

# Computational Data Analysis

## Machine Learning

**Yao Xie, Ph.D.**

*Associate Professor*

Harold R. and Mary Anne Nash Early Career Professor

H. Milton Stewart School of Industrial and Systems  
Engineering

Support Vector Machine (SVM)



# Main types approaches to design classifiers

- Bayes rule + assumption for  $p(x|y = 1)$ 
  - Assume  $p(x|y = 1)$  is Gaussian
  - Assume  $p(x|y = 1)$  is fully factorized
- Use geometric intuitions
  - k-nearest neighbor classifier
  - Support vector machine
- Directly go for the decision boundary  $h(x) = -\ln \frac{q_i(x)}{q_j(x)}$ 
  - Logistic regression
  - Neural networks

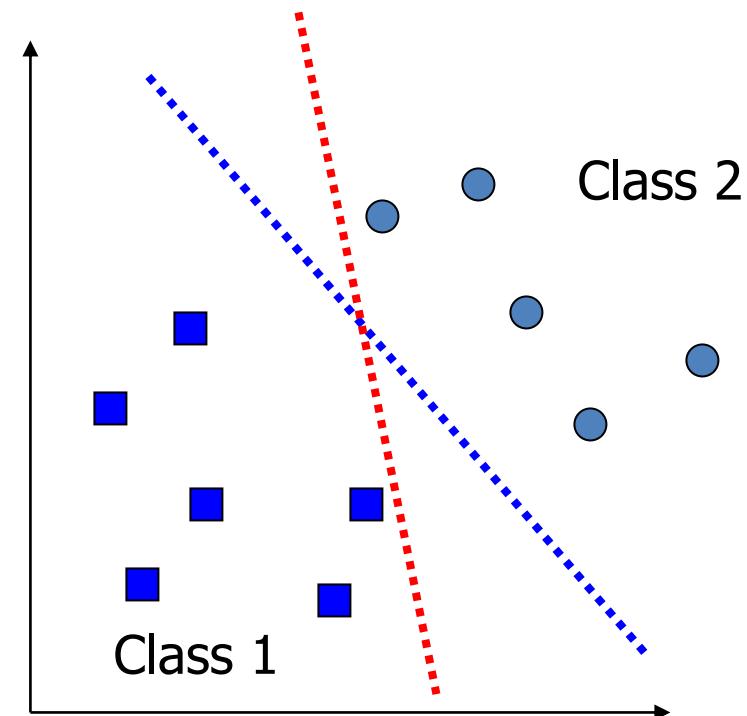
# Which decision boundary is better?

Suppose the training samples are linearly separable

We can find a decision boundary which gives zero **training** error

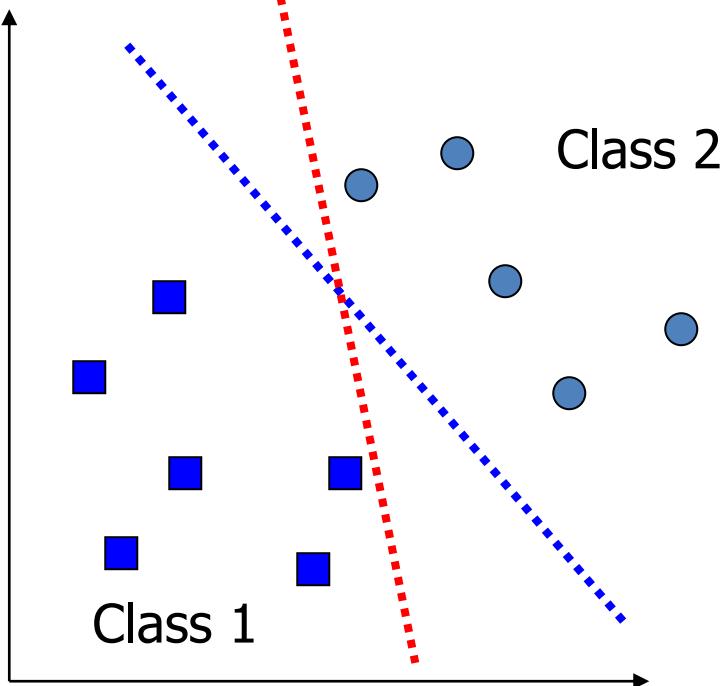
But there are many such decision boundaries

Which one is better?



# Compare two decision boundaries

Generalization error: Suppose we perturb data (new test data), which boundary is more susceptible to error?



# Support vector machine

In machine learning, support-vector machines (SVMs, also support-vector networks) are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis.

An SVM maximize the "gap" of training samples from the decision boundary to be as wide as possible.

Cortes, Corinna; Vapnik, Vladimir N. (1995).  
"Support-vector networks". Machine Learning.  
20 (3): 273–297.



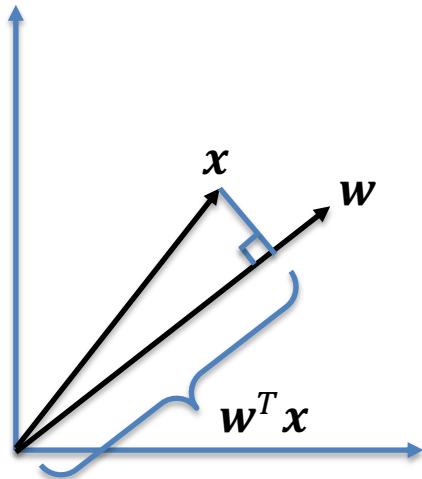
Cortes, Corinna



Vladimir Naumovich Vapnik

# Background: Projection in vector space

Projection  $\mathbf{w}^T \mathbf{x}$



Projection

$$\mathbf{w}^T \mathbf{x} = \left( \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x} \cdot \|\mathbf{w}\| = \mathbf{w}'^T \mathbf{x} \cdot \|\mathbf{w}\|$$

First project along a direction  $\mathbf{w}'$

Then stretch by magnitude  $\|\mathbf{w}\|$

# Background: Hyperplane

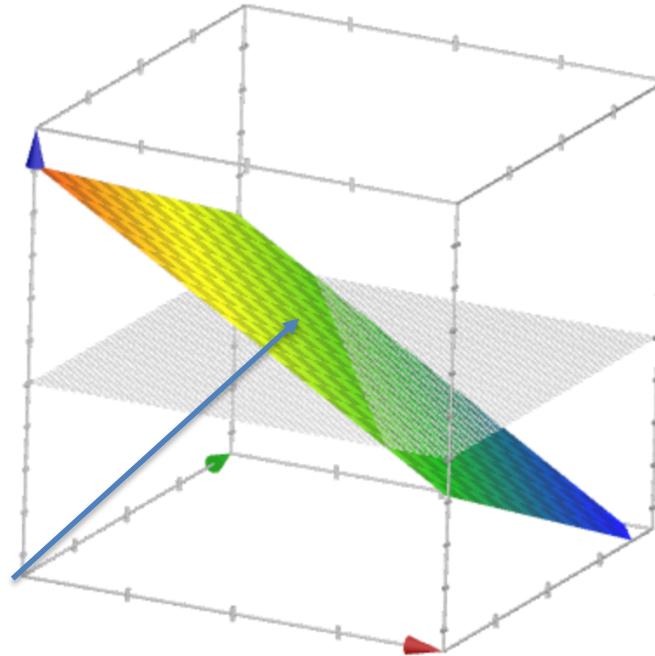
Hyperplane representation  $w$

$$z = -3x - 2y + 1$$

$$3x + 2y + z - 1 = 0$$

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

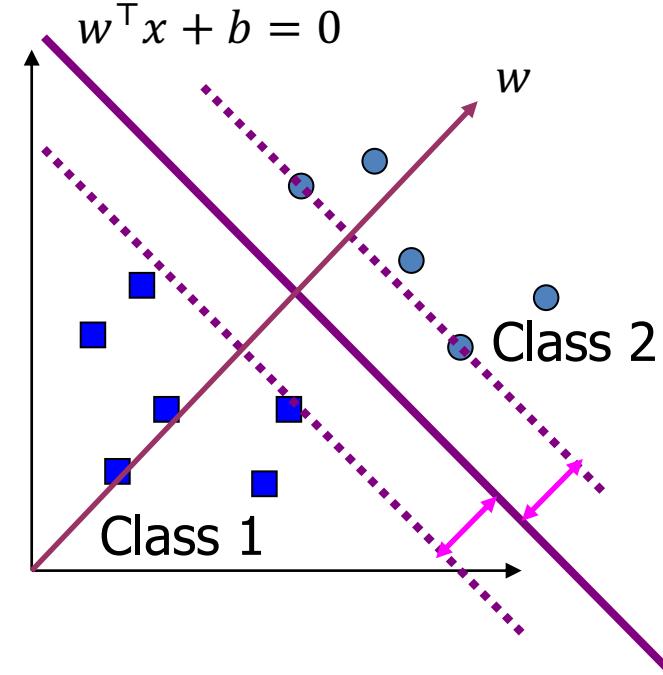
$$w = [3, 2, 1]^T, x = [x, y, z]^T, b = -1$$



	From	To
x	-10.0000	10.0000
y	-10.0000	10.0000
z	-44.9785	59.6919

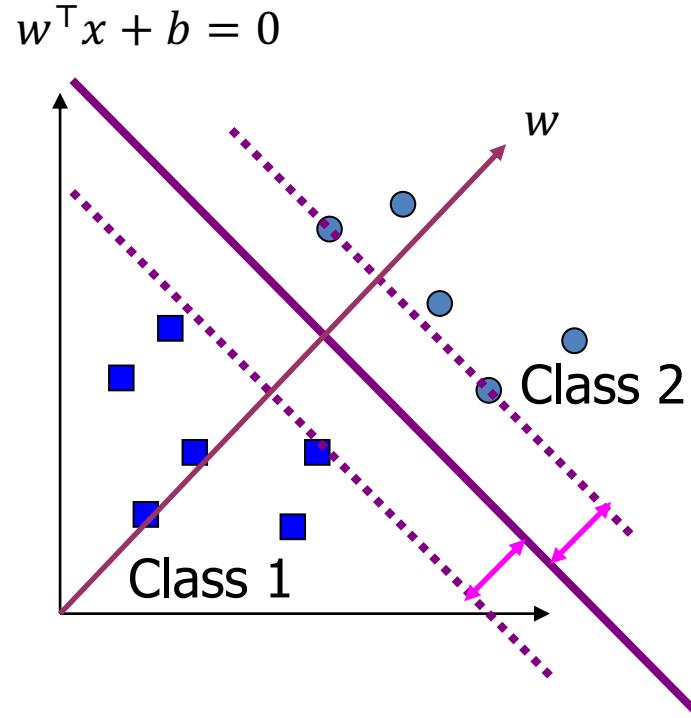
# Geometric interpretation of a classifier

- Parameterizing decision boundary as:  $w^T x + b = 0$ 
  - $w$  denotes a **vector orthogonal** to the decision boundary
  - $b$  is a **scalar offset** term
- Dash lines are parallel to decision boundary and they just hit the data points



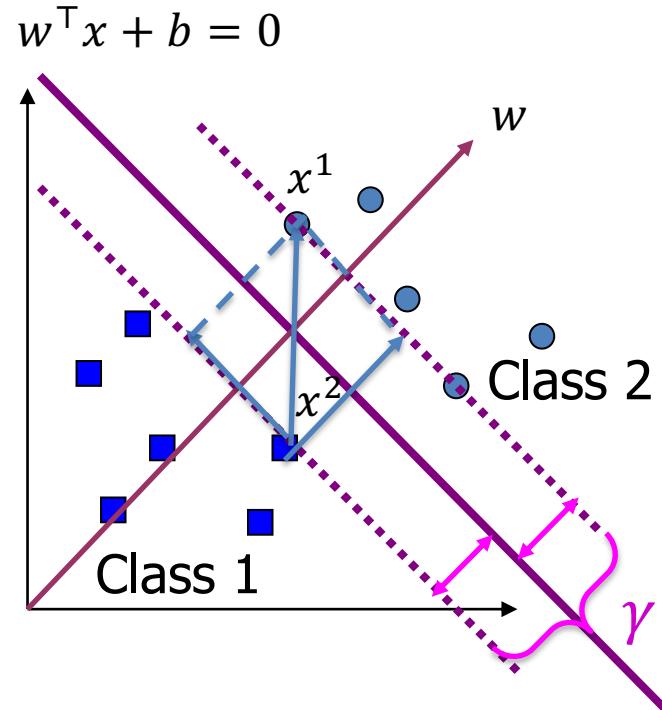
# Constraints on data points: correct separation

- Constraints on data points
  - For all  $x$  in class 2,  $y = 1$   
 $w^T x + b \geq c$
  - For all  $x$  in class 1,  $y = -1$   
 $w^T x + b \leq -c$
- Or more compactly  
 $(w^T x + b)y \geq c$



# Classifier margin

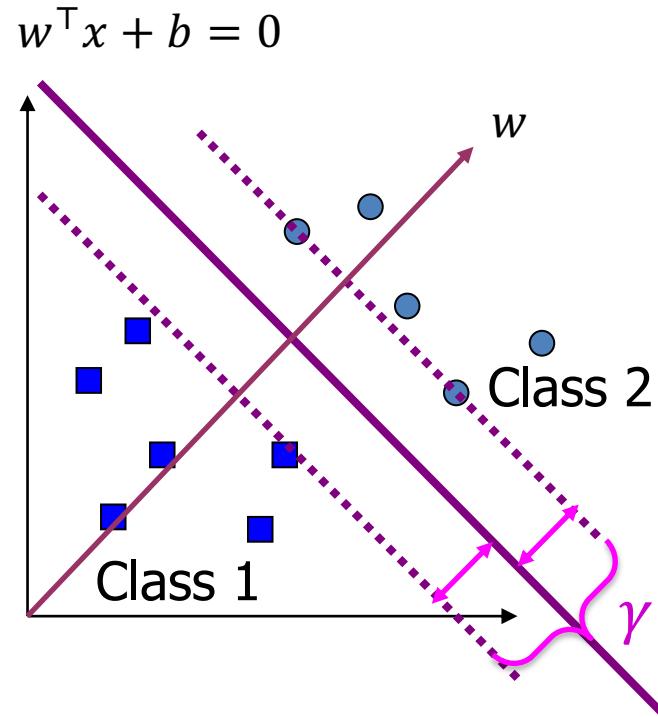
- Pick two data points  $x^1$  and  $x^2$  which are on each dash line respectively
- The **unnormalized** margin is  $w^\top(x^1 - x^2) := 2c$
- On the other hand, we can show  $w^\top(x^1 - x^2) = \gamma \|w\|$
- The margin is  $\gamma = \frac{2c}{\|w\|}$



# Maximum margin classifier

Find decision boundary  $w$  as far from data point as possible

$$\max_{w,b} \gamma = \frac{2c}{||w||}$$
$$s.t. y^i(w^\top x^i + b) \geq c, \forall i$$



# Equivalent form

$$\begin{aligned} & \max_{w,b} \frac{2c}{\|w\|} \\ \text{subject to } & y^i(w^\top x^i + b) \geq c, \forall i \end{aligned}$$

- Note that the magnitude of  $c$  merely scales  $w$  and  $b$ , and does not change the relative goodness of different classifiers
- Set  $c = 1$  (and drop the 2) to get a cleaner problem

$$\begin{aligned} & \max_{w,b} \frac{1}{\|w\|} \\ s.t. \quad & y^i(w^\top x^i + b) \geq 1, \forall i \end{aligned}$$

# Support vector machines

- A constrained convex quadratic programming problem

$$\begin{aligned} & \min_{w,b} \|w\|^2 \\ \text{subject to } & y^i(w^\top x^i + b) \geq 1, \forall i \end{aligned}$$

- Only a few of the constraints are relevant → **support vectors**
- Kernel methods are introduced for nonlinear classification problem

Solving constrained  
optimization problem

# Lagrangian Duality

The primal problem

$$\begin{aligned} & \min_w f(w) \\ & \text{s.t. } g_i(w) \leq 0, i = 1, \dots, k \\ & \quad h_i(w) = 0, i = 1, \dots, l \end{aligned}$$

The Lagrangian function

$$L(w, \alpha, \beta) = f(w) + \sum_i^k \alpha_i g_i(w) + \sum_i^l \beta_i h_i(w)$$

$\alpha_i \geq 0$ , and  $\beta_i$  are called the Lagrangian multipliers

# The KKT conditions

If there exists some saddle point of  $L$ , then the saddle point satisfies the following "Karush-Kuhn-Tucker" (KKT) conditions:

$$\frac{\partial L}{\partial w} = 0$$

$$\frac{\partial L}{\partial \alpha} = 0$$

$$\frac{\partial L}{\partial \beta} = 0$$

$$g_i(w) \leq 0$$

$$h_i(w) = 0$$

$$\alpha_i \geq 0$$

$$\alpha_i g_i(w) = 0$$

# Dual problem of support vector machines

$$\begin{aligned} & \min_{w,b} \|w\|^2 \\ & s.t. y^i(w^\top x^i + b) \geq 1, \forall i \end{aligned}$$

Convert to standard form

$$\begin{aligned} & \min_{w,b} \frac{1}{2} w^\top w \\ & s.t. 1 - y^i(w^\top x^i + b) \leq 0, \forall i \end{aligned}$$

The Lagrangian function

$$L(w, b, \alpha) = \frac{1}{2} w^\top w + \sum_{i=1}^m \alpha_i (1 - y^i(w^\top x^i + b))$$

# Deriving the dual problem

$$L(w, b, \alpha) = \frac{1}{2} w^\top w + \sum_{i=1}^m \alpha_i (1 - y^i (w^\top x^i + b))$$

Taking derivative and set to zero

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^m \alpha_i y^i x^i = 0$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i y^i = 0$$

# Plug back relation of w and b

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} \left( \sum_{i=1}^m \alpha_i y^i x^i \right)^\top \left( \sum_{j=1}^m \alpha_j y^j x^j \right) \\ &\quad + \sum_{i=1}^m \alpha_i \left( 1 - y^i \left( \left( \sum_{j=1}^m \alpha_j y^j x^j \right)^\top x^i + b \right) \right) \end{aligned}$$

After simplification (plug in the optimal  $w^*$  and  $b^*$ ):

$$L(w^*, b^*, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^i y^j (x^{i^\top} x^j)$$

# The dual problem of SVM

$$\begin{aligned} \max_{\alpha} L(w^*, b^*, \alpha) &= \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m \alpha_i \alpha_j y^i y^j (x^{i^\top} x^j) \\ s.t. \alpha_i &\geq 0, i = 1, \dots, m \\ \sum_i^m \alpha_i y^i &= 0 \end{aligned}$$

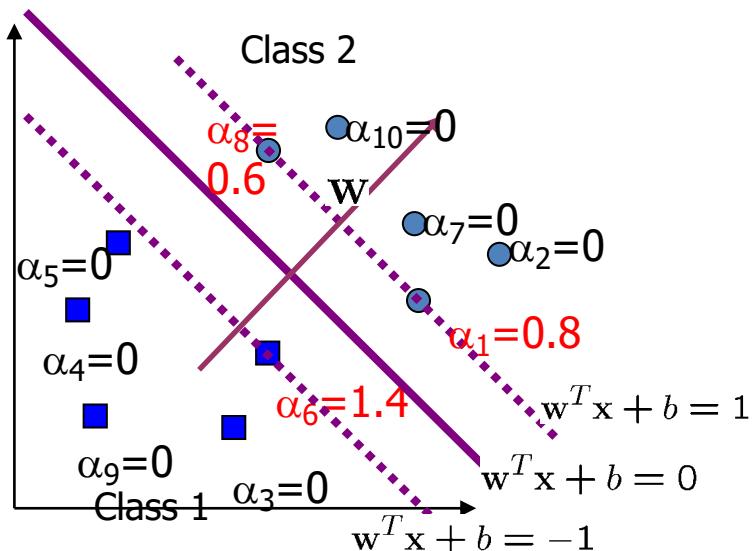
- This is a constrained **quadratic programming**
- Nice and convex, and global maximum can be found
- Very similar to the dual of minimum enclosing ball problem

# Properties of SVM

# Support vectors

Note that the KKT condition  $\alpha_i g_i(w) = 0$ :  $\alpha_i \left(1 - y^i(w^\top x^i + b)\right) = 0$

- For data points with  $\left(1 - y^i(w^\top x^i + b)\right) < 0, \alpha_i = 0$
- For data points with  $\left(1 - y^i(w^\top x^i + b)\right) = 0, \alpha_i \geq 0$



Call the training data points whose  $\alpha_i$ 's are nonzero the **support vectors (SV)**

# Computing b and obtain the classifier

- Pick any data point with  $\alpha_i > 0$ , solve for  $b$  with

$$1 - y^i(w^\top x^i + b) = 0$$

- One KKT condition:  $\frac{\partial L}{\partial w} = 0$

$$w = \sum_{i=1}^m \alpha_i y^i x^i$$

- For a new test point  $z$

- Compute

$$w^\top z + b = \sum_{i \in \text{support vectors}} \alpha_i y^i (x^{i^T} z) + b$$

- Classify  $z$  as class 1 if the  $w^\top z + b > 0$ , and class 2 otherwise

# Interpretation of support vector machines

The optimal  $w$  is a linear combination of a small number of data points. This “sparse” representation can be viewed as data compression

To compute the weights  $\alpha_i$ , and to use support vector machines we need to specify only the inner products (or kernel) between the examples  $x^{i^T} x^j$

We make decisions by comparing each new example  $z$  with only the support vectors:

$$y^* = \text{sign} \left( \sum_{i \in \text{support vectors}} \alpha_i y^i (x^{i^T} z) + b \right)$$

# Comparison

	Bayes classifier (Gaussian)	KNN	Logistic regression	SVM
Number of parameters	$O(n^2)$	1 (non-parametric)	$O(n)$	$O(n)$
Robustness to model mismatch	No	Yes	relatively	relatively
Noisy prediction?	No	Yes	Yes (when logistic function value is close to 0.5)	No

# Extensions of SVM

- Kernelized SVM: In addition to performing linear classification, SVMs can perform a **non-linear classification (non-linear decision boundary)** by kernelized SVM implicitly mapping their inputs into high-dimensional feature spaces.
- Soft-margin SVM: allow a few data points to be mis-classified.
- Multi-class SVM

# Kernelized SVM

$$\max_{\alpha} L(w^*, b^*, \alpha) = \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m \alpha_i \alpha_j y^i y^j (x^{i^\top} x^j)$$

s.t.  $\alpha_i \geq 0, i = 1, \dots, m$

$$\sum_i^m \alpha_i y^i = 0$$

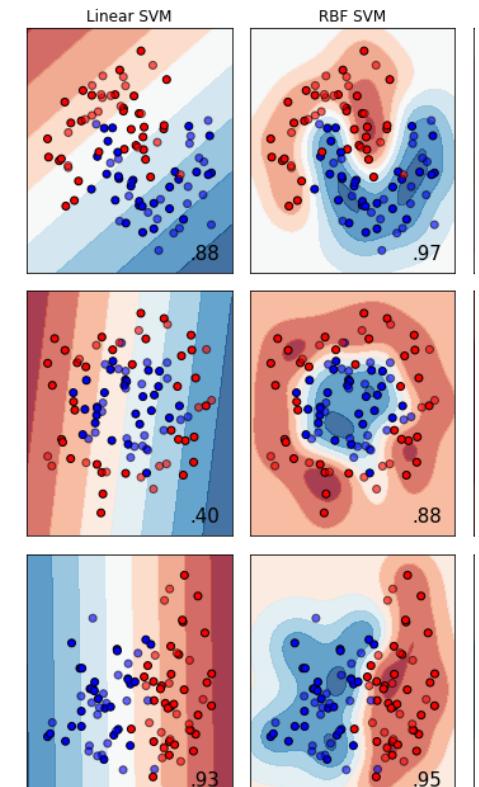
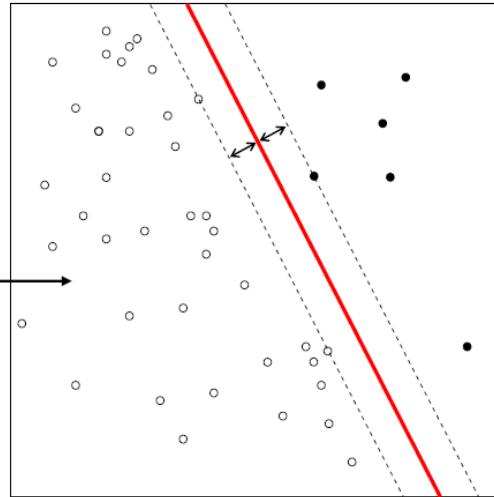
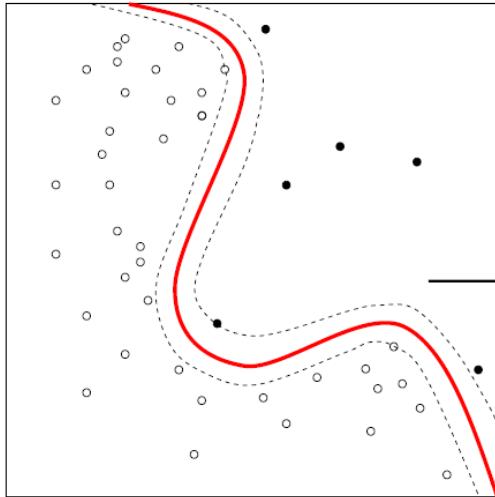
e.g. Gaussian RBF kernel

*Replaced by  $k(x^i, x^j)$*

- Implicit map data to non-linear feature space

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

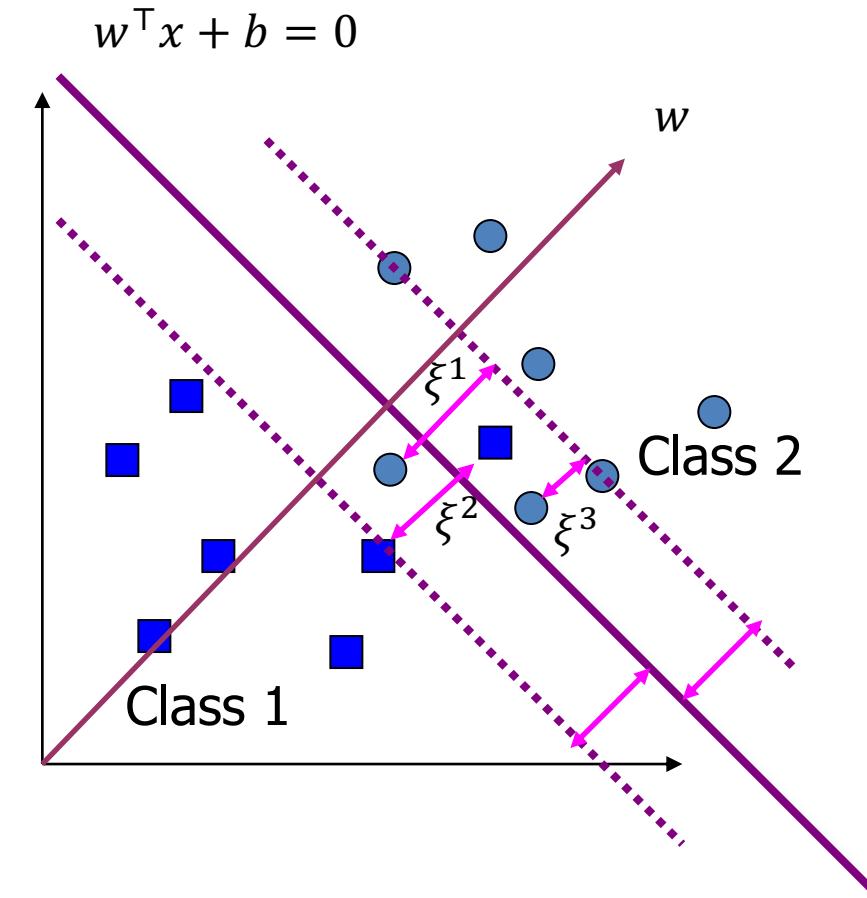
# Implicitly map data to non-linear feature space



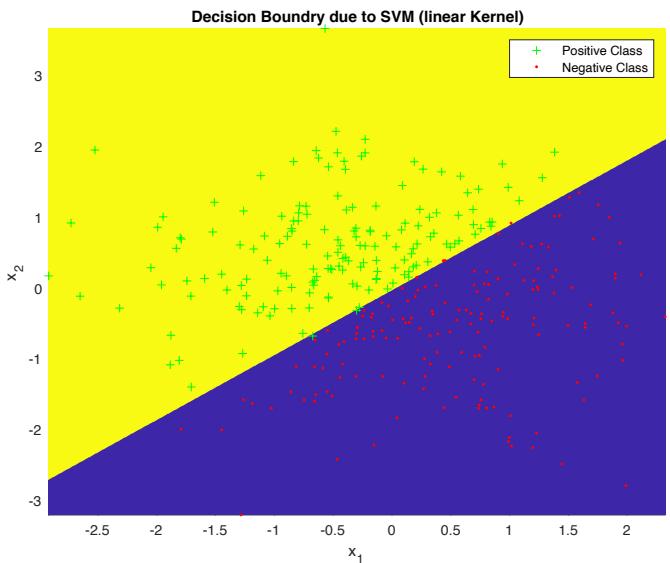
# Soft-margin SVM

- What if the data is not linearly separable?
- We will allow a few points to violate the hard margin constraint
- $(w^T x + b)y \geq 1 - \xi$

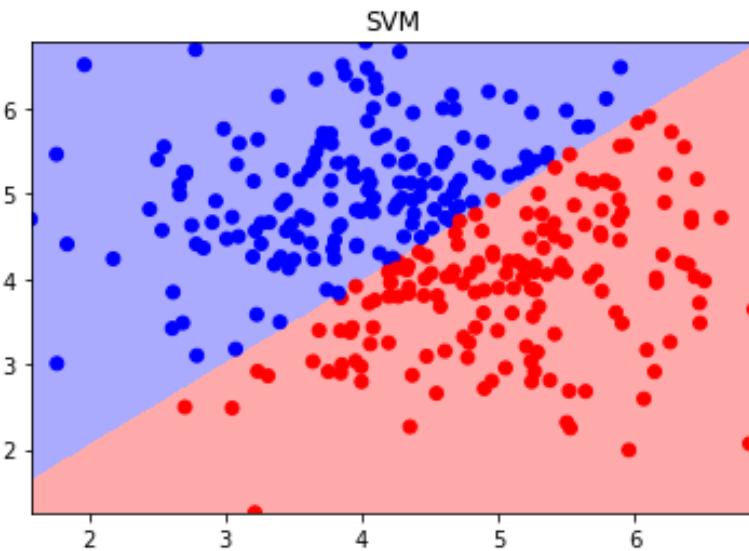
$$\begin{aligned} & \min_{w,b,\xi} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi^i \\ \text{s.t. } & y^i(w^T x^i + b) \geq 1 - \xi^i, \xi^i \geq 0, \forall i \end{aligned}$$



# Demo



MATLAB



Python

# Comparison with logistic regression

- Logistic regression focuses on maximizing the probability of the data.
- SVM tries to find the separating hyperplane that maximizes the distance of the closest points to the margin
- Which one to use?
  - SVM typically works better for “clearly” linearly/nonlinearly separable classes
  - Logistic regression can work better for classes are not separable “in the middle”, due to its probabilistic formulation
- In practice, you should try both and compare.

# Scikit-learn – python library for machine learning



Hand-written digits classification example using SVM

<http://scikit-learn.org/stable/tutorial/basic/tutorial.html#introduction>

[http://scikit-learn.org/stable/auto\\_examples/classification/plot\\_digits\\_classification.html](http://scikit-learn.org/stable/auto_examples/classification/plot_digits_classification.html)

# Demo: comparison of classification algorithms

