

Support Vector Machine (SVM)

CSE 5334 Data Mining
Spring 2020

Won Hwa Kim

(Slides courtesy of Mark Craven and David Page Jr. at UW - Madison)



Goals for this lecture

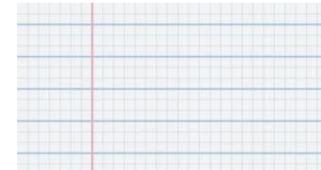
you should understand the following concepts

- the margin
- slack variables
- the linear support vector machine
- hinge loss
- non-linear SVMs
- the kernel trick
- the primal and dual formulations of SVM learning
- support vectors
- the kernel matrix
- valid kernels
- polynomial kernel
- Gaussian kernel
- string kernels
- support vector regression

Four key SVM ideas

- **Maximize the margin**

don't choose just *any* separating hyperplane



- **Penalize misclassified examples**

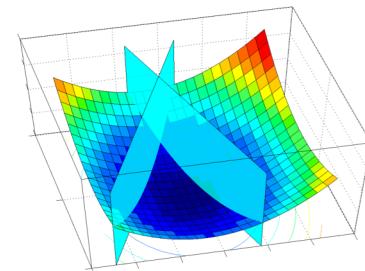
use soft constraints and slack variables



- **Use optimization methods to find model**

linear programming

quadratic programming



- **Use kernels to represent nonlinear functions and handle complex instances (sequences, trees, graphs, etc.)**



Some key vector concepts

the *dot product* between two vectors w and x is defined as:

$$w \cdot x = w^T x = \sum_i w_i x_i$$

for example

$$\begin{bmatrix} 1 \\ 3 \\ -5 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ -2 \\ -1 \end{bmatrix} = (1)(4) + (3)(-2) + (-5)(-1) = 3$$

the 2-norm (Euclidean length) of a vector x is defined as:

$$\|x\|_2 = \sqrt{\sum_i |x_i|^2}$$

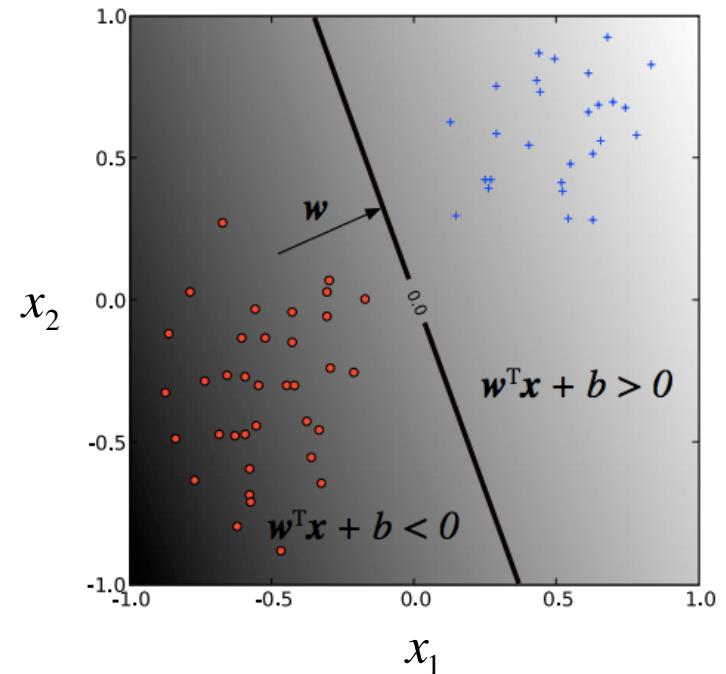
Linear separator learning revisited

suppose we encode our classes as $\{-1, +1\}$ and consider a linear classifier

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \left(\sum_{i=1}^n w_i x_i \right) + b > 0 \\ -1 & \text{otherwise} \end{cases}$$

an instance $\langle \mathbf{x}, y \rangle$ will be classified correctly if

$$y(\mathbf{w}^\top \mathbf{x} + b) > 0$$



Large margin classification

- Given a training set that is linearly separable, there are infinitely many hyperplanes that could separate the positive/negative instances.
 - Which one should we choose?
- In SVM learning, we find the hyperplane that maximizes the margin.

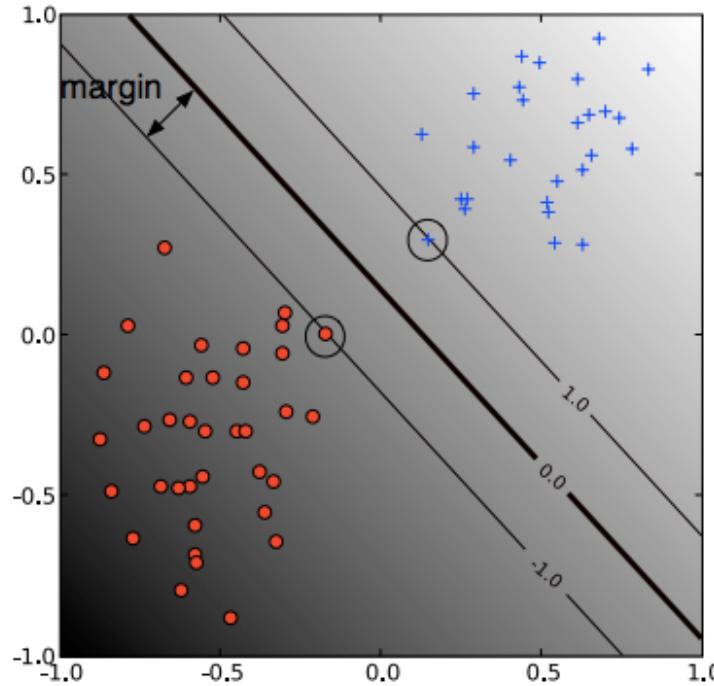
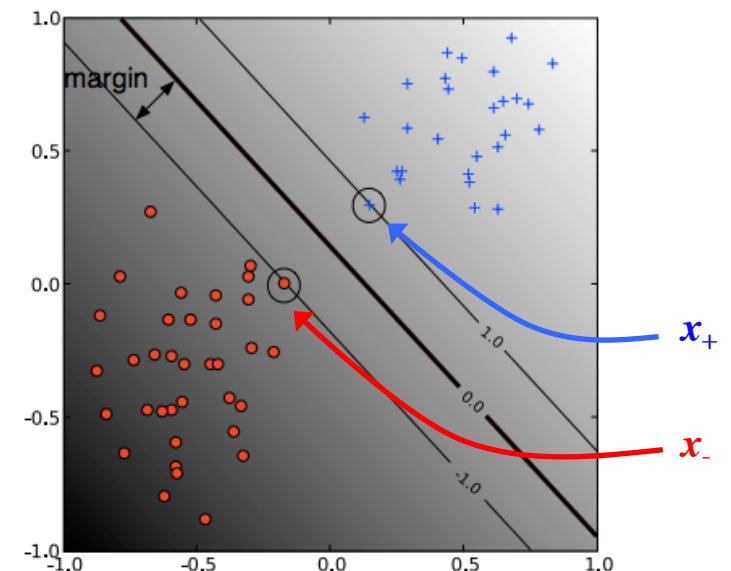


Figure from Ben-Hur & Weston,
Methods in Molecular Biology 2010

The margin

- suppose we learn a hyperplane h given a training set D
- let x_+ denote the closest instance to the hyperplane among positive instances, and similarly for x_- and negative instances
- since $\mathbf{w}^T \mathbf{x} + b = 0$ and $c(\mathbf{w}^T \mathbf{x} + b) = 0$ define the same hyperplane, we have the freedom to choose the normalization of \mathbf{w}
- choose normalization such that $\mathbf{w}^T \mathbf{x}_+ + b = 1$ and $\mathbf{w}^T \mathbf{x}_- + b = -1$ for positive and negative support vectors respectively
- then the margin is given by

$$\begin{aligned}\text{margin}_D(h) &= \frac{1}{2} \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \cdot (\mathbf{x}_+ - \mathbf{x}_-) \\ &= \frac{\mathbf{w}^T (\mathbf{x}_+ - \mathbf{x}_-)}{\|\mathbf{w}\|_2} \\ &= \frac{1}{\|\mathbf{w}\|_2}\end{aligned}$$



The hard-margin SVM

- given a training set $D = \{ \langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle \}$
- we can frame the goal of maximizing the margin as a constrained optimization problem

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2$$

subject to constraints : $y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1$
for $i = 1, \dots, m$

adjust these parameters

to minimize this

correctly classify $\mathbf{x}^{(i)}$ with room to spare

- and use standard algorithms to find an optimal solution to this problem

The soft-margin SVM

[Cortes & Vapnik, *Machine Learning* 1995]

- if the training instances are not linearly separable, the previous formulation will fail
- we can adjust our approach by using *slack variables* (denoted by ξ) to tolerate errors

$$\underset{\mathbf{w}, b, \xi^{(1)} \dots \xi^{(m)}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi^{(i)}$$



subject to constraints : $y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}$

$$\xi^{(i)} \geq 0$$

for $i = 1, \dots, m$

- C determines the relative importance of maximizing margin vs. minimizing slack

The effect of C in a soft-margin SVM

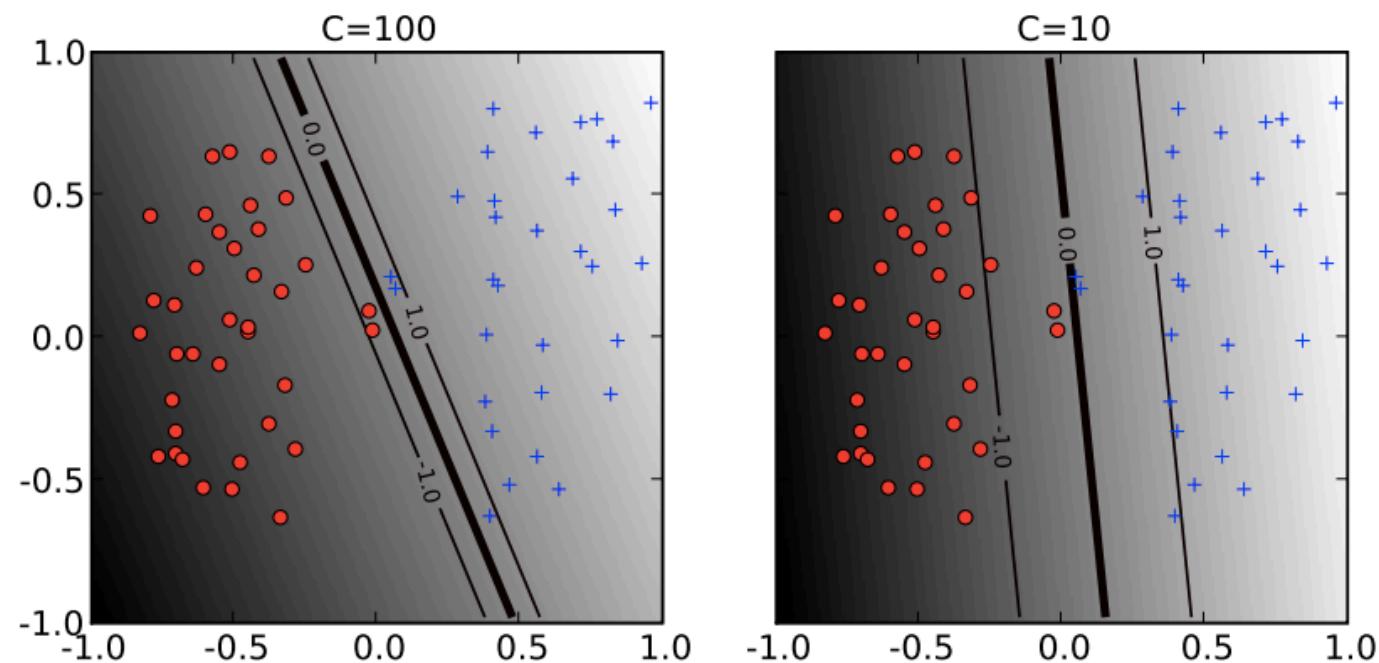
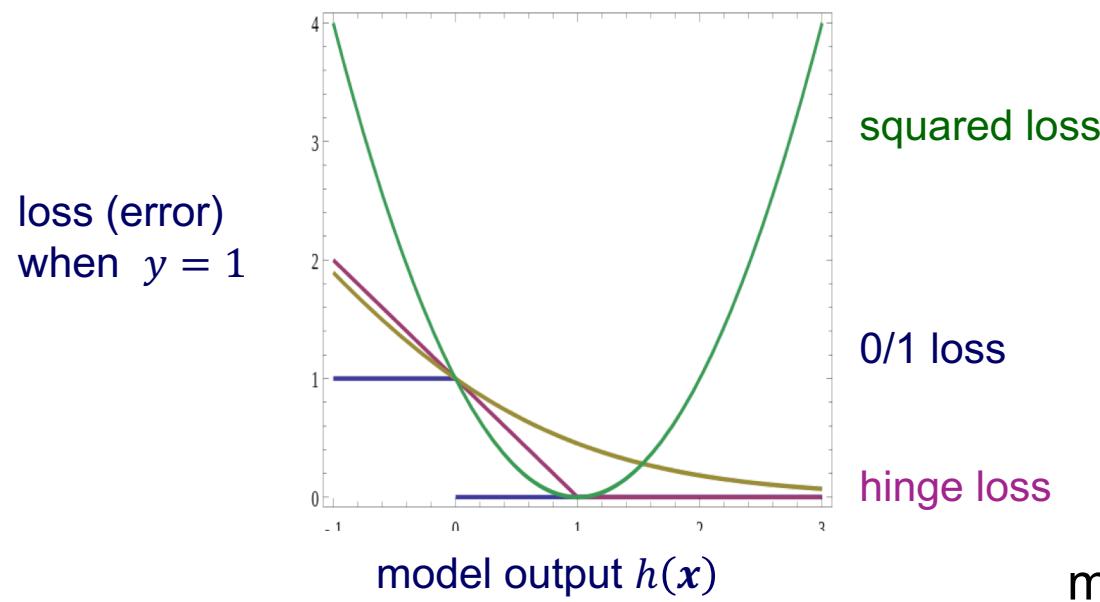


Figure from Ben-Hur & Weston,
Methods in Molecular Biology 2010

Hinge loss

- when we covered neural nets, we talked about minimizing squared loss and cross-entropy loss
- soft-margin SVMs minimize *hinge loss*



$$\underset{\mathbf{w}, b, \xi^{(1)} \dots \xi^{(m)}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi^{(i)}$$

Non-linear classifiers

- What if a linear separator is not an appropriate decision boundary for a given task?
- For any data set, there exists a mapping ϕ to a higher-dimensional space such that the data is linearly separable

$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_k(\mathbf{x}))$$

- Example: mapping to quadratic space

$$\mathbf{x} = \langle x_1, x_2 \rangle \quad \text{suppose } \mathbf{x} \text{ is represented by 2 features}$$

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$$

- now try to find a linear separator in this space

Non-linear classifiers

- for the linear case, our discriminant function was given by

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \\ -1 & \text{otherwise} \end{cases}$$

- for the nonlinear case, it can be expressed as

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \hat{\mathbf{w}} \cdot \phi(\mathbf{x}) + b > 0 \\ -1 & \text{otherwise} \end{cases}$$

where $\hat{\mathbf{w}}$ is a higher dimensional vector

SVM with polynomial kernels

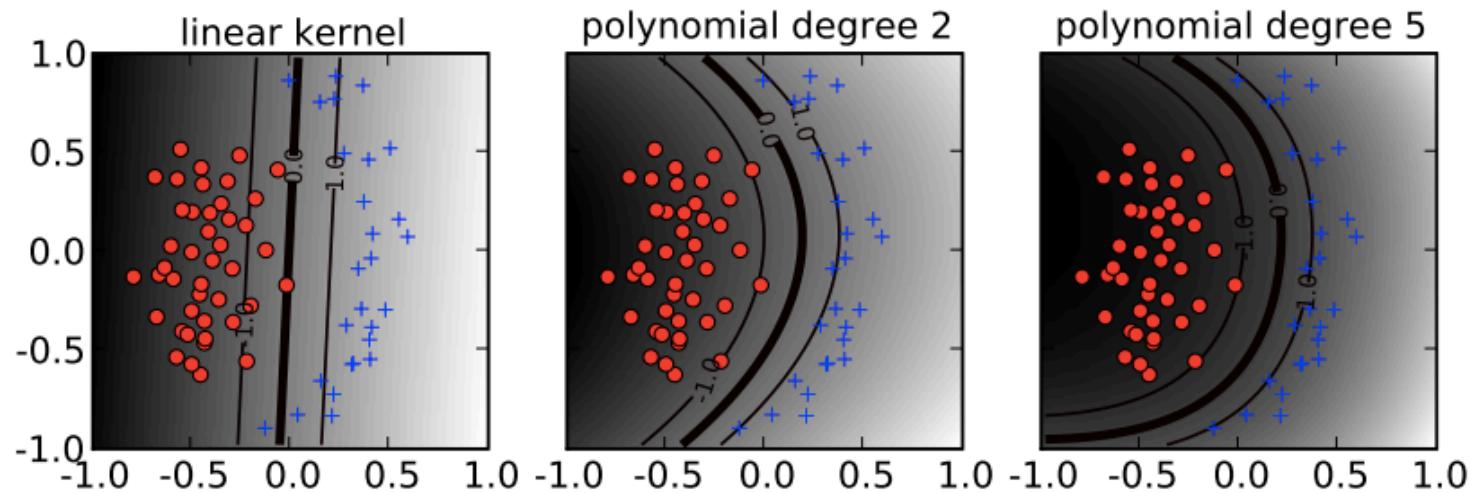


Figure from Ben-Hur & Weston,
Methods in Molecular Biology 2010

The kernel trick

- explicitly computing this nonlinear mapping does not scale well
- a dot product between two higher-dimensional mappings can sometimes be implemented by a *kernel function*
- example: quadratic kernel

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x} \cdot \mathbf{z} + 1)^2 \\ &= (x_1 z_1 + x_2 z_2 + 1)^2 \\ &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 + 2x_1 z_1 + 2x_2 z_2 + 1 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1) \cdot \\ &\quad (z_1^2, \sqrt{2}z_1 z_2, z_2^2, \sqrt{2}z_1, \sqrt{2}z_2, 1) \\ &= \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) \end{aligned}$$

The kernel trick

- thus we can use a kernel to compute the dot product without explicitly mapping the instances to a higher-dimensional space

$$k(x, z) = (x \cdot z + 1)^2 = \phi(x) \cdot \phi(z)$$

But why is the kernel trick helpful?

The kernel trick

- given a training set $D = \{ \langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle \}$
- suppose the weight vector can be represented as a linear combination of the training instances

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \mathbf{x}^{(i)}$$

- then we can represent a linear SVM as

$$\sum_{i=1}^m \alpha_i \mathbf{x}^{(i)} \cdot \mathbf{x} + b$$

- and a nonlinear SVM as

$$\begin{aligned} & \sum_{i=1}^m \alpha_i \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^m \alpha_i k(\mathbf{x}^{(i)}, \mathbf{x}) + b \end{aligned}$$

The *primal* and *dual* form of the hard-margin SVM

primal

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2$$

subject to constraints: $y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1$

for $i = 1, \dots, m$

dual

$$\underset{\alpha_1, \dots, \alpha_m}{\text{maximize}} \quad \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^m \alpha_j \alpha_k y^{(j)} y^{(k)} (\mathbf{x}^{(j)} \cdot \mathbf{x}^{(k)})$$

subject to constraints: $\alpha_i \geq 0 \quad \text{for } i = 1, \dots, m$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

The *dual*/form with a kernel (hard margin version)

primal

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2$$

subject to constraints: $y^{(i)}(\mathbf{w}^\top \phi(\mathbf{x}^{(i)}) + b) \geq 1$

for $i = 1, \dots, m$

dual

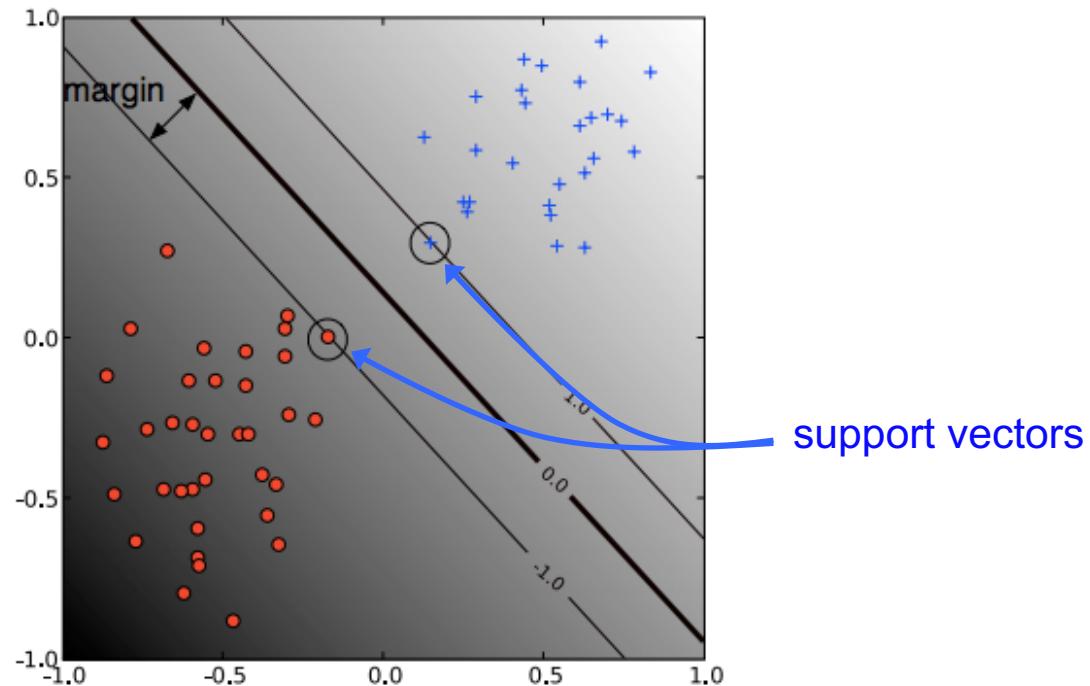
$$\underset{\alpha_1, \dots, \alpha_m}{\text{maximize}} \quad \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^m \alpha_j \alpha_k y^{(j)} y^{(k)} k(\mathbf{x}^{(j)}, \mathbf{x}^{(k)})$$

subject to constraints: $\alpha_i \geq 0 \quad \text{for } i = 1, \dots, m$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Support vectors

- the solution to the dual formulation is a sparse linear combination of the training instances
- those instances having $\alpha_i > 0$ are called *support vectors* – they lie on the margin boundary
- the solution wouldn't change if all the instances with $\alpha_i = 0$ were deleted



The kernel matrix

- the kernel matrix (a.k.a. Gram matrix) represents pairwise similarities for instances in the training set

$$\begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & k(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \cdots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(m)}) \\ k(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & \ddots & & \\ \vdots & & & \\ k(\mathbf{x}^{(m)}, \mathbf{x}^{(1)}) & & & k(\mathbf{x}^{(m)}, \mathbf{x}^{(m)}) \end{bmatrix}$$

- it represents the information about the training set that is provided as input to the optimization process

Some common kernels

- polynomial of degree d

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^d$$

- polynomial of degree up to d

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z} + 1)^d$$

- radial basis function (RBF) (a.k.a. Gaussian)

$$k(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$$

The RBF kernel

- the feature mapping ϕ for the RBF kernel is infinite dimensional!
- recall that $k(x, z) = \phi(x) \cdot \phi(z)$

$$k(x, z) = \exp\left(-\frac{1}{2}\|x - z\|^2\right) \text{ for } \gamma = \frac{1}{2}$$

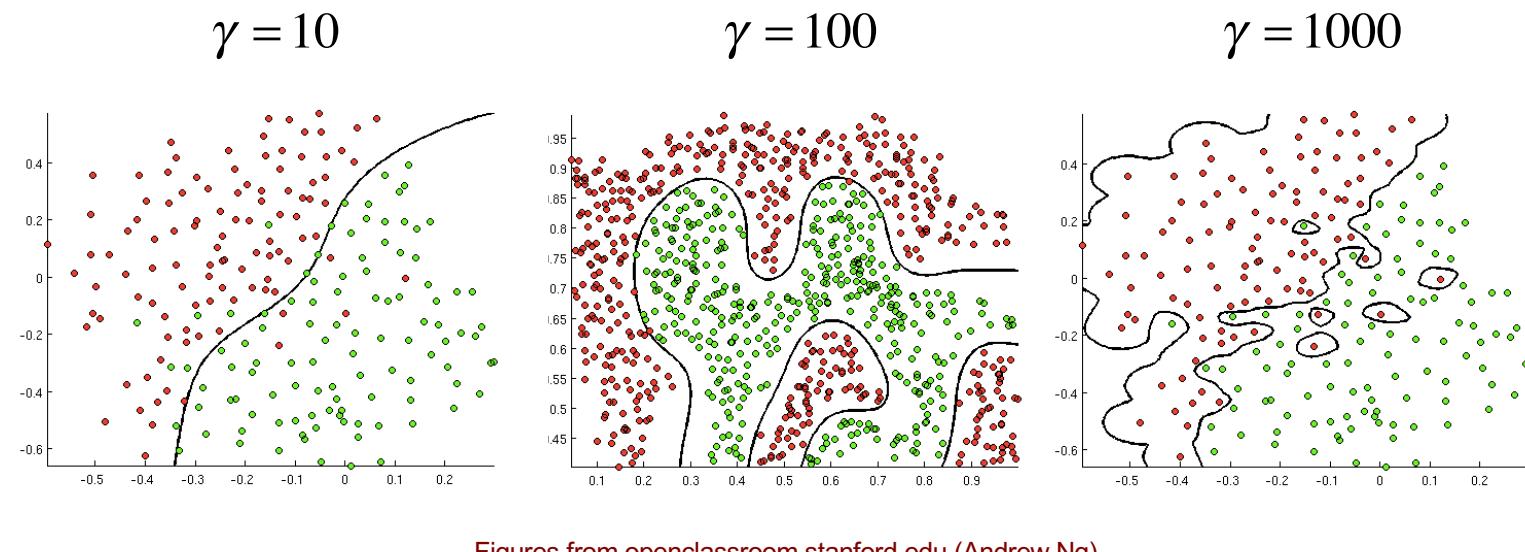
$$= \exp\left(-\frac{1}{2}\|x\|^2\right)\exp\left(-\frac{1}{2}\|z\|^2\right)\exp(x \cdot z)$$

$$= \exp\left(-\frac{1}{2}\|x\|^2\right)\exp\left(-\frac{1}{2}\|z\|^2\right)\left(\sum_{n=0}^{\infty} \frac{(x \cdot z)^n}{n!}\right)$$



from the Taylor series
expansion of $\exp(x \cdot z)$

The RBF kernel illustrated



What makes a valid kernel?

- $k(x, z)$ is a valid kernel if there is some ϕ such that

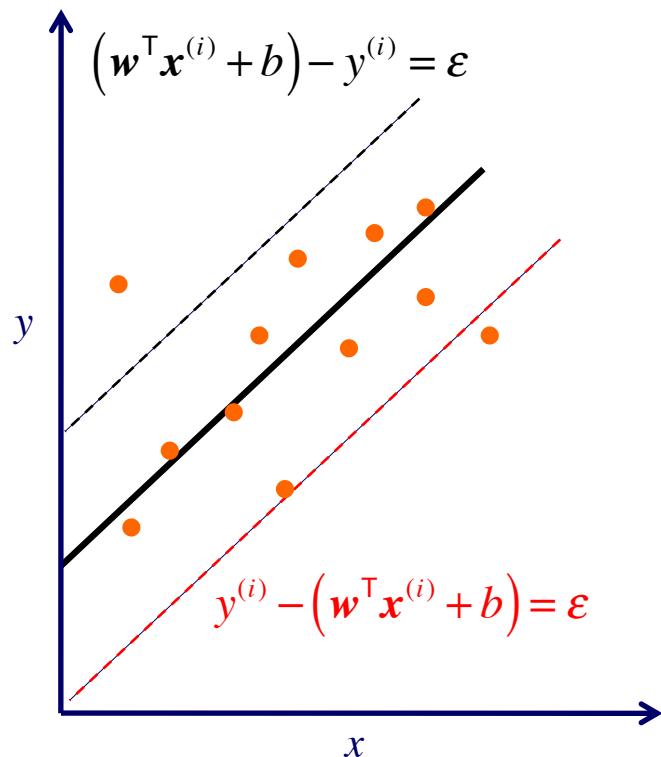
$$k(x, z) = \phi(x) \cdot \phi(z)$$

- this holds for a symmetric function $k(x, z)$ if and only if the kernel matrix K is positive semidefinite for any training set (Mercer's theorem)

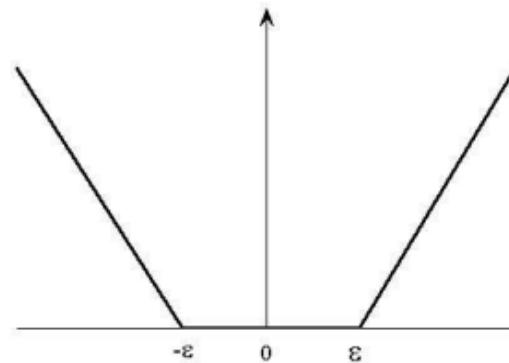
definition of positive semidefinite (p.s.d): $\forall v : v^\top K v \geq 0$

Support vector regression

- the SVM idea can also be applied in regression tasks
- an ε -insensitive error function specifies that a training instance is well explained if the model's prediction is within ε of $y^{(i)}$



ε -insensitive error function



Support vector regression

$$\underset{\mathbf{w}, b, \xi^{(1)} \dots \xi^{(m)}, \hat{\xi}^{(1)} \dots \hat{\xi}^{(m)}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m (\xi^{(i)} + \hat{\xi}^{(i)})$$

$$\text{subject to constraints: } (\mathbf{w}^\top \mathbf{x}^{(i)} + b) - y^{(i)} \leq \varepsilon + \xi^{(i)}$$

$$y^{(i)} - (\mathbf{w}^\top \mathbf{x}^{(i)} + b) \leq \varepsilon + \hat{\xi}^{(i)}$$

$$\xi^{(i)}, \hat{\xi}^{(i)} \geq 0$$

for $i = 1, \dots, m$

slack variables allow predictions
for some training instances to be
off by more than ε

Comments...

- we can find solutions that are globally optimal (maximize the margin)
 - because the learning task is framed as a convex optimization problem
 - no local minima, in contrast to multi-layer neural nets
- there are two formulations of the optimization: *primal* and *dual*
 - dual represents classifier decision in terms of support vectors
 - dual enables the use of kernel functions
- we can use a wide range of optimization methods to learn SVM
 - standard quadratic programming solvers
 - SMO [Platt, 1999]
 - linear programming solvers for some formulations
 - etc.