

SVMs: Duality and Kernel Trick

Machine Learning - 10601

Geoff Gordon, Miroslav Dudík

[partly based on slides of Ziv-Bar Joseph]

<http://www.cs.cmu.edu/~ggordon/10601/>

November 18, 2009

SVMs as quadratic programs

Two optimization problems: For the separable and non separable cases

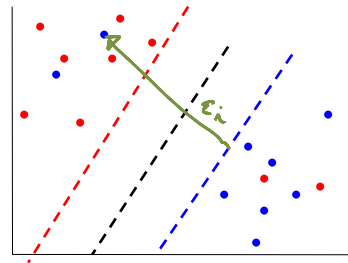
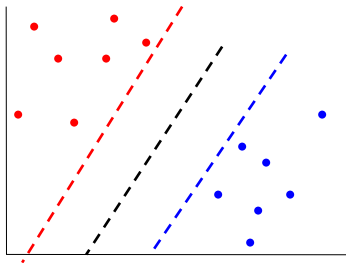
$$\min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$$

$$\underline{\underline{(\mathbf{w}^T \mathbf{x}_i + b) y_i \geq 1}}$$

$$\min_{\mathbf{w}, b, \varepsilon_i} \frac{\mathbf{w}^T \mathbf{w}}{2} + \sum_{i=1}^n C \varepsilon_i$$

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0$$



Dual for separable case

$$\min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$$

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \geq 0 \quad (\alpha_i)$$

\triangle # constraints
 $=$ # data pts
 $=$ # dual vars

Dual for separable case

$$\min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$$

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \geq 0 \quad (\alpha_i)$$

Lagrangian

$$\underline{\underline{L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{\mathbf{w}^T \mathbf{w}}{2} - \sum_i \alpha_i ((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1)}}$$

Dual for separable case

$$\min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$$

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \geq 0 \quad (\alpha_i)$$

EQUIVALENT

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} L(\mathbf{w}, b, \alpha)$$

=

$$\max_{\alpha \geq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$$

Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{\mathbf{w}^T \mathbf{w}}{2} - \sum_i \alpha_i ((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1)$$

constraint in \mathbf{w}, b
constraint in α

Dual for separable case

$$\min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$$

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \geq 0 \quad (\alpha_i)$$

$$\min_{\mathbf{w}, b} \max_{\alpha \geq 0} L(\mathbf{w}, b, \alpha)$$

optimality of α_i

for all i :

$$\alpha_i \geq 0$$

$$\alpha_i ((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1) = 0$$

MULT. COMBINE

Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{\mathbf{w}^T \mathbf{w}}{2} - \sum_i \alpha_i ((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1)$$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$$

optimality of \mathbf{w} and b

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i y_i$$

$$\sum_i \alpha_i y_i = 0$$

$$\left. \begin{array}{l} \mathbf{w} = \sum_i \alpha_i \mathbf{x}_i y_i \\ \sum_i \alpha_i y_i = 0 \end{array} \right\} \frac{\partial L}{\partial \mathbf{w}} = 0, \frac{\partial L}{\partial b} = 0$$

KARUSH
- KUHN-TUCKER
KKT
conditions

Dual for separable case

Dual formulation

$$\max_{\mathbf{a} \geq 0} \min_{\mathbf{w}, b} \left(\frac{\mathbf{w}^T \mathbf{w}}{2} - \sum_i \alpha_i \left((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \right) \right)$$

Optimality conditions (KKT conditions)

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i \mathbf{x}_i y_i \\ \sum_i \alpha_i y_i &= 0 \\ \alpha_i &\geq 0 \\ \alpha_i \left((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \right) &= 0 \end{aligned}$$

Dual for separable case

Dual formulation

$$\max_{\mathbf{a} \geq 0} \min_{\mathbf{w}, b} \left(\frac{\mathbf{w}^T \mathbf{w}}{2} - \sum_i \alpha_i \left((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \right) \right)$$

Optimality conditions (KKT conditions)

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i \mathbf{x}_i y_i \\ \sum_i \alpha_i y_i &= 0 \\ \alpha_i &\geq 0 \\ \alpha_i \left((\mathbf{w}^T \mathbf{x}_i + b) y_i - 1 \right) &= 0 \end{aligned}$$

$$\max_{\mathbf{a}} \left(\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right)$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

OPTIMALITY
OF \mathbf{w}, b

lower
bound

Dual for separable case

Dual formulation $\frac{1}{2} \alpha^T R \alpha$

$$\max_{\alpha} \left(\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right)$$

$\sum_i \alpha_i y_i = 0$

$\alpha_i \geq 0 \quad \forall i$

VARIABLES = # data pts
 # CONSTR = # data pts + 1
 =
 NEW POINT: x
 $y = \text{sign}(w^T x + b)$

Optimality conditions
 (KKT conditions)

$$w = \sum_i \alpha_i \mathbf{x}_i y_i$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

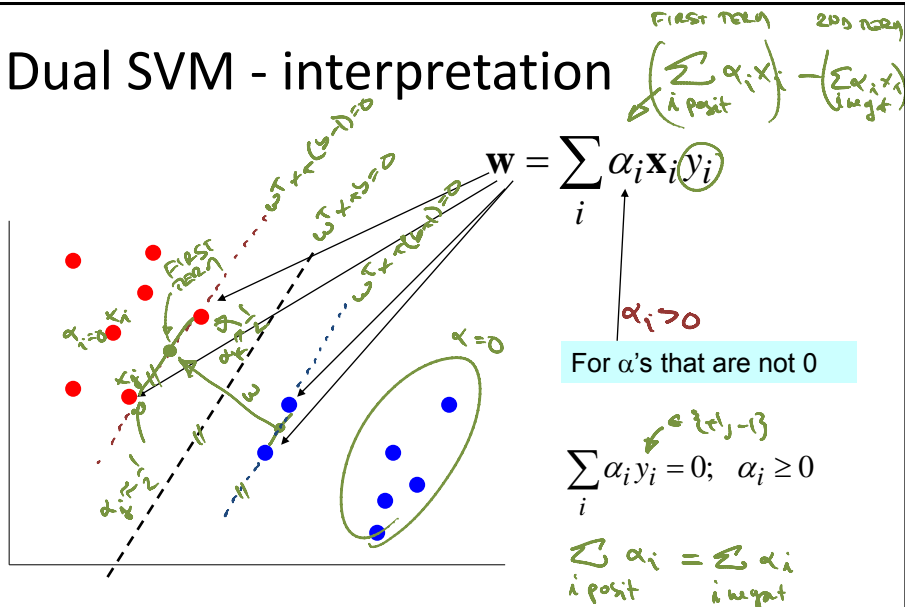
$$\alpha_i (w^T \mathbf{x}_i + b) y_i - 1 = 0$$

SOME $\alpha_i > 0$

$w^T \mathbf{x}_i + b = 1$

$b = 1 - w^T \mathbf{x}_i$

Dual SVM - interpretation



Dual SVM for linearly separable case

Our dual target function: $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$

ALL WE NEED ARE DOT PRODUCTS

$$\sum_i \alpha_i y_i = 0$$

Dot product across all pairs of training samples

$$\alpha_i \geq 0 \quad \forall i$$

Dot product with all training samples

To evaluate a new sample \mathbf{x} we need to compute:

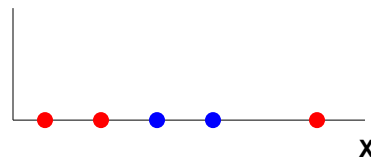
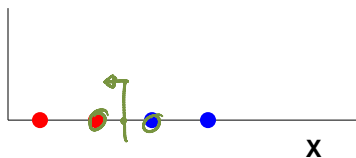
$$\underline{\mathbf{w}}^T \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

This might be too much work!
(e.g. when lifting \mathbf{x} into high dimensions)

Classifying in 1-d

Can an SVM correctly classify this data?

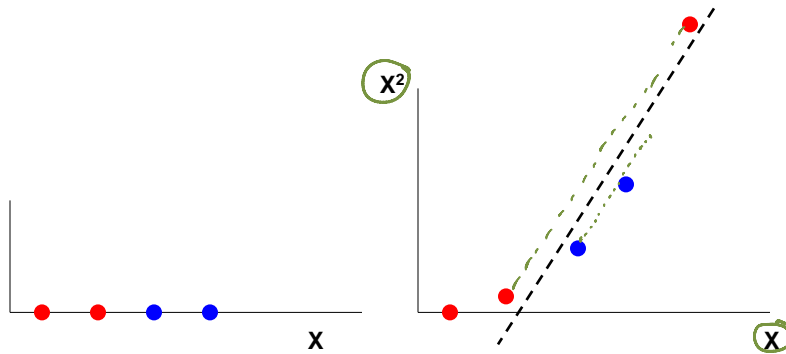
What about this?



Classifying in 1-d

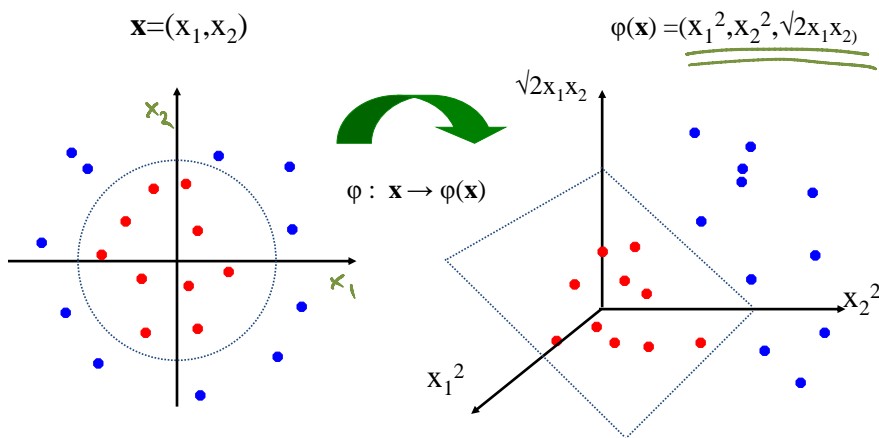
Can an SVM correctly classify this data?

And now?



Non-linear SVDs in 2-d

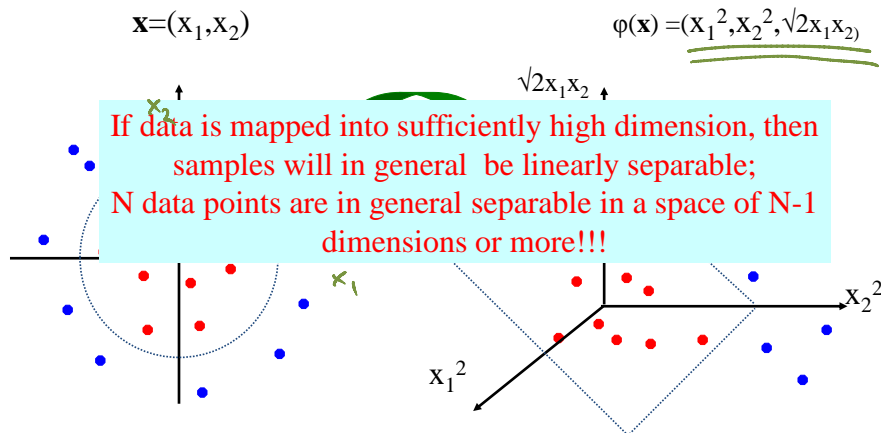
- The original input space (\mathbf{x}) can be mapped to some higher-dimensional feature space ($\varphi(\mathbf{x})$) where the training set is separable:



This slide is courtesy of www.iro.umontreal.ca/~pift6080/documents/papers/svm_tutorial.ppt

Non-linear SVDs in 2-d

- The original input space (\mathbf{x}) can be mapped to some higher-dimensional feature space ($\phi(\mathbf{x})$) where the training set is separable:

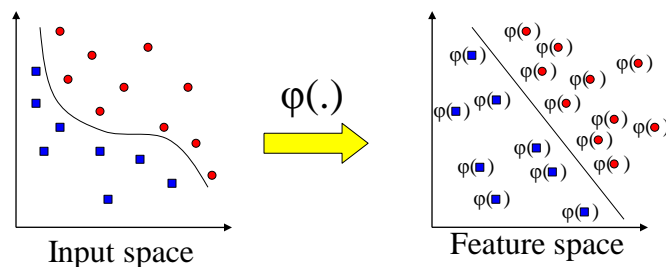


This slide is courtesy of www.iro.umontreal.ca/~pift6080/documents/papers/svm_tutorial.ppt

Transformation of Inputs

- Possible problems
 - High computation burden due to high-dimensionality
 - Many more parameters *possibly overfitting*
- SVM solves these two issues simultaneously
 - “Kernel tricks” for efficient computation
 - Dual formulation only assigns parameters to samples, not features

MARLINS: solve stat. problem



Polynomials of degree two

- While working in higher dimensions is beneficial, it also increases our running time because of the dot product computation

$$\max_{\alpha} \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

- However, there is a neat trick we can use

- consider all quadratic terms for $x_1, x_2 \dots x_m$

m is the number of features in each vector

The $\sqrt{2}$ term will become clear in the next slide

$$\phi(x) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \vdots \\ \sqrt{2}x_m \\ x_1^2 \\ \vdots \\ x_m^2 \\ \sqrt{2}x_1x_2 \\ \vdots \\ \sqrt{2}x_{m-1}x_m \end{pmatrix}$$

$m+1$ linear terms

m quadratic terms

$m(m-1)/2$ pairwise terms

$\sim m^2$

Dot product for polynomials of degree two

How many operations do we need for the dot product?

$$\phi(x)\phi(z) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \vdots \\ \sqrt{2}x_m \\ x_1^2 \\ \vdots \\ x_m^2 \\ \sqrt{2}x_1x_2 \\ \vdots \\ \sqrt{2}x_{m-1}x_m \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \sqrt{2}z_1 \\ \vdots \\ \sqrt{2}z_m \\ z_1^2 \\ \vdots \\ z_m^2 \\ \sqrt{2}z_1z_2 \\ \vdots \\ \sqrt{2}z_{m-1}z_m \end{pmatrix}$$

$$= \sum_i 2x_i z_i + \sum_i x_i^2 z_i^2 + \sum_{i,j=i+1} 2x_i x_j z_i z_j + 1$$

m m $m(m-1)/2$ $\sim m^2$

Polynomials of degree d in m variables

$$x_1, \dots, x_m$$

$$\varphi(x) = \left\{ \begin{array}{l} 1 \\ x_1 \\ \vdots \\ x_m \\ \text{quad. terms} \end{array} \right\} \sim m^2$$

$$\left\{ \begin{array}{l} \text{cubic terms} \\ \vdots \\ \text{d-th degree} \end{array} \right\} \sim m^d$$

Polynomials of degree d in m variables, n data points

Original formulation

$$\text{Min } (w^T w)/2$$

$$(w^T \varphi(x_i) + b) y_i \geq 1$$

\leftarrow dimensionality of φ is $O(m^d)$
constraints = n

Dual formulation

$$\max_{\alpha} \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j y_i y_j \varphi(x_i) \varphi(x_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

\leftarrow dimension of α is n
constraints = $n+1$

The kernel trick

How many operations do we need for the dot product?

$$\underbrace{\varphi(x)\varphi(z)}_{K(x,z)} = \underbrace{\sum_i 2x_i z_i}_m + \underbrace{\sum_i x_i^2 z_i^2}_m + \underbrace{\sum_{i,j=i+1} 2x_i x_j z_i z_j + 1}_{m(m-1)/2} \approx m^2 \quad \dots = \sim m^d$$

There's **structure** to this dot product... we can do this faster!

$$\begin{aligned} (xz+1)^2 &= (xz)^2 + 2(xz) + 1 \\ &= \left(\sum_i x_i z_i\right)^2 + \sum_i 2x_i z_i + 1 \\ &= \sum_i 2x_i z_i + \sum_i x_i^2 z_i^2 + \sum_{i,j=i+1} 2x_i x_j z_i z_j + 1 \end{aligned}$$

We only need **m** operations!

$\varphi(x) = \begin{pmatrix} \frac{1}{\sqrt{2}} \cdot x \\ \frac{1}{\sqrt{2}} \cdot x^2 \\ \vdots \end{pmatrix}$
 $(x \cdot z + 1)^d = \binom{d}{0} (x \cdot z)^0 + \binom{d}{1} (x \cdot z)^1 + \dots + \binom{d}{d} (x \cdot z)^d$

The kernel trick

How many operations do we need for the dot product?

$$\underbrace{\varphi(x)\varphi(z)}_{K(x,z)} = \underbrace{\sum_i 2x_i z_i}_m + \underbrace{\sum_i x_i^2 z_i^2}_m + \underbrace{\sum_{i,j=i+1} 2x_i x_j z_i z_j + 1}_{m(m-1)/2} \approx m^2$$

There's **structure** to this dot product... we can do this faster!

$$\begin{aligned} (xz+1)^2 &= (xz)^2 + 2(xz) + 1 \\ &= \left(\sum_i x_i z_i\right)^2 + \sum_i 2x_i z_i + 1 \\ &= \sum_i 2x_i z_i + \sum_i x_i^2 z_i^2 + \sum_{i,j=i+1} 2x_i x_j z_i z_j + 1 \end{aligned}$$

We only need **m** operations!

Note that to evaluate a new sample we are also using dot products so we save there as well

Where we are

[lifting into poly's of degree d]

Our dual target function:

$$\max_{\alpha_i} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

mn^2 operations to evaluate all coefficients

$O(m)$ operations per (i,j)

To evaluate a new sample \mathbf{x} we need to compute:

$$\mathbf{w}^T \phi(\mathbf{x}) + b = \sum_i \alpha_i y_i \phi(\mathbf{x}_i) \phi(\mathbf{x}) + b$$

mr operations where r is the number of support vectors ($\alpha_i > 0$)

$\phi(\mathbf{x}_i; \mathbf{x})$

Other kernels

• Beyond polynomials there are other very high dimensional basis functions that can be made practical by finding the right kernel function

- Radial-Basis Function:

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

- kernel functions for discrete objects (graphs, strings, etc.)

Kernels are measures of similarity

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

- LARGEST IF $\mathbf{x} = \mathbf{z}$
- DROPS AS $\|\mathbf{x} - \mathbf{z}\|$ INCREASES

Decision rule for a new sample \mathbf{x} :

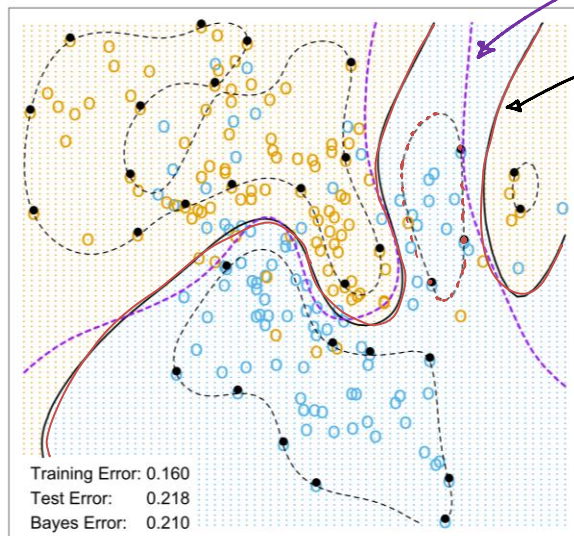
$$\mathbf{w}^T \phi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

[CAVEATS]
KERNELS MUST
DEFINE DOT
PRODUCTS

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x})$$

WE CARE
ABOUT THE SIGN

SVM - Radial Kernel



OPTIMAL
DEC. BOUNDARY

SVM
DEC. BOUNDARY

$$K(\mathbf{x}, \mathbf{z}) = \exp(\dots)$$

This slide is courtesy of Hastie-Tibshirani-Friedman, 2nd ed.

Dual formulation for non-separable case

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \alpha_i \geq 0 \quad \forall i$$

DEFOLE:
C → ∞

The only difference is that the α_i 's are now bounded

To evaluate a new sample \mathbf{x} we need to compute:

$$\mathbf{w}^T \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \mathbf{x} + b$$

Why do SVMs work?

• If we are using huge features spaces (with kernels) how come we are not overfitting the data?

- We maximize margin!

→ SIMPLE SOLUTIONS

- We minimize loss + regularization

$$C \left(\sum_i \text{hinge loss} \right) + \text{reg.}$$

↑
TRAIN ERR

Software

- A list of SVM implementation can be found at <http://www.kernel-machines.org/software.html>
- Some implementation (such as LIBSVM) can handle multi-class classification
- SVMLight is among one of the earliest implementation of SVM
- Several Matlab toolboxes for SVM are also available

Applications of SVMs

- Bioinformatics
- Machine Vision
- Text Categorization
- Ranking (e.g., Google searches)
- Handwritten Character Recognition
- Time series analysis

→ Lots of very successful applications!!!

Handwritten digit recognition



3-nearest-neighbor = 2.4% error

400–300–10 unit MLP = 1.6% error

LeNet: 768–192–30–10 unit MLP = 0.9% error

Current best (kernel machines, vision algorithms) \approx 0.6% error

Important points

- Difference between regression classifiers and SVMs'
- Maximum margin principle
- Target function for SVMs
- Linearly separable and non separable cases
- Dual formulation of SVMs
- Kernel trick and computational complexity