non-zero: support vectors. The sum is therefore only to be over the support vectors)

- $\mathbf{b} = y_K(1 - \xi_K) - \mathbf{w}_K^T x_K = y_K(1 - \xi_K) - \sum_{i=1}^{N} \alpha_i y_i \boxed{x_i^T x_K}$ with $K = \arg \max_i \alpha_i$

*Note:* Classification:

$$\mathbf{f}(\mathbf{x}_t) = \mathbf{w}^T\mathbf{x}_t + \mathbf{b} = \sum_{i=1}^{N} \alpha_i y_i \boxed{x_i^T x_t} + b$$

Relies on an ***inner product*** between the test point $x_t$ and the support vectors $x_i$

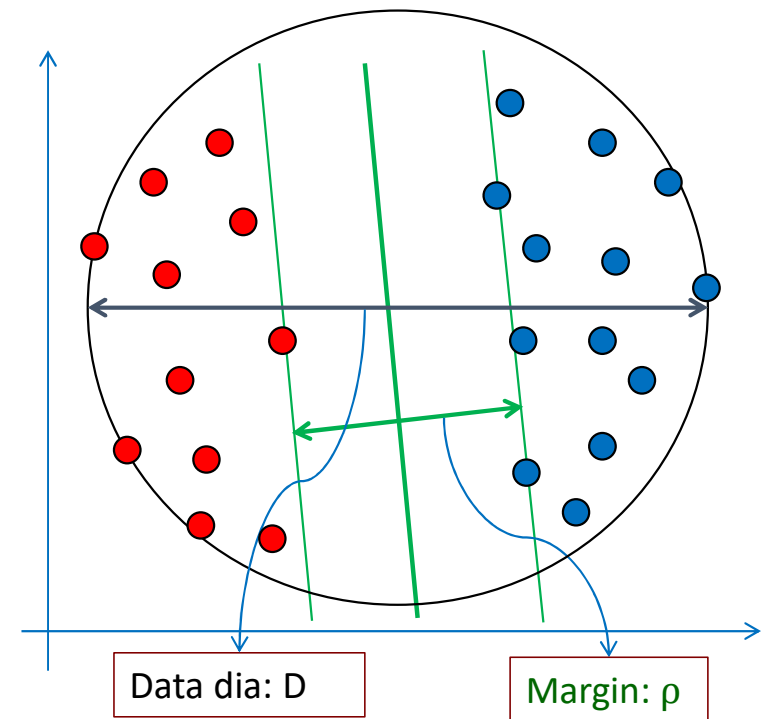# Theoretical Justification for Maximum Margin

- V.N. Vapnik and A.Ya. Červonenkis quantified complexity

  - Higher the VC-dimension ($h$), more complex the classifier

  - Bound on Expected Loss:

$$R_{tst}(\alpha) = R_{trn}(\alpha) + f(h, N)$$

$$h \leq min\left\{d, \left\lceil\frac{1}{m^2}\right\rceil\right\} + 1,$$

    where $d$ is the dimensionality. $m = \rho/D$ is the relative margin, with $\rho$ as the margin, and $D$ as the diameter of the smallest sphere that encloses all of the training examples.

  - Implication: If $\rho$ /D is high, VC dimension is low and expected error is low, regardless of the dimensionality $\boldsymbol{d}$

- Complexity of the classifier is kept small regardless of dimensionality.
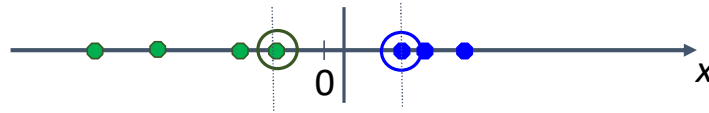


Data dia: D

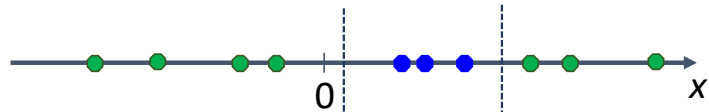Margin: ρ

# Linear SVMs: Overview

- Convex Optimization guarantees the global optimum.

- Support vectors are automatically identified.

- Use of max margin optimizes test accuracy

  - One of the best classifiers, given a feature representation

  - SVMs works well even with fewer training samples.

    - Note: Minimizes overfitting

  - Does not scale to huge datasets

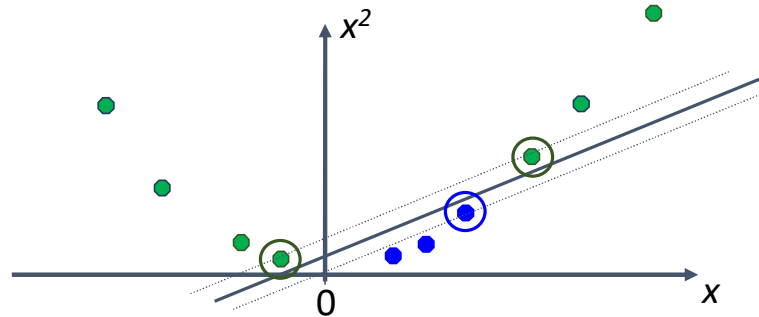    - GD based approaches work better in such cases.

# Non-linear SVMs

- Datasets that are linearly separable with some noise work out great with Linear SVMs:

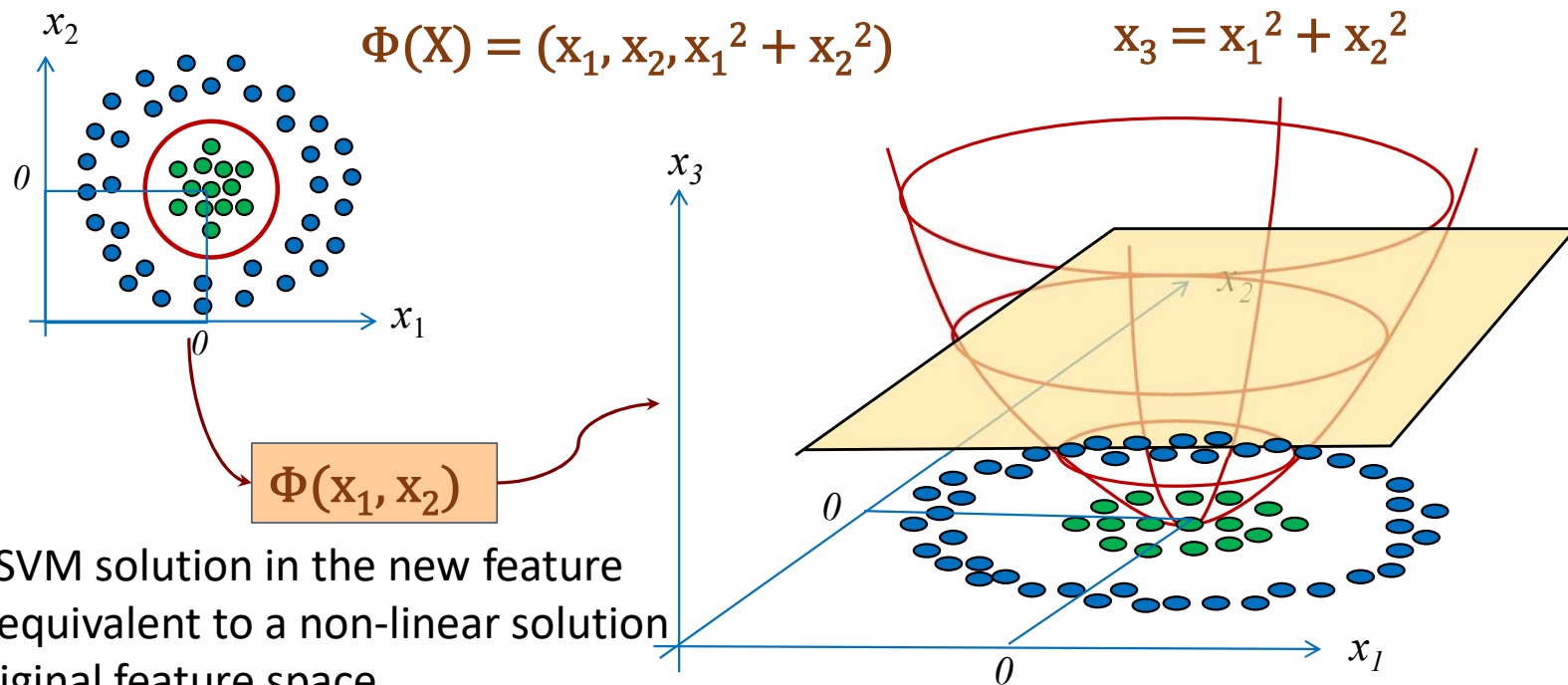- But what should be done if the dataset is just not linearly separable.

- Apply a non-linear transformation, to the feature space such that the samples become linearly separable

$$\Phi_k = (x_k, x_k^2)$$

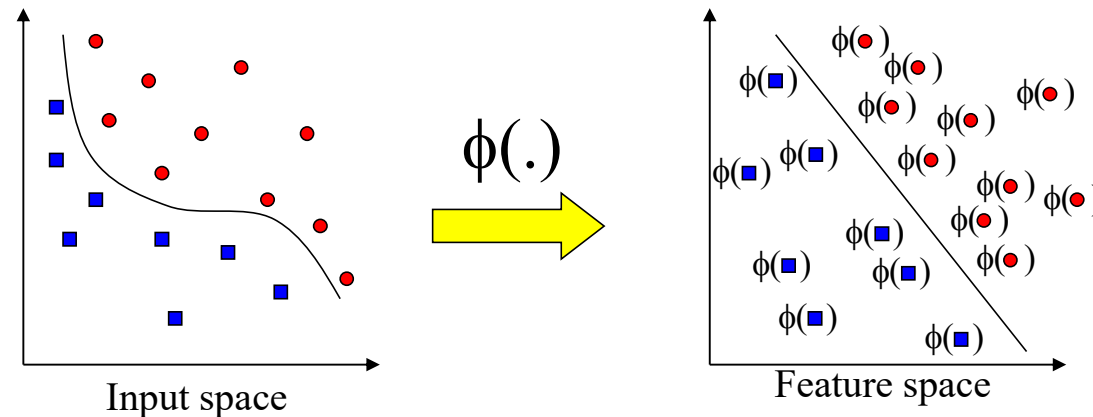# Circular Boundary

$$\Phi(X) = (x_1, x_2, x_1^2 + x_2^2)$$

$$x_3 = x_1^2 + x_2^2$$

$\Phi(x_1, x_2)$

- A linear SVM solution in the new feature space is equivalent to a non-linear solution in the original feature space.

- The mapping, $\Phi(X)$, is however unknown. Depends on the distribution of points in the feature space. It can be a complex mapping in general

# Transforming the Data



$$\phi(.)$$

Input space

Feature space

Note: feature space is of higher dimension
than the input space in practice

- Computation in the feature space can be costly because it is high dimensional
  - The feature space is typically infinite-dimensional!
- Question: Can we find the SVM solution without knowing $\Phi(\mathbf{X})$?
- The kernel trick comes to rescue

# Quadratic Basis Function

- Let there be a mapping from $\mathbf{x} \to \mathbf{\Phi(x)}$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Choose $\mathbf{\Phi(x)}$ to be pairwise monomial terms, $\mathbf{\Phi(x)} \in \mathbb{R}^{m^2}$, with $x \in \mathbb{R}^m$

$$\mathbf{x} \to \mathbf{\Phi(x)} \to \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \to \mathbf{\Phi(x)} \to \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ \dots \end{bmatrix}$$

$$\mathbf{\Phi(x)} \to \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \quad \mathbf{\Phi(z)} \to \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix} \quad \Phi^T(x)\Phi(z) = \begin{bmatrix} x_1 x_1 & x_1 x_2 & x_1 x_3 & \dots x_3 x_3 \end{bmatrix} \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix} = \sum_{i=1}^{3} \sum_{j=1}^{3} (x_i x_j)(z_i z_j)$$

# SVM Solution

1. Maximize: $Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \boxed{\mathbf{x}_i^{\mathrm{T}} \mathbf{x}_j}$,

$$\boldsymbol{\Phi}^T(x_i)\boldsymbol{\Phi}(x_j)$$

subject to $0 \leq \alpha_i \leq C \;\; \forall i, \quad \sum_{i=1}^{N} \alpha_i y_i = 0$

- $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i = \sum_{i=1}^{N} \alpha_i y_i \boldsymbol{\Phi}(x_i)$

*Note:* Classification:

$$\mathbf{f}(\mathbf{x}_t) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_t + \mathbf{b} = \sum_{i=1}^{N} \alpha_i y_i \boxed{\mathbf{x}_i^{\mathrm{T}} \mathbf{x}_t} + b$$

$$\boldsymbol{\Phi}^T(x_i)\boldsymbol{\Phi}(x_j)$$

- We need to do $\frac{N(N+1)}{2} \approx \frac{N^2}{2}$ dot products to calculate $\boldsymbol{\Phi}^T(x_i)\boldsymbol{\Phi}(x_j)$

- Each dot product further requires $m^2$ calculations

- The whole calculation will cost $\frac{N^2 \, m^2}{2}$

# Quadratic Basis Function

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \Phi(\mathbf{x}) \rightarrow \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \quad \Phi(\mathbf{z}) \rightarrow \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$$

Choose $\mathbf{\Phi}(\mathbf{x})$ to be pairwise monomial terms,
$\mathbf{\Phi}(\mathbf{x}) \in \mathbb{R}^{m^2}$, with $x \in \mathbb{R}^m$

$$\Phi^T(x)\Phi(z) = [x_1 x_1 \; x_1 x_2 \; x_1 x_3 \; \dots x_3 x_3] \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$$

$$= \sum_{i=1}^{3}\sum_{j=1}^{3}(x_i x_j)(z_i z_j)$$

$O(m^2)$

- Just out of casual, innocent, interest, let's look at another function of $x$ and $z$:

$$(x^T z)^2 = \left(\sum_{i=1}^{3} x_i z_i\right)\left(\sum_{j=1}^{3} x_j z_j\right) = \sum_{i=1}^{3}\sum_{j=1}^{3} x_i z_i x_j z_j = \sum_{i=1}^{3}\sum_{j=1}^{3}(x_i x_j)(z_i z_j)$$

Both are same, but $(x^T z)^2$ is only $O(m)$ to compute

# Quadratic Form

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

$$m = 3$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \mathbf{\Phi(x)} \rightarrow \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \sqrt{2}x_2x_3 \end{bmatrix} = O(10) = O\left(2m + \frac{m(m-1)}{2} + 1\right) \approx O\left(\frac{m^2}{2}\right)$$

$$\mathbf{\Phi(x).\,\Phi(z)} = 1 + 2\sum_{i=1}^{m} x_i z_i + \sum_{i=1}^{m} x_i^2 z_i^2 + \sum_{i=1}^{m}\sum_{j=i+1}^{m} 2x_i x_j z_i z_j$$

$$(x.z + 1)^2 = (x.z)^2 + 2\,x.z + 1 = \left(\sum_{i=1}^{m} x_i z_i\right)^2 + 2\sum_{i=1}^{m} x_i z_i + 1 = \sum_{i=1}^{m}\sum_{j=1}^{m} x_i z_i x_j z_j + 2\sum_{i=1}^{m} x_i z_i + 1$$

$$= 1 + 2\sum_{i=1}^{m} x_i z_i + \sum_{i=1}^{m} x_i^2 z_i^2 + \sum_{i=1}^{m}\sum_{j=i+1}^{m} 2x_i x_j z_i z_j \qquad O(m)$$

# For a Cubic Kernel

- If the original Space is 3-dimensional:

$$\mathbf{K}(\mathbf{X}, \mathbf{Z}) = (\mathbf{X} \cdot \mathbf{Z})^3 = (x_1 z_1 + x_2 z_2 + x_3 z_3)^3$$

- Equivalent to working in a 10-dimensional space
- Kernel:   5(3+2):$\times$ and 2:$+$
- $\Phi(X) = 13$: $\times$,   $\Phi(Z) = 13$: $\times$,
- $\Phi(X)\Phi(Z) = 10$: $\times$, 9: $+$
- Total 36(13+13+10):$\times$ and 9:$+$

$$\Phi(\mathbf{X}) = \Phi\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} x_1^3 \\ x_2^3 \\ x_3^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_1^2 x_3 \\ x_1 x_3^2 \\ x_2^2 x_3 \\ x_2 x_3^2 \\ x_1 x_2 x_3 \end{bmatrix}$$

# Higher Order Polynomials

**Kernel Function**

| Polynomial | $\Phi(x)$ | Cost to build $Q_{kl}$ matrix traditionally | Cost if 100 features | $\Phi(x). \Phi(z)$ $K(x, z)$ | Cost to build $Q_{kl}$ matrix sneakily | Cost if 100 features |
|---|---|---|---|---|---|---|
| Quadratic | All $m^2/2$ terms up to degree 2 | $m^2 N^2 /4$ | $2{,}500\ N^2$ | $(x.z+1)^2$ | $m N^2 / 2$ | $50\ N^2$ |
| Cubic | All $m^3/6$ terms up to degree 3 | $m^3 N^2 /12$ | $83{,}000\ N^2$ | $(x.z+1)^3$ | $m N^2 / 2$ | $50\ N^2$ |
| Quartic | All $m^4/24$ terms up to degree 4 | $m^4 N^2 /48$ | $1{,}960{,}000\ N^2$ | $(x.z+1)^4$ | $m N^2 / 2$ | $50\ N^2$ |

Slide inspired from, Andrew W. Moore

# SVM Solution: Kernel function

1. Maximize:

$$Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \, \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$$

2. $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i;$

3. $\mathbf{b} = \mathbf{1} - \mathbf{w}^{\mathrm{T}} \mathbf{x}_{s+} = \mathbf{1} - \sum_{i=1}^{N} \alpha_i y_i \mathbf{K}(\mathbf{x}_i, \mathbf{x}_{s+})$
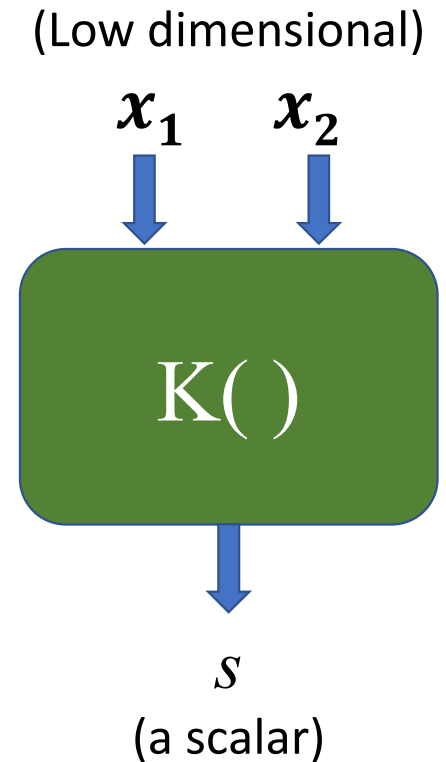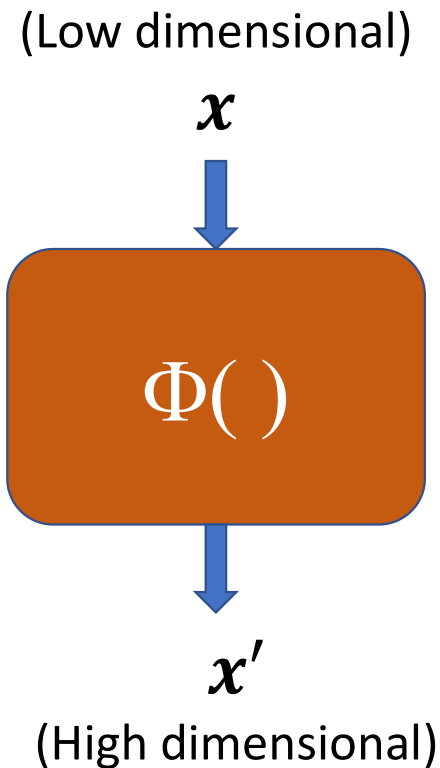
> Do we know the $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$.?
>
> Let us compare $\Phi$ and $\mathbf{K}()$.

*Note:* Classification:

$$\mathbf{g}(\mathbf{x}_t) = \mathbf{w}^{\mathrm{T}} \mathbf{x}_t + \mathbf{b} = \sum_{i=1}^{N} \alpha_i y_i \, \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) + b$$

# Comparing compare $\Phi()$ and $\mathbf{K}()$

- $\Phi(\mathbf{x})$ is a complex non-linear mapping of $\mathbf{x}$ into a high-dimensional space.

- $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j)$ is a simple function that measures the similarity between two vectors.

- If we know the $\mathbf{K}$ function, we do not need $\Phi$.

- A kernel function *implicitly* maps data to a high-dimensional space (without the need to compute each $\phi(\mathbf{x})$ explicitly).

(Low dimensional)

$x$

$\Phi(\ )$

$x'$

(High dimensional)

(Low dimensional)

$x_1 \quad x_2$

$\mathrm{K}(\ )$

$s$

(a scalar)

# What Functions are Kernels?

- For some functions $K(x_i, x_j)$ checking that $K(x_i, xj) = \Phi^T(x_i)\Phi(x_j)$ can be cumbersome.

- Mercer's theorem determines which functions can be used as a kernel function:

  ***Every semi-positive definite symmetric function is a kernel***

- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

**K =**

| | | | | |
|---|---|---|---|---|
| $K(x_1,x_1)$ | $K(x_1,x_2)$ | $K(x_1,x_3)$ | ... | $K(x_1,x_n)$ |
| $K(x_2,x_1)$ | $K(x_2,x_2)$ | $K(x_2,x_3)$ | | $K(x_2,x_n)$ |
| | | | | |
| ... | ... | ... | ... | ... |
| $K(x_n,x_1)$ | $K(x_n,x_2)$ | $K(x_n,x_3)$ | ... | $K(x_n,x_n)$ |

✓ $K(x_i, x_j)$ measures the similarity or proximity between two data points $x_i$ and $x_j$ in the input space

# Examples of Kernel Functions

- Linear: $K(x_i, x_j) = x_i^T x_j$

  - Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ is $x$ itself

- Polynomial of power $p$: $K(x_i, x_j) = (1 + x_i^T x_j)^p$

  - Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ has $^{d+p}_p C$ dimensions

- Gaussian (radial-basis function): $K(x_i, x_j) = e^{-\frac{\left(\|x_i - x_j\|\right)^2}{2\sigma^2}}$

  - Mapping $\Phi: x \rightarrow \Phi(x)$, where $\Phi(x)$ is *infinite-dimensional*: every point is mapped to *a function* (a Gaussian); combination of functions for support vectors is the separator.

# Example

- Suppose we have 5 one-dimensional data points
    - $x_1 = 1, x_2 = 2, x_3 = 4, x_4 = 5, x_5 = 6$, with 1, 2, 6 as class 1 and 4, 5 as class 2 $\Rightarrow y_1 = 1, y_2 = 1, y_3 = -1, y_4 = -1, y_5 = 1$

- We use the polynomial kernel of degree 2
    - $K(x, y) = (xy + 1)^2$
    - C is set to 100

- We first find $\alpha_i$ $(i = 1, \ldots, 5)$ by

$$\text{max.} \sum_{i=1}^{5} \alpha_i - \frac{1}{2} \sum_{i=1}^{5} \sum_{j=1}^{5} \alpha_i \alpha_j y_i y_j \, (x_i x_j + 1)^2$$

$$\text{subject to } 0 \leq \alpha_i \leq 100, \sum_{i=1}^{5} \alpha_i y_i = 0$$

# Example

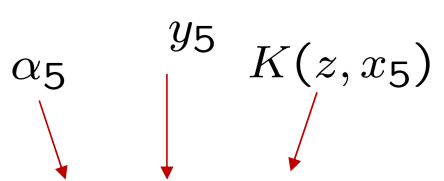- By using a QP solver, we get
  - $\alpha_1 = 0, \alpha_2 = 2.5, \alpha_3 = 0, \alpha_4 = 7.333, \alpha_5 = 4.833$
  - Note that the constraints are indeed satisfied
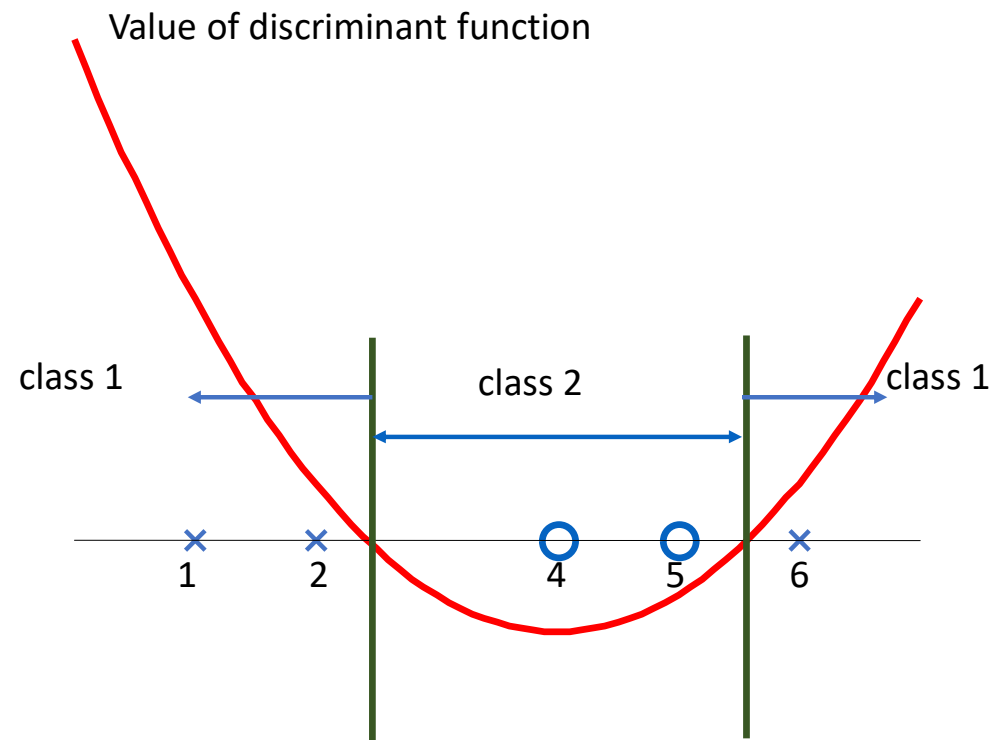  - The support vectors are $\{x_2 = 2, x_4 = 5, x_5 = 6\}$
- The discriminant function is

$$\alpha_5 \qquad y_5 \quad K(z, x_5)$$

$$\mathbf{f}(\mathbf{z}) = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{z} + b = \mathbf{2.5}\,(\mathbf{1})(\mathbf{2z}+\mathbf{1})^2 + \mathbf{7.333}\,(-\mathbf{1})(\mathbf{5z}+\mathbf{1})^2 + \mathbf{4.833}\,(\mathbf{1})(\mathbf{6z}+\mathbf{1})^2 + \mathbf{b}$$

$$= \mathbf{0.6667}\,\mathbf{z}^2 - \mathbf{5.333}\,\mathbf{z} + \mathbf{b}$$

- $b$ is recovered by solving f(2)=1 or by f(5)=-1 or by f(6)=1, as $x_2$ and $x_5$ lie on the line $\mathbf{w}^T\mathbf{x} + \mathbf{b} = \mathbf{1}$ and $x_4$ lies on the line $\mathbf{w}^T\mathbf{x} + \mathbf{b} = -\mathbf{1}$
- All three give b=9 $\Rightarrow \mathbf{f}(\mathbf{z}) = \mathbf{0.6667}\,\mathbf{z}^2 - \mathbf{5.333}\,\mathbf{z} + \mathbf{9}$
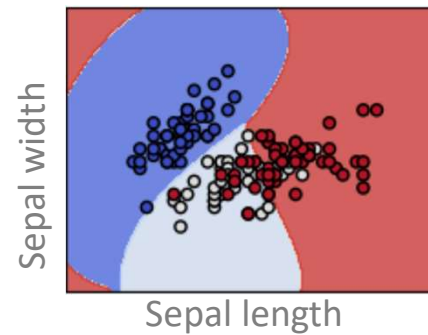
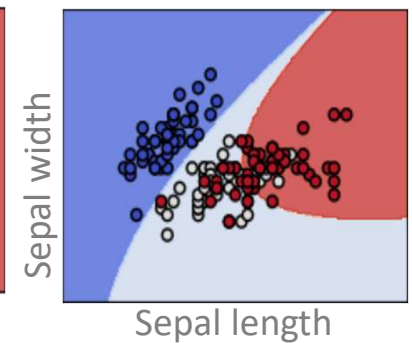# Example

# Scikit Learn Implementation

- **libsvm** based: sklearn.svm.SVC
- Several Kernels: linear, polynomial, rbf, sigmoid, custom.
- Other HyperParameters: C, kernel params
- Usage:

```
>>>  from sklearn import svm
>>>  X = [[0, 0], [1, 1]]
>>>  y = [0, 1]
>>>  clf = svm.SVC(kernel='rbf')
>>>  clf.fit(X, y)
>>>  clf.predict([[2., 2.]])
```

SVC: RBF Kernel          SVC: Poly (3) Kernel



```
>>>  # get support vectors
>>>  clf.support_vectors_
>>>  # get support vector indices
>>>  clf.support_
```

https://scikit-learn.org/stable/modules/svm.html

# Summary

- Linear SVMs generalize well, but cannot separate non-linear data
- If features can be transformed appropriately, SVMs can learn non-linear boundaries.
- How do we find the feature transformation?
  - Use some popular Kernel functions.
- Kernels (nonlinear) SVMs are also good at generalization and can deal with non-linear data.