

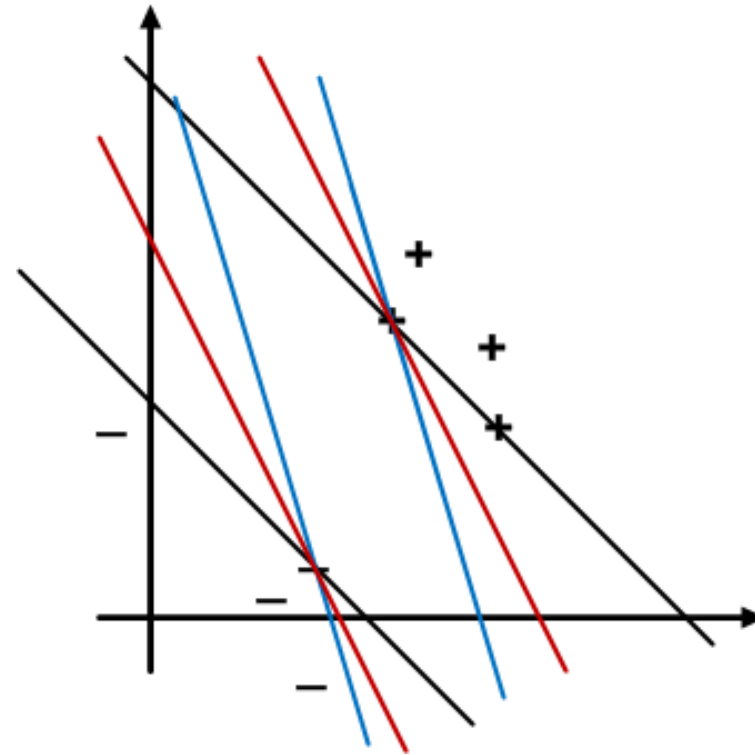
Support Vector Machines (SVMs)

- Vladimir Vapnik

$$h: X \rightarrow \{-1, +1\}$$

- Widest street approach
- Maximum margin classifier

- Which street is the best classifier: blue, red or black?



Support Vector Machines (cont.)

- Vladimir Vapnik

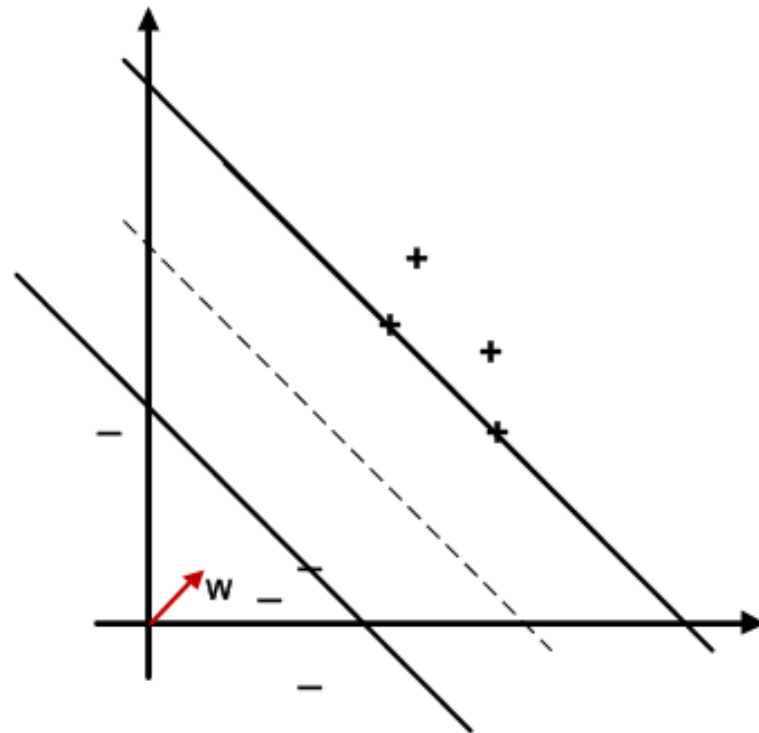
$$h: X \rightarrow \{-1, +1\}$$

- The separating hyperplane can be described as follows:

$$\mathbf{w} \cdot \mathbf{x} + b = 0^{**}$$

- Vector \mathbf{w} needs to be perpendicular to the street (why?)

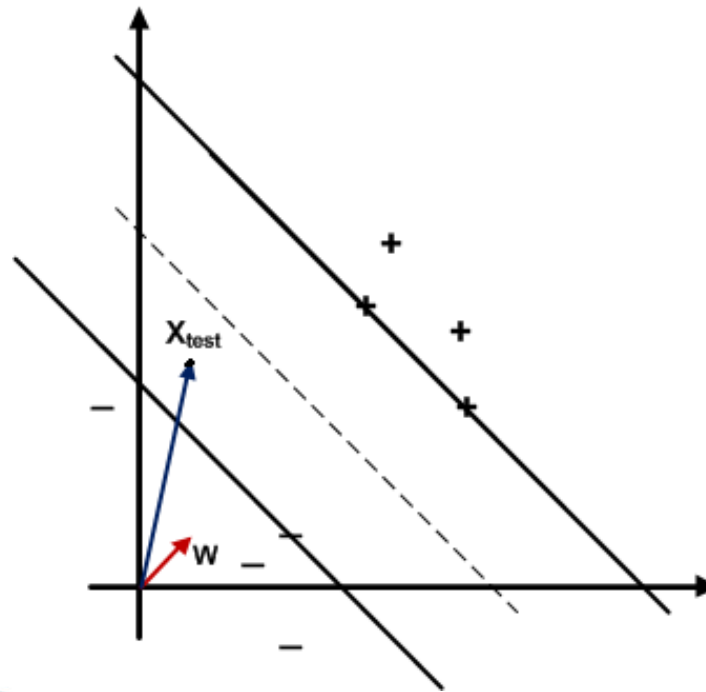
$$^{**} \mathbf{w}^T \mathbf{x} + b = 0$$



SVM Classifier

If $(\mathbf{w} \cdot \mathbf{x}_{test} \geq c)$ Then Class is +

If $(\mathbf{w} \cdot \mathbf{x}_{test} + b \geq 0)$ Then Class is +, s.t. $b = -c$

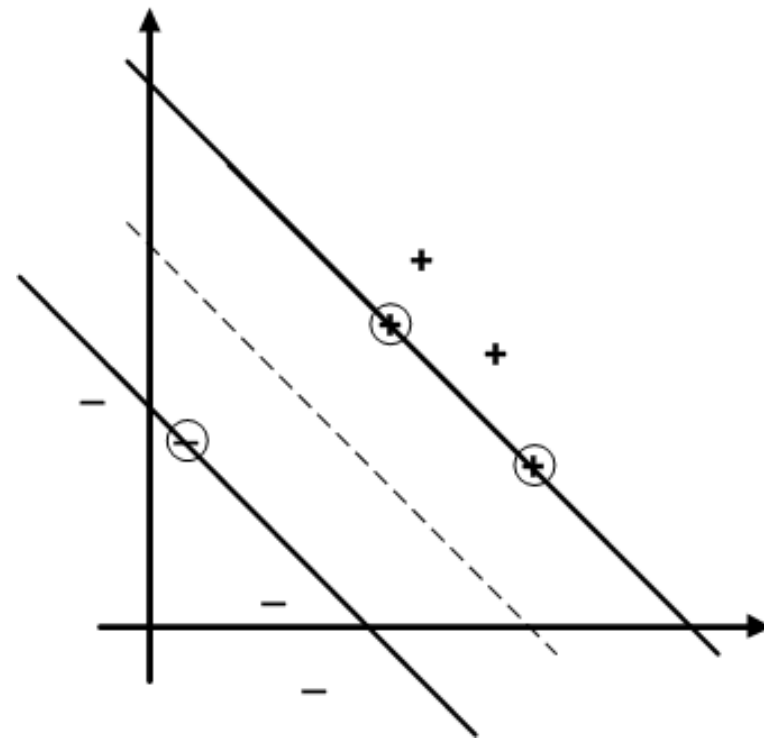


Constraints

□ Let assume the following constraints:

$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_+ + b \geq 1 & (y = +1) \\ \mathbf{w} \cdot \mathbf{x}_- + b \leq -1 & (y = -1) \end{cases}$$

$$\Rightarrow y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0, \forall i$$



Constraints (cont.)

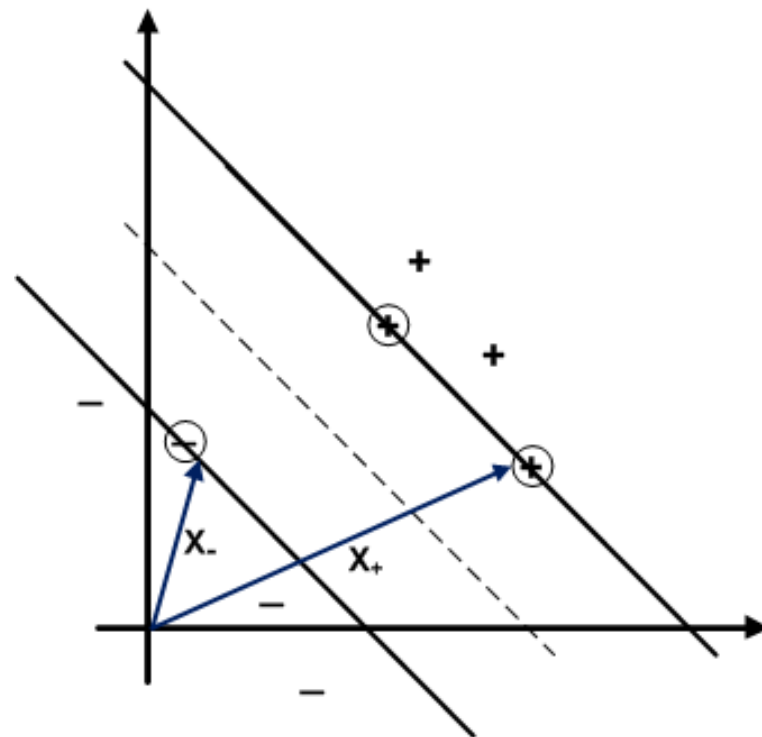
□ Let assume the following constraints:

$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_+ + b \geq 1 & (y = +1) \\ \mathbf{w} \cdot \mathbf{x}_- + b \leq -1 & (y = -1) \end{cases}$$

$$\Rightarrow y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0, \forall i$$

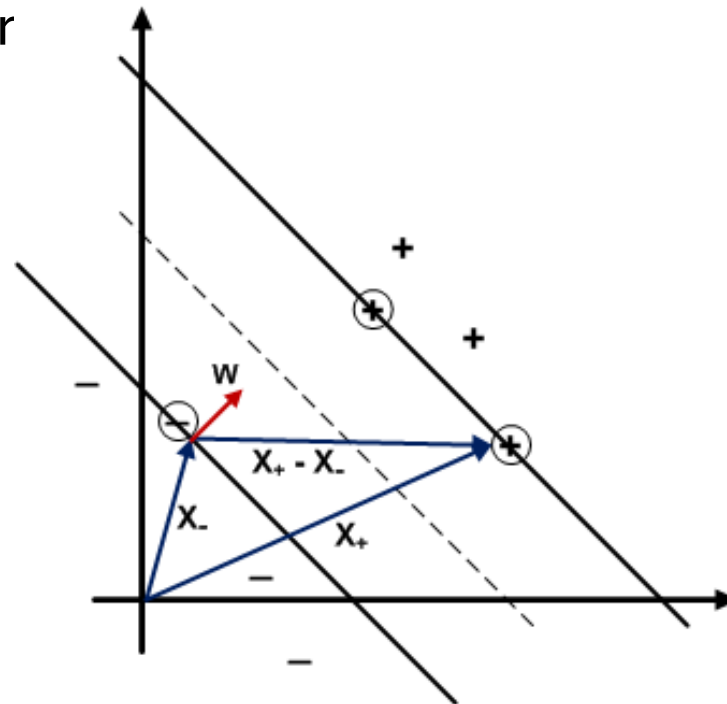
$$\boxed{y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 = 0}$$

IFF \mathbf{x}_i is a support vector



Width of the Street

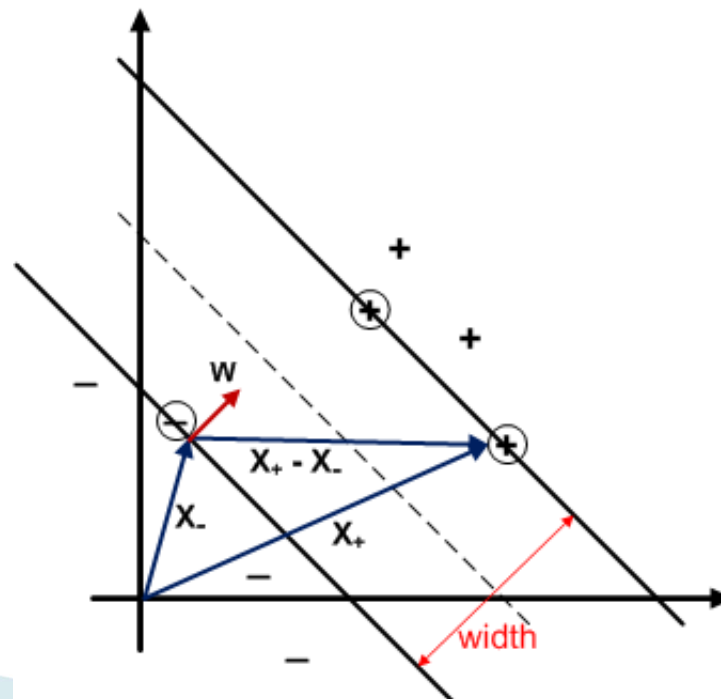
- Recall that vector w is perpendicular to the street.
- If we project vector $(x_+ - x_-)$ to vector w , the width of the street is obtained.
 - Vector w has to be a unit vector



Width of the Street (cont.)

$$Width = (x_+ - x_-) \cdot \frac{w}{||w||} = (w \cdot x_+ - w \cdot x_-) \cdot \frac{1}{||w||} = \frac{2}{||w||}$$

- Note that $w \cdot x_+ + b = 1$ and $w \cdot x_- + b = -1$



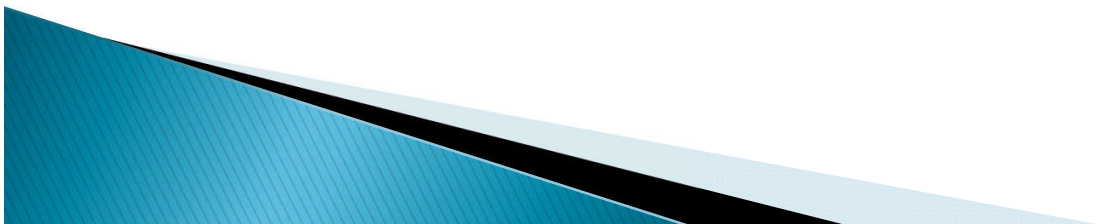
Constrained Optimization Problem

$$Width = (x_+ - x_-) \cdot \frac{w}{||w||} = (w \cdot x_+ - w \cdot x_-) \cdot \frac{1}{||w||} = \frac{2}{||w||}$$

- The goal is to maximize $\frac{2}{||w||}$, or to minimize $||w||$,
subject to: $y_i(w \cdot x_i + b) - 1 \geq 0$

↓

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



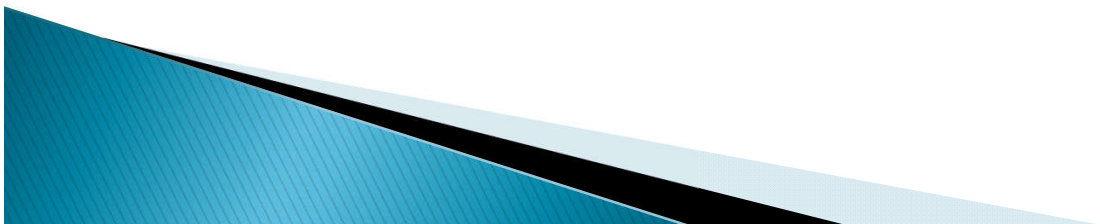
Lagrange Multipliers

- The goal is to minimize L , w.r.t. \mathbf{w}, b & maximize L , w.r.t. each α_i

$$L(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1]$$

subject to: $\alpha_i \geq 0, \forall i$

- This quadratic optimization problem is known as the dual problem.



Dual Problem

- The goal is to minimize L , w.r.t. \mathbf{w}, b and maximize L , w.r.t. each α_i

$$L(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1]$$

$$\begin{cases} \frac{\partial L(\mathbf{w}, b)}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 \rightarrow \boxed{\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i} \\ \frac{\partial L(\mathbf{w}, b)}{\partial b} = -\sum_i \alpha_i y_i = 0 \rightarrow \boxed{\sum_i \alpha_i y_i = 0} \end{cases}$$

- The Representer Theorem states that the solution \mathbf{w} can always be written as a linear combination of the training data.

Dual Problem (cont.)

- If we substitute \mathbf{w} and b into the Lagrange multipliers formula:

$$L(\boldsymbol{\alpha}) = \frac{1}{2} (\sum_i \alpha_i y_i \mathbf{x}_i) \cdot (\sum_j \alpha_j y_j \mathbf{x}_j) - (\sum_i \alpha_i y_i \mathbf{x}_i) \cdot (\sum_j \alpha_j y_j \mathbf{x}_j) - \sum_i \alpha_i y_i b + \sum_i \alpha_i$$

$$L(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_{i=1}^m \alpha_i, \forall i$$

subject to: $\alpha_i \geq 0, \forall i$

- This quadratic problem over α_i is known as the dual problem.
- It is shown that the problem space is convex, so it doesn't get stuck in local minimum/maximum.

Dual Problem (cont.)

- Quadratic optimization problem

$$L(\alpha) = -\frac{1}{2} \sum_{\substack{i \text{ is a SV} \\ j \text{ is a SV}}} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_{i \text{ is a SV}} \alpha_i$$

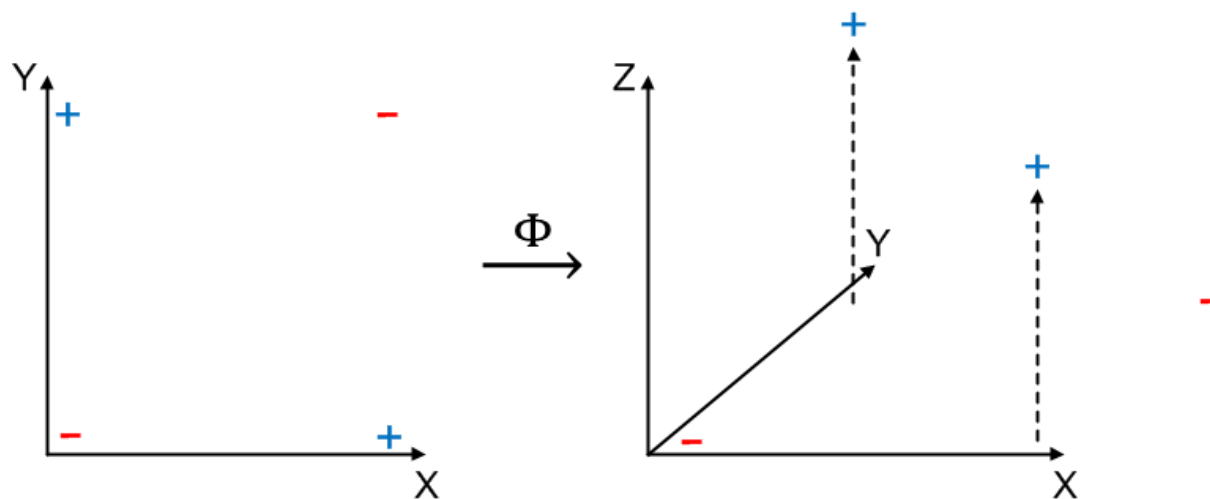
subject to: $\alpha_i > 0, \forall i$

- Classification Rule

$$\sum_{i \text{ is a SV}} \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_{test} + b \geq 0 \Rightarrow \text{Class is +}$$

Kernel Trick Intuition

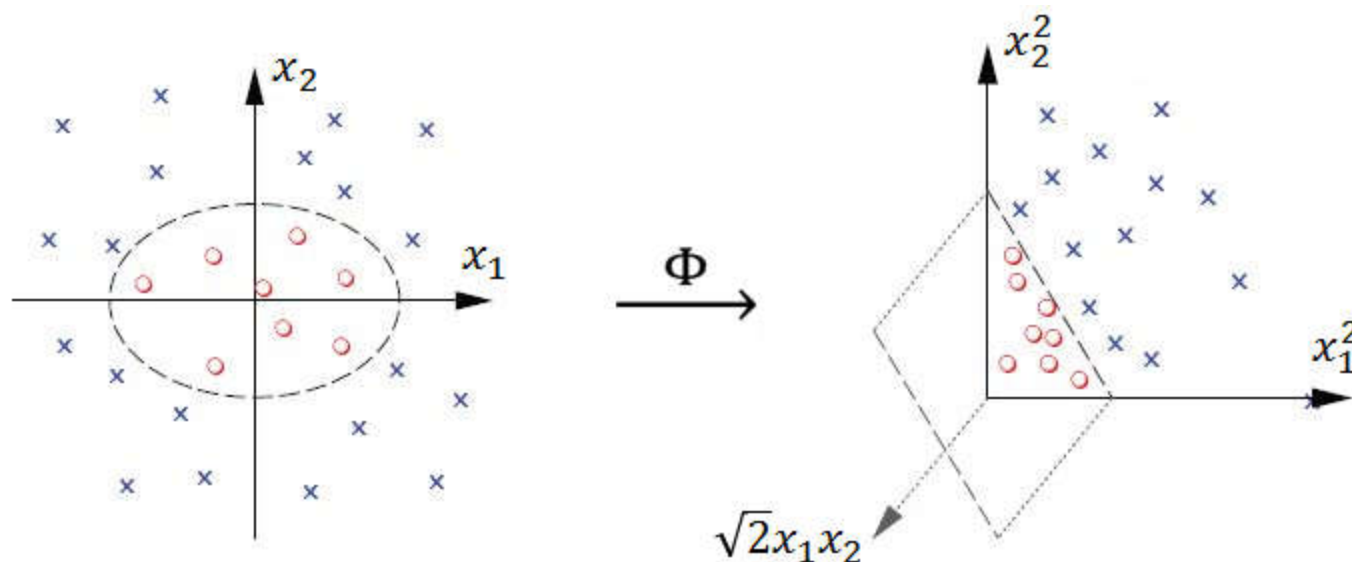
- ❑ SVM Solution for linearly inseparable problems, such as XOR
 - Kernel Trick: using a linear classifier to solve a non-linear problem.



Kernel Trick (cont.)

- Higher dimensional feature space – example:

$$(x_1, x_2) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$



Kernel Trick – Formal

$$\begin{cases} L(\alpha) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_i \alpha_i \\ h(\mathbf{x}_{test}) = \text{sgn}(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_{test} + b) \end{cases}$$

1

□ SVM Transformation:

$$L(\alpha) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) + \sum_i \alpha_i$$

2

Kernel Trick – Formal (cont.)

□ SVM Transformation:

$$L(\alpha) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) + \sum_i \alpha_i \quad \text{2}$$

□ Kernel Trick (to avoid expensive data transformations)

$$L(\alpha) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \boxed{K(\mathbf{x}_i, \mathbf{x}_j)} + \sum_i \alpha_i \quad \text{3}$$

Kernel Trick – Formal (cont.)

$$\Phi: X \rightarrow \mathbb{Z}$$

$$L(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_i \alpha_i$$

- If $K(\mathbf{x}_i, \mathbf{x}_j)$ is an inner product in some space, we are good!

Kernel Function Properties

- Function $K(\mathbf{x}, \mathbf{x}')$ is a valid kernel if:
 - It computes an inner product in some space \mathbb{Z} .
 - We just need to know that space \mathbb{Z} exists!
 - It is symmetric / commutative, i.e. $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$.
 - It should (preferably) be positive semi-definite, i.e. satisfy Mercer's theorem.



Kernel Types

□ The most frequently used kernel types:

- Linear: $(\mathbf{x}_i \cdot \mathbf{x}_j + c)$

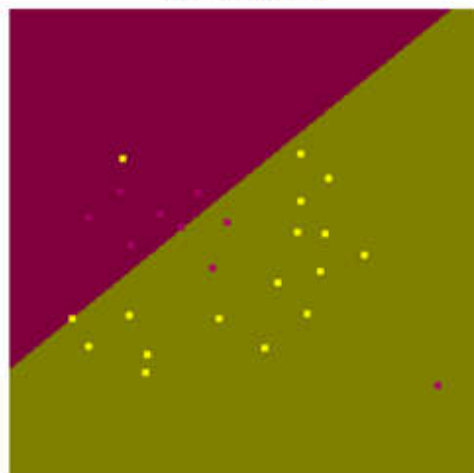
- If $c = 0$, it is homogenous.

- Polynomial: $(\alpha \mathbf{x}_i \cdot \mathbf{x}_j + c)^d$, subject to: $d > 1$

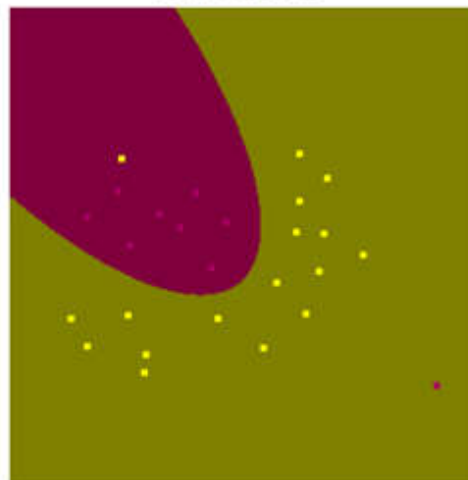
- Gaussian RBF: $\exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) = \exp\left(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$,
subject to: $\gamma = \frac{1}{2\sigma^2}$

Kernel Types

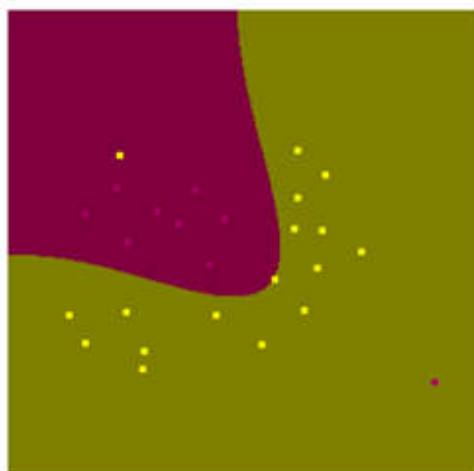
Linear kernel



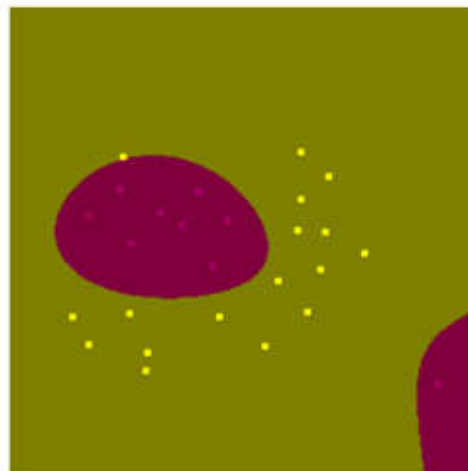
Poly degree=2



Poly degree=4



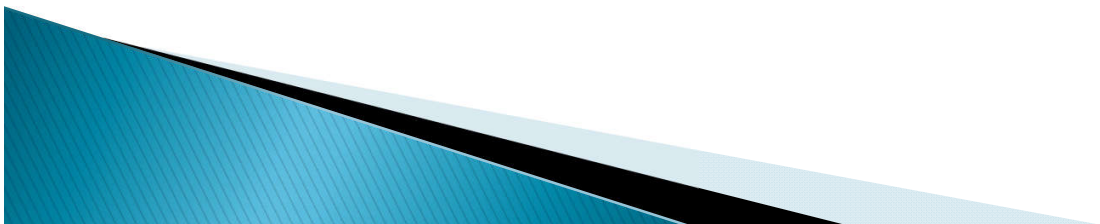
RBF kernel



Kernel Choice

- ❑ time of SVM learning: linear < poly < rbf
- ❑ ability to fit any data: linear < poly < rbf
- ❑ risk of overfitting: linear < poly < rbf
- ❑ risk of underfitting: rbf < poly < linear
- ❑ number of hyper-parameters: linear < rbf < poly

- ❑ So which one to choose? [1]
 - Ockham's razor



Other Kernels

□ List of other well-known kernel functions:

- Exponential Kernel
- Laplacian Kernel
- ANOVA Kernel
- Hyperbolic Tangent (Sigmoid) Kernel
- Rational Quadratic Kernel
- Multiquadric Kernel
- Inverse Multiquadric Kernel
- Circular Kernel
- Spherical Kernel
- Power Kernel
- Log Kernel
- Spline Kernel
- B-Spline Kernel
- Bessel Kernel
- Cauchy Kernel
- Chi-Square Kernel
- Histogram Intersection Kernel
- T-Student Kernel
- Bayesian Kernel
- Wavelet Kernel

References

1. Support–vector networks, Corinna Cortes & Vladimir Vapnik, Machine Learning 20, 273–297, 1995.
2. Lecture on Support Vector Machines, Patrick Winston, Massachusetts Institute of Technology, 2010.
3. Learning from data, Yaser Abu–Mostafa, Malik Magdon–Ismail, Hsuan–Tien Lin, 2012.

