### Supervised Learning Methods

- k-nearest-neighbors
- Decision trees
- Neural networks
- Naïve Bayes
- Support vector machines (SVM)

## **Support Vector Machines**

Chapter 18.9 and the optional paper "Support vector machines" by M. Hearst, ed., 1998

Acknowledgments: These slides combine and modify ones provided by Andrew Moore (CMU), Carla Gomes (Cornell), Mingyue Tan (UBC), Jerry Zhu (Wisconsin), Glenn Fung (Wisconsin), and Olvi Mangasarian (Wisconsin)

13

cloud

ice

land

snow

water



15

12

### What are Support Vector Machines (SVMs) Used For?

- Classification
- Regression and data-fitting
- Supervised and unsupervised learning



16



## Linear Classifiers (aka Linear Discriminant Functions)

### • Definition:

A function that is a linear combination of the components of the input (column vector)  $\mathbf{x}$ :

$$f(\mathbf{x}) = \sum_{j=1}^{d} w_j x_j + b = \mathbf{w}^{\mathrm{T}} \mathbf{x} + b$$

where  ${\bf w}$  is the weight (column vector) and b is the bias

• A **2-class classifier** then uses the rule:

Decide class  $c_1$  if  $f(\mathbf{x}) \ge 0$  and class  $c_2$  if  $f(\mathbf{x}) < 0$ or, equivalently, decide  $c_1$  if  $\mathbf{w}^{\mathsf{T}}\mathbf{x} \ge -b$  and  $c_2$  otherwise













Maximum Margin  $f(\mathbf{x}, \mathbf{w}, b) = sign(\mathbf{w}^{\mathsf{T}} \mathbf{x} + b)$  denotes +1 The maximum denotes -1 margin linear classifier is the linear classifier Support Vector with the maximum are those data margin points that the margin pushes This is the against simplest kind of SVM (Called an LSVM) Linear SVM

























Given a new point **x**, we can classify it by

• Computing score:  $\mathbf{w}^{\mathsf{T}}\mathbf{x} + b$ 

Score > t: yes

Score < -t: no

Else: don't know

- Deciding class based on whether < 0 or > 0
- If desired, can set confidence threshold t



# SVM as Constrained Optimization

- Unknowns: w, b
- Objective function: maximize the margin:  $M = 2 / ||\mathbf{w}||$
- Equivalent to **minimizing**  $||\mathbf{w}||$  or  $||\mathbf{w}||^2 = \mathbf{w}^T \mathbf{w}$
- *N* training points:  $(\mathbf{x}_k, y_k), y_k = +1$  or -1
- Subject to each training point correctly classified, i.e.,

subject to  $y_k(\mathbf{w}^{\mathsf{T}}\mathbf{x}_k + b) \ge 1$  for all k < N constraints

This is a **quadratic optimization problem** (**QP**), which can be solved efficiently

42

### SVMs: More than Two Classes

- SVMs can only handle two-class problems
- *k*-class problem: Split the task into *k* **binary** tasks and learn *k* SVMs:
  - Class 1 vs. the rest (classes 2 k)
  - Class 2 vs. the rest (classes 1, 3 k)
  - ...
  - Class k vs. the rest
- Pick the class that puts the point *farthest into its positive region*







### Non Linearly-Separable Data

Approach 1: Allow a few points on the wrong side (**slack variables**)

"Soft Margin Classification"























Approach 2: Map data to a higher dimensional space, and then do linear classification there (called the kernel trick)



64



![](_page_11_Figure_1.jpeg)

![](_page_11_Figure_2.jpeg)

![](_page_11_Figure_3.jpeg)

![](_page_12_Figure_0.jpeg)

![](_page_12_Figure_2.jpeg)

![](_page_12_Figure_3.jpeg)

### **Improving Efficiency**

- Time complexity of the original optimization formulation depends on the dimensionality, *k*, of z (*k* >> *d*)
- We can convert the optimization problem into an equivalent form, called the *Dual Form*, with time complexity *O*(*N*<sup>3</sup>) that depends on *N*, the number of training examples, *not k*
- Dual Form will also allow us to rewrite the mapping functions in Φ in terms of "kernel functions" instead

![](_page_13_Figure_0.jpeg)

### Algorithm

- Compute N x N matrix Q by computing y<sub>i</sub> y<sub>j</sub> (x<sub>i</sub><sup>T</sup> x<sub>j</sub>) between all pairs of training examples
- Solve the optimization problem to compute α<sub>i</sub>

for *i* = 1, ..., N

- Each non-zero *α<sub>i</sub>* indicates that example **x**<sub>i</sub> is a support vector
- Compute  $\boldsymbol{w}$  and  $\boldsymbol{b}$
- Then classify test example **x** with:

$$f(\mathbf{x}) = sign(\mathbf{w}^{\mathsf{T}} \mathbf{x} - \mathbf{b})$$

![](_page_13_Figure_10.jpeg)

![](_page_13_Figure_12.jpeg)

![](_page_14_Figure_0.jpeg)

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

$$\Phi(x_i)^{\mathrm{T}} \cdot \Phi(x_j)$$
  

$$= x_{i_1}^2 x_{j_1}^2 + \sqrt{2}x_{i_1} x_{i_2} \sqrt{2}x_{j_1} x_{j_2} + x_{i_2}^2 x_{j_2}^2$$
  

$$= x_{i_1}^2 x_{j_1}^2 + 2x_{i_1} x_{i_2} x_{j_1} x_{j_2} + x_{i_2}^2 x_{j_2}^2$$
  

$$= (x_{i_1} x_{j_1} + x_{i_2} x_{j_2})^2$$
  

$$= (x_i^{\mathrm{T}} \cdot x_j)^2$$

Dual formulation of the optimization problem depends on the input data only in dot products of the form: Φ(**x**<sub>i</sub>)<sup>T</sup> · Φ(**x**<sub>j</sub>) where **x**<sub>i</sub> and **x**<sub>j</sub> are two examples
We can compute these dot products efficiently for certain types of Φ's where *k*(**x**<sub>i</sub>, **x**<sub>i</sub>) = Φ(**x**<sub>i</sub>)<sup>T</sup> · Φ(**x**<sub>j</sub>)

• Example:

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

$$\Phi(\mathbf{x}_i)^{\mathsf{T}} \cdot \Phi(\mathbf{x}_j) = (\mathbf{x}_i^{\mathsf{T}} \cdot \mathbf{x}_j)^2 = K(\mathbf{x}_i, \mathbf{x}_j)$$

Since the data *only* appears as dot products, we do *not* need to map the data to higher dimensional space (using Φ(**x**) because we can use the kernel function *K* instead

82

### **Kernel Functions**

- A kernel, *K*(**x**<sub>i</sub>, **x**<sub>j</sub>), is a dot product in *some* feature space
- A kernel function is a function that can be applied to pairs of input examples to evaluate dot products in some corresponding (possibly infinite dimensional) feature space
- We do *not* need to compute  $\Phi$  explicitly

![](_page_15_Figure_0.jpeg)

### Some Commonly Used Kernels

- Linear kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j$
- Quadratic kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j + 1)^2$
- Polynomial of degree *d* kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^d$
- Radial-Basis Function (Gaussian) kernel:  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-||\mathbf{x}_i \cdot \mathbf{x}_j||^2 / \sigma^2)$
- Many possible kernels; picking a good one is tricky
- Hacking with SVMs: create various kernels, hope their space  $\Phi$  is meaningful, plug them into SVM, pick one with good classification accuracy
- Kernel usually combined with slack variables because no guarantee of linear separability in new space

86

### Algorithm

- Compute N x N matrix Q by computing
   y<sub>i</sub> y<sub>i</sub> K(x<sub>i</sub>, x<sub>j</sub>) between all pairs of training points
- Solve optimization problem to compute α<sub>i</sub>
  - for *i* = 1, ..., *N*
  - Each non-zero *α<sub>i</sub>* indicates that example **x**<sub>i</sub> is a support vector
- Compute **w** and b
- Classify test example **x** using:

$$f(\mathbf{x}) = sign(\mathbf{w}^{\mathsf{T}} \mathbf{x} - b)$$

### Applications of SVMs

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

 $\rightarrow$  Lots of very successful applications!

![](_page_16_Figure_0.jpeg)

![](_page_16_Figure_2.jpeg)

### Example Application: The Federalist Papers Dispute

- Written in 1787-1788 by Alexander Hamilton, John Jay, and James Madison to persuade the citizens of New York to ratify the U.S. Constitution
- Papers consisted of short essays, 900 to 3500 words in length
- Authorship of 12 of those papers have been in dispute (Madison or Hamilton); these papers are referred to as the disputed Federalist papers

70-Word Dictionary									
1	a	15	do	29	is	43	or	57	this
2	all	16	down	30	it	44	our	58	to
3	also	17	even	31	its	45	shall	59	up
4	an	18	every	32	may	46	should	60	upon
5	and	19	for	33	more	47	so	61	was
6	any	20	from	34	must	48	some	62	were
7	are	21	had	35	$_{my}$	49	such	63	what
8	as	22	has	36	no	50	than	64	when
9	at	23	have	37	not	51	that	65	which
10	be	24	her	38	now	52	the	66	who
11	been	25	his	39	of	53	their	67	will
12	but	26	if	40	on	54	then	68	with
13	by	27	in	41	one	55	there	69	would
14	can	28	into	42	only	56	things	70	your

### Feature Selection for Classifying the Disputed Federalist Papers

- Apply the SVM algorithm for feature selection to:
  - Train on the 106 Federalist papers with known authors
  - Find a classification hyperplane (LSVM) that uses as few words as possible
- Use the hyperplane to classify the 12 disputed papers

97

![](_page_17_Figure_6.jpeg)

### Hyperplane Classifier Using 3 Words

• A hyperplane depending on three words was found:

0.537*to* + 24.663*upon* + 2.953*would* = 66.616

• All disputed papers ended up on the Madison side of the plane

98

### Summary

- Learning linear functions
  - Pick *separating hyperplane* that maximizes margin
  - Separating plane defined in terms of support vectors (small number of training examples) *only*
- Learning non-linear functions
  - Project examples into a higher dimensional space
  - Use kernel functions for efficiency
- Generally avoids overfitting problem
- Global optimization method; no local optima
- Can be expensive to apply, especially for multi-class problems
- Biggest Drawback: The choice of kernel function
  - There is no "set-in-stone" theory for choosing a kernel function for any given problem
  - Once a kernel function is chosen, there is only ONE modifiable parameter, the error penalty  ${\ensuremath{\mathcal{C}}}$

### Software

- A list of SVM implementations can be found at http://www.kernelmachines.org/software.html
- Some implementations (such as LIBSVM) can handle multi-class classification
- SVMLight is one of the earliest and most frequently used implementations of SVMs
- Several Matlab toolboxes for SVMs are available