

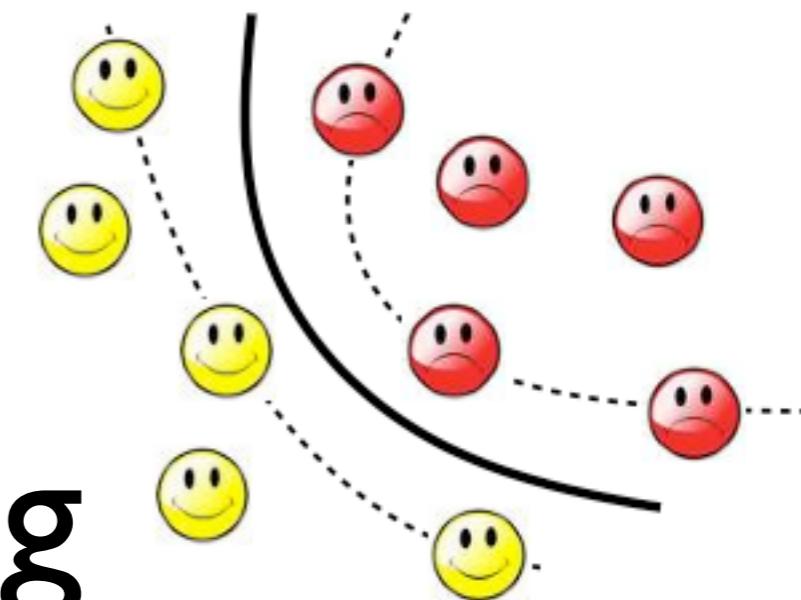


THE UNIVERSITY OF  
SYDNEY

# Machine Learning and Data Mining (COMP 5318)

Logistic Regression and SVMs

Tongliang Liu

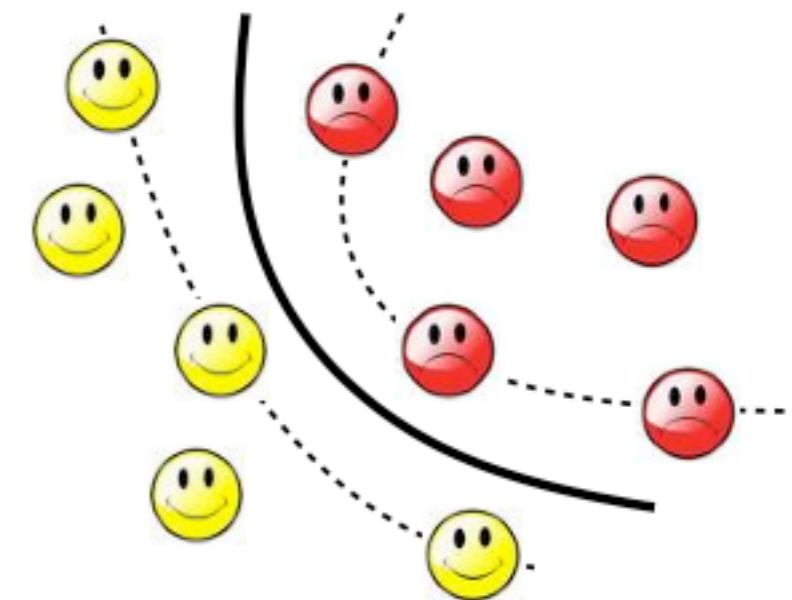




THE UNIVERSITY OF  
SYDNEY

# Announcements

- This lecture is based on:
  - Murphy's book: Chapters 8, 14
  - Bishop's book: Chapters 4, 7
  - Ullman's book: Chapter 12



# Quick Review



# Supervised learning

- Learn a mapping function  $f$  from  $\mathbf{x}$  to  $y$

$$y = f(\mathbf{x})$$

- If  $y \in \{1, 2, \dots, C\}$  the problem is called classification
- If  $y \in \mathbb{R}$  the problem is called regression



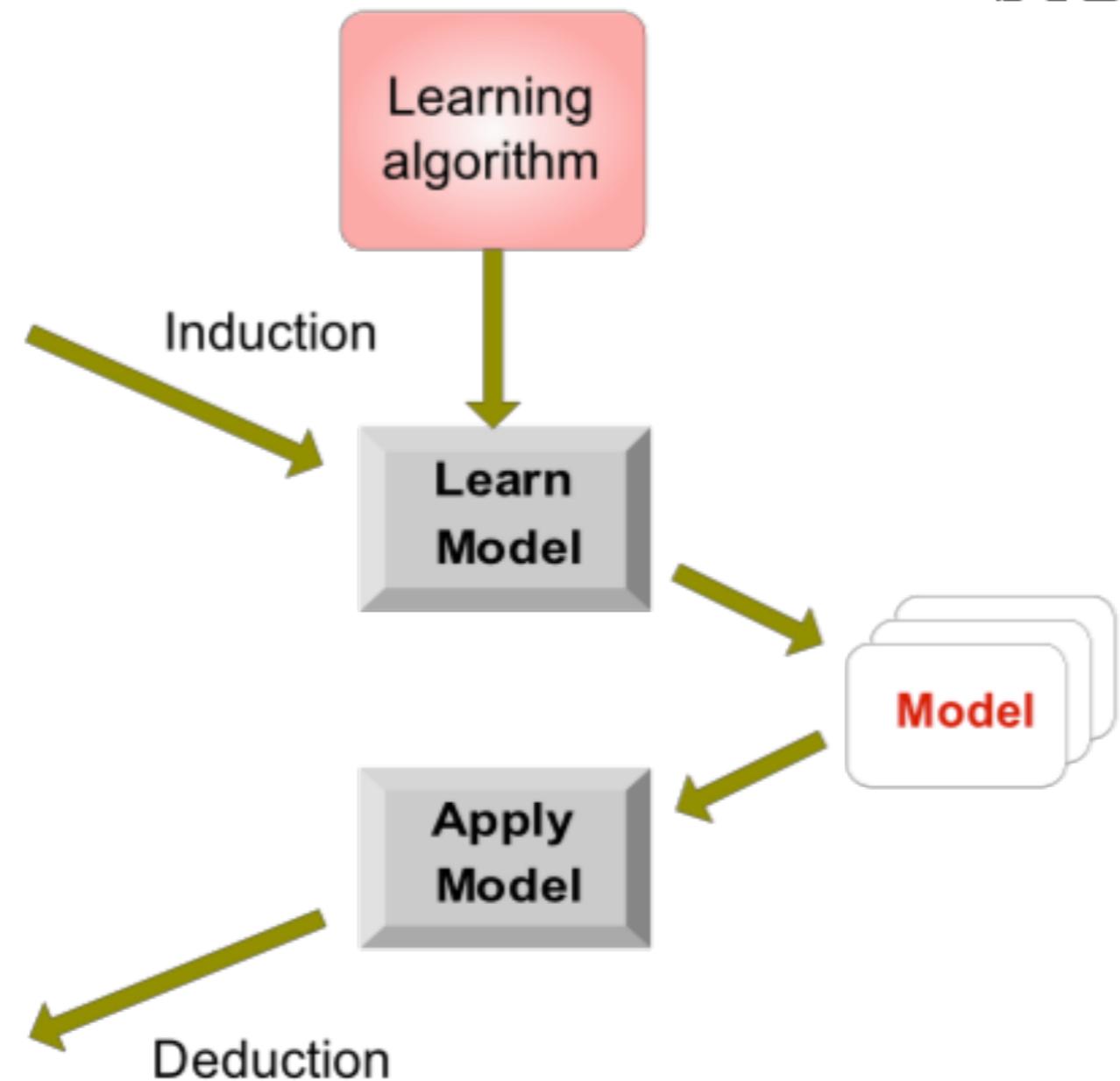
# Illustrating Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

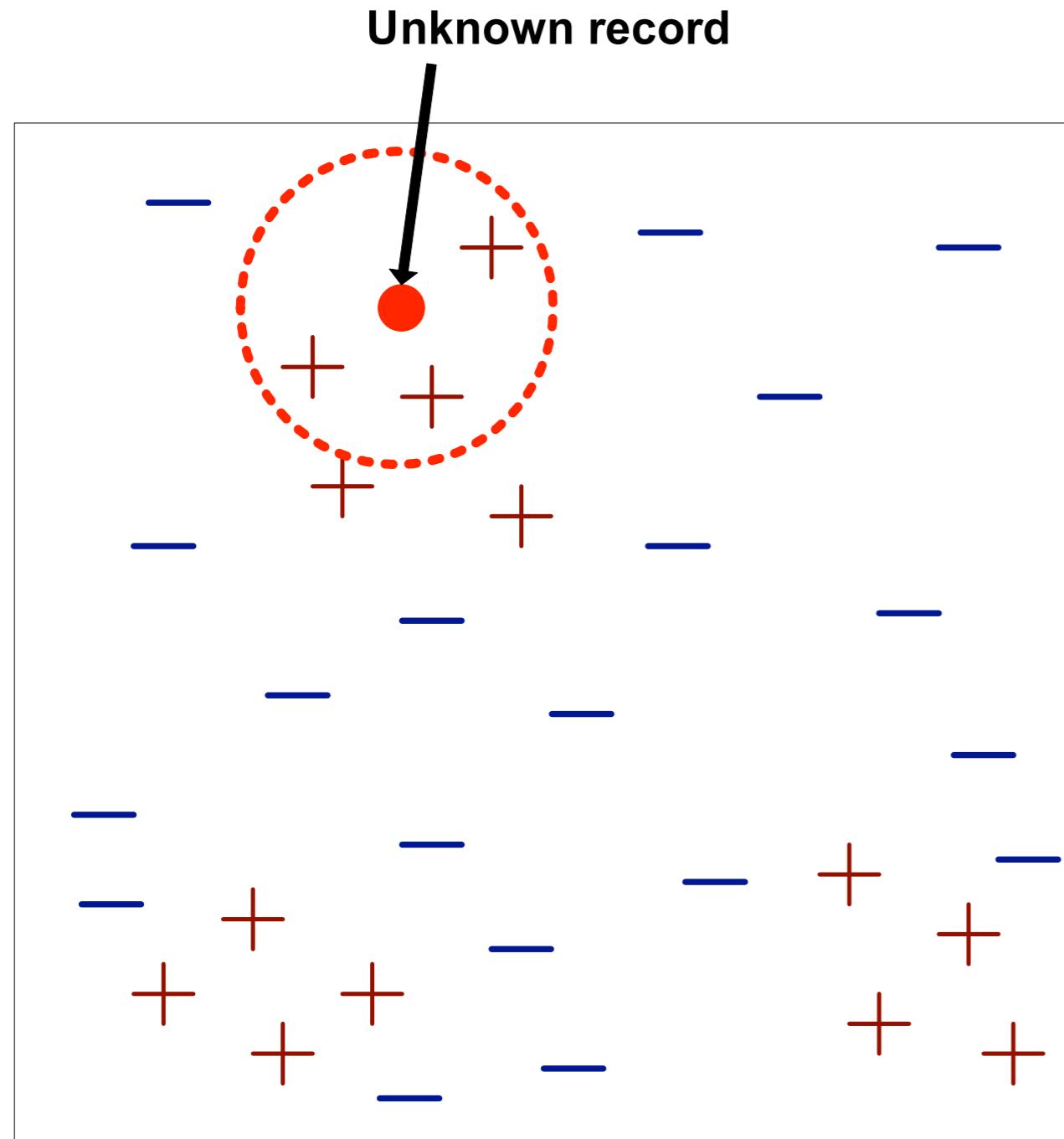
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set





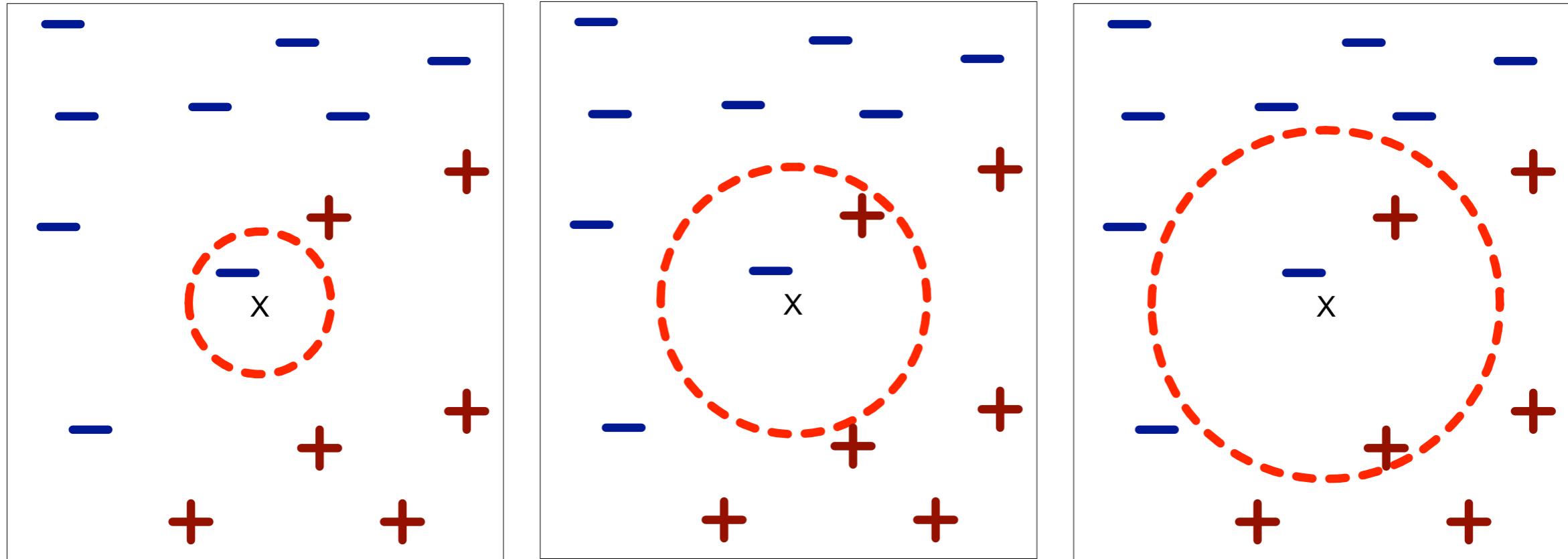
# Nearest Neighbour Classifiers



- Requires three things
  - The set of stored records
  - Distance Metric to compute distance between records
  - The value of  $k$ , the number of nearest neighbours to retrieve
- To classify an unknown record:
  - Compute distance to other training records
  - Identify  $k$  nearest neighbours
  - Use class labels of nearest neighbours to determine the class label of unknown record (e.g., by taking majority vote)



# Definition of Nearest Neighbour



(a) 1-nearest neighbor

(b) 2-nearest neighbor

(c) 3-nearest neighbor

K-nearest neighbours of a record  $x$  are data points that have the  $k$  smallest distance to  $x$



# Bayesian Classifiers

- Approach:
  - compute the posterior probability  $P(C | A_1, A_2, \dots, A_n)$  for all values of  $C$  using the Bayes' theorem

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- Choose value of  $C$  that maximises  $P(C | A_1, A_2, \dots, A_n)$
- Equivalent to choosing value of  $C$  that maximises  $P(A_1, A_2, \dots, A_n | C) P(C)$
- How to estimate  $P(A_1, A_2, \dots, A_n | C)$ ?



# Naïve Bayes Classifier

- Assume independence among attributes  $A_i$  when class is given:
  - $P(A_1, A_2, \dots, A_n | C) = P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$
  - Can estimate  $P(A_i | C_j)$  for all  $A_i$  and  $C_j$ .
  - New point is classified to  $C_j$  if  $P(C_j) \prod P(A_i | C_j)$  is maximal.

# How to Estimate Probabilities from Data?



THE UNIVERSITY OF  
SYDNEY

#	Refund	Status	Salary	Class
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class:  $P(C) = N_c/N$

- e.g.,  $P(No) = 7/10$ ,  
 $P(Yes) = 3/10$

- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_{C^k}$$

- where  $|A_{ik}|$  is number of instances having attribute  $A_i$  and belongs to class  $C_k$

- Examples:

$$P(\text{Status}=\text{Married} | \text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes} | \text{Yes}) = 0$$

# Example of Naïve Bayes Classifier



THE UNIVERSITY OF  
SYDNEY

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

A: attributes

M: mammals

N: non-mammals

$$P(A|M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A|N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A|M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A|N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$$P(A|M)P(M) > P(A|N)P(N)$$

=> Mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?



# Metrics for Performance Evaluation

- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix:

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a	b
	Class>No	c	d

a: TP (true positive)  
b: FN (false negative)  
c: FP (false positive)  
d: TN (true negative)



# Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a + c}$$

- a: TP (true positive)
- b: FN (false negative)
- c: FP (false positive)
- d: TN (true negative)

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

- Precision is biased towards C(Yes | Yes) & C(Yes | No)
- Recall is biased towards C(Yes | Yes) & C(No | Yes)
- F-measure is biased towards all except C(No | No)

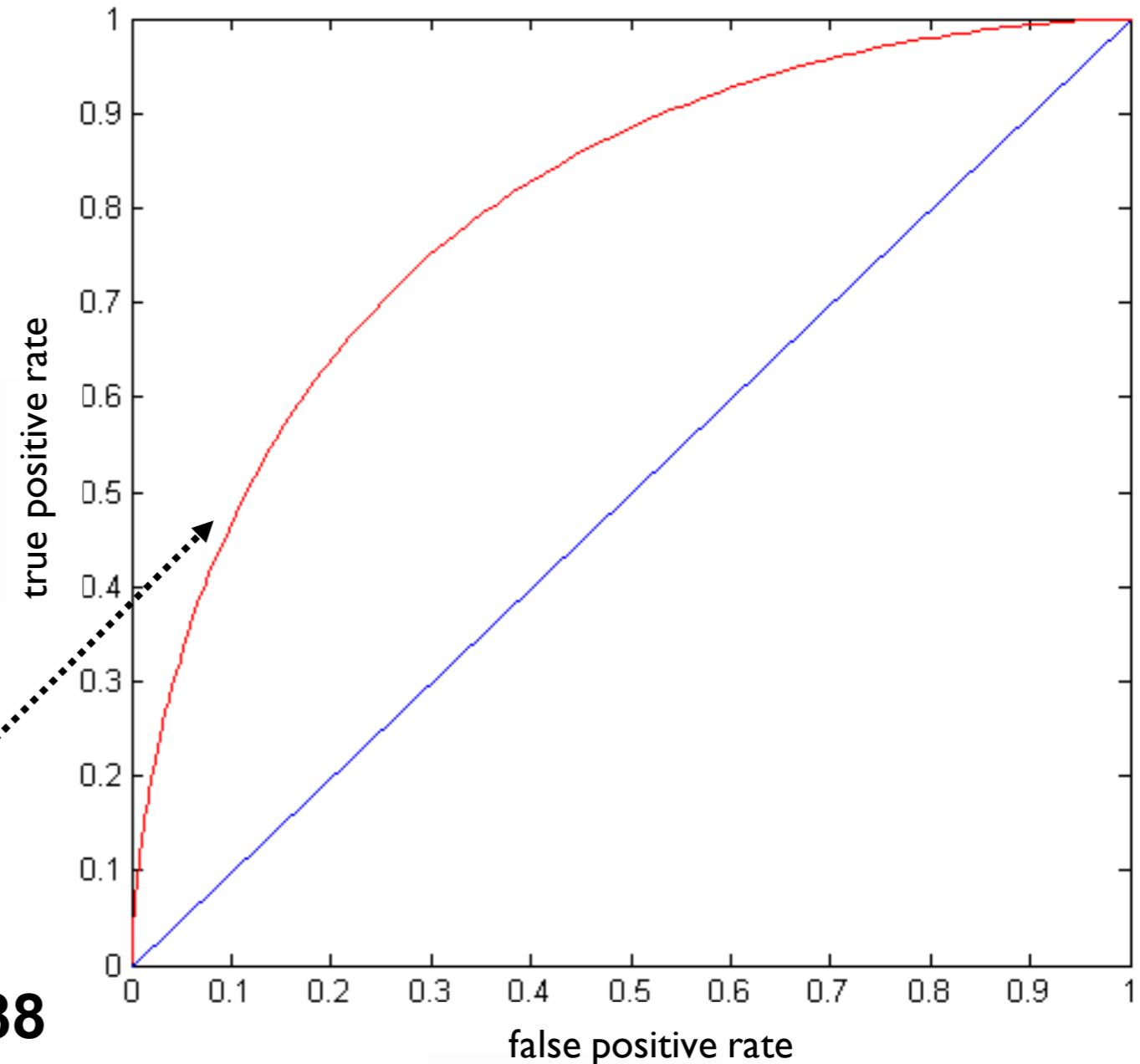
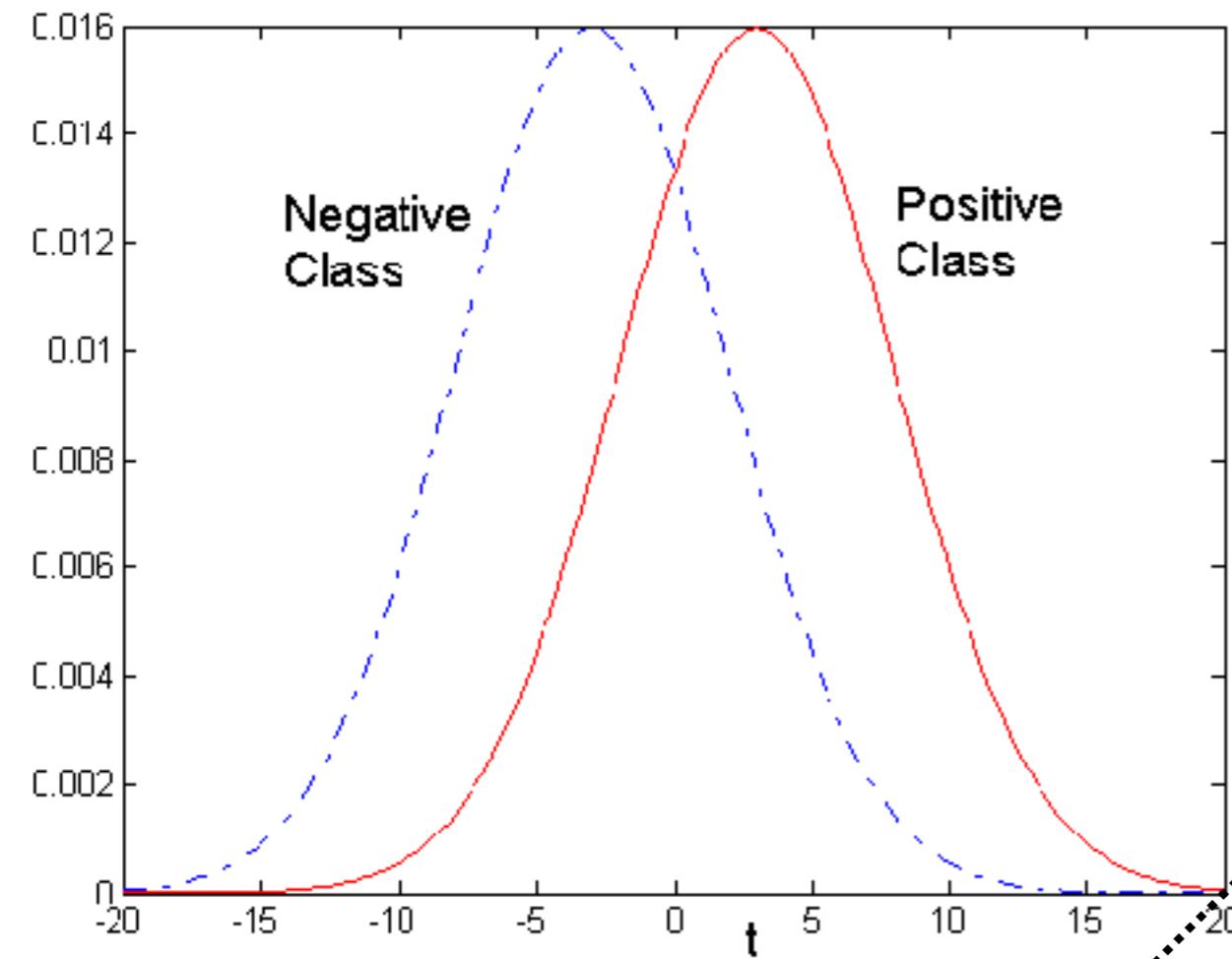
$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$



# ROC Curve

- TP rate,  $TPR = TP/(TP+FN)$
- FP rate,  $FPR = FP/(FP + TN)$

- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at  $x > t$  is classified as positive



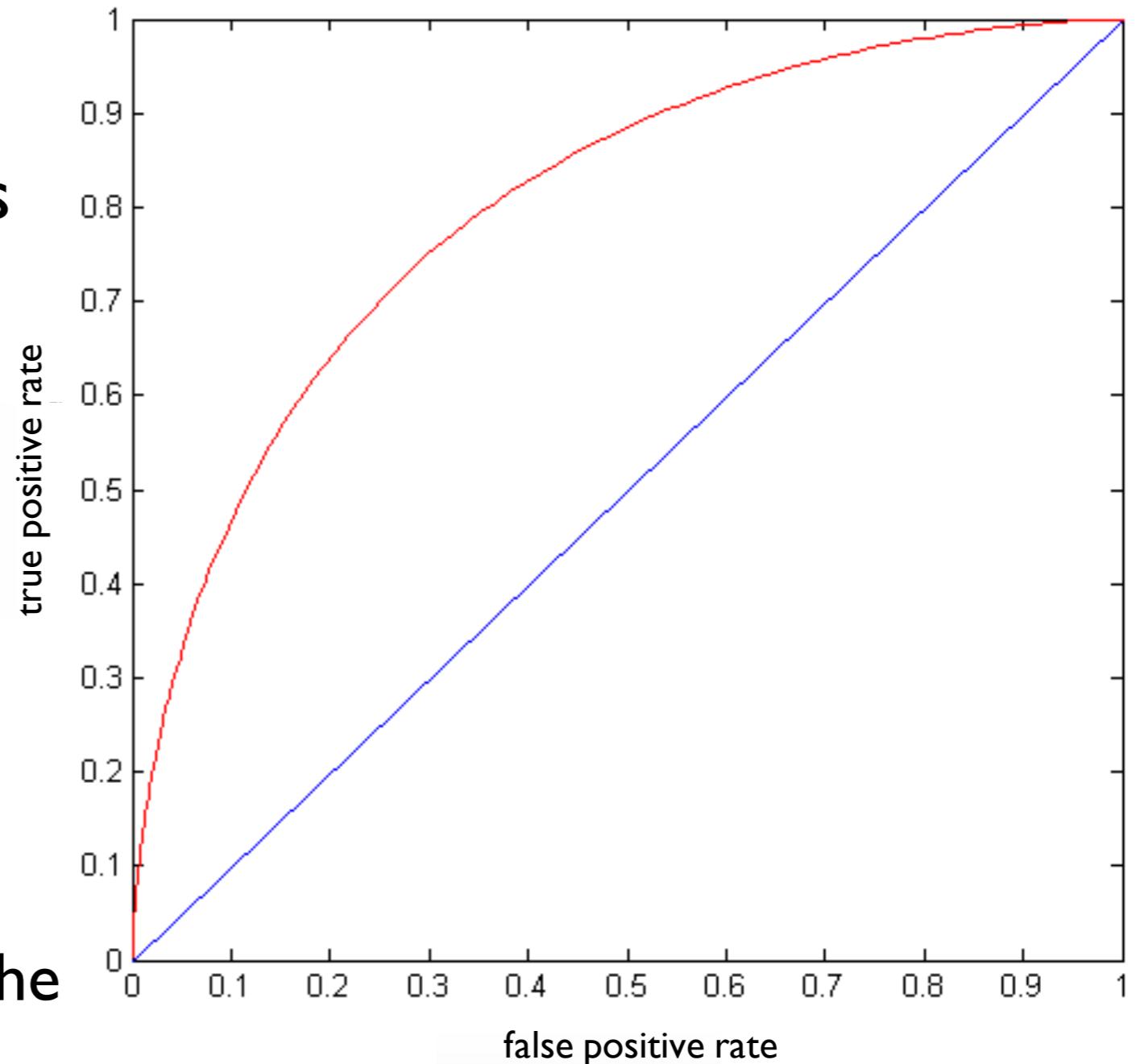
At threshold  $t$ :

$TP=0.5$ ,  $FN=0.5$ ,  $FP=0.12$ ,  $FN=0.88$

# ROC Curve

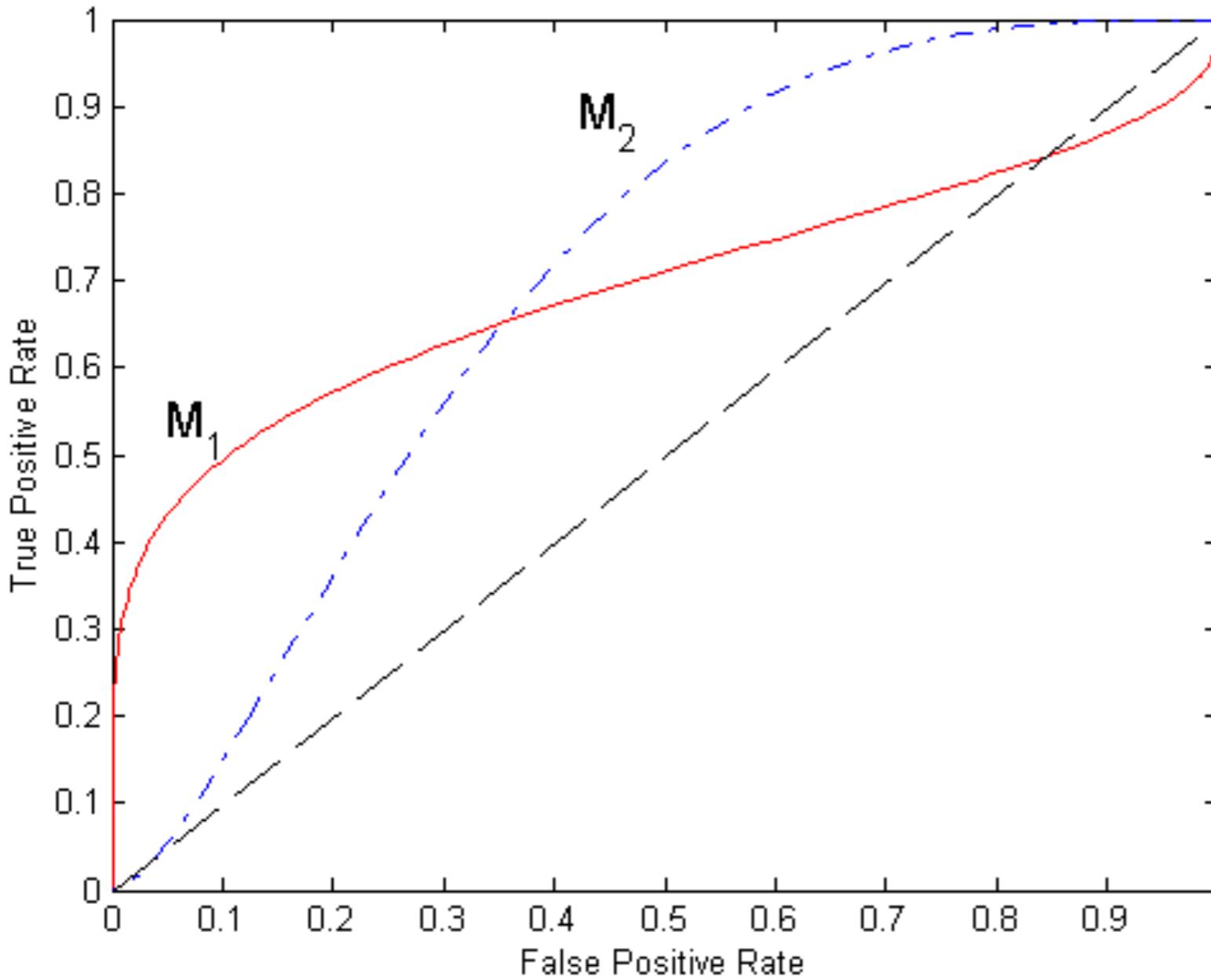
(TP,FP):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal
- Diagonal line:
  - Random guessing
  - Below diagonal line:
    - prediction is opposite of the true class





# Using ROC for Model Comparison



- No model consistently outperform the other
- $M_1$  is better for small FPR
- $M_2$  is better for large FPR
- Area Under the ROC curve
  - Ideal:
    - Area = 1
  - Random guess:
    - Area = 0.5



# How to construct a ROC curve

Instance	$P(+ A)$	True Class
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

- Use classifier that produces posterior probability for each test instance  $P(+|A)$
- Sort the instances according to  $P(+|A)$  in decreasing order
- Apply threshold at each unique value of  $P(+|A)$
- Count the number of TP, FP, TN, FN at each threshold
- TP rate,  $TPR = TP/(TP+FN)$
- FP rate,  $FPR = FP/(FP + TN)$

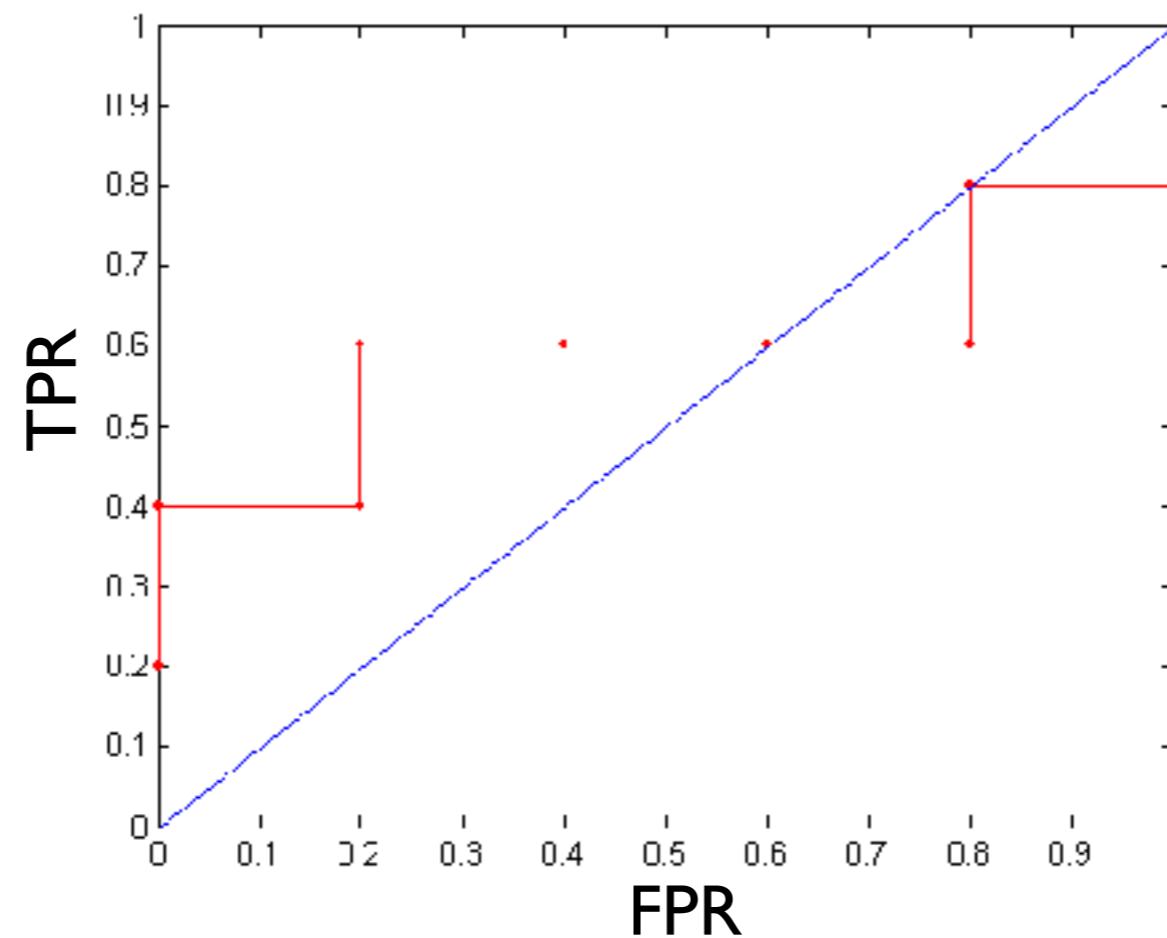


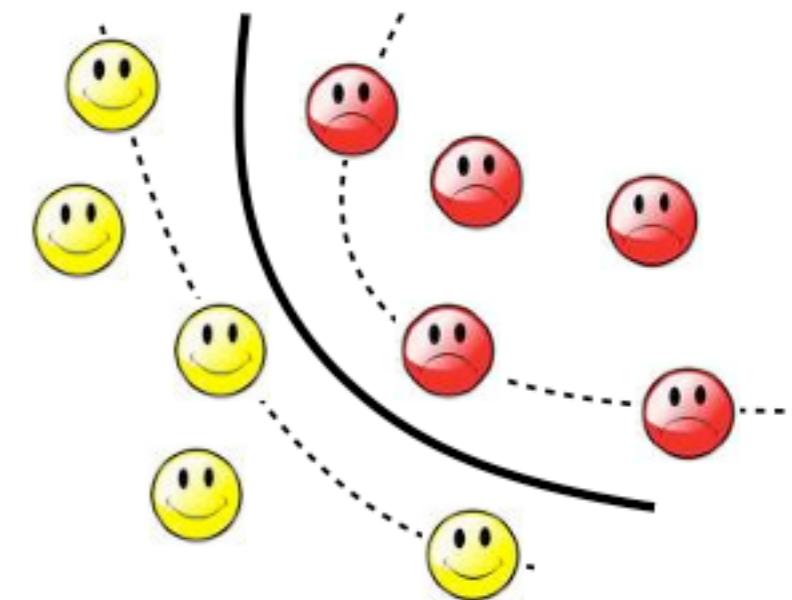
# How to construct a ROC curve

THE UNIVERSITY OF  
SYDNEY

Class	+	-	+	-	-	-	+	-	+	+	+	
Threshold >=	0.25	0.43	0.53	0.76	0.85	0.85	0.85	0.87	0.93	0.95	1.00	
TP	5	4	4	3	3	3	3	2	2	1	0	
FP	5	5	4	4	3	2	1	1	0	0	0	
TN	0	0	1	1	2	3	4	4	5	5	5	
FN	0	1	1	2	2	2	2	3	3	4	5	
TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0	
FPR	1	1	0.8	0.8	0.6	0.4	0.2	0.2	0	0	0	

ROC Curve:

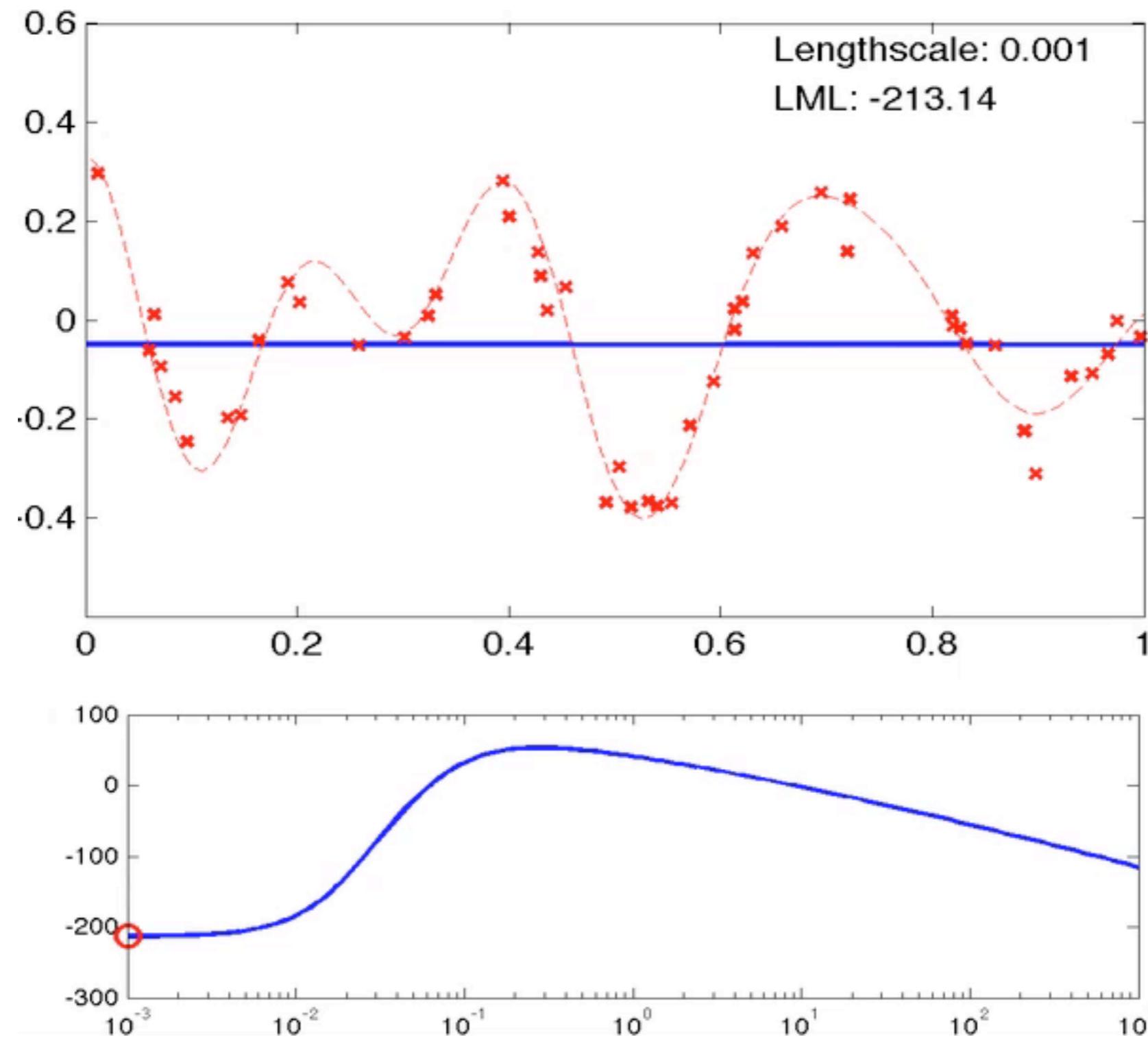




# Classification II



# Learning as optimisation





# Loss functions

- Squared error, 0-1 Loss

$$L(y, \hat{y}) = (y - \hat{y})^2$$

$$L(y, \hat{y}) = I(y \neq \hat{y})$$

- Minimise risk, (expected risk, empirical risk)

$$R(\hat{f}) = E_{\mathbf{x}, y} L(f(\mathbf{x}), \hat{f}(\mathbf{x}))$$

$$\hat{R}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(\mathbf{x}_i))$$



# Loss for density estimation

- Suppose output is  $\hat{p}(y|\mathbf{x})$ , truth is  $p(y|\mathbf{x})$
- Use KL (Kullback-Leibler) divergence

$$L(p(y|\mathbf{x}), \hat{p}(y|\mathbf{x})) = KL(p(y|\mathbf{x}), \hat{p}(y|\mathbf{x})) = \sum_y p(y|\mathbf{x}) \log \frac{p(y|\mathbf{x})}{\hat{p}(y|\mathbf{x})}$$

- Risk is expected negative log likelihood

$$R(\hat{p}) = -E_{\mathbf{x}} \sum_y p(y|\mathbf{x}) \log \hat{p}(y|\mathbf{x}) = -E_{\mathbf{x}, y} \log \hat{p}(y|\mathbf{x})$$



# Estimation vs Inference

- Learning as optimisation (frequentist): Given  $D$ , choose  $\hat{f}$  to approximate  $f$  as closely as possible, so as to minimise (future) expected risk
- Usually compute parameter estimate  $\hat{\theta}$
- Learning as inference (Bayesian): Given  $D$ , compute posterior over functions  $p(f|D)$
- Or posterior over parameters

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

- Decision theory demonstrates that one of the best ways to minimise frequentist risk is to be Bayesian.



# Generative vs Discriminative

- Generative approach:
  - Model  $p(y, \mathbf{x}) = p(\mathbf{x}|y)p(y)$
  - Use Bayes' theorem  $p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$
- Discriminative approach:
  - Model  $p(y|\mathbf{x})$  directly



# Naïve Bayes Classifier

Generative model:

$$p(\mathbf{x}, \mathcal{C}_k) = p(\mathcal{C}_k)p(\mathbf{x}_n | \mathcal{C}_k) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k)}{\sum_{\mathcal{C}_j} p(\mathbf{x} | \mathcal{C}_j)p(\mathcal{C}_j)}$$

class posterior

class conditional density

class prior

normalising constant

```
graph TD; A[class posterior] --> B[p(x|C_k)p(C_k)]; C[class conditional density] --> B; D[class prior] --> B; E[normalising constant] --> F[sum_{C_j} p(x|C_j)p(C_j)]
```



THE UNIVERSITY OF  
**SYDNEY**

# Logistic Regression

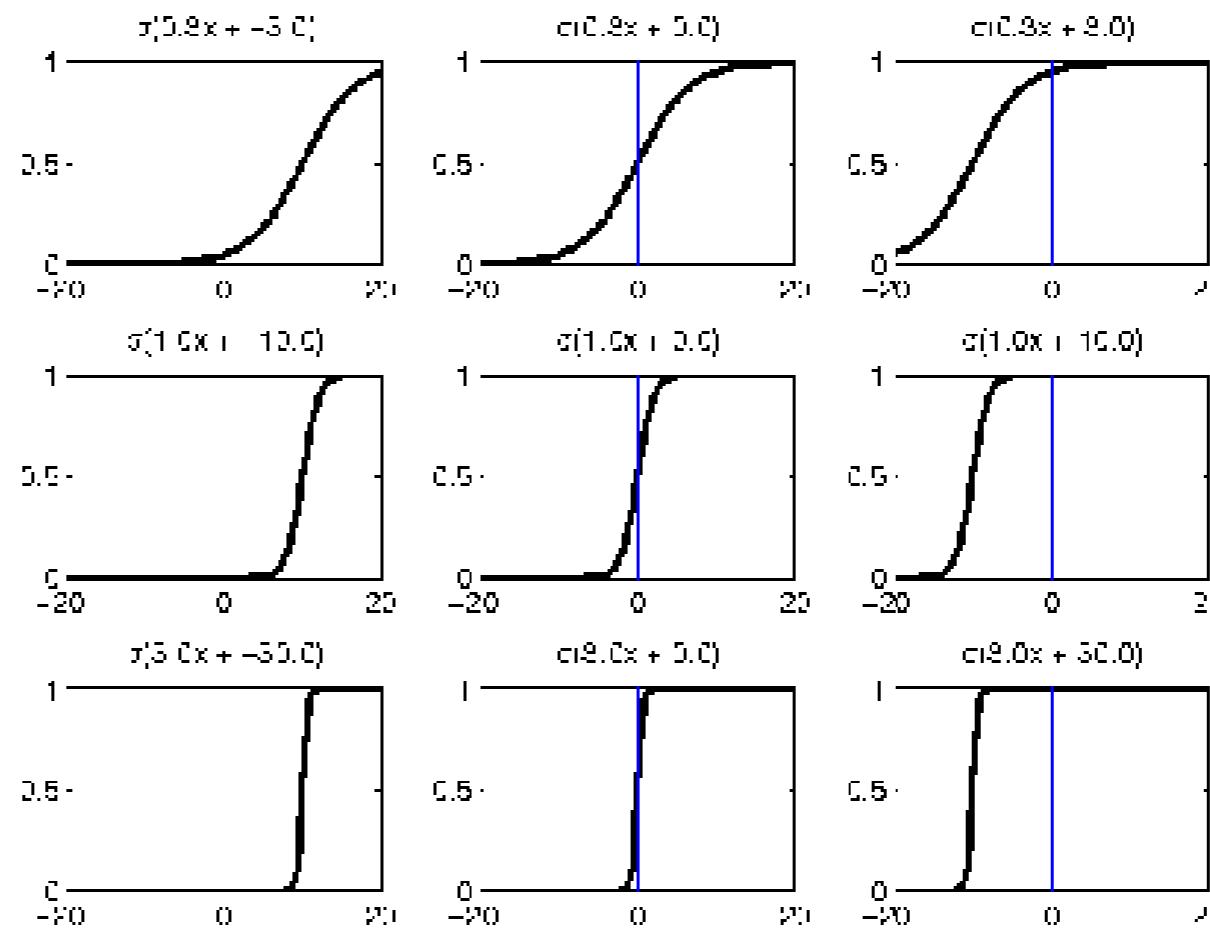


# Logistic Regression

- Discriminative model for binary classification

$$\begin{aligned} p(y|\mathbf{x}, \mathbf{w}) &= \text{Ber}(y|\sigma(\eta)) = \sigma(\eta)^y(1 - \sigma(\eta))^{1-y} \\ \eta &= \mathbf{w}^T \mathbf{x} \end{aligned}$$

$$\sigma(\eta) \stackrel{\text{def}}{=} \frac{1}{1 + \exp(-\eta)} = \frac{e^\eta}{e^\eta + 1}$$

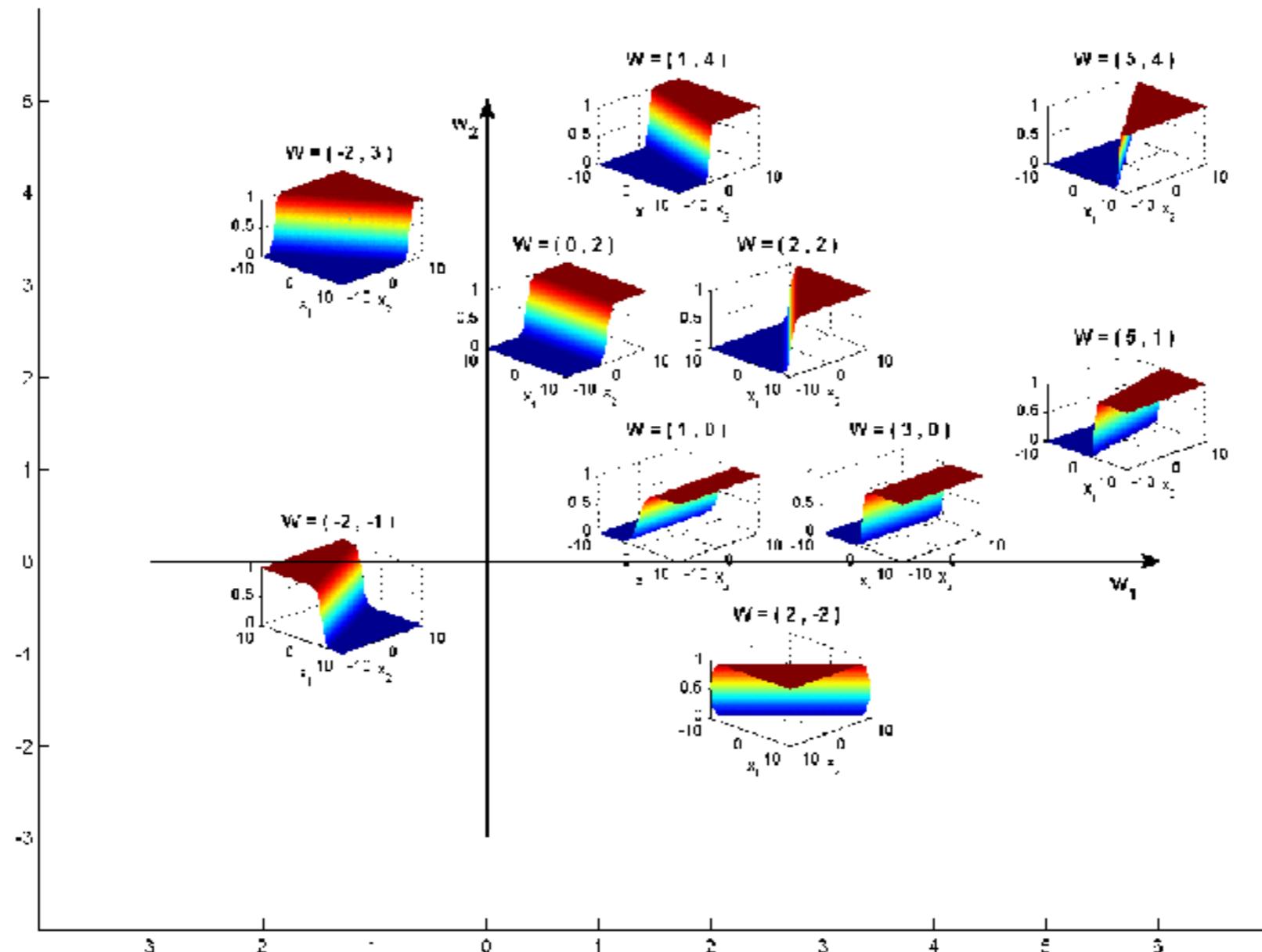


Sigmoid  
or  
Logistic  
function

# Decision boundary

- Logistic Regression in 2D

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$





# Logistic Regression

- Assumes a parametric form for directly estimating  $P(Y | X)$ . For binary concepts, this is:

$$P(Y = 0|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

$$\begin{aligned} P(Y = 1|X) &= 1 - P(Y = 0|X) \\ &= \frac{\exp(w_0 + \sum_{i=1}^n w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \end{aligned}$$

- Equivalent to a one-layer backpropagation neural net.
- Logistic regression is the source of the sigmoid function used in backpropagation.
- Objective function for training is somewhat different.



# Logistic Regression as a Log-Linear Model

- Logistic regression is basically a linear model, which is demonstrated by taking logs.

$$\text{Assign label } Y = 0 \text{ iff } 1 < \frac{P(Y = 0 | X)}{P(Y = 1 | X)}$$
$$1 < \exp(w_0 + \sum_{i=1}^n w_i X_i)$$
$$0 < w_0 + \sum_{i=1}^n w_i X_i$$

or equivalently  $w_0 > \sum_{i=1}^n -w_i X_i$

- Also called a **maximum entropy model (MaxEnt)** because it can be shown that standard training for logistic regression gives the distribution with maximum entropy that is consistent with the training data.



# Logistic Regression Training

- Weights are set during training to maximise the **conditional data likelihood** :

$$W \leftarrow \operatorname{argmax}_W \prod_{d \in D} P(Y^d | X^d, W)$$

where  $D$  is the set of training examples and  $Y^d$  and  $X^d$  denote, respectively, the values of  $Y$  and  $X$  for example  $d$ .

- Equivalently viewed as maximising the **conditional log likelihood** (CLL)

$$W \leftarrow \operatorname{argmax}_W \sum_{d \in D} \ln P(Y^d | X^d, W)$$



# Logistic Regression Training

- Like neural-nets, can use standard gradient descent to find the parameters (weights) that optimise the CLL objective function.
- Many other more advanced training methods are possible to speed up convergence.
  - Conjugate gradient
  - Generalised Iterative Scaling (GIS)
  - Modified Iterative Scaling (MIS)
  - Limited-memory quasi-Newton (L-BFGS)
  - Stochastic gradient descent



# Logistic Regression Training

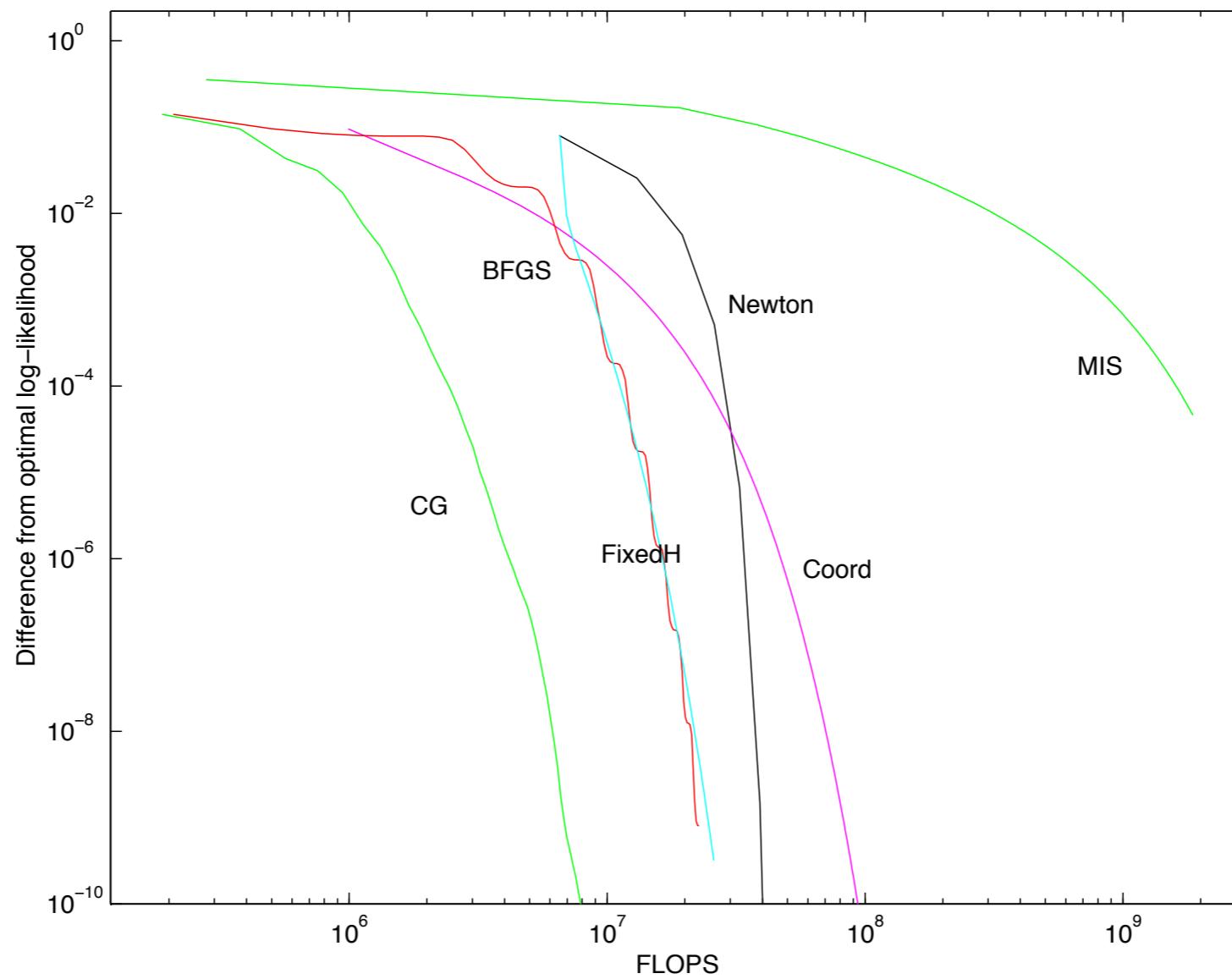


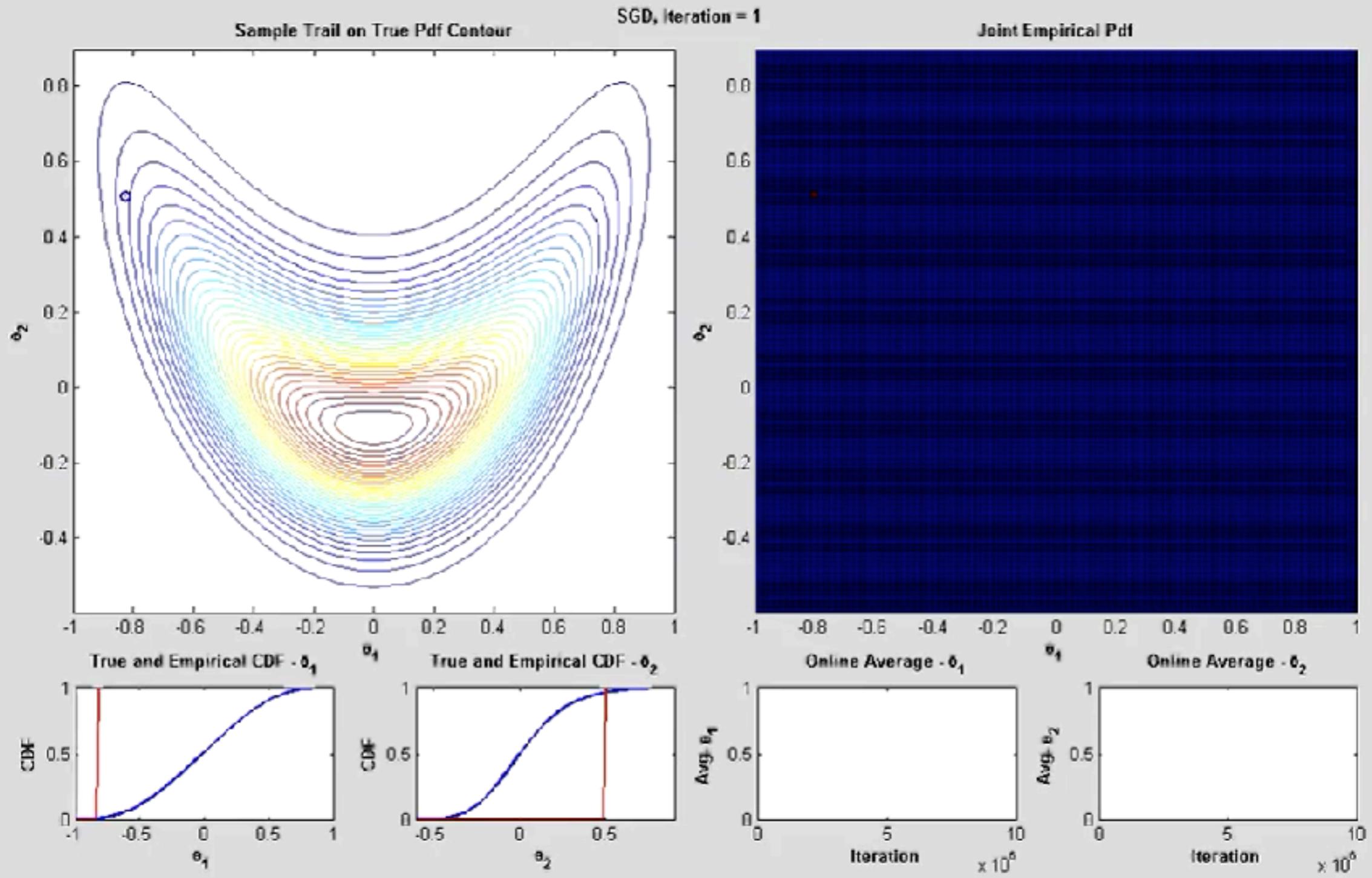
Figure 1: Cost vs. performance of six logistic regression algorithms. The dataset had 300 points in 100 dimensions. “CG” is conjugate gradient (section 4), “Coord” is coordinate-wise Newton, “FixedH” is Fixed-Hessian, and “MIS” is modified iterative scaling. CG also has the lowest actual time in Matlab.

See: “A comparison of numerical optimizers for logistic regression” by Thomas Minka, 2003.



THE UNIVERSITY OF  
SYDNEY

# Gradient Descent





# Preventing Overfitting in Logistic Regression

- To prevent overfitting, one can use **regularisation** (a.k.a. smoothing) by penalising large weights by changing the training objective:

$$W \leftarrow \operatorname{argmax}_W \sum_{d \in D} \ln P(Y^d | X^d, W) - \frac{\lambda}{2} \|W\|^2$$

Where  $\lambda$  is a constant that determines the amount of smoothing

- This can be shown to be equivalent to assuming a Gaussian prior for  $W$  with zero mean and a variance related to  $1/\lambda$ .



THE UNIVERSITY OF  
SYDNEY

# Multinomial Logistic Regression

- Logistic regression can be generalised to multi-class problems (where  $Y$  has a multinomial distribution).
- Effectively constructs a linear classifier for each category.

# Relation Between NB and Logistic Regression



THE UNIVERSITY OF  
SYDNEY

- Naïve Bayes with Gaussian distributions for features (GNB), can be shown to give the same functional form for the conditional distribution  $P(Y|X)$ .
- But converse is not true, so Logistic Regression makes a weaker assumption.
- Logistic regression is a **discriminative** rather than generative model, since it models the conditional distribution  $P(Y|X)$  and directly attempts to fit the training data for predicting  $Y$  from  $X$ . Does not specify a full joint distribution.

# Relation Between NB and Logistic Regression (continued)

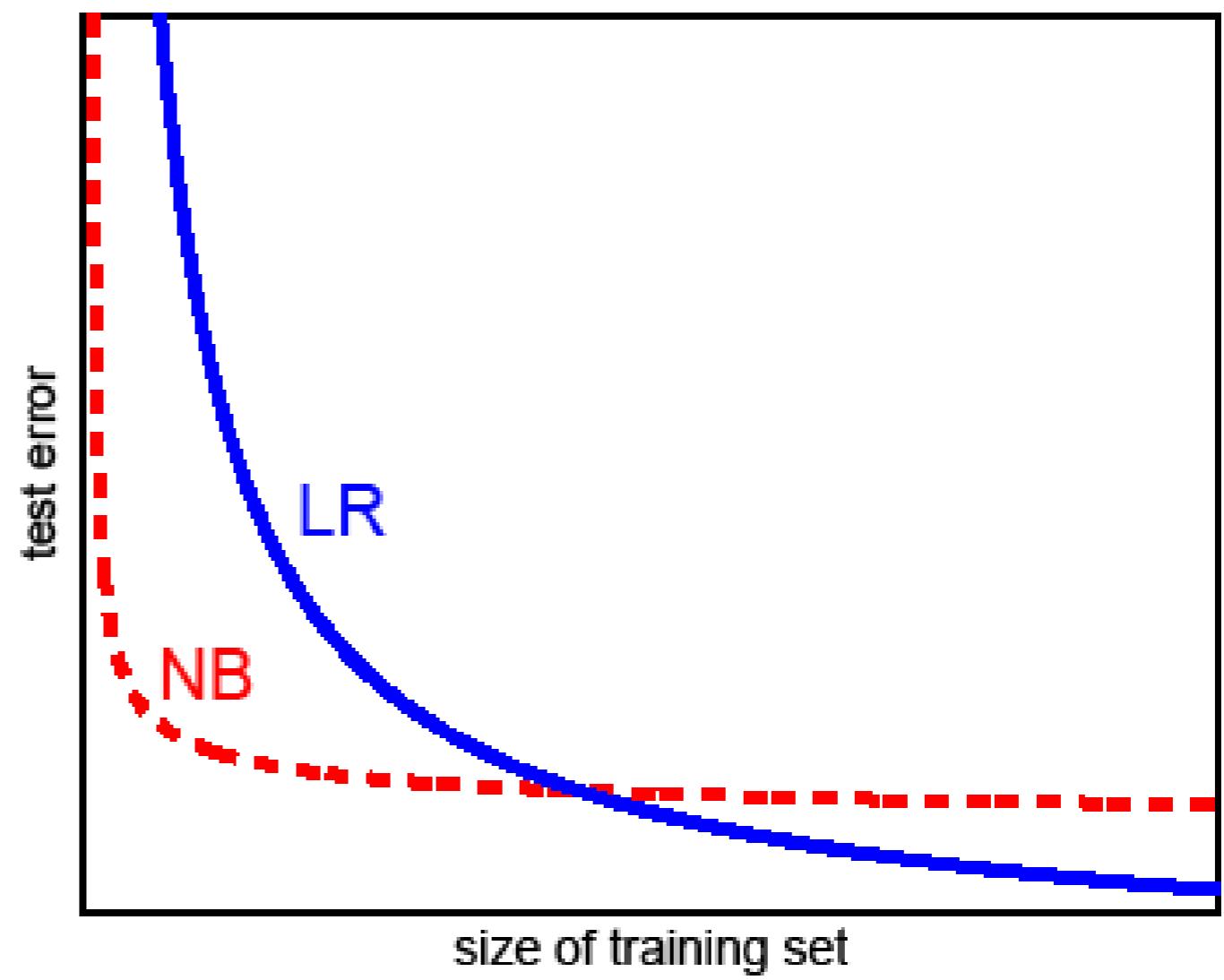
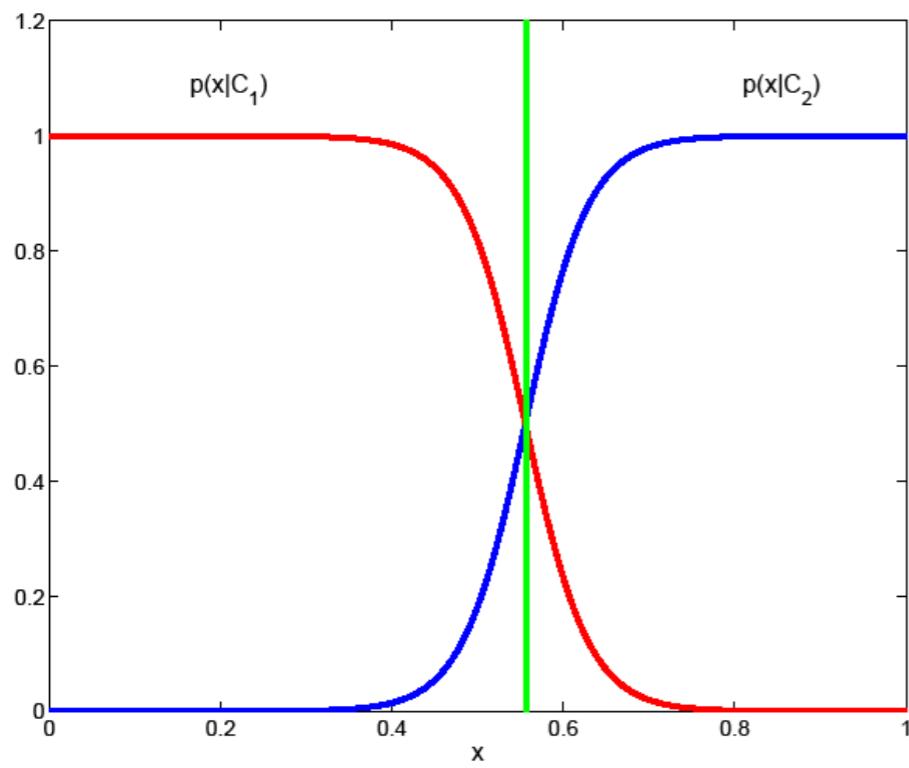
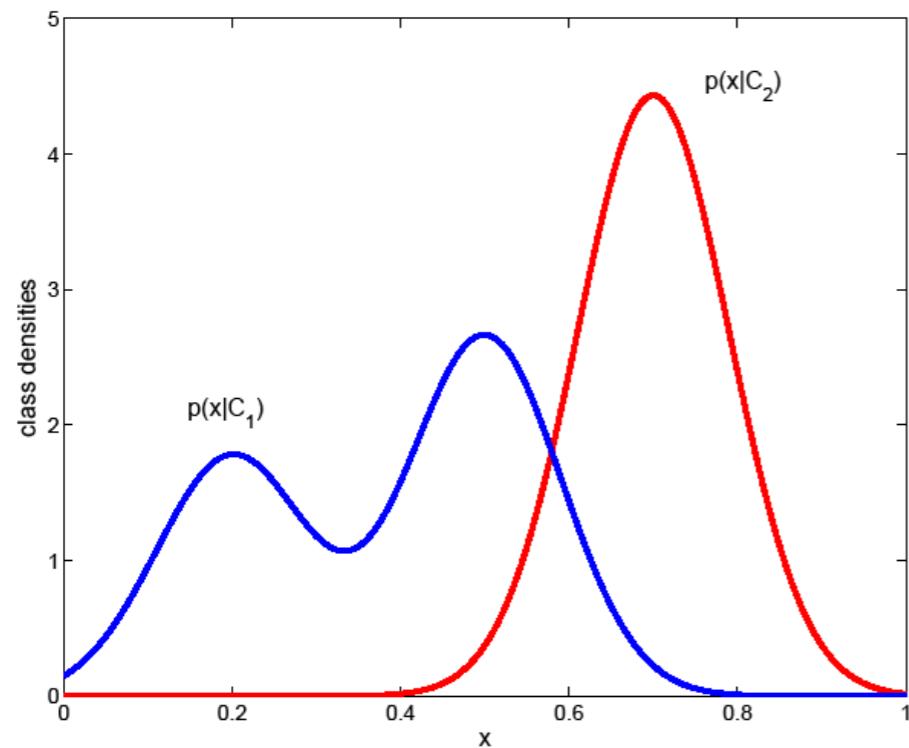


THE UNIVERSITY OF  
SYDNEY

- When conditional independence is violated, logistic regression gives better generalisation if it is given sufficient training data.
- GNB converges to accurate parameter estimates faster ( $O(\log n)$  examples for  $n$  features) compared to Logistic Regression ( $O(n)$  examples).
- Experimentally, GNB is better when training data is scarce, logistic regression is better when it is plentiful.



# Logistic Regression vs NB





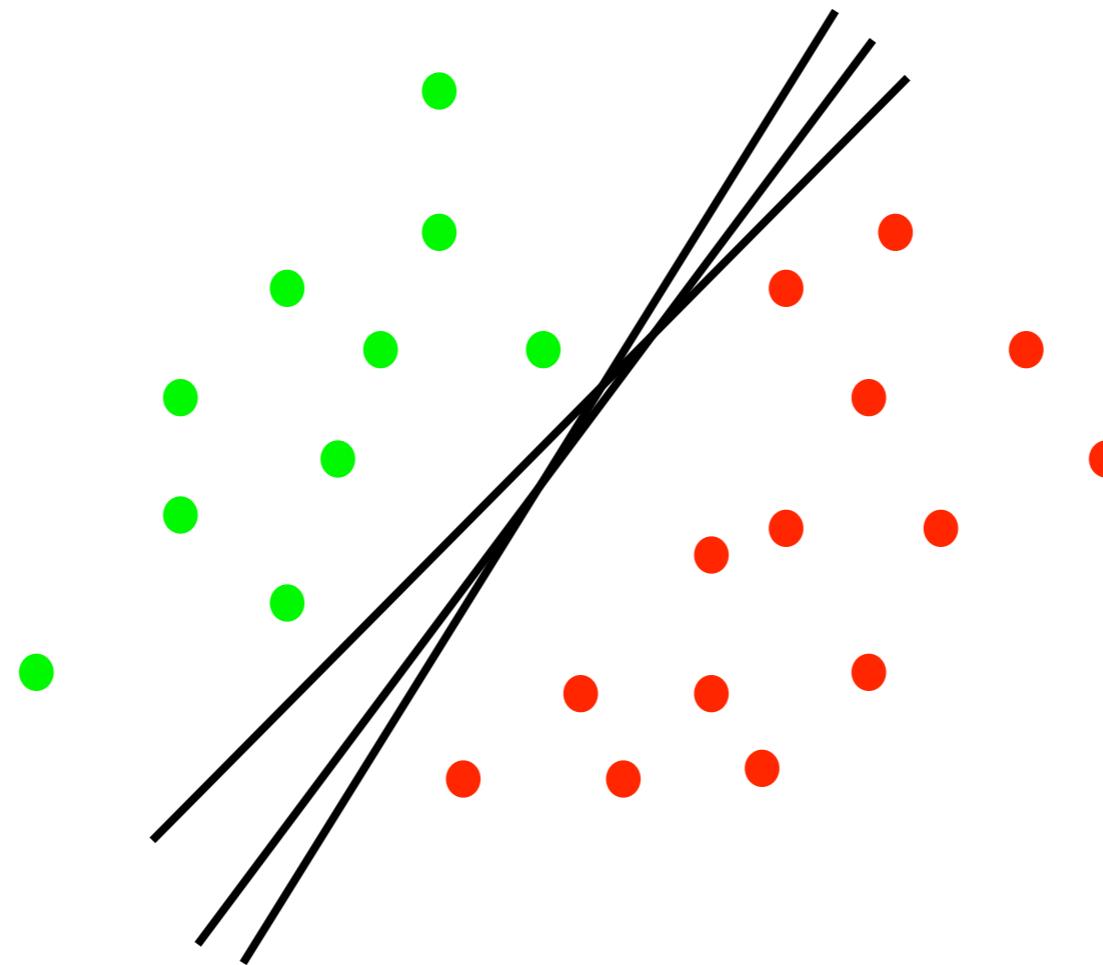
THE UNIVERSITY OF  
**SYDNEY**

# Support Vector Machines



THE UNIVERSITY OF  
SYDNEY

# Linear classifiers

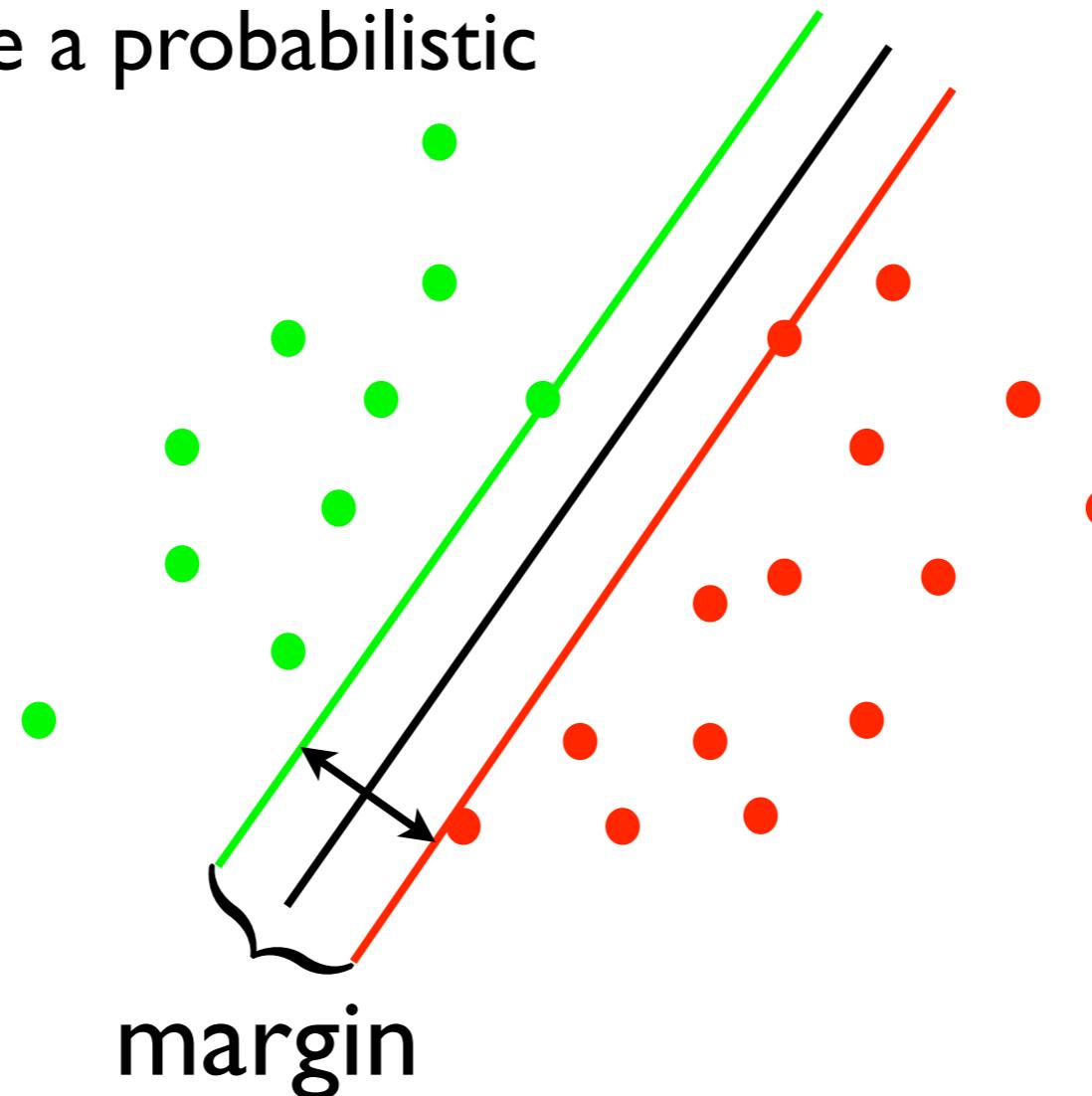


Which line is better?

# Maximum margin

## Things to note

- The boundary is defined by only a few points (support vectors)
- Difficult to define a probabilistic explanation





# Linear Basis Function Models

Generally

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

where  $\phi_j(\mathbf{x})$  are known as *basis functions*.

Typically,  $\phi_0(\mathbf{x}) = 1$ , so that  $w_0$  acts as a bias.

In the simplest case, we use linear basis functions :

$$\phi_d(\mathbf{x}) = x_d$$



# Maximum margin classifiers (I)

Given a training set

Inputs  $\mathbf{x}_1, \dots, \mathbf{x}_N$

Targets

$t_1, \dots, t_N$  where  $t_n \in \{-1, 1\}$

Linear classifier

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

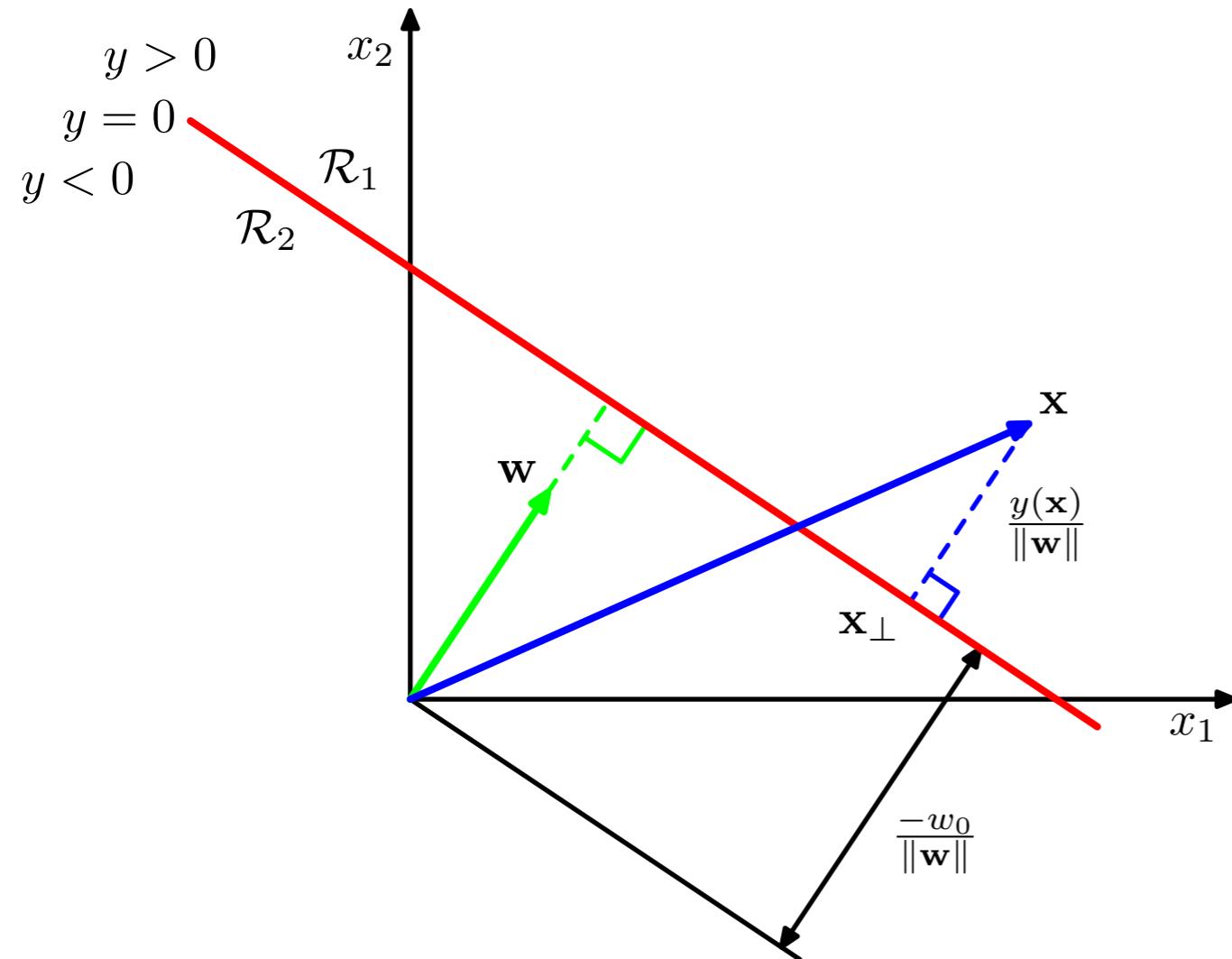
Output sign of  $y(\mathbf{x})$

so that  $t_n y(\mathbf{x}_n) > 0$

for all points in the training set



# Maximum margin classifiers(2)



Assuming  $\phi(\mathbf{x}) = \mathbf{x}$

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

For points in the boundary

$$y(\mathbf{x}) = 0 \quad \text{and}$$

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

For arbitrary  $\mathbf{x}$

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad \text{and} \quad r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$



# Maximum margin classifiers(3)

For all points correctly classified  $t_n y(\mathbf{x}_n) > 0$  ,  $b = \omega_0$

and 
$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

Maximum margin is found by

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

However, this is too difficult to optimise !



# Maximum margin classifiers(4)

Note that if  $\mathbf{w} \rightarrow \kappa \mathbf{w}$      $t_n y(\mathbf{x}_n)/\|\mathbf{w}\|$  is unchanged  
 $b \rightarrow \kappa b$

Setting  $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$  for the support vectors:

## Quadratic programming

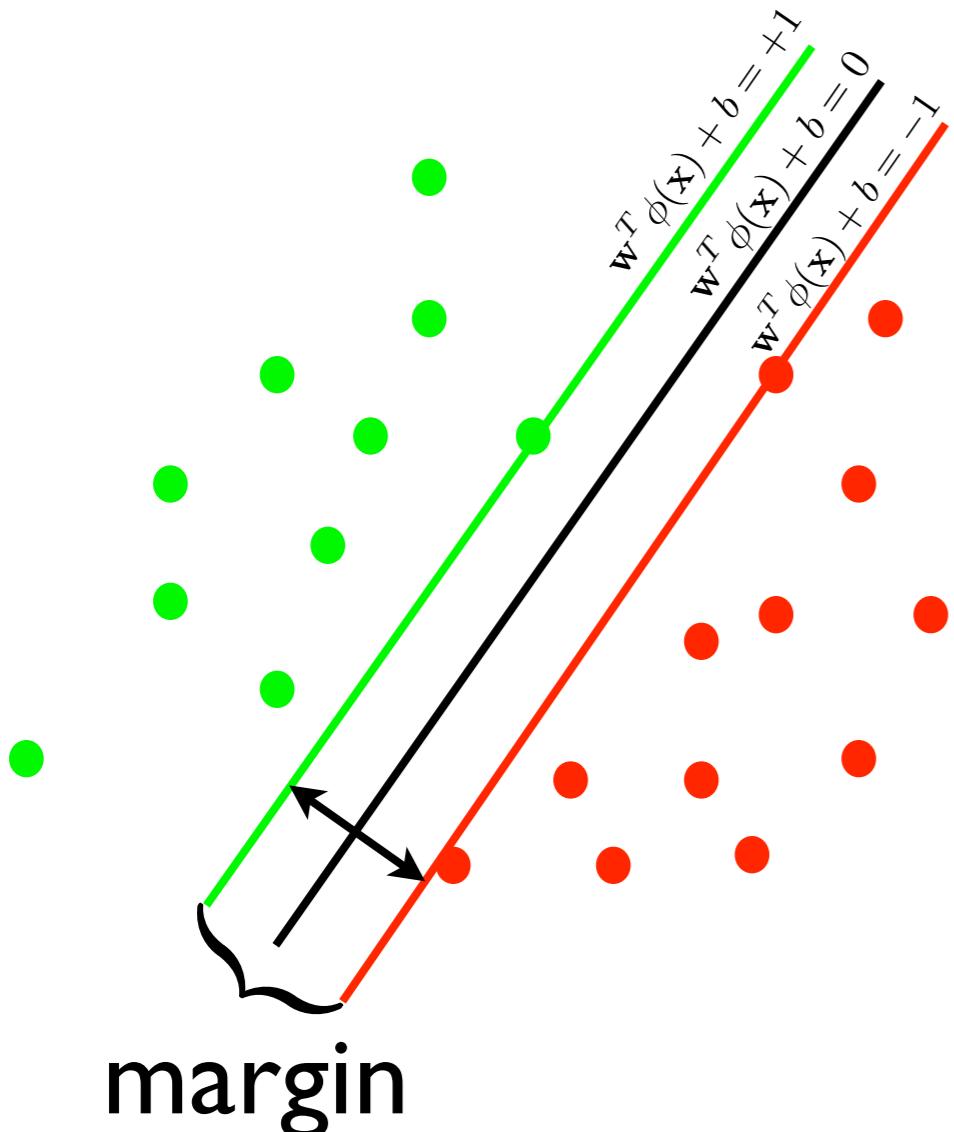
$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

s.t.

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N$$



# Support Vector Machines



## Quadratic programming

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

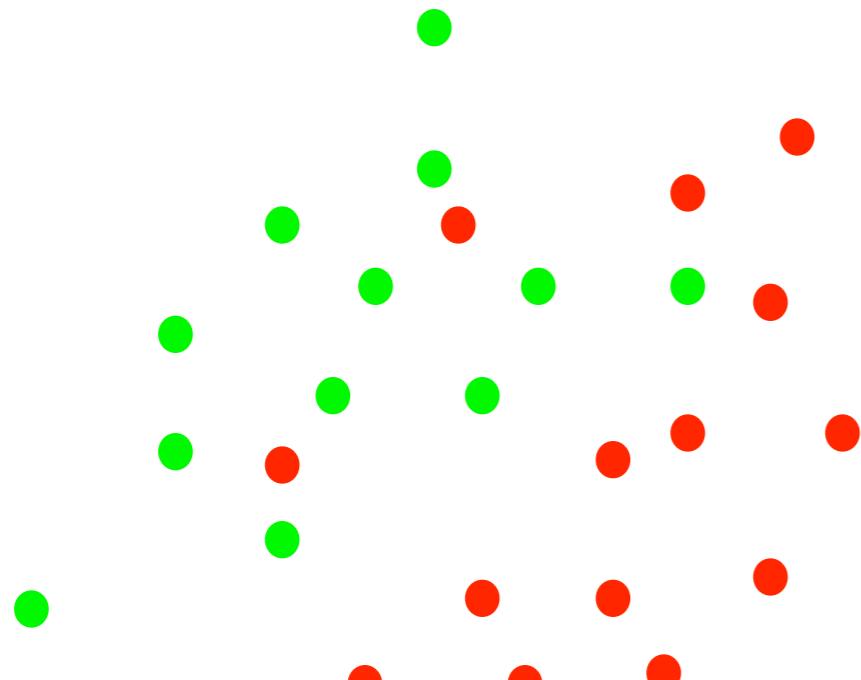
s.t.

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N$$

- Solve efficiently by quadratic programming (QP)
- Hyperplane defined by support vectors



# What happens if the data is not linearly separable?



## Quadratic programming

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

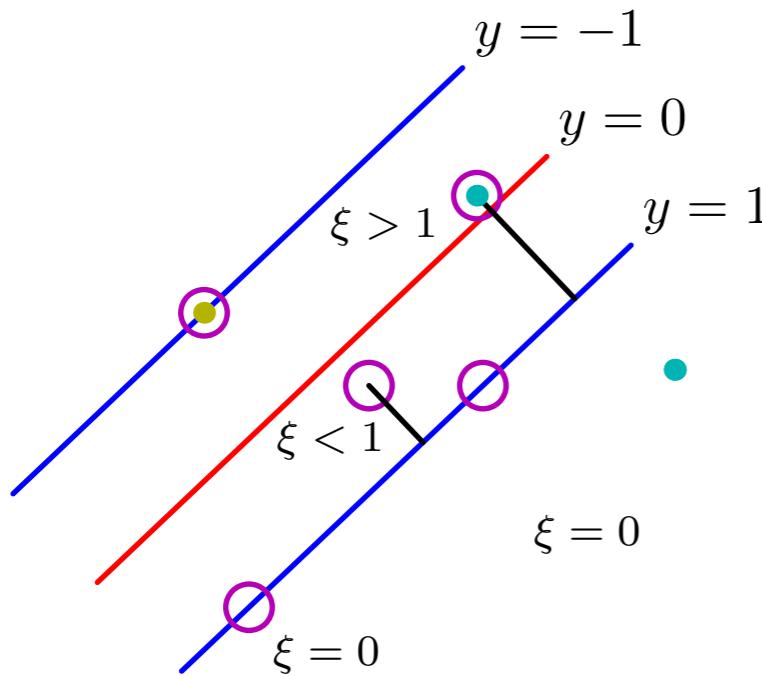
s.t.

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N$$

If margin  $\geq 1$ , don't care  
If margin  $< 1$ , pay linear penalty



# Soft Margin Classification



## Quadratic programming

$$\arg \min_{\mathbf{w}, b} C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

s.t.

$$t_n y(\mathbf{x}_n) \geqslant 1 - \xi_n$$

$$\xi_n \geqslant 0$$

$$n = 1, \dots, N$$



# The Dual SVM formulation

## Quadratic programming

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

s.t.

$$t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N$$

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

Kernel trick

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

## Lagrangian function

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\}$$

setting the derivatives w.r.t.  
 $\mathbf{w}$  and  $b$  equal to zero:

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

$$0 = \sum_{n=1}^N a_n t_n.$$

$$\begin{aligned} a_n &\geq 0, & n = 1, \dots, N, \\ \text{s.t. } \sum_{n=1}^N a_n t_n &= 0. \end{aligned}$$



# Kernel SVMs

- Solve the QP problem

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

with respect to  $\mathbf{a}$  subject to the constraints

$$a_n \geq 0, \quad n = 1, \dots, N,$$

$$\sum_{n=1}^N a_n t_n = 0.$$

- Support vectors satisfy

$$t_n \left( \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1$$

- For query points compute

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b.$$



# Think about kernels not features

## Polynomial kernel

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

## Gaussian kernel

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

Infinite dimensional feature space!





# What Functions are Kernels?

- For some functions  $K(\mathbf{x}_i, \mathbf{x}_j)$  checking that  $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  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:

$K =$

$K(\mathbf{x}_1, \mathbf{x}_1)$	$K(\mathbf{x}_1, \mathbf{x}_2)$	$K(\mathbf{x}_1, \mathbf{x}_3)$	...	$K(\mathbf{x}_1, \mathbf{x}_n)$
$K(\mathbf{x}_2, \mathbf{x}_1)$	$K(\mathbf{x}_2, \mathbf{x}_2)$	$K(\mathbf{x}_2, \mathbf{x}_3)$		$K(\mathbf{x}_2, \mathbf{x}_n)$
...	...	...	...	...
$K(\mathbf{x}_n, \mathbf{x}_1)$	$K(\mathbf{x}_n, \mathbf{x}_2)$	$K(\mathbf{x}_n, \mathbf{x}_3)$	...	$K(\mathbf{x}_n, \mathbf{x}_n)$



# Examples of Kernel Functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$ 
  - Mapping  $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$ , where  $\phi(\mathbf{x})$  is  $\mathbf{x}$  itself
- Polynomial of power  $p$ :  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$ 
  - Mapping  $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$ , where  $\phi(\mathbf{x})$  has  $\binom{d+p}{p}$  dimensions
- Gaussian (radial-basis function):  $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}}$ 
  - Mapping  $\Phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$ , where  $\phi(\mathbf{x})$  is *infinite-dimensional*: every point is mapped to *a function* (a Gaussian); combination of functions for support vectors is the separator.
  - Higher-dimensional space still has *intrinsic* dimensionality  $d$  (the mapping is not *onto*), but linear separators in it correspond to *non-linear* separators in original space.



# Overfitting

- Huge feature space with kernels, what about overfitting???
- Maximising margin leads to sparse set of support vectors
- Some interesting theory says that SVMs search for simple hypothesis with large margin
- Often robust to overfitting



# Summary

- SVM advantages
  - Efficient QP learning
  - Sparse
- SVM disadvantages
  - Not probabilistic
  - Do not directly handle multi-class problems
  - QP might struggle in very large datasets