

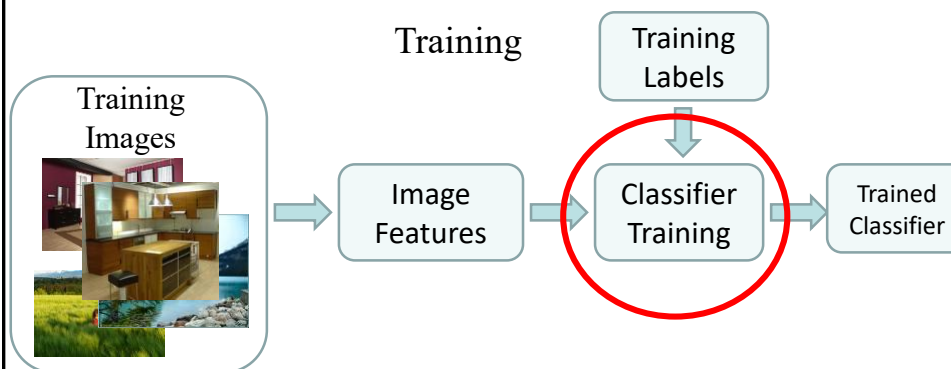
# ECS797P - Machine Learning for Visual Data Analysis

*Tao Xiang*

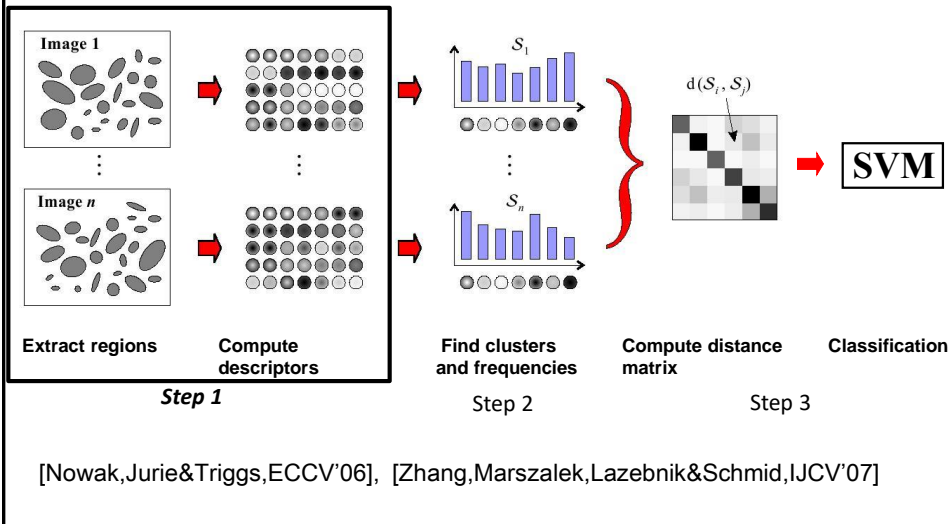
*t.xiang@qmul.ac.uk*

*Lecture 2*

## Part 2: Classifiers

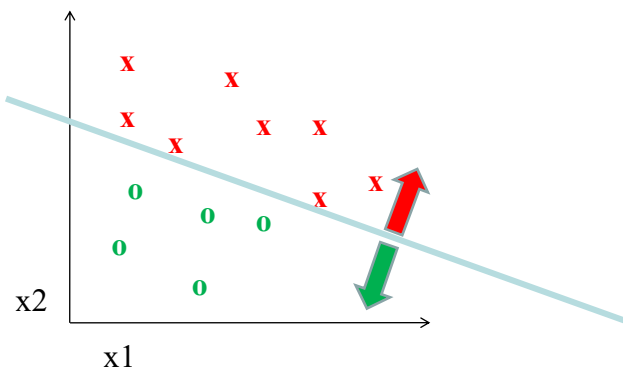


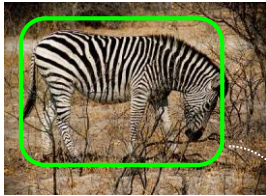
# Bag-of-features for image classification



## Learning a classifier

Given some set of features with corresponding labels, learn a function to predict the labels from the features





## Classifier-based methods



### Many classifiers to choose from

- SVM
- Neural networks
- Naïve Bayes
- Bayesian network
- Logistic regression
- Randomized Forests
- Boosted Decision Trees
- K-nearest neighbor
- RBMs
- Etc.

Which is the best one?

## The perfect classification algorithm

- Objective function: encodes the right loss for the problem
- Parameterization: makes assumptions that fit the problem
- Regularization: right level of regularization for amount of training data
- Training algorithm: can find parameters that maximize objective on training set
- Inference algorithm: can solve for objective function in evaluation

## Generative vs. Discriminative Classifiers

### Generative

- Training
  - Models the data and the labels
  - Assume (or learn) probability distribution and dependency structure
  - Can impose priors
- Testing
  - $P(y=1, x) / P(y=0, x) > t$ ?
- Examples
  - Foreground/background GMM
  - Naïve Bayes classifier
  - Bayesian network

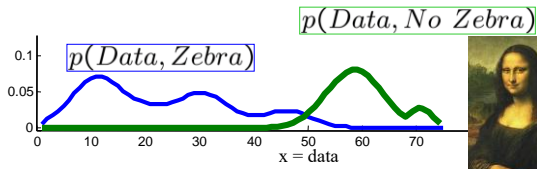
### Discriminative

- Training
  - Learn to directly predict the labels from the data
  - Assume form of boundary
  - Margin maximization or parameter regularization
- Testing
  - $f(x) > t$  ; e.g.,  $w^T x > t$
- Examples
  - Logistic regression
  - SVM
  - Boosted decision trees

## Discriminative vs. generative

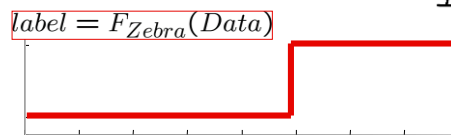
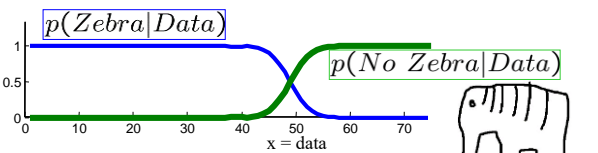
- Generative model

(The artist)

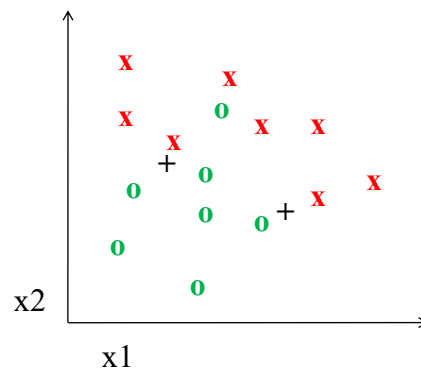


- Discriminative model

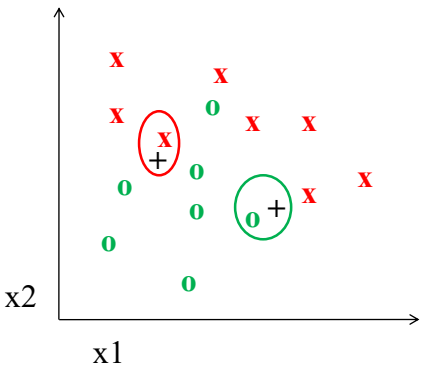
(The lousy painter)



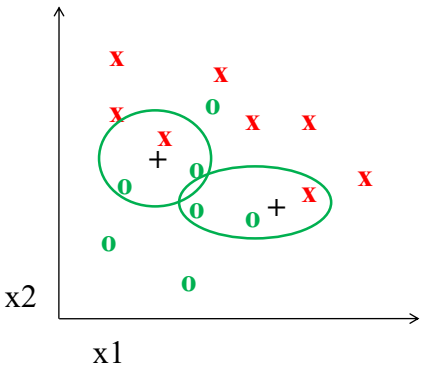
## K-nearest neighbor



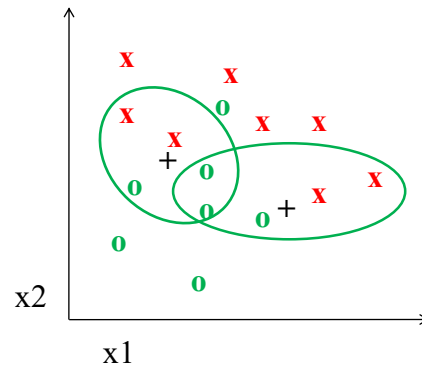
# 1-nearest neighbor



# 3-nearest neighbor



## 5-nearest neighbor



What is the parameterization? The regularization? The training algorithm? The inference?

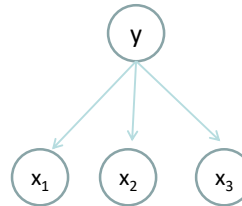
Is K-NN generative or discriminative?

## Using K-NN

- Simple, and a good one to try first
- With infinite examples, 1-NN provably has error that is at most twice Bayes optimal error

# Naïve Bayes

- Objective
- Parameterization
- Regularization
- Training
- Inference



[Naïve Bayes Classifier wikipedia link](#)

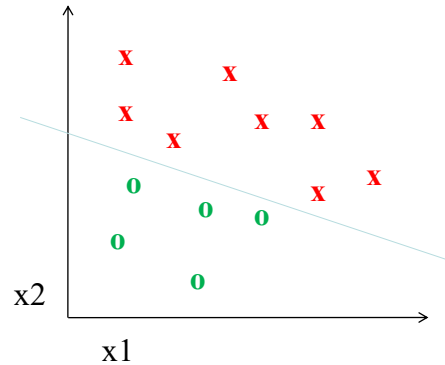
## Using Naïve Bayes

- Simple thing to try for categorical data
- Very fast to train/test



## Classifiers: Logistic Regression

- Objective
- Parameterization
- Regularization
- Training
- Inference



[Logistic Regression wikipedia link](#)

[A document on how to use Logistic Regression as a classifier](#)

## Using Logistic Regression

- Quick, simple classifier (try it first)
- Use L2 or L1 regularization
  - L1 does feature selection and is robust to irrelevant features but slower to train

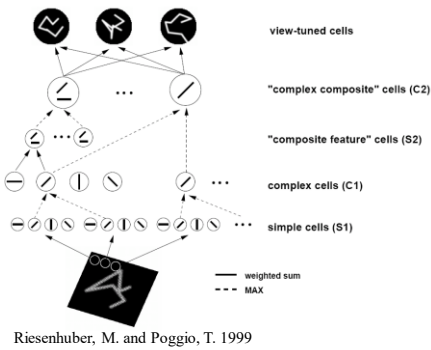
# Classifier: Neural Networks

Fukushima's Neocognitron, 1980

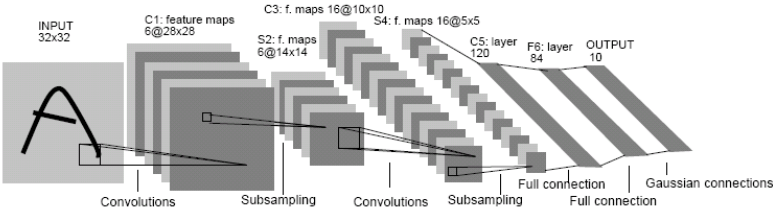
Rowley, Baluja, Kanade 1998

LeCun, Bottou, Bengio, Haffner 1998

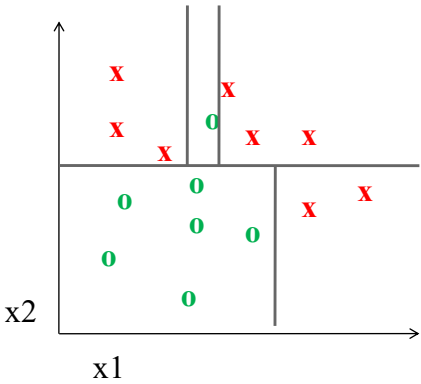
Serre et al. 2005



LeNet convolutional architecture (LeCun 1998)



# Classifiers: Decision Trees



[Decision Tree wikipedia link](#)

# Ensemble Methods: Boosting

## Discrete AdaBoost(Freund & Schapire 1996b)

1. Start with weights  $w_i = 1/N$ ,  $i = 1, \dots, N$ .
2. Repeat for  $m = 1, 2, \dots, M$ :
  - (a) Fit the classifier  $f_m(x) \in \{-1, 1\}$  using weights  $w_i$  on the training data.
  - (b) Compute  $\text{err}_m = E_w[1_{(y \neq f_m(x))}]$ ,  $c_m = \log((1 - \text{err}_m)/\text{err}_m)$ .
  - (c) Set  $w_i \leftarrow w_i \exp[c_m \cdot 1_{(y_i \neq f_m(x_i))}]$ ,  $i = 1, 2, \dots, N$ , and renormalize so that  $\sum_i w_i = 1$ .
3. Output the classifier  $\text{sign}[\sum_{m=1}^M c_m f_m(x)]$

figure from Friedman et al. 2000

## Classifier: Boosting

### Viola & Jones 2001

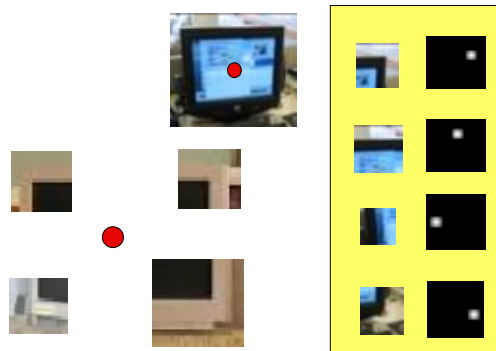
Haar features via Integral Image  
Cascade  
Real-time performance

.....



### Torralla et al., 2004

Part-based Boosting  
Each weak classifier is a part  
Part location modeled by  
offset mask



## Two ways to think about classifiers

1. What is the objective? What are the parameters? How are the parameters learned? How is the learning regularized? How is inference performed?
2. How is the data modeled? How is similarity defined? What is the shape of the boundary?

Comparison			
	Learning Objective	Training	Inference
Naïve Bayes	$\text{maximize } \sum_i \left[ \sum_j \log P(x_{ij}   y_i; \theta_j) \right] + \log P(y_i; \theta_0)$	$\theta_{kj} = \frac{\sum_i \delta(x_{ij} = 1 \wedge y_i = k) + r}{\sum_i \delta(y_i = k) + Kr}$	$\theta_1^T \mathbf{x} + \theta_0 > 0$ <p>assuming <math>\mathbf{x}</math> in <math>\{0, 1\}</math></p> <p>where <math>\theta_{kj} = \log \frac{P(x_{ij} = 1   y = 1)}{P(x_{ij} = 1   y = 0)}</math>  <math>\theta_{0j} = \log \frac{P(x_{ij} = 0   y = 1)}{P(x_{ij} = 0   y = 0)}</math></p>
Logistic Regression	$\text{maximize } \sum_i \log(P(y_i   \mathbf{x}, \boldsymbol{\theta})) + \lambda \ \boldsymbol{\theta}\ $ <p>where <math>P(y_i   \mathbf{x}, \boldsymbol{\theta}) = 1 / (1 + \exp(-y_i \boldsymbol{\theta}^T \mathbf{x}))</math></p>	Gradient ascent	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Linear SVM	$\text{minimize } \lambda \sum_i \xi_i + \frac{1}{2} \ \boldsymbol{\theta}\ ^2$ <p>such that <math>y_i \boldsymbol{\theta}^T \mathbf{x} \geq 1 - \xi_i \quad \forall i</math></p>	Linear programming	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Kernelized SVM	complicated to write	Quadratic programming	$\sum_i y_i \alpha_i K(\hat{\mathbf{x}}_i, \mathbf{x}) > 0$
Nearest Neighbor	most similar features $\rightarrow$ same label	Record data	$y_i$ <p>where <math>i = \operatorname{argmin}_i K(\hat{\mathbf{x}}_i, \mathbf{x})</math></p>

## What to remember about classifiers

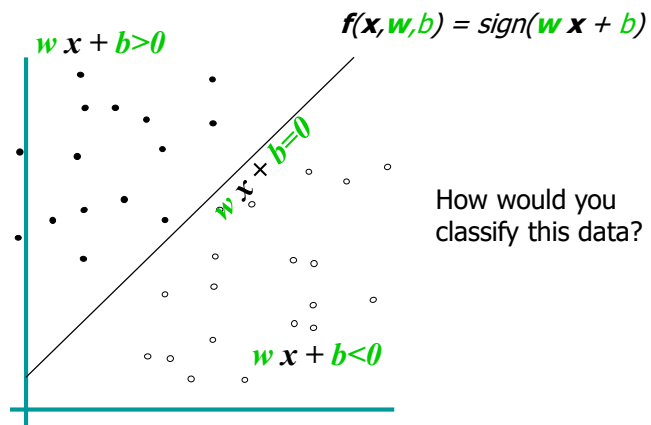
- No free lunch: machine learning algorithms are tools, not dogmas
- Try simple classifiers first
- Better to have smart features and simple classifiers than simple features and smart classifiers
- Use increasingly powerful classifiers with more training data (bias-variance tradeoff)

## **SUPPORT VECTOR MACHINES**

## Linear Classifiers



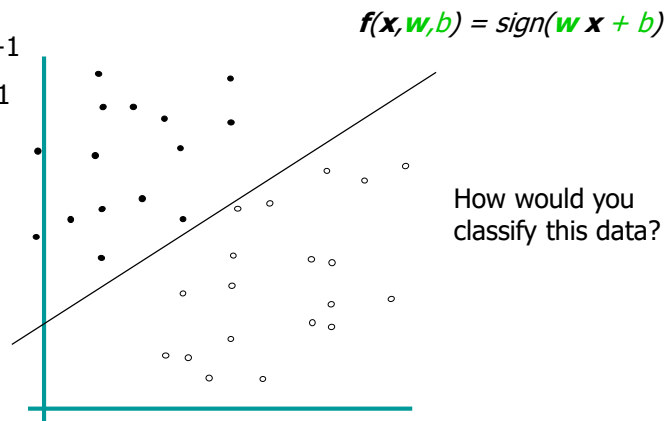
- denotes +1
- denotes -1



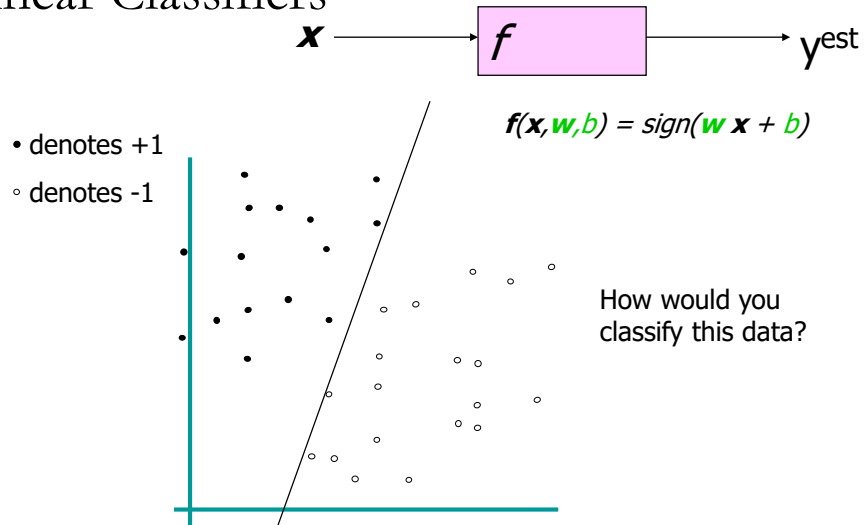
## Linear Classifiers



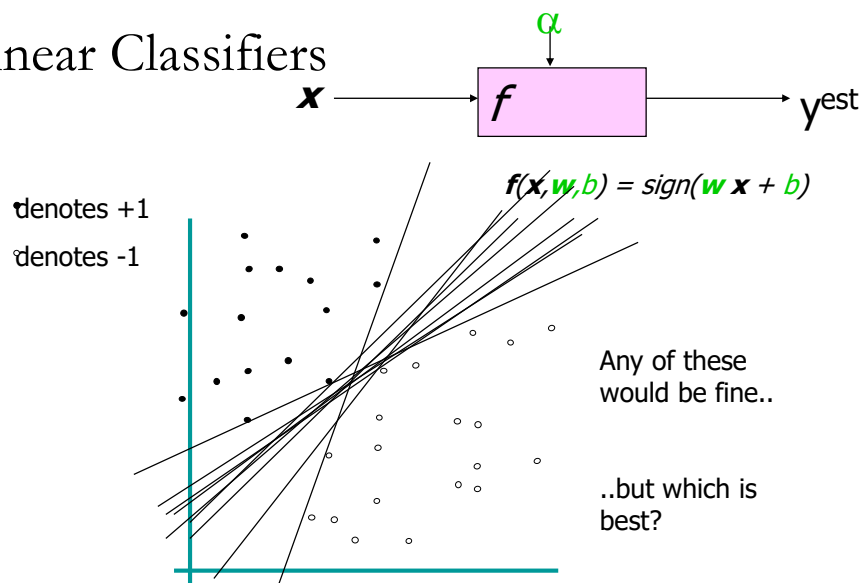
- denotes +1
- denotes -1

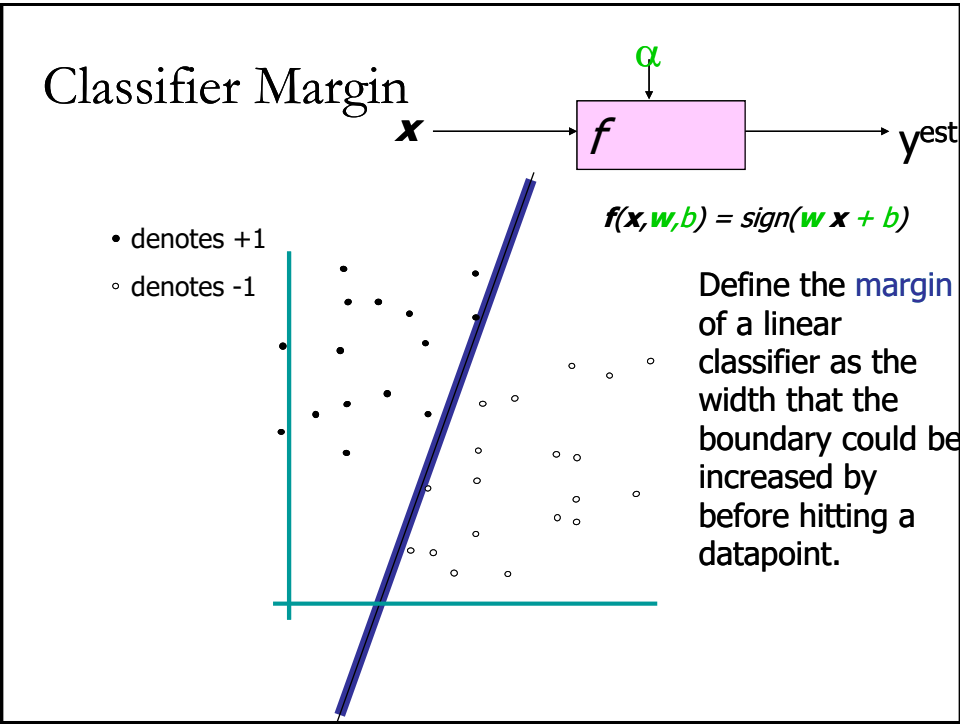
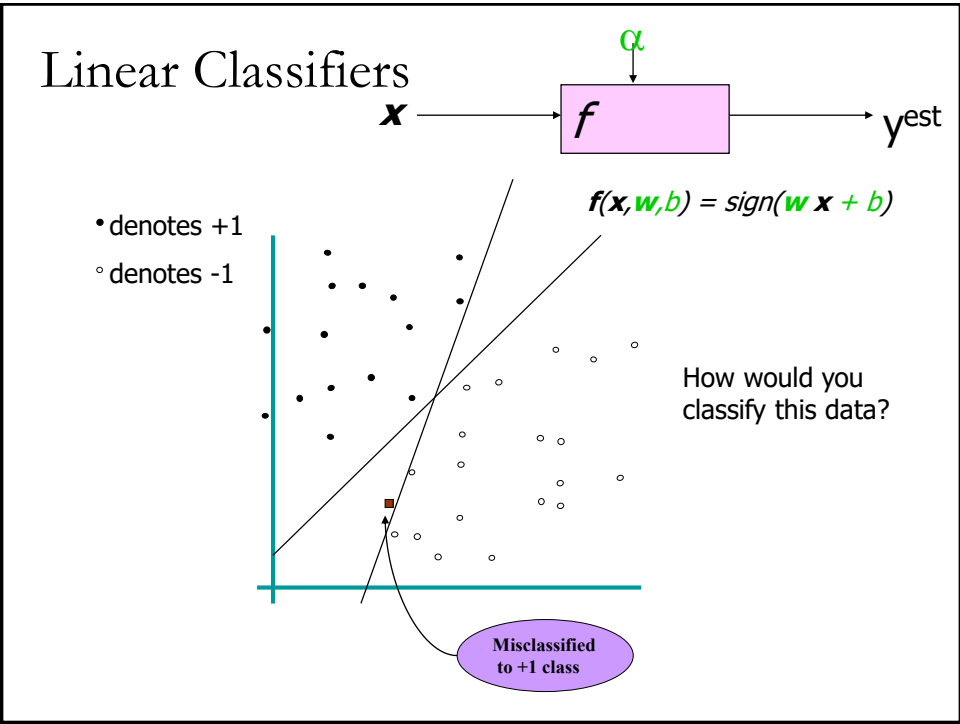


## Linear Classifiers

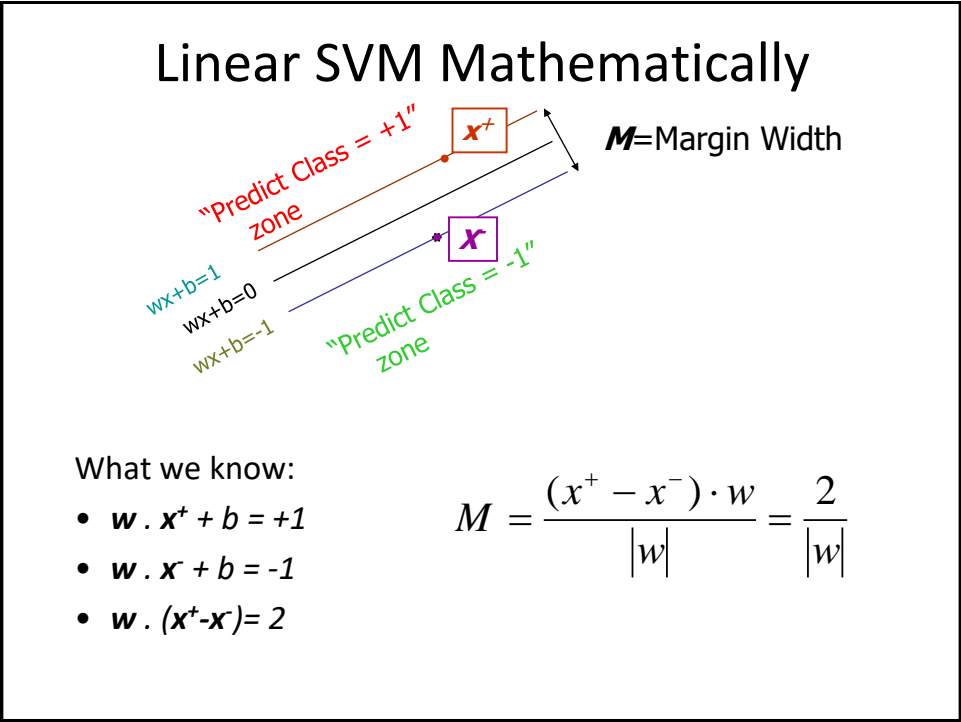
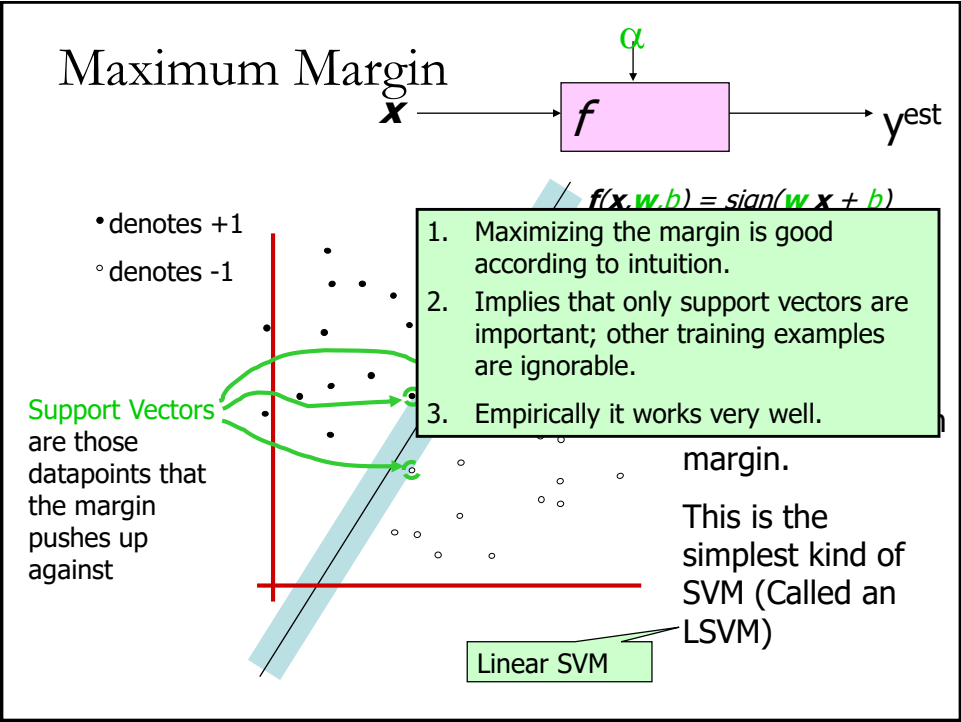


## Linear Classifiers









## Linear SVM Mathematically

- Goal: 1) Correctly classify all training data

$$\begin{aligned}
 & wx_i + b \geq 1 \text{ if } y_i = +1 \\
 & wx_i + b \leq -1 \text{ if } y_i = -1 \\
 & y_i(wx_i + b) \geq 1 \text{ for all } i
 \end{aligned}
 \left. \vphantom{\begin{aligned} & wx_i + b \geq 1 \text{ if } y_i = +1 \\ & wx_i + b \leq -1 \text{ if } y_i = -1 \end{aligned}} \right\} \curvearrowright$$

- 2) Maximize the Margin  
same as minimize

$$M = \frac{2}{|w|}$$

$$\frac{1}{2} w^T w$$

- We can formulate a Quadratic Optimization Problem and solve for  $w$  and  $b$

$$\begin{aligned}
 & \text{Minimize } \Phi(w) = \frac{1}{2} w^T w \\
 & \text{subject to } y_i(wx_i + b) \geq 1 \quad \forall i
 \end{aligned}$$

## Solving the Optimization Problem

Find  $w$  and  $b$  such that  
 $\Phi(w) = \frac{1}{2} w^T w$  is minimized;  
 and for all  $\{(\mathbf{x}_i, y_i)\}$ :  $y_i(w^T \mathbf{x}_i + b) \geq 1$

- Need to optimize a *quadratic* function subject to *linear* constraints.
- Quadratic optimization problems are a well-known class of mathematical programming problems, and many (rather intricate) algorithms exist for solving them.
- The solution involves constructing a *dual problem* where a *Lagrange multiplier*  $\alpha_i$  is associated with every constraint in the primary problem:

Find  $\alpha_1 \dots \alpha_N$  such that  
 $Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and  
 (1)  $\sum \alpha_i y_i = 0$   
 (2)  $\alpha_i \geq 0$  for all  $\alpha_i$

## The Optimization Problem Solution

- The solution has the form:

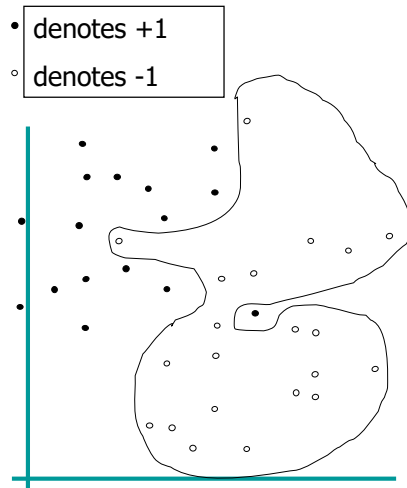
$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_k - \mathbf{w}^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ such that } \alpha_k \neq 0$$

- Each non-zero  $\alpha_i$  indicates that corresponding  $\mathbf{x}_i$  is a support vector.
- Then the classifying function will have the form:

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point  $\mathbf{x}$  and the support vectors  $\mathbf{x}_i$  – we will return to this later.
- Also keep in mind that solving the optimization problem involved computing the inner products  $\mathbf{x}_i^T \mathbf{x}_j$  between all pairs of training points.

## Dataset with noise

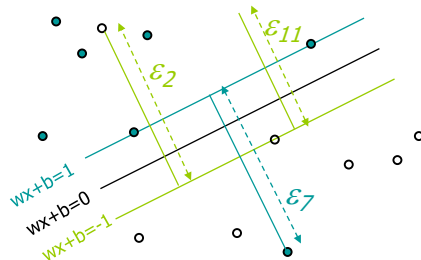


- **Hard Margin:** So far we require all data points be classified correctly
  - No training error
- **What if the training set is noisy?**
  - **Solution 1:** use very powerful kernels

Overfitting?

# Soft Margin Classification

*Slack variables  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.*



What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \varepsilon_k$$

## Hard Margin v.s. Soft Margin

- **The old formulation:**

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$   
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

- **The new formulation incorporating slack variables:**

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$   
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$  and  $\xi_i \geq 0$  for all  $i$

- **Parameter  $C$  can be viewed as a way to control overfitting.**

## Linear SVMs: Overview

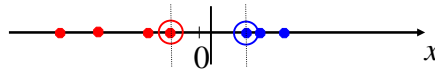
- The classifier is a *separating hyperplane*.
- Most “important” training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points  $\mathbf{x}_i$  are support vectors with non-zero Lagrangian multipliers  $\alpha_i$ .
- Both in the dual formulation of the problem and in the solution training points appear only inside dot products:

Find  $\alpha_1 \dots \alpha_N$  such that  
 $Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and  
 (1)  $\sum \alpha_i y_i = 0$   
 (2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$

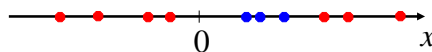
$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

## Non-linear SVMs

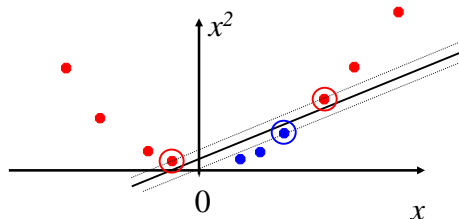
- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?

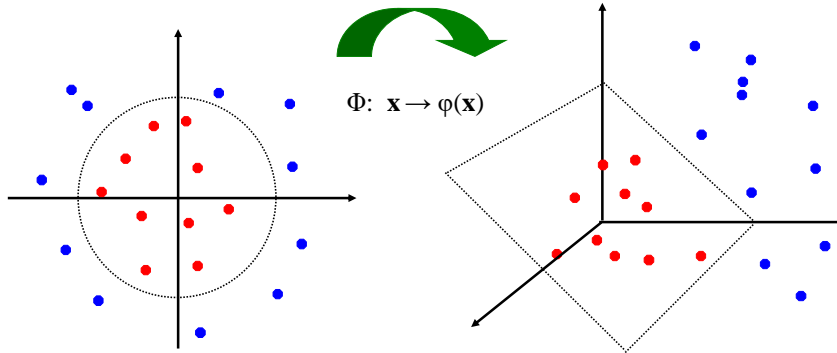


- How about... mapping data to a higher-dimensional space:



## Non-linear SVMs: Feature spaces

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



## The “Kernel Trick”

- The linear classifier relies on dot product between vectors  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- If every data point is mapped into high-dimensional space via some transformation

$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$ , the dot product becomes:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$$

- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.

## What Functions are Kernels?

- For some functions  $K(\mathbf{x}_i, \mathbf{x}_j)$  checking that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \text{ can be cumbersome.}$$

- Mercer's theorem:

*Every semi-positive definite symmetric function is a kernel*

- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

$$K = \begin{array}{|c|c|c|c|c|} \hline K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & K(\mathbf{x}_1, \mathbf{x}_3) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \hline K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & K(\mathbf{x}_2, \mathbf{x}_3) & & K(\mathbf{x}_2, \mathbf{x}_N) \\ \hline \dots & \dots & \dots & \dots & \dots \\ \hline K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & K(\mathbf{x}_N, \mathbf{x}_3) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \\ \hline \end{array}$$

## Examples of Kernel Functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial of power  $p$ :  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
- Gaussian (radial-basis function network):

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

- Sigmoid:  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T \mathbf{x}_j + \beta_1)$

## Non-linear SVMs Mathematically

- **Dual problem formulation:**

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j)$  is maximized and

(1)  $\sum \alpha_i y_i = 0$

(2)  $\alpha_i \geq 0$  for all  $\alpha_i$

- **The solution is:**

$$f(x) = \sum \alpha_i y_i K(x_i, x) + b$$

- **Optimization techniques for finding  $\alpha_i$ 's remain the same!**

## Nonlinear SVM - Overview

- SVM locates a separating hyperplane in the feature space and classify points in that space
- It does not need to represent the space explicitly, simply by defining a kernel function
- The kernel function plays the role of the dot product in the feature space.



## Properties of SVM

- **Flexibility in choosing a similarity function**
- **Sparseness of solution when dealing with large data sets**
  - only support vectors are used to specify the separating hyperplane
- **Ability to handle large feature spaces**
  - complexity does not depend on the dimensionality of the feature space
- **Overfitting can be controlled by soft margin approach**
- **Nice math property:** a simple convex optimization problem which is guaranteed to converge to a single global solution

## Some Issues

- **Choice of kernel**
  - Gaussian or polynomial kernel is default
  - if ineffective, more elaborate kernels are needed
  - domain experts can give assistance in formulating appropriate similarity measures
- **Choice of kernel parameters**
  - e.g.  $\sigma$  in Gaussian kernel
  - $\sigma$  is the distance between closest points with different classifications
  - In the absence of reliable criteria, applications rely on the use of a validation set or cross-validation to set such parameters.

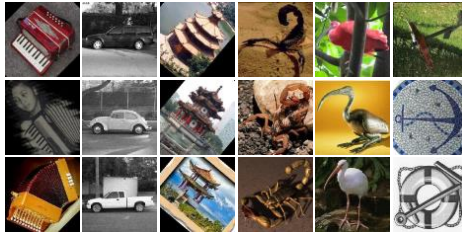
## References

- **An excellent tutorial on VC-dimension and Support Vector Machines:**  
C.J.C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):955-974, 1998.
- **The VC/SRM/SVM Bible:**  
Statistical Learning Theory by Vladimir Vapnik, Wiley-Interscience; 1998
- **The most widely used code library:**  
<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- **A good link:**  
<http://www.kernel-machines.org/>

## Summary: SVMs for image classification

1. Pick an image representation (in our case, bag of features)
2. Pick a kernel function for that representation
3. Compute the matrix of kernel values between every pair of training examples
4. Feed the kernel matrix into your favorite SVM solver to obtain support vectors and weights
5. At test time: compute kernel values for your test example and each support vector, and combine them with the learned weights to get the value of the decision function

## Category classification – CalTech101

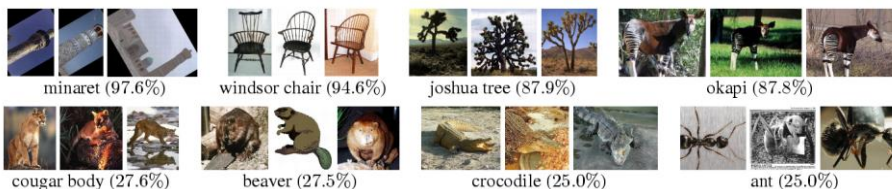


L	Single-level	Pyramid
0(1x1)	41.2±1.2	
1(2x2)	55.9±0.9	57.0 ±0.8
2(4x4)	63.6±0.9	64.6 ±0.8
3(8x8)	60.3±0.9	64.6 ±0.7

Bag-of-features approach by Zhang et al.'07: 54 %

## CalTech101

### Easiest and hardest classes



- Sources of difficulty:
  - Lack of texture
  - Camouflage
  - Thin, articulated limbs
  - Highly deformable shape

# Evaluation of image classification

- Averaged classification accuracy
- Confusion matrix

Actual class	Predicted class		
	Cat	Dog	Rabbit
	Cat	5	3
	Dog	2	3
	Rabbit	0	2

