

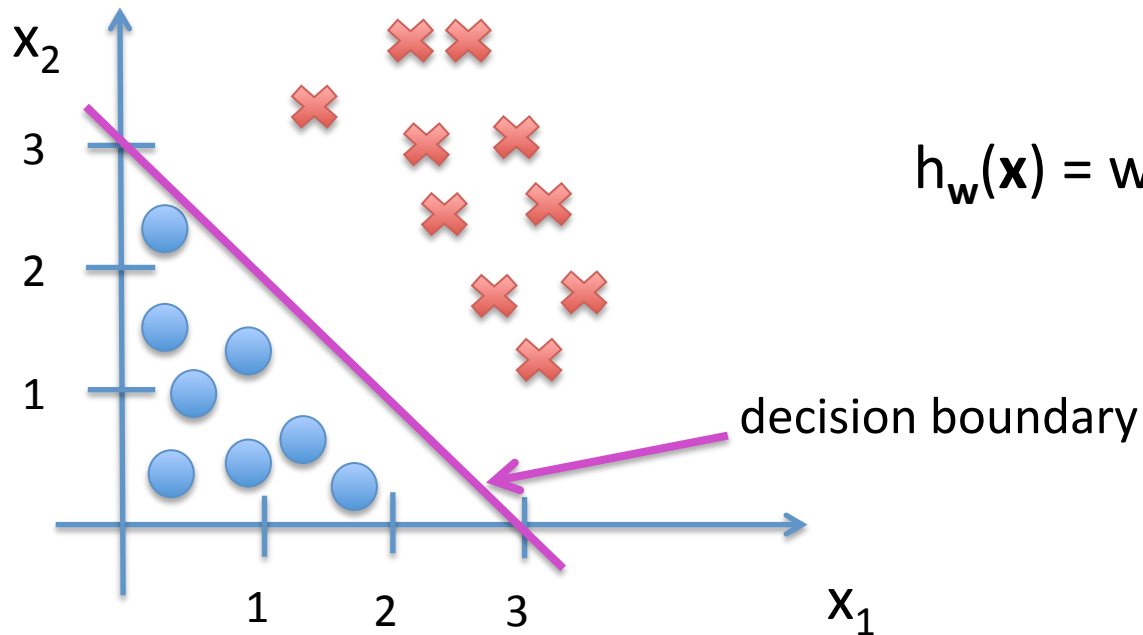
Machine Learning: Support Vector Machines (SVMs), Part 2

Slides adapted from David Sontag, who adapted from Luke Zettlemoyer, Vibhav Gogate, and Carlos Guestrin

Kernel Trick

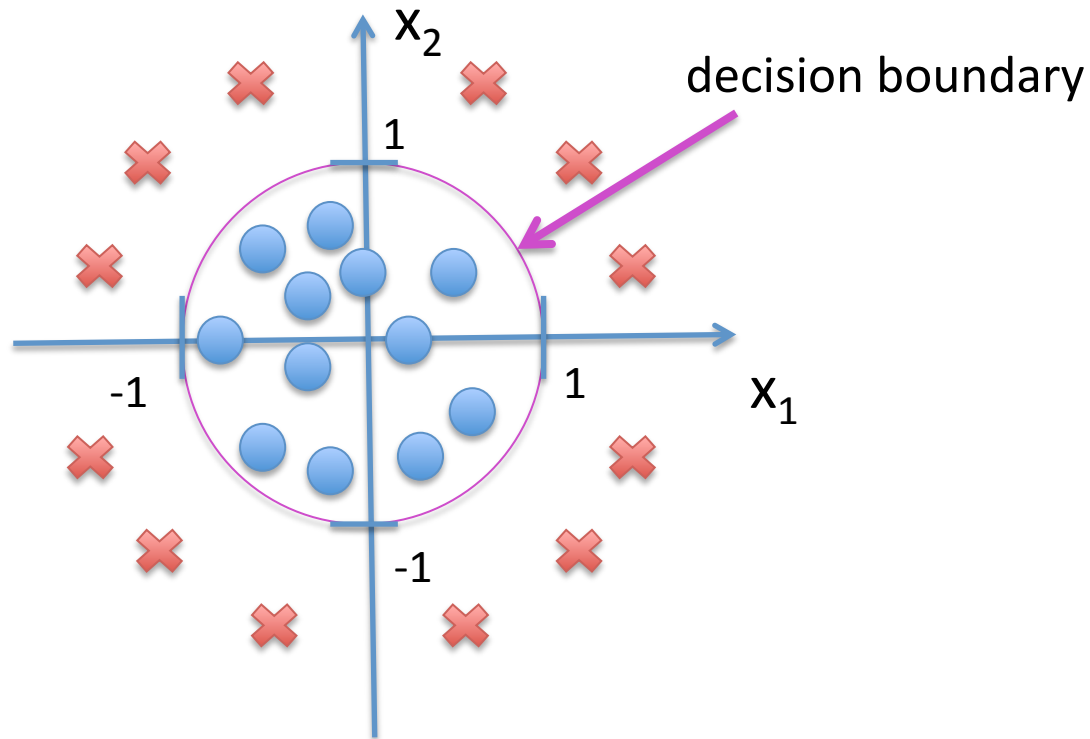
- Allows for non-linear decision boundaries
- Allows you to dramatically increase number of features without significantly increasing storage or computation
- Mathematical: Uses duality theory from convex optimization

Decision Boundary



Predict “ $y=1$ ” if $-3+x_1+x_2 \geq 0$

Non-linear decision boundaries



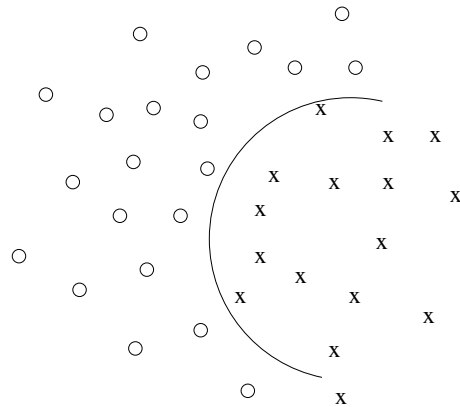
$$\mathbf{w} \in \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2$$

Predict "y=1" if $-1 + x_1^2 + x_2^2 \geq 0$

Non-linear decision region

$$b + w_1x_1 + w_2x_2 + \dots + w_nx_n + w_{n+1}x_1x_2 + w_{n+2}x_1x_3 + \dots + w_qe^{x_1} > 0$$



Non-linear separator in the **original x-space**

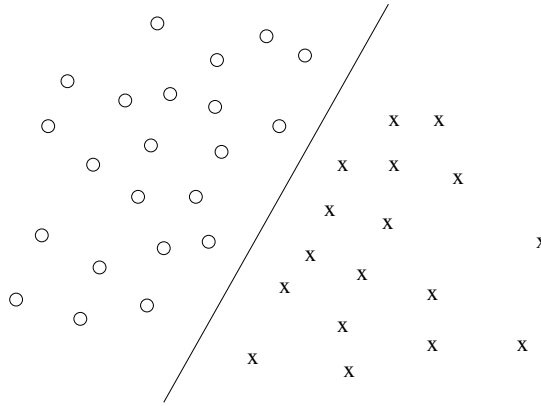
- Can potentially reduce errors
- But becomes much more difficult to solve non-linear optimization problem

Consider using features that are functions of features

- For example, in addition to current features x_1, x_2, \dots, x_n , use features like $x_1x_2, x_1x_3, \dots, e^{x_1}, e^{x_2}$.
- New set of features:
 - $\varphi_1 = x_1, \varphi_2 = x_2, \varphi_n = x_n, \varphi_{n+1} = x_1x_2, \varphi_{n+2} = x_1x_3, \dots$
- New feature vector
 - $\boldsymbol{\varphi} = (x_1, x_2, \dots, x_n, x_1x_2, x_1x_3, \dots, e^{x_1}, e^{x_2}, \dots)$
- From original data $x^{(1)}, x^{(2)}, \dots, x^{(m)}$, can create new data $\boldsymbol{\varphi}^{(1)}, \boldsymbol{\varphi}^{(2)}, \dots, \boldsymbol{\varphi}^{(m)}$

Decision Region in ϕ space

$$b + w_1 \phi_1 + \dots + w_n \phi_n + w_{n+1} \phi_{n+1} + w_{n+2} \phi_{n+2} + \dots + w_Q \phi_Q > 0$$



Linear separator in the **feature ϕ -space**

- This is now linear! Can use standard SVM methodology.

Minimize _{\mathbf{w}, b, ξ} $\mathbf{w} \bullet \mathbf{w} + C \sum_i \xi^{(i)}$

subject to

$(\mathbf{w} \bullet \boldsymbol{\phi}(\mathbf{x}^{(i)}) + b)y^{(i)} \geq 1 - \xi^{(i)}$ for all i

- But \mathbf{w} and $\boldsymbol{\phi}$ are now $Q \gg N$ dimensional.

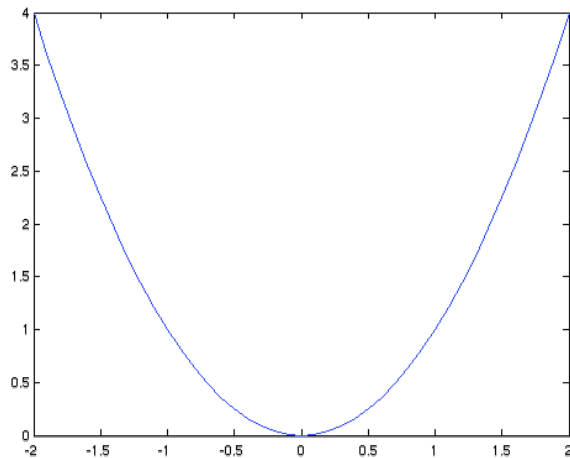
What's Next!

- One of the most interesting and exciting advances in machine learning: Kernel trick
- Basic idea: for some classes of non-linear features, dimensionality doesn't increase
- But first, a detour
 - Constrained optimization!

Constrained optimization

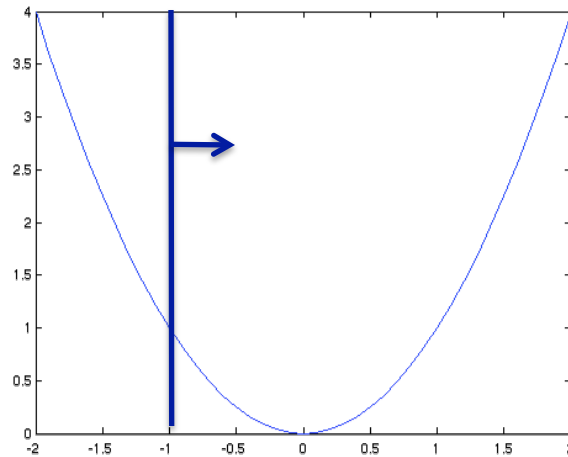
$$\begin{array}{ll} \min_x & x^2 \\ \text{s.t.} & x \geq b \end{array}$$

No Constraint



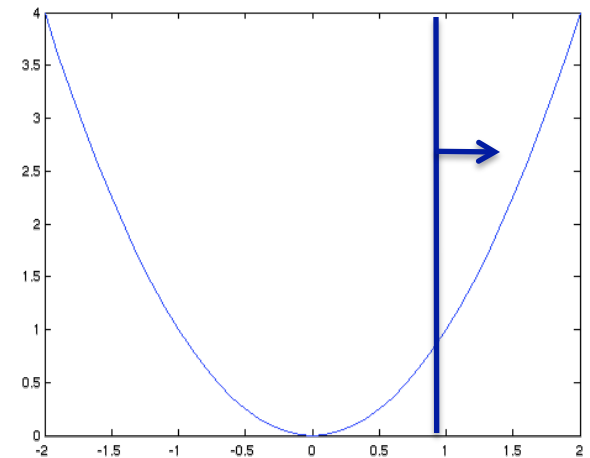
$$x^* = 0$$

$x \geq -1$



$$x^* = 0$$

$x \geq 1$

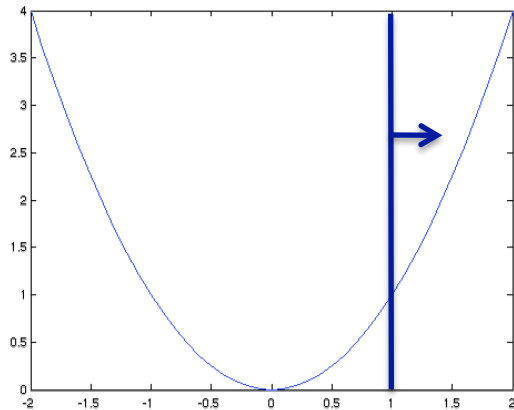


$$x^* = 1$$

How do we solve with constraints?

→ Lagrange Multipliers!!!

Lagrange multipliers – Dual variables



$$\begin{array}{ll} \min_x & x^2 \\ \text{s.t.} & x \geq b \end{array}$$

Add Lagrange multiplier

Rewrite
Constraint

Introduce Lagrangian (objective):

$$L(x, \alpha) = x^2 - \alpha(x - b)$$

We will solve:

$$\begin{array}{ll} \min_x & \max_{\alpha} L(x, \alpha) \\ \text{s.t.} & \underline{\alpha} \geq 0 \end{array}$$

Add new
constraint

Can easily show that if x^*, α^* optimal for min-max problem, then x^* is optimal for original problem.

Dual SVM derivation (hard margin SVM)

Original optimization problem:

Minimize $\frac{1}{2} \mathbf{w} \bullet \mathbf{w}$

s.t. $(\mathbf{w} \bullet \mathbf{x}^{(i)} + b) y^{(i)} \geq 1$ for all i

Rewrite
constraints

One Lagrange multiplier
per example

Lagrangian:

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w} \bullet \mathbf{w} - \sum_i \alpha_i [(\mathbf{w} \bullet \mathbf{x}^{(i)} + b) y^{(i)} - 1]$$
$$\alpha_i \geq 0 \text{ for all } i.$$

Our goal now is to solve:

$$\begin{aligned} & \min_{\mathbf{w}, b} \max_{\alpha \geq 0} L(\mathbf{w}, \alpha) \\ &= \max_{\alpha \geq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, \alpha) \end{aligned}$$

Equivalent dual problem

$$\max_{\alpha \geq 0} \min_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w} \bullet \mathbf{w} - \sum_i \alpha_i [(\mathbf{w} \bullet \mathbf{x}^{(i)} + b) y^{(i)} - 1]$$

- Take the partial derivative of $L(\mathbf{w}, \alpha)$ w.r.t w_1, \dots, w_n, b and set to zero:

$$\mathbf{w} = \sum_i \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

$$\sum_i \alpha_i y^{(i)} = 0$$

- Substitute back in and simplify

$\begin{aligned} \text{(dual)} \quad & \max_{\alpha \geq 0} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)} \bullet \mathbf{x}^{(j)} \\ & \text{s.t. } \sum_i \alpha_i y^{(i)} = 0 \end{aligned}$

From dual to primal

$$\begin{aligned} \text{(dual)} \quad & \max_{\alpha \geq 0} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)} \bullet \mathbf{x}^{(j)} \\ & \text{s.t. } \sum_i \alpha_i y^{(i)} = 0 \end{aligned}$$

- After solving dual, can get optimal primal solution:

$$\mathbf{w} = \sum_i \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

$b =$ see Andrew Ng notes

Predict with $y \leftarrow \text{sign} [\mathbf{w} \bullet \mathbf{x} + b]$

$$= \text{sign} [\sum_i \alpha_i y^{(i)} (\mathbf{x} \bullet \mathbf{x}^{(i)}) + b]$$

Features of features

$$\phi(x) = \begin{pmatrix} x^{(1)} \\ \dots \\ x^{(n)} \\ x^{(1)}x^{(2)} \\ x^{(1)}x^{(3)} \\ \dots \\ e^{x^{(1)}} \\ \dots \end{pmatrix}$$

Feature vector: $\phi(\mathbf{x}^{(i)})$

Weight vector: \mathbf{w}

Feature vector and weight
vector super-high dimensional

Minimize $\|\mathbf{w}\|^2$

subject to:

$$y_t(\mathbf{w} \bullet \phi(\mathbf{x}^{(i)})) + b \geq 1 \text{ for all } i$$

Consider Dual

$$\phi(x) = \begin{pmatrix} x^{(1)} \\ \dots \\ x^{(n)} \\ x^{(1)}x^{(2)} \\ x^{(1)}x^{(3)} \\ \dots \\ e^{x^{(1)}} \\ \dots \end{pmatrix}$$

Maximize over $\alpha \geq 0$:

$$\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) := \Phi(\mathbf{x}^{(i)}) \bullet \Phi(\mathbf{x}^{(j)})$$

$$\text{s.t. } \sum_i \alpha_i y^{(i)} = 0$$

$$\mathbf{w} = \sum_i \alpha_i y^{(i)} \Phi(\mathbf{x}^{(i)})$$

Still looks hard, since we still have to deal with huge-dimensional $\Phi(\mathbf{x}^{(i)})$

Kernel trick: for some $\Phi(\mathbf{x})$, can evaluate kernel $K(\mathbf{x}, \mathbf{z})$ directly without evaluating $\Phi(\mathbf{x})$!

Example

Consider following feature of feature vector:

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1x_1 \\ x_1x_2 \\ x_1x_3 \\ x_2x_1 \\ x_2x_2 \\ x_2x_3 \\ x_3x_1 \\ x_3x_2 \\ x_3x_3 \end{bmatrix}.$$

$$\begin{aligned} \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) &= (\mathbf{x} \cdot \mathbf{z})^2 \\ &:= K(\mathbf{x}, \mathbf{z}) \end{aligned}$$

In this example, we increased the number of features from $O(n)$ to $O(n^2)$; but the dual is:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j (\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)})^2 \\ \text{s.t.} \quad & \sum_i \alpha_i y^{(i)} = 0 \end{aligned}$$

$$y \leftarrow \text{sign} \left[\sum_i \alpha_i y^{(i)} (\mathbf{x} \cdot \mathbf{x}^{(i)})^2 + b \right]$$

Only have to deal with feature vectors of length n ! Thus, obtain order of magnitude of new features & nonlinear decision boundary for free!

Story summary

- Suppose you want to use lots of non-linear features $\phi(\mathbf{x})$. Computationally impossible.
- But suppose $K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z})$ can be evaluated without explicitly evaluating $\phi(\mathbf{x})$.
- Then solving the dual for the α_i 's and then doing classification is tractable.
- Research problem: find “kernels” $K(\mathbf{x}, \mathbf{z})$ that are easy to evaluate and can be expressed as $K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z})$

Quadratic kernel

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z} + c)^2 = \left(\sum_{j=1}^n x^{(j)} z^{(j)} + c \right) \left(\sum_{\ell=1}^n x^{(\ell)} z^{(\ell)} + c \right) \\&= \sum_{j=1}^n \sum_{\ell=1}^n x^{(j)} x^{(\ell)} z^{(j)} z^{(\ell)} + 2c \sum_{j=1}^n x^{(j)} z^{(j)} + c^2 \\&= \sum_{j,\ell=1}^n (x^{(j)} x^{(\ell)}) (z^{(j)} z^{(\ell)}) + \sum_{j=1}^n (\sqrt{2c} x^{(j)}) (\sqrt{2c} z^{(j)}) + c^2,\end{aligned}$$

Feature mapping given by:

$$\Phi(\mathbf{x}) = [x^{(1)2}, x^{(1)}x^{(2)}, \dots, x^{(3)2}, \sqrt{2c}x^{(1)}, \sqrt{2c}x^{(2)}, \sqrt{2c}x^{(3)}, c]$$

Common kernels

- Polynomials of degree exactly d

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

- Polynomials of degree up to d

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

- Gaussian kernels

$$K(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|_2^2}{2\sigma^2}\right)$$

- And many others!

Soft SVM with kernels

Maximize:

$$\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

$$\text{s.t. } \sum_i \alpha_i y^{(i)} = 0$$

$$0 \leq \alpha_i \leq C \quad \text{for } i$$

$$K(\mathbf{x}, \mathbf{z}) := \boldsymbol{\Phi}(\mathbf{x}) \bullet \boldsymbol{\Phi}(\mathbf{z})$$

- Almost the same as hard-margin dual
- Once again, tractable for kernels that are easy to evaluate. Never compute features explicitly!!!
- $O(m^2)$ computation in size of dataset m to compute objective
 - much work on speeding up