

NONLINEAR CLASSIFICATION AND REGRESSION

Nonlinear Classification and Regression: Outline

2

Probability & Bayesian Inference

- Multi-Layer Perceptrons
 - The Back-Propagation Learning Algorithm
- Generalized Linear Models
 - Radial Basis Function Networks
 - Sparse Kernel Machines
 - Nonlinear SVMs and the Kernel Trick
 - Relevance Vector Machines

Nonlinear Classification and Regression: Outline

3

Probability & Bayesian Inference

□ Multi-Layer Perceptrons

- The Back-Propagation Learning Algorithm

□ Generalized Linear Models

- Radial Basis Function Networks

- Sparse Kernel Machines

- Nonlinear SVMs and the Kernel Trick

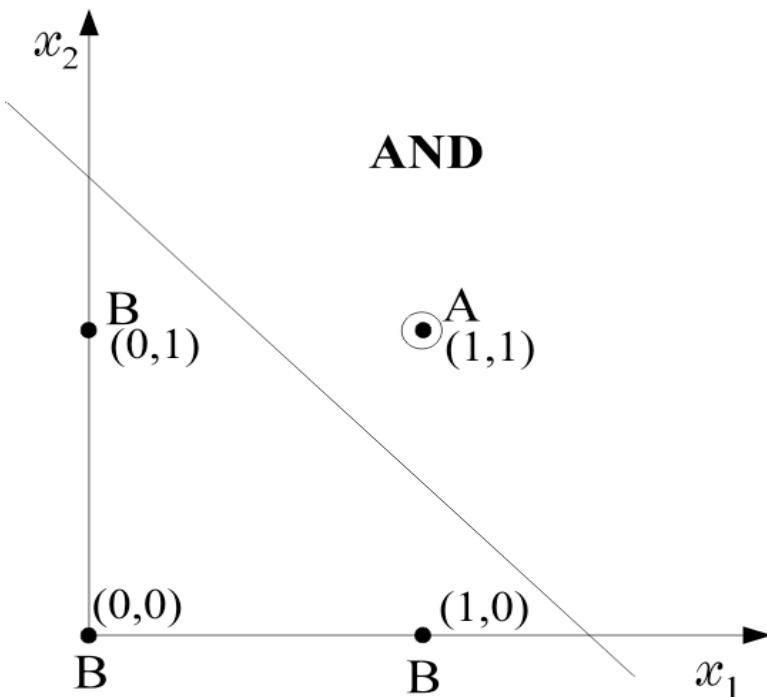
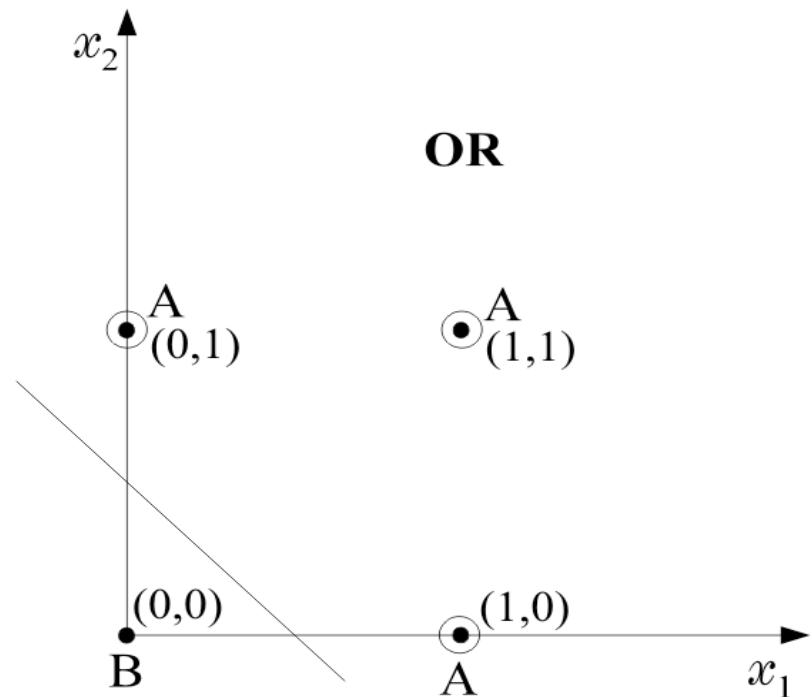
- Relevance Vector Machines

Implementing Logical Relations

4

Probability & Bayesian Inference

- AND and OR operations are linearly separable problems

**AND****OR**

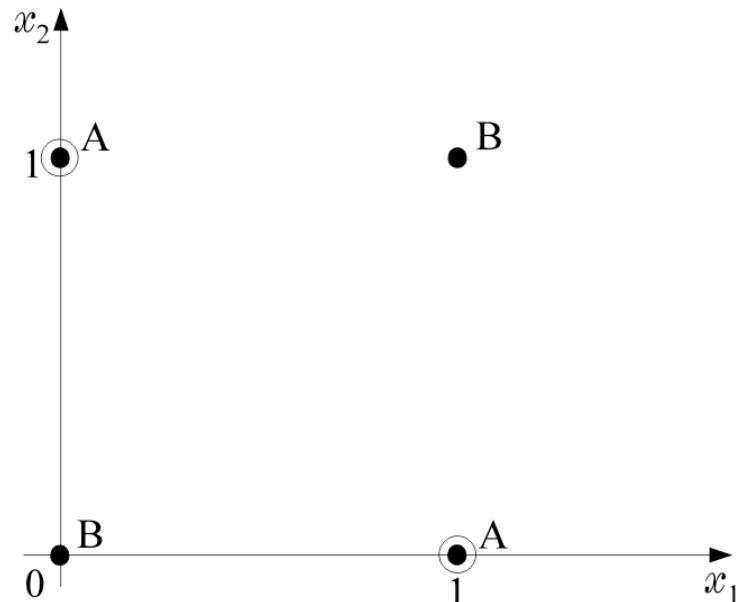
The XOR Problem

5

Probability & Bayesian Inference

- XOR is not linearly separable.

| x_1 | x_2 | XOR | Class |
|-------------------------|-------------------------|------------|--------------|
| 0 | 0 | 0 | B |
| 0 | 1 | 1 | A |
| 1 | 0 | 1 | A |
| 1 | 1 | 0 | B |



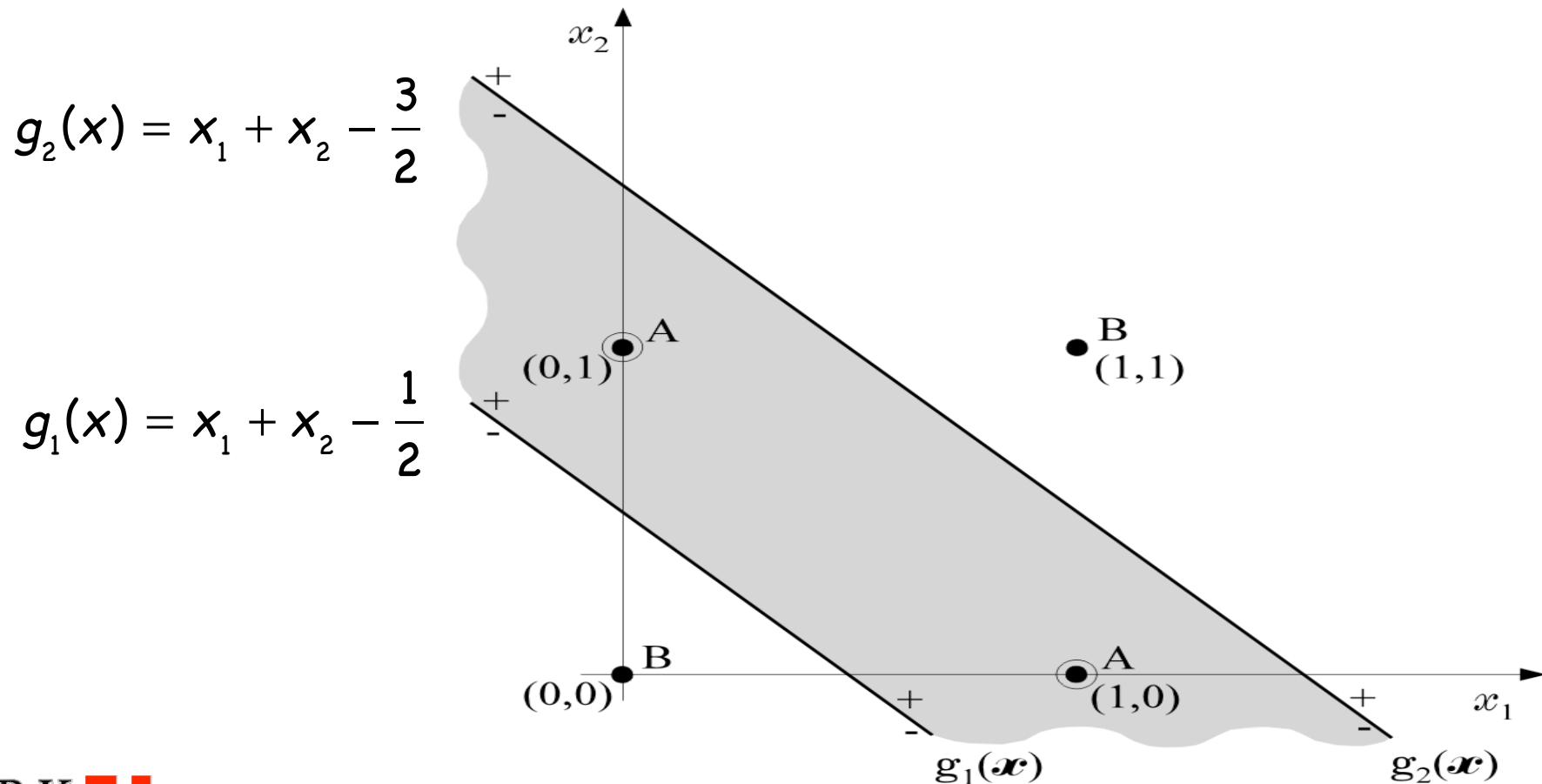
- How can we use linear classifiers to solve this problem?

Combining two linear classifiers

6

Probability & Bayesian Inference

- Idea: use a logical combination of two linear classifiers.



Combining two linear classifiers

7

Probability & Bayesian Inference

Let $f(x)$ be the unit step activation function:

$$f(x) = 0, \quad x < 0$$

$$f(x) = 1, \quad x \geq 0$$

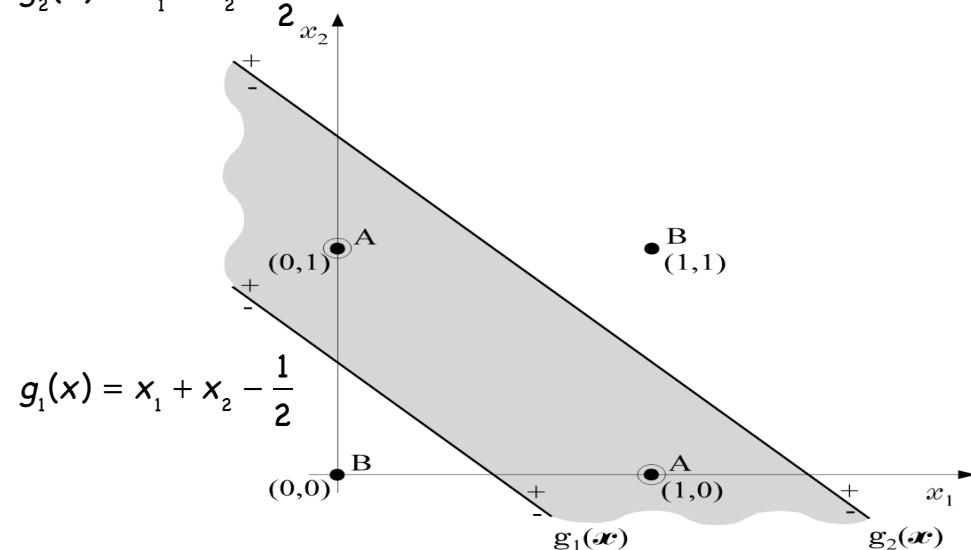
Observe that the classification problem is then solved by

$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

$$g_2(x) = x_1 + x_2 - \frac{3}{2}$$

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$



Combining two linear classifiers

8

Probability & Bayesian Inference

- This calculation can be implemented sequentially:
 1. Compute y_1 and y_2 from x_1 and x_2 .
 2. Compute the decision from y_1 and y_2 .
- Each layer in the sequence consists of one or more linear classifications.
- This is therefore a two-layer perceptron.

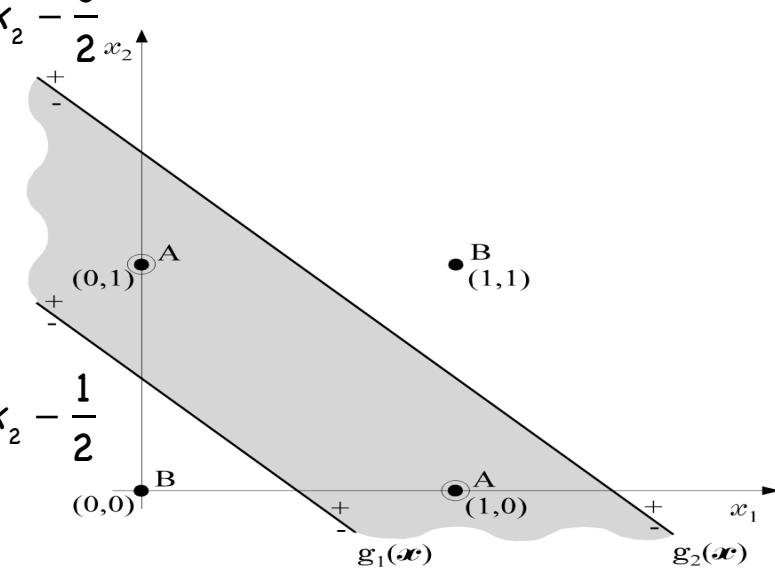
$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$

$$g_2(x) = x_1 + x_2 - \frac{3}{2}$$

$$g_1(x) = x_1 + x_2 - \frac{1}{2}$$



The Two-Layer Perceptron

9

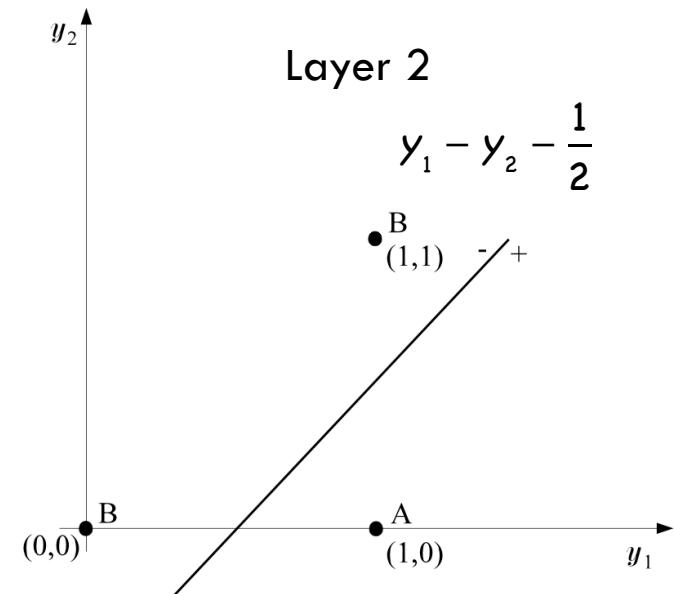
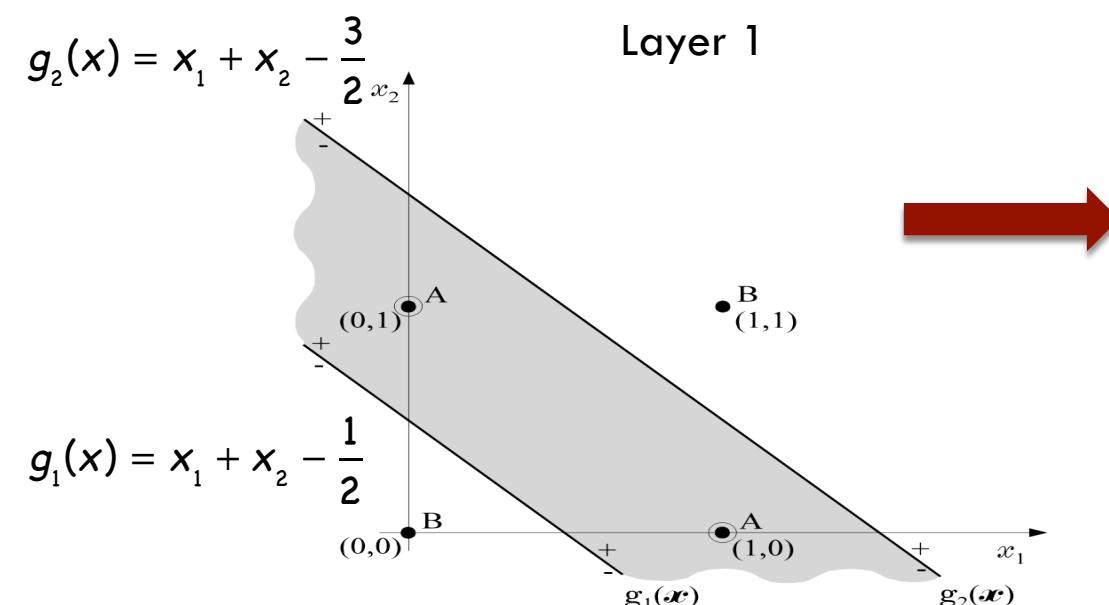
Probability & Bayesian Inference

| Layer 1 | | | | Layer 2 |
|---------|-------|-------|-------|---------|
| x_1 | x_2 | y_1 | y_2 | |
| 0 | 0 | 0(-) | 0(-) | B(0) |
| 0 | 1 | 1(+) | 0(-) | A(1) |
| 1 | 0 | 1(+) | 0(-) | A(1) |
| 1 | 1 | 1(+) | 1(+) | B(0) |

$$f\left(y_1 - y_2 - \frac{1}{2}\right)$$

where

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$



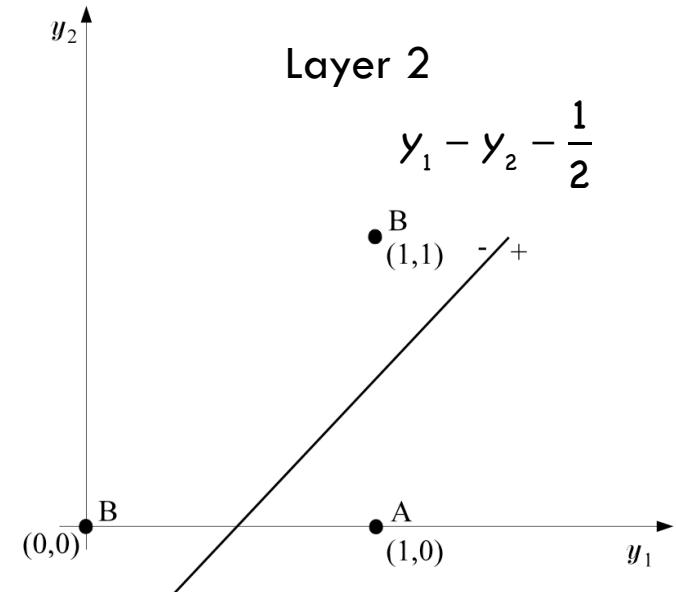
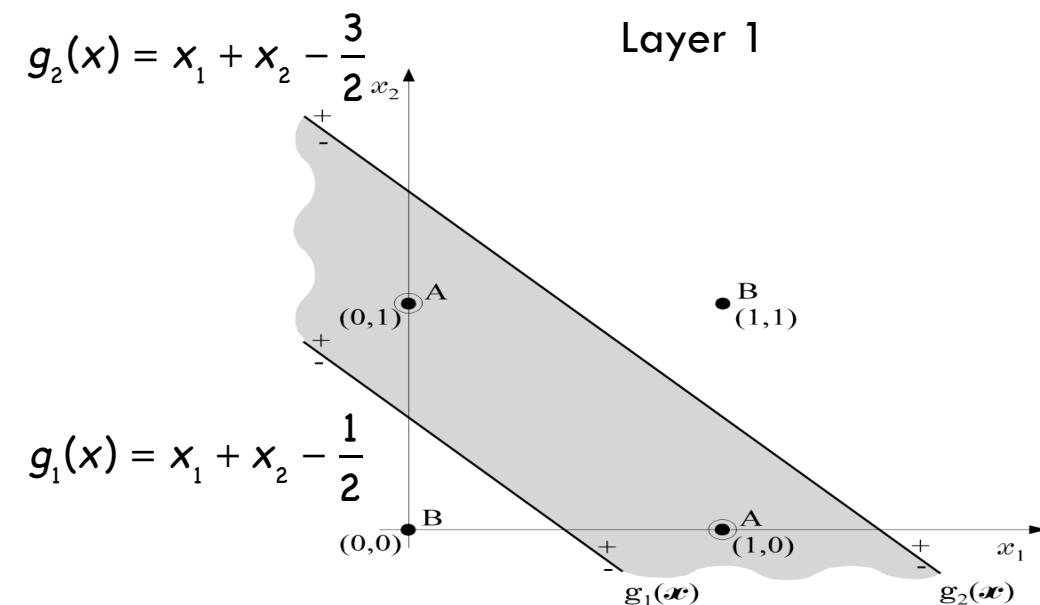
The Two-Layer Perceptron

10

Probability & Bayesian Inference

- The first layer performs a nonlinear mapping that makes the data linearly separable.

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$



The Two-Layer Perceptron Architecture

11

Probability & Bayesian Inference

Input Layer

$$x_1 \diamond \xrightarrow{1}$$

1

$$x_2 \diamond \xrightarrow{1}$$

1

Hidden Layer

$$g_1(x) = x_1 + x_2 - \frac{1}{2}$$

$$\text{circle} \diamond -\frac{1}{2}$$

$$g_2(x) = x_1 + x_2 - \frac{3}{2}$$

$$\text{circle} \diamond -\frac{3}{2}$$

Output Layer

$$y_1 - y_2 - \frac{1}{2}$$

$$\text{circle} \diamond$$

$$-\frac{1}{2}$$

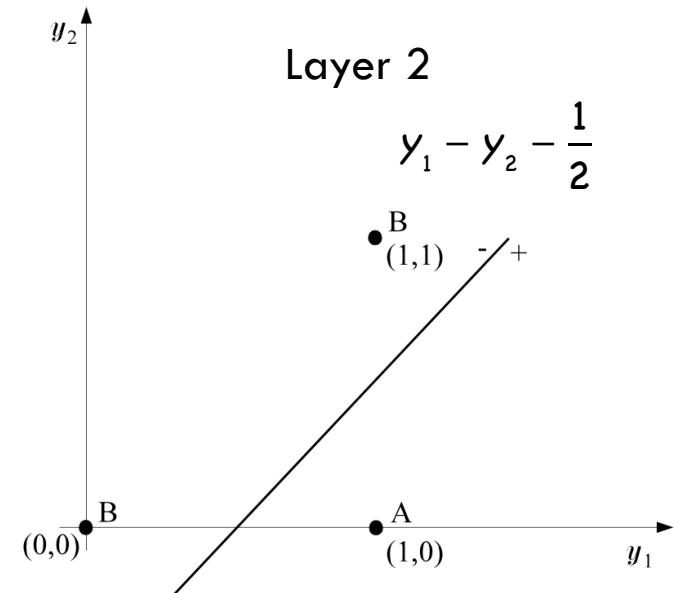
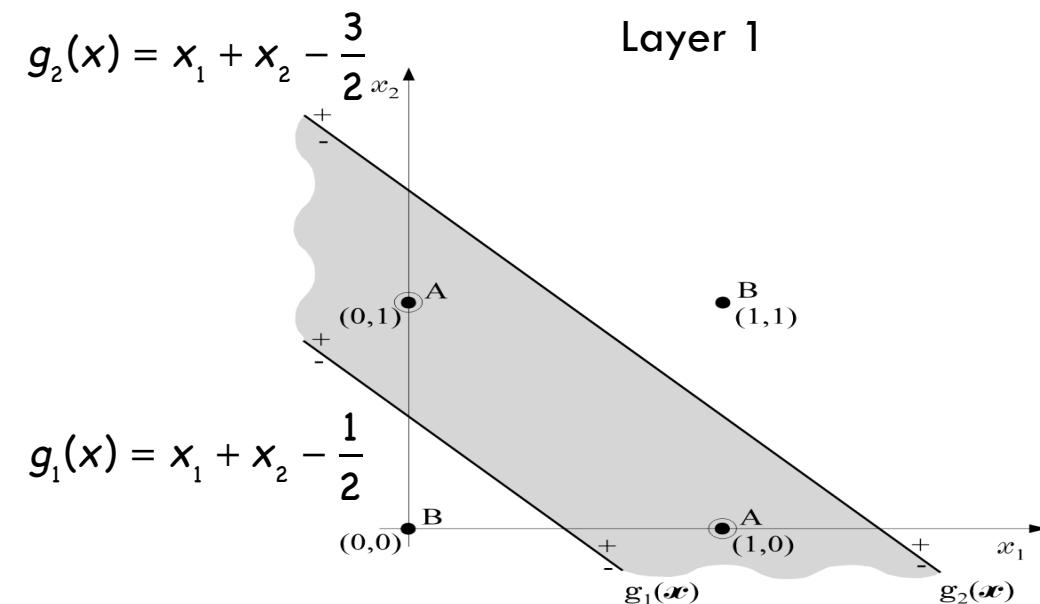
The Two-Layer Perceptron

12

Probability & Bayesian Inference

- Note that the hidden layer maps the plane onto the vertices of a unit square.

$$y_1 = f(g_1(x)) \text{ and } y_2 = f(g_2(x))$$

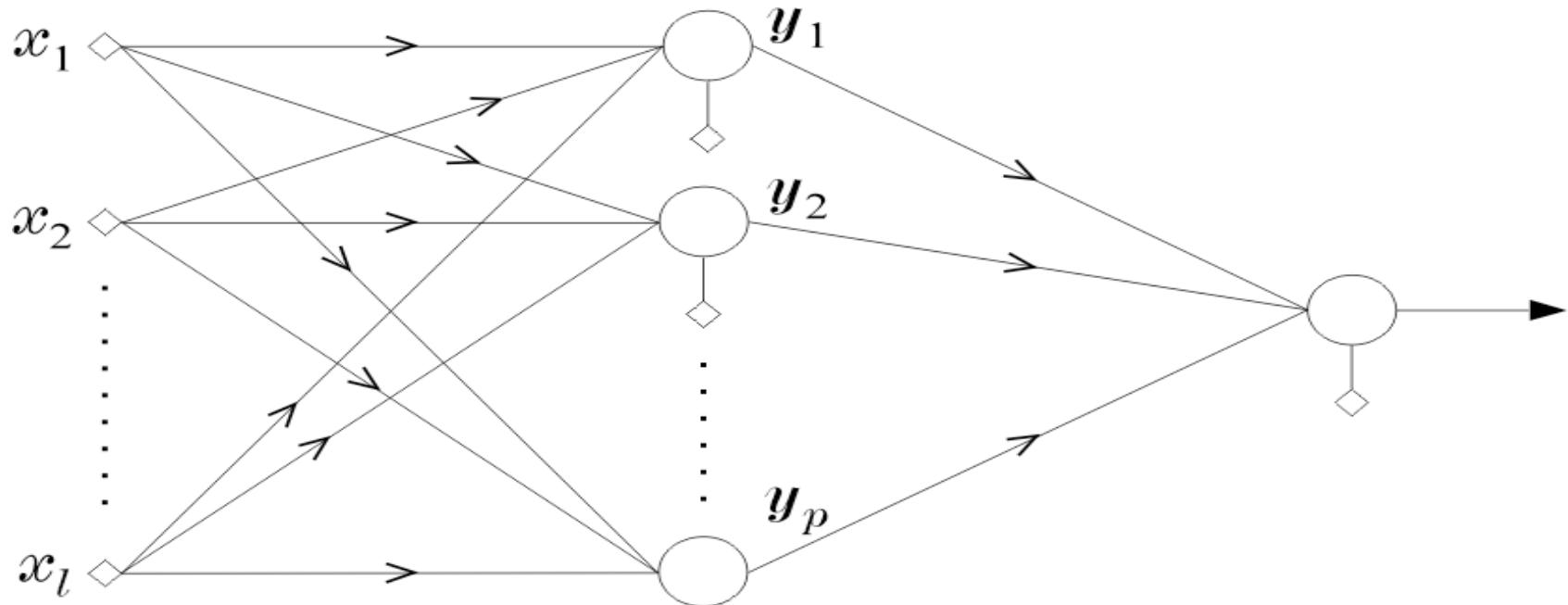


Higher Dimensions

13

Probability & Bayesian Inference

- Each hidden unit realizes a hyperplane discriminant function.
- The output of each hidden unit is 0 or 1 depending upon the location of the input vector relative to the hyperplane.



$$\underline{x} \in R^l$$

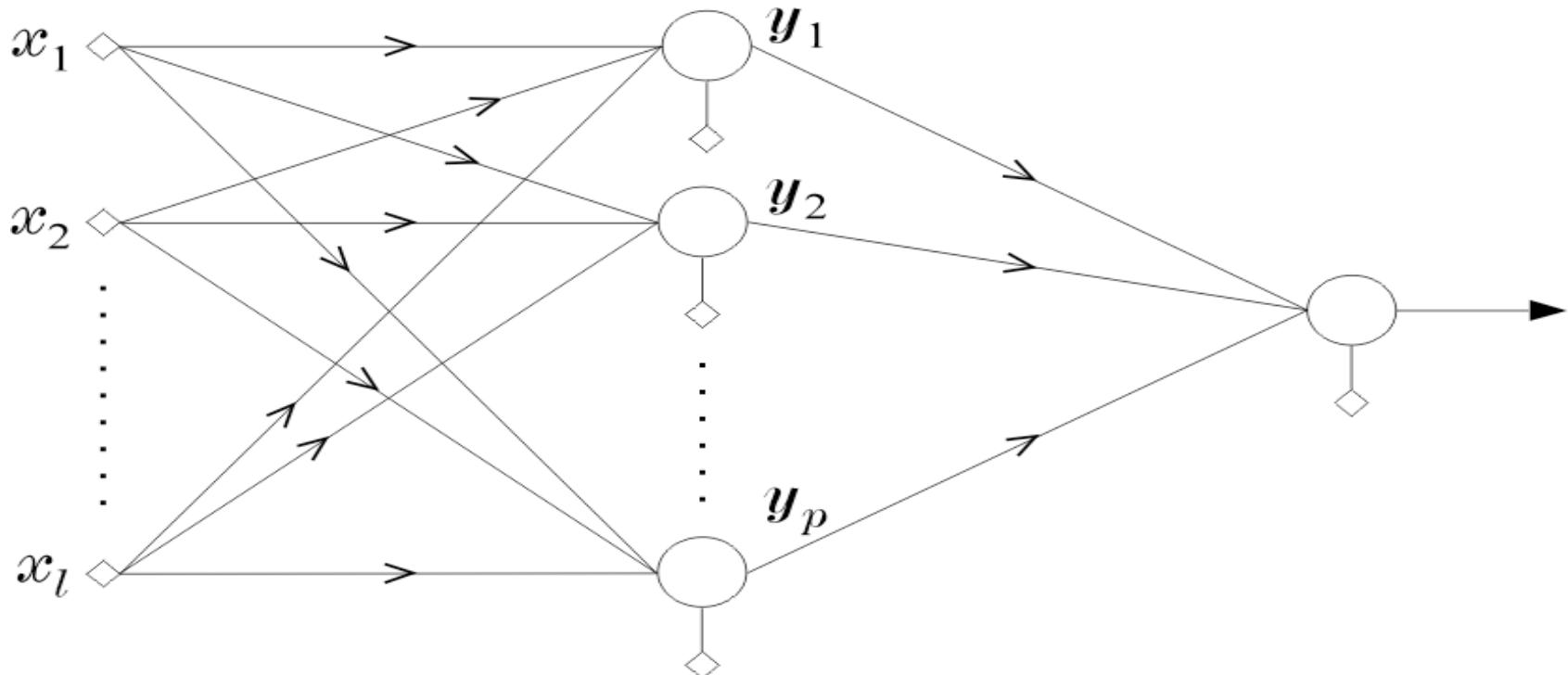
$$\underline{x} \rightarrow \underline{y} = [y_1, \dots, y_p]^T, y_i \in \{0, 1\} \quad i = 1, 2, \dots, p$$

Higher Dimensions

14

Probability & Bayesian Inference

- Together, the hidden units map the input onto the vertices of a p -dimensional unit hypercube.



$$\underline{x} \in R^l$$

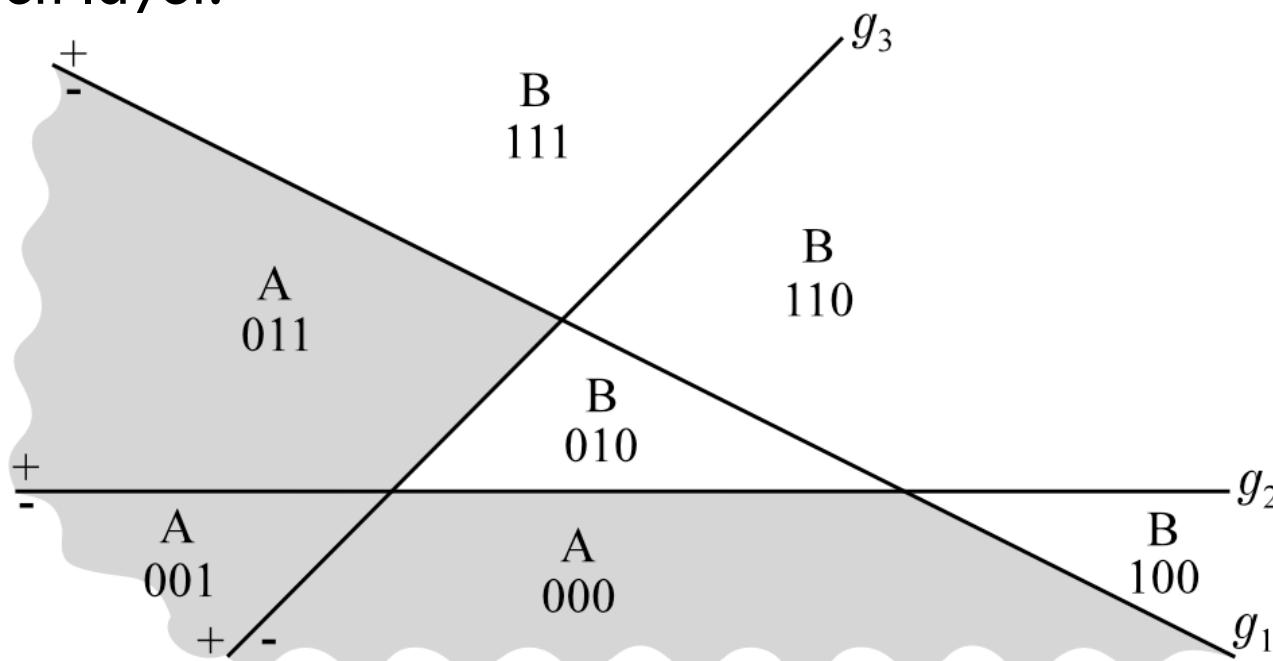
$$\underline{x} \rightarrow \underline{y} = [y_1, \dots, y_p]^T, y_i \in \{0, 1\} \quad i = 1, 2, \dots, p$$

Two-Layer Perceptron

15

Probability & Bayesian Inference

- These p hyperplanes partition the l -dimensional input space into polyhedral regions
- Each region corresponds to a different vertex of the p -dimensional hypercube represented by the outputs of the hidden layer.

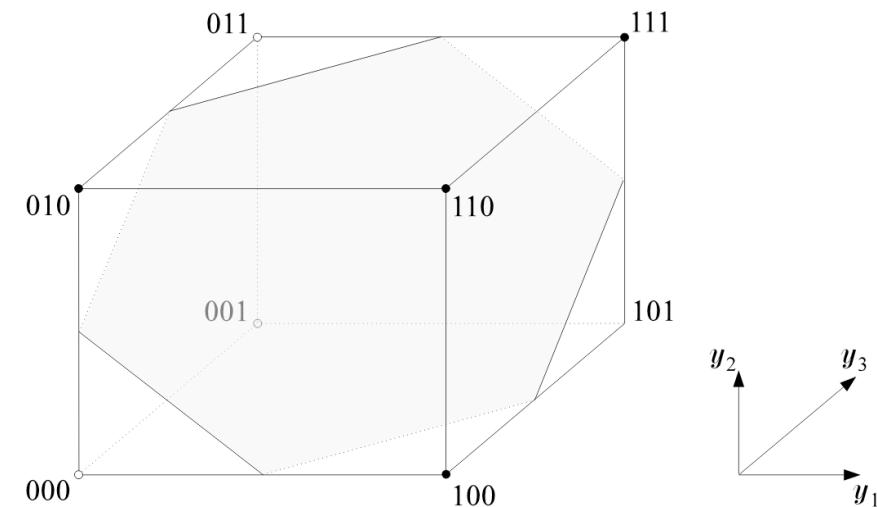
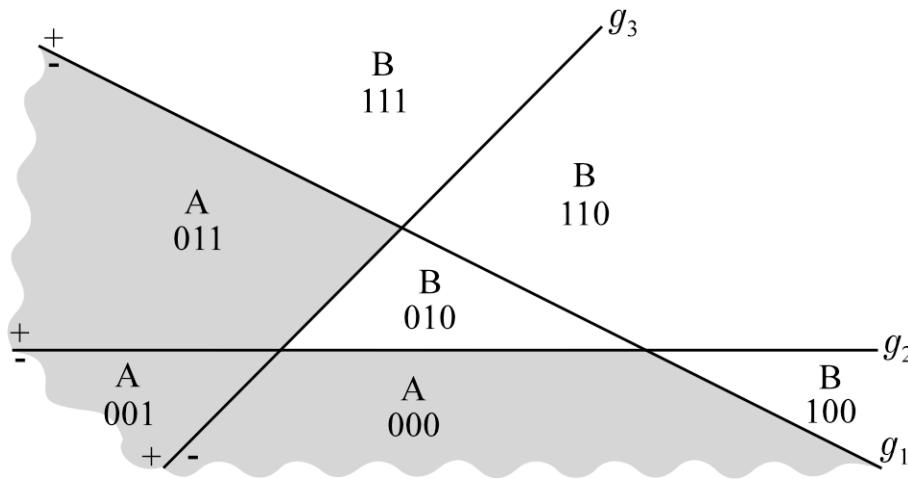


Two-Layer Perceptron

16

Probability & Bayesian Inference

- In this example, the vertex $(0, 0, 1)$ corresponds to the region of the input space where:
 - $g_1(x) < 0$
 - $g_2(x) < 0$
 - $g_3(x) > 0$



Limitations of a Two-Layer Perceptron

17

Probability & Bayesian Inference

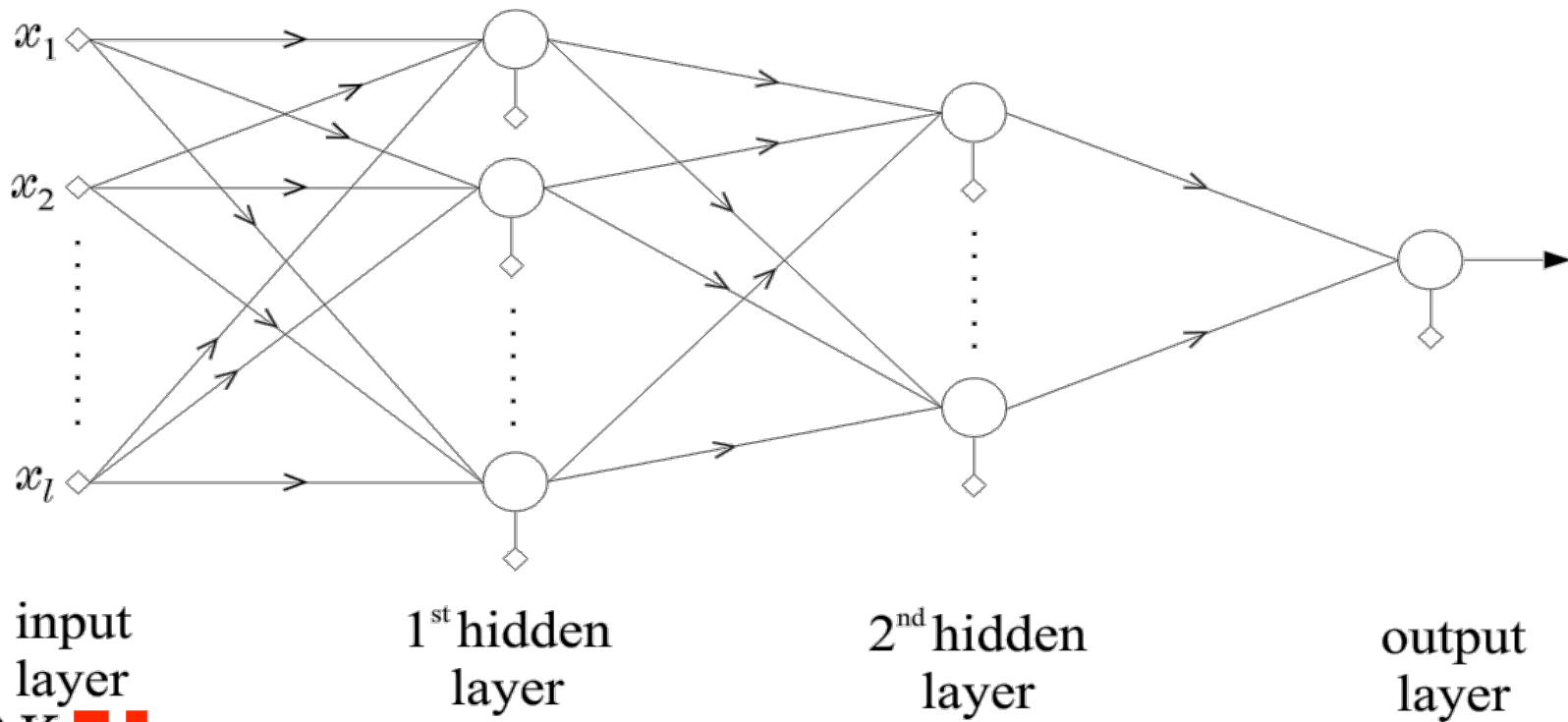
- The output neuron realizes a hyperplane in the transformed space that partitions the p vertices into two sets.
- Thus, the two layer perceptron has the capability to classify vectors into **classes that consist of unions of polyhedral regions.**
- But **NOT ANY** union. It depends on the relative position of the corresponding vertices.
- How can we solve this problem?

The Three-Layer Perceptron

18

Probability & Bayesian Inference

- Suppose that Class A consists of the union of K polyhedra in the input space.
- Use K neurons in the 2nd hidden layer.
- Train each to classify one region as positive, the rest negative.
- Now use an output neuron that implements the OR function.

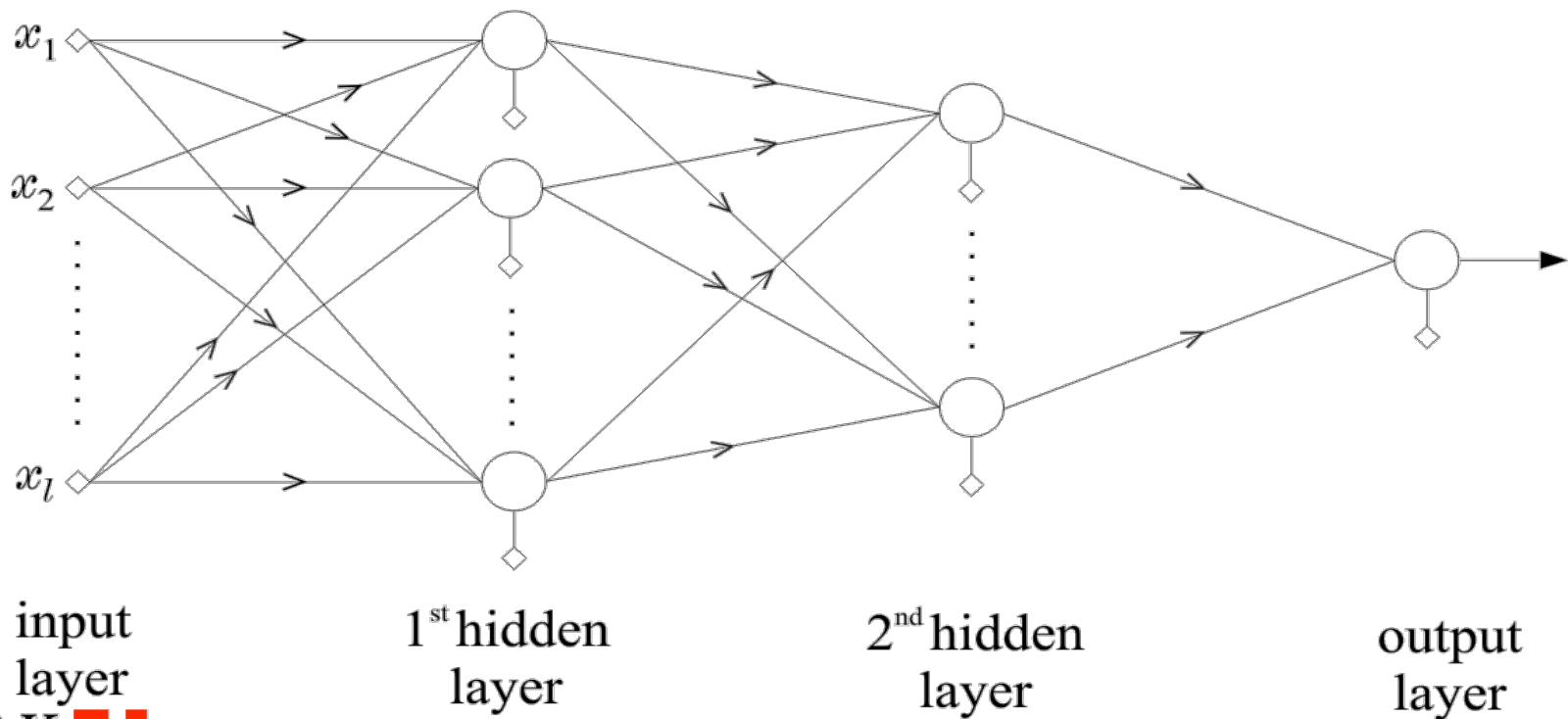


The Three-Layer Perceptron

19

Probability & Bayesian Inference

- Thus the three-layer perceptron can separate classes resulting from any union of polyhedral regions in the input space.

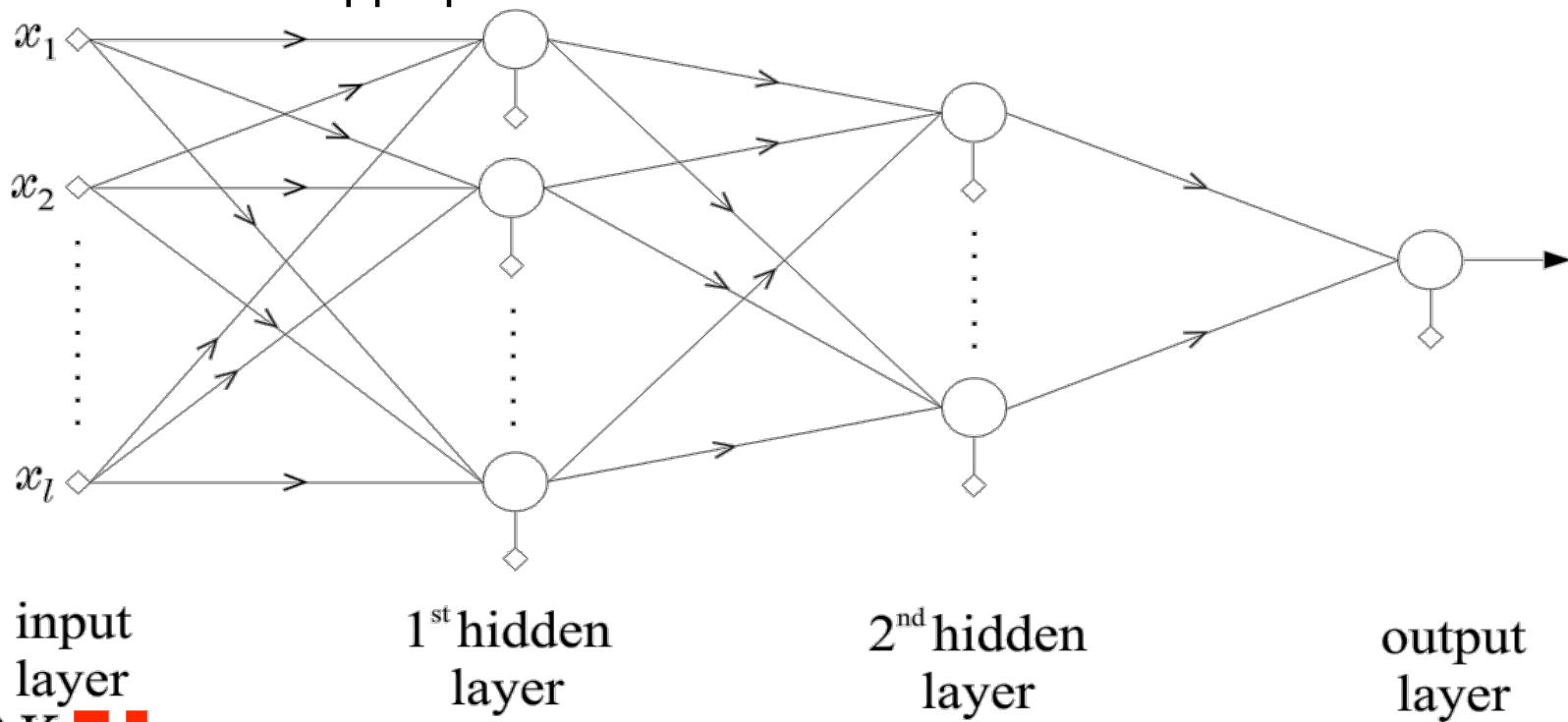


The Three-Layer Perceptron

20

Probability & Bayesian Inference

- The first layer of the network forms the **hyperplanes** in the input space.
- The second layer of the network forms the **polyhedral regions** of the input space
- The third layer forms the appropriate **unions of these regions** and maps each to the appropriate class.



Learning Parameters

Training Data

22

Probability & Bayesian Inference

- The training data consist of N input-output pairs:

$$(\mathbf{y}(i), \mathbf{x}(i)), \quad i \in 1, \dots, N$$

where

$$\mathbf{y}(i) = \left[y_1(i), \dots, y_{k_L}(i) \right]^t$$

and

$$\mathbf{x}(i) = \left[x_1(i), \dots, x_{k_O}(i) \right]^t$$

Choosing an Activation Function

23

Probability & Bayesian Inference

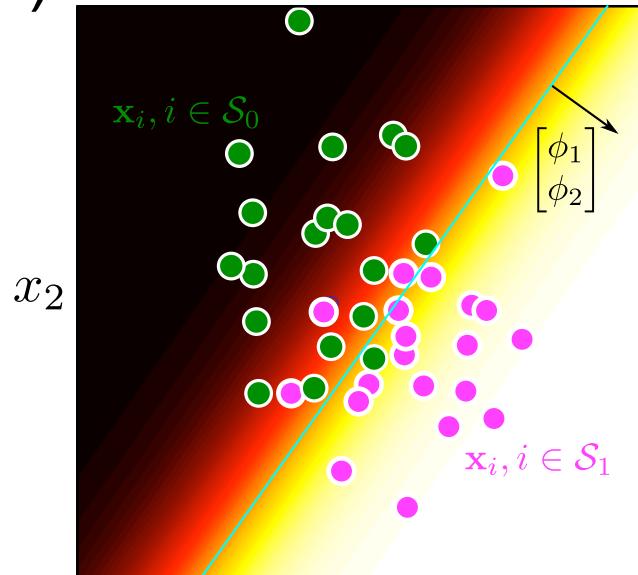
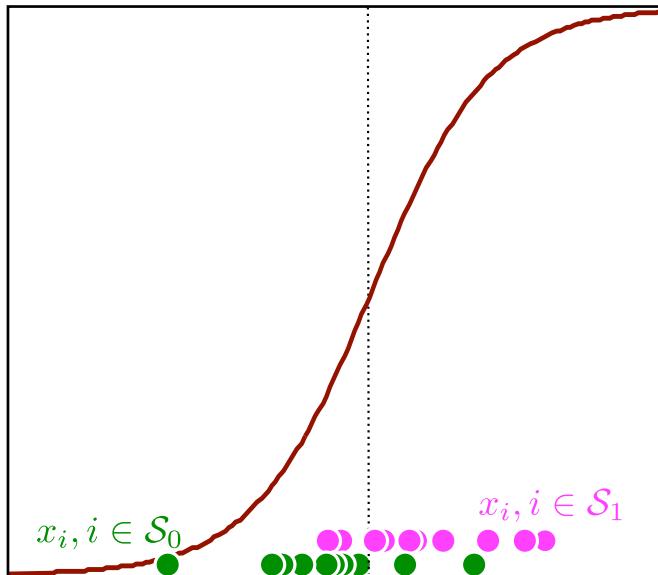
- The unit step activation function means that the error rate of the network is a discontinuous function of the weights.
- This makes it difficult to learn optimal weights by minimizing the error.
- To fix this problem, we need to use a smooth activation function.
- A popular choice is the sigmoid function we used for logistic regression:

Smooth Activation Function

24

Probability & Bayesian Inference

$$f(a) = \frac{1}{1 + \exp(-a)}$$



$$\mathbf{w}^t \phi$$

 x_1

Output: Two Classes

25

Probability & Bayesian Inference

- For a binary classification problem, there is a single output node with activation function given by

$$f(a) = \frac{1}{1 + \exp(-a)}$$

- Since the output is constrained to lie between 0 and 1, it can be interpreted as the probability of the input vector belonging to Class 1.

Output: $K > 2$ Classes

26

Probability & Bayesian Inference

- For a K -class problem, we use K outputs, and the softmax function given by

$$y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

- Since the outputs are constrained to lie between 0 and 1, and sum to 1, y_k can be interpreted as the probability that the input vector belongs to Class K .

Non-Convex

27

Probability & Bayesian Inference

- Now each layer of our multi-layer perceptron is a logistic regressor.
- Recall that optimizing the weights in logistic regression results in a convex optimization problem.
- Unfortunately the cascading of logistic regressors in the multi-layer perceptron makes the problem non-convex.
- This makes it difficult to determine an exact solution.
- Instead, we typically use **gradient descent** to find a locally optimal solution to the weights.
- The specific learning algorithm is called the **backpropagation** algorithm.

End of Lecture

Nov 21, 2011

Nonlinear Classification and Regression: Outline

29

Probability & Bayesian Inference

- Multi-Layer Perceptrons
 - **The Back-Propagation Learning Algorithm**
- Generalized Linear Models
 - Radial Basis Function Networks
 - Sparse Kernel Machines
 - Nonlinear SVMs and the Kernel Trick
 - Relevance Vector Machines

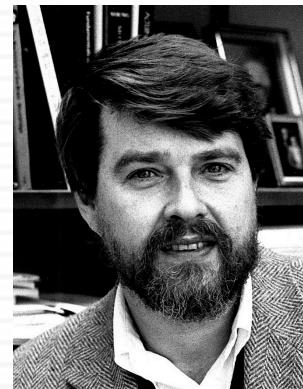
The Backpropagation Algorithm

Paul J. Werbos. Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences. PhD thesis, Harvard University, 1974

Rumelhart, David E.; Hinton, Geoffrey E., Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". *Nature* **323** (6088): 533–536.



Werbos



Rumelhart



Hinton

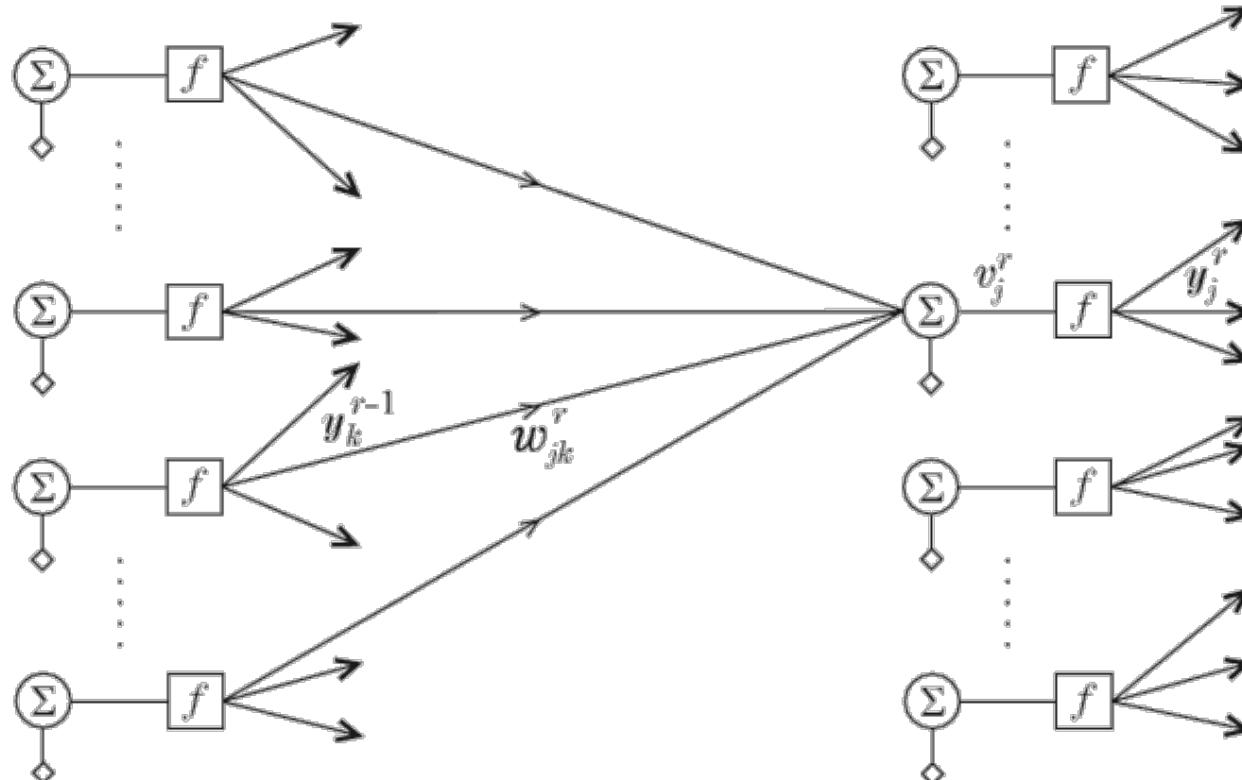
Notation

31

Probability & Bayesian Inference

- Assume a network with L layers

- k_0 nodes in the input layer.
- k_r nodes in the r 'th layer.



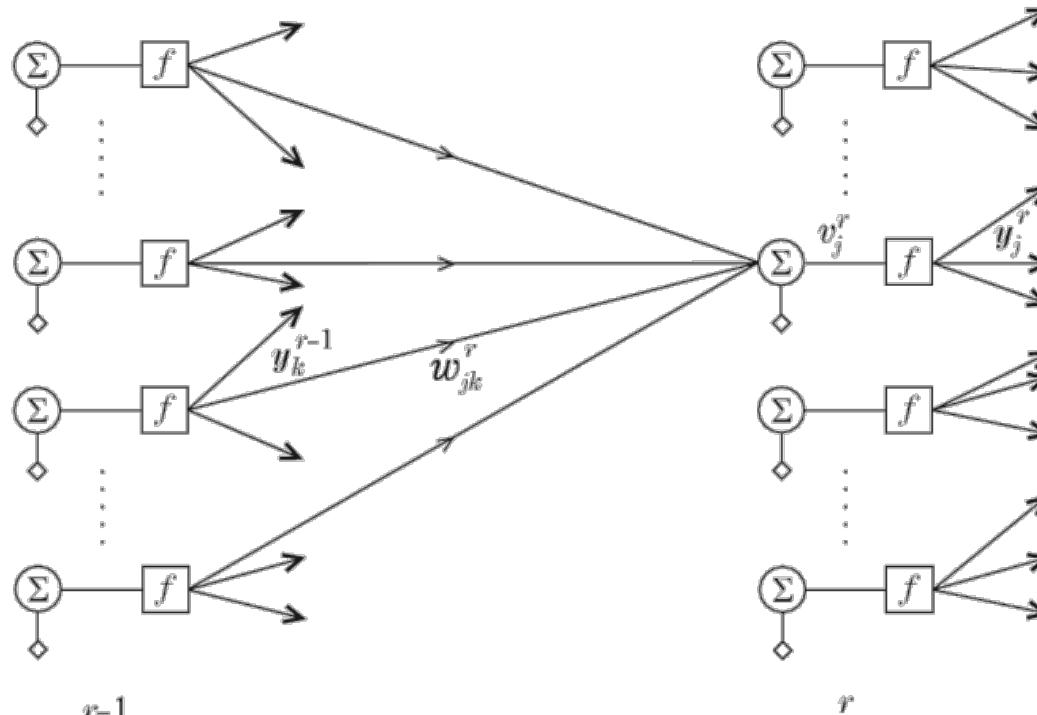
Notation

32

Probability & Bayesian Inference

Let y_k^{r-1} be the output of the k th neuron of layer $r - 1$.

Let w_{jk}^r be the weight of the synapse on the j th neuron of layer r from the k th neuron of layer $r - 1$.

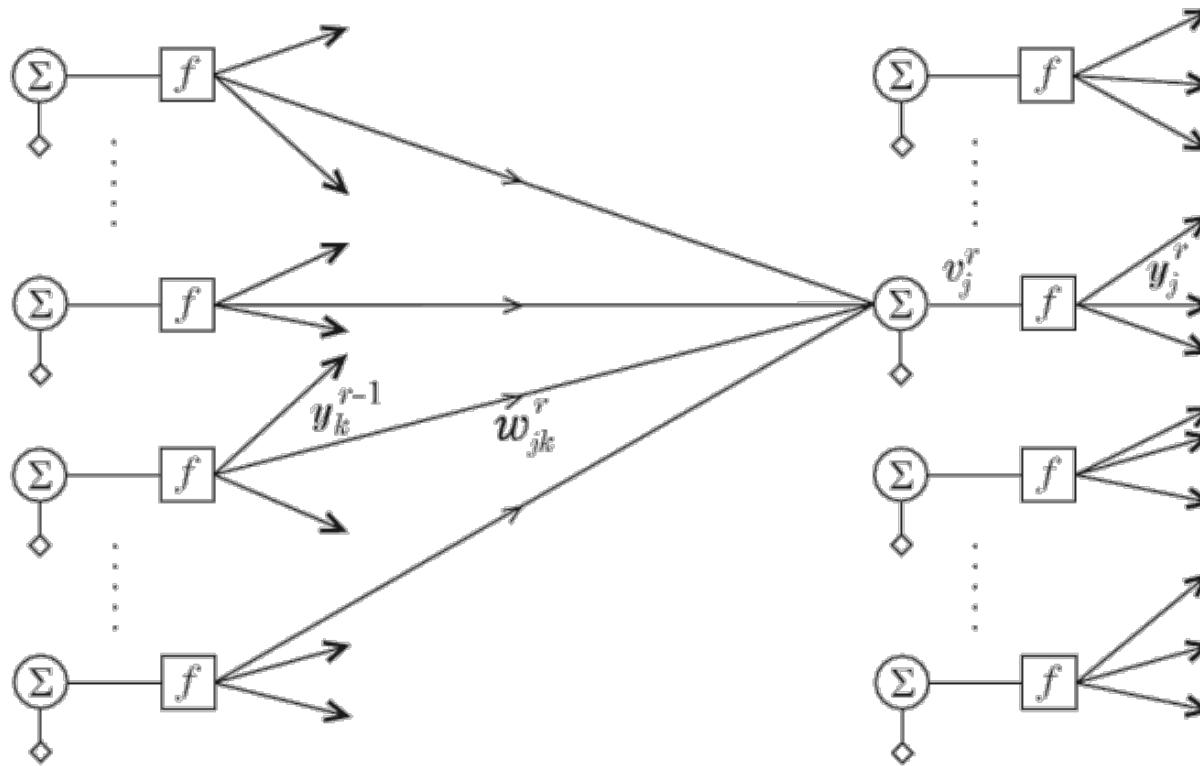


Input

33

Probability & Bayesian Inference

$$y_k^0(i) = x_k(i), \ k = 1, \dots, k_0$$

 $r-1$ r

Notation

34

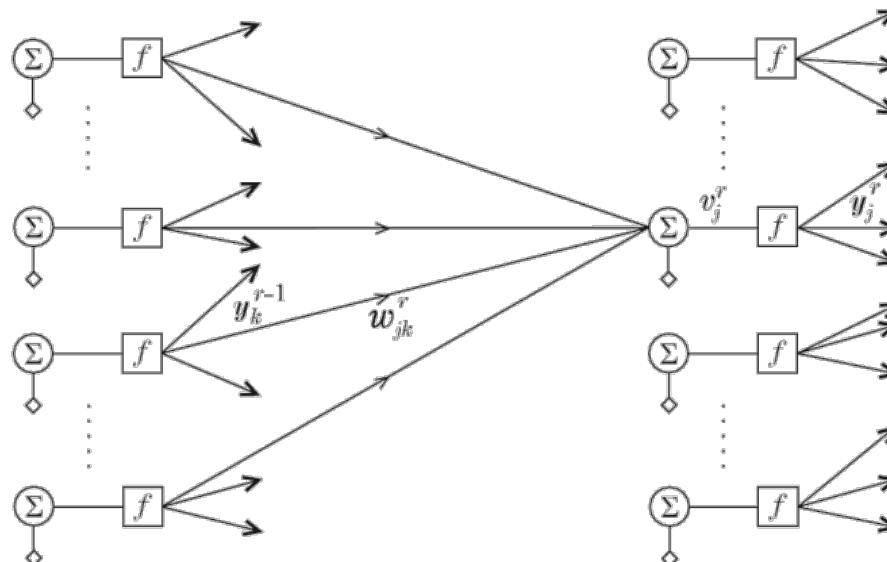
Probability & Bayesian Inference

Let v_j^r be the total input to the j th neuron of layer r :

$$v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i) = \sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)$$

where we define $y_0^r(i) = +1$ to incorporate the bias term.

Then $y_j^r(i) = f(v_j^r(i)) = f\left(\sum_{k=0}^{k_{r-1}} w_{jk}^r y_k^{r-1}(i)\right)$



Cost Function

35

Probability & Bayesian Inference

- A common cost function is the squared error:

$$\mathcal{J} = \sum_{i=1}^N \varepsilon(i)$$

where $\varepsilon(i) \triangleq \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (y_m(i) - \hat{y}_m(i))^2$

and

$\hat{y}_m(i) = y_k^r(i)$ is the output of the network.

Cost Function

36

Probability & Bayesian Inference

- To summarize, the error for input i is given by

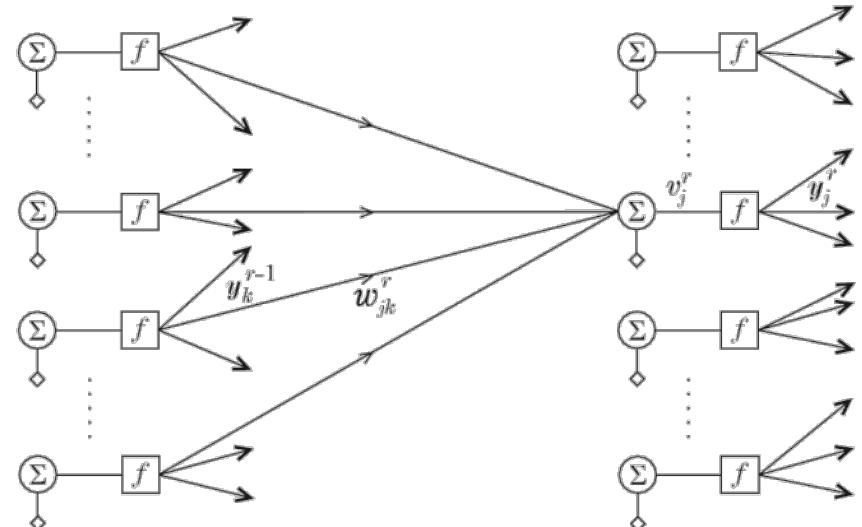
$$\varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

where $\hat{y}_m(i) = y_k^r(i)$ is the output of the output layer
and each layer is related to the previous layer through

$$y_j^r(i) = f(v_j^r(i))$$

and

$$v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i)$$



Gradient Descent

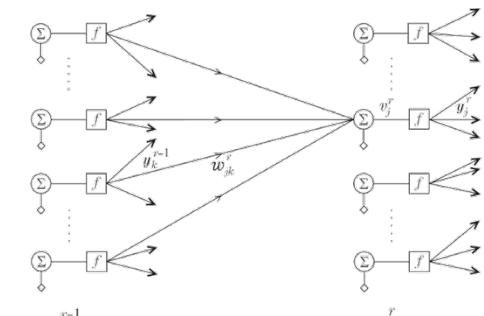
37

Probability & Bayesian Inference

$$\varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

- Gradient descent starts with an initial guess at the weights over all layers of the network.
- We then use these weights to compute the network output $\hat{y}(i)$ for each input vector $x(i)$ in the training data.
- This allows us to calculate the error $\varepsilon(i)$ for each of these inputs.
- Then, in order to minimize this error, we incrementally update the weights in the negative gradient direction:

$$w_j^r(\text{new}) = w_j^r(\text{old}) - \mu \frac{\partial J}{\partial w_j^r} = w_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial w_j^r}$$



Gradient Descent

38

Probability & Bayesian Inference

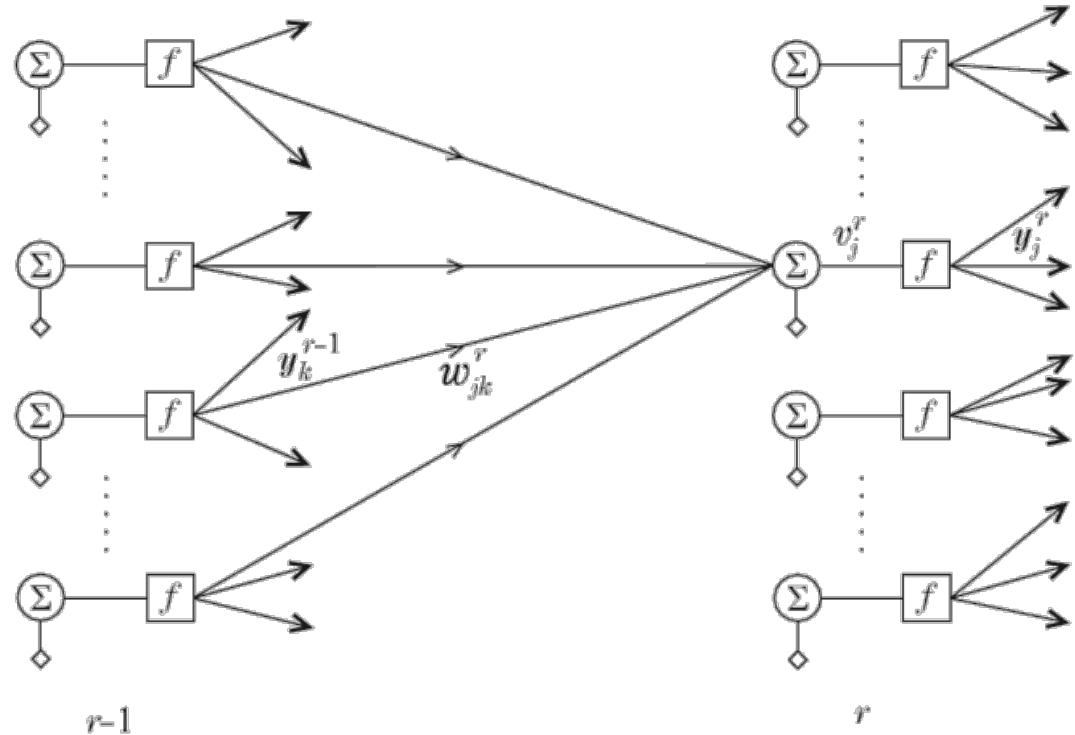
□ Since $v_j^r(i) = (\mathbf{w}_j^r)^t \mathbf{y}^{r-1}(i)$,

the influence of the j th weight of the r th layer on the error can be expressed as:

$$\begin{aligned}\frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} &= \frac{\partial \varepsilon(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r} \\ &= \delta_j^r(i) \mathbf{y}^{r-1}(i)\end{aligned}$$

where

$$\delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$



$r-1$

r

Gradient Descent

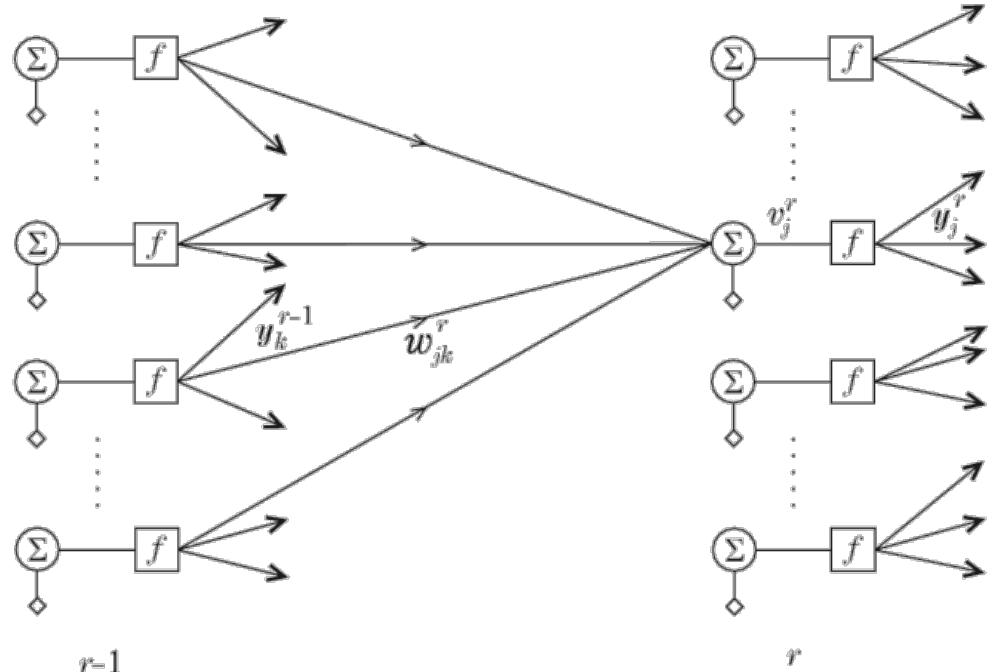
39

Probability & Bayesian Inference

$$\frac{\partial \varepsilon(i)}{\partial w_j^r} = \delta_j^r(i) y^{r-1}(i),$$

where

$$\delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$



For an intermediate layer r ,
we cannot compute $\delta_j^r(i)$ directly.

However, $\delta_j^r(i)$ can be computed inductively,
starting from the output layer.

Backpropagation: The Output Layer

40

Probability & Bayesian Inference

$$\frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i), \text{ where } \delta_j^r(i) \triangleq \frac{\partial \varepsilon(i)}{\partial v_j^r(i)}$$

$$\text{and } \varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

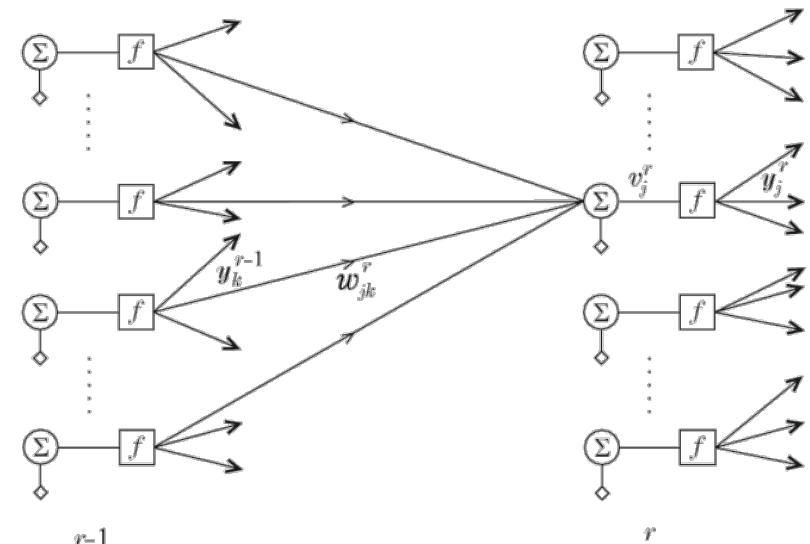
Recall that $\hat{y}_m(i) = y_j^L(i) = f(v_j^L(i))$

Thus at the output layer we have

$$\delta_j^L(i) = \frac{\partial \varepsilon(i)}{\partial v_j^L(i)} = \frac{\partial \varepsilon(i)}{\partial e_j^L(i)} \frac{\partial e_j^L(i)}{\partial v_j^L(i)} = e_j^L(i) f'(v_j^L(i))$$

$$f(a) = \frac{1}{1 + \exp(-a)} \rightarrow f'(a) = f(a)(1 - f(a))$$

$$\boxed{\delta_j^L(i) = e_j^L(i) f(v_j^L(i)) (1 - f(v_j^L(i)))}$$



Backpropagation: Hidden Layers

41

Probability & Bayesian Inference

- Observe that the dependence of the error on the total input to a neuron in a previous layer can be expressed in terms of the dependence on the total input of neurons in the following layer:

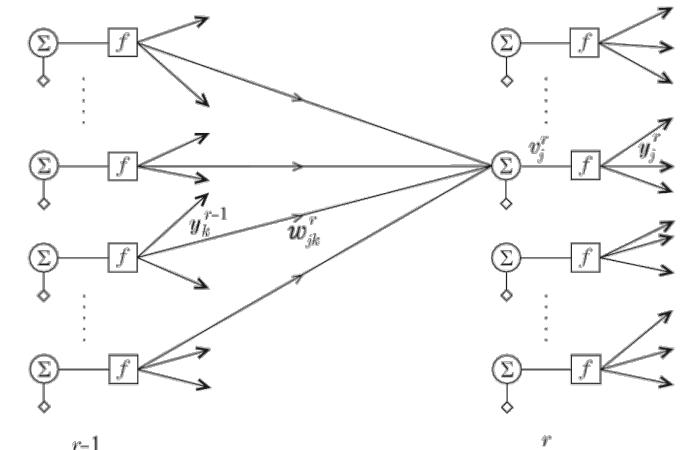
$$\delta_j^{r-1}(i) = \frac{\partial \varepsilon(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^k \frac{\partial \varepsilon(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^k \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)}$$

where $v_k^r(i) = \sum_{m=0}^{k-1} w_{km}^r y_m^r(i) = \sum_{m=0}^{k-1} w_{km}^r f(v_m^{r-1}(i))$

Thus we have $\frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = w_{kj}^r f'(v_j^{r-1}(i))$

and so $\delta_j^{r-1}(i) = \frac{\partial \varepsilon(i)}{\partial v_j^{r-1}(i)} = f'(v_j^{r-1}(i)) \sum_{k=1}^k \delta_k^r(i) w_{kj}^r = f(v_j^L(i)) (1 - f(v_j^L(i))) \sum_{k=1}^k \delta_k^r(i) w_{kj}^r$

Thus once the $\delta_k^r(i)$ are determined they can be propagated backward to calculate $\delta_j^{r-1}(i)$ using this inductive formula.



Backpropagation: Summary of Algorithm

42

Probability & Bayesian Inference

1. Initialization

- Initialize all weights with small random values

2. Forward Pass

- For each input vector, run the network in the forward direction, calculating:

$$v_j^r(i) = (\mathbf{w}_j^r)^T \mathbf{y}^{r-1}(i); \quad \mathbf{y}_j^r(i) = f(v_j^r(i))$$

$$\text{and finally } \varepsilon(i) = \frac{1}{2} \sum_{m=1}^{k_L} (e_m(i))^2 = \frac{1}{2} \sum_{m=1}^{k_L} (\hat{y}_m(i) - y_m(i))^2$$

3. Backward Pass

- Starting with the output layer, use our inductive formula to compute the $\delta_j^{r-1}(i)$:

- Output Layer (Base Case): $\delta_j^L(i) = e_j^L(i) f'(v_j^L(i))$

- Hidden Layers (Inductive Case): $\delta_j^{r-1}(i) = f'(v_j^{r-1}(i)) \sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r$

4. Update Weights

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

Batch vs Online Learning

43

Probability & Bayesian Inference

- As described, on each iteration backprop updates the weights based upon all of the training data.
This is called **batch learning**.

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \sum_{i=1}^N \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

- An alternative is to update the weights after each training input has been processed by the network, based only upon the error for that input. This is called **online learning**.

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} \quad \text{where } \frac{\partial \varepsilon(i)}{\partial \mathbf{w}_j^r} = \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

Batch vs Online Learning

44

Probability & Bayesian Inference

- One advantage of batch learning is that averaging over all inputs when updating the weights should lead to smoother convergence.
- On the other hand, the randomness associated with online learning might help to prevent convergence toward a local minimum.
- Changing the order of presentation of the inputs from epoch to epoch may also improve results.

Remarks

45

Probability & Bayesian Inference

- Local Minima
 - The objective function is in general non-convex, and so the solution may not be globally optimal.
- Stopping Criterion
 - Typically stop when the change in weights or the change in the error function falls below a threshold.
- Learning Rate
 - The speed and reliability of convergence depends on the learning rate μ .

Nonlinear Classification and Regression: Outline

46

Probability & Bayesian Inference

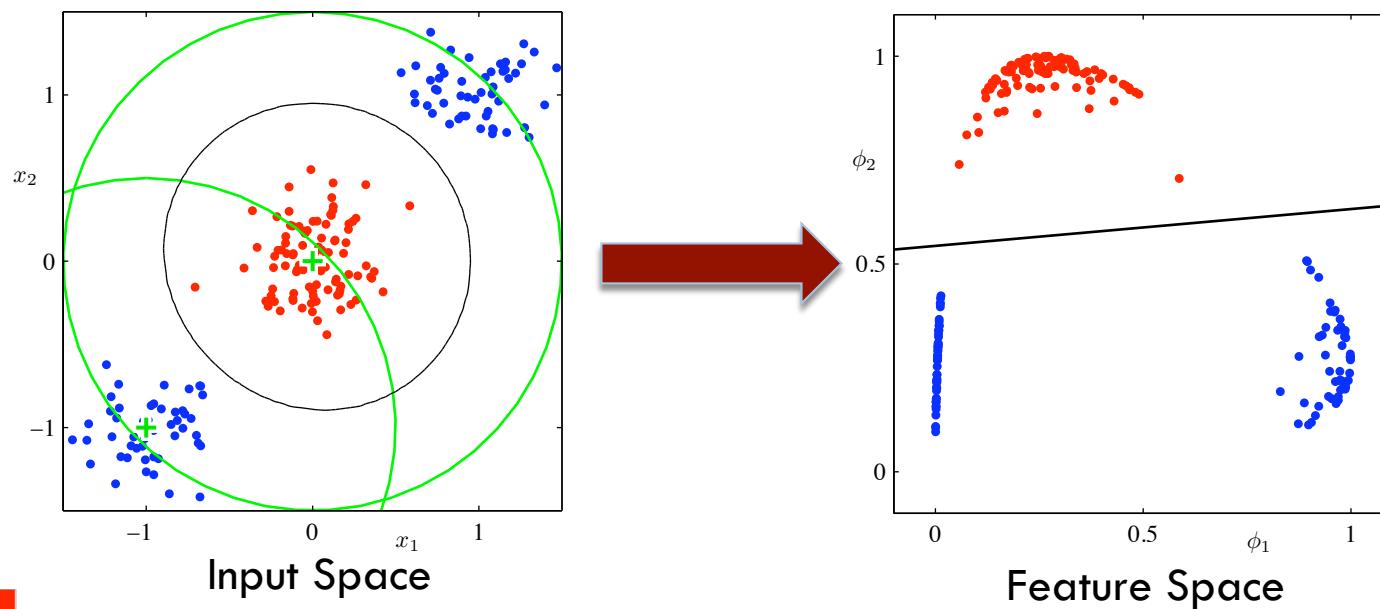
- Multi-Layer Perceptrons
 - The Back-Propagation Learning Algorithm
- Generalized Linear Models
 - Radial Basis Function Networks
 - Sparse Kernel Machines
 - Nonlinear SVMs and the Kernel Trick
 - Relevance Vector Machines

Generalizing Linear Classifiers

47

Probability & Bayesian Inference

- One way of tackling problems that are not linearly separable is to transform the input in a nonlinear fashion prior to applying a linear classifier.
- The result is that decision boundaries that are linear in the resulting feature space may be highly nonlinear in the original input space.



Nonlinear Basis Function Models

48

Probability & Bayesian Inference

- 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.

Nonlinear basis functions for classification

49

Probability & Bayesian Inference

- In the context of classification, the discriminant function in the feature space becomes:

$$g(y(x)) = w_0 + \sum_{i=1}^M w_i y_i(x) = w_0 + \sum_{i=1}^M w_i \phi_i(x)$$

- This formulation can be thought of as an input space approximation of the true separating discriminant function $g(x)$ using a set of interpolation functions $\phi_i(x)$.

Dimensionality

50

Probability & Bayesian Inference

- The dimensionality M of the feature space may be less than, equal to, or greater than the dimensionality D of the original input space.
 - $M < D$: This may result in a factoring out of irrelevant dimensions, reduction in the number of model parameters, and resulting improvement in generalization (reduced overlearning).
 - $M > D$: Problems that are not linearly separable in the input space may become separable in the feature space, and the probability of linear separability generally increases with the dimensionality of the feature space. Thus choosing $M \gg D$ helps to make the problem linearly separable.

Cover's Theorem

51

Probability & Bayesian Inference

- “A complex pattern-classification problem, cast in a high-dimensional space nonlinearly, is more likely to be linearly separable than in a low-dimensional space, provided that the space is not densely populated.”
— Cover, T.M. , *Geometrical and Statistical properties of systems of linear inequalities with applications in pattern recognition.*, 1965

Example

Nonlinear Classification and Regression: Outline

52

Probability & Bayesian Inference

- Multi-Layer Perceptrons
 - The Back-Propagation Learning Algorithm
- Generalized Linear Models
 - Radial Basis Function Networks
 - Sparse Kernel Machines
 - Nonlinear SVMs and the Kernel Trick
 - Relevance Vector Machines

Radial Basis Functions

53

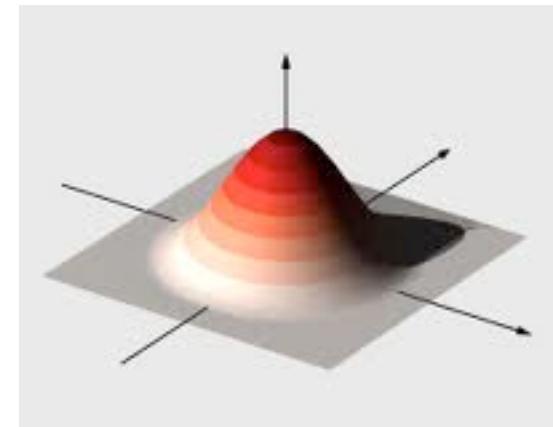
Probability & Bayesian Inference

- Consider interpolation functions (kernels) of the form

$$\phi_i(\|x - \mu_i\|)$$

- In other words, the feature value depends only upon the Euclidean distance to a ‘centre point’ in the input space.
- A commonly used RBF is the isotropic Gaussian:

$$\phi_i(x) = \exp\left(-\frac{1}{2\sigma_i^2}\|x - \mu_i\|^2\right)$$



Relation to KDE

- We can use Gaussian RBFs to approximate the discriminant function $g(x)$:

$$g(y(x)) = w_0 + \sum_{i=1}^M w_i y_i(x) = w_0 + \sum_{i=1}^M w_i \phi_i(x)$$

- where

$$\phi_i(x) = \exp\left(-\frac{1}{2\sigma_i^2}\|x - \mu_i\|^2\right)$$

- This is reminiscent of kernel density estimation, where we approximated probability densities as a normalized sum of Gaussian kernels.

Relation to KDE

55

Probability & Bayesian Inference

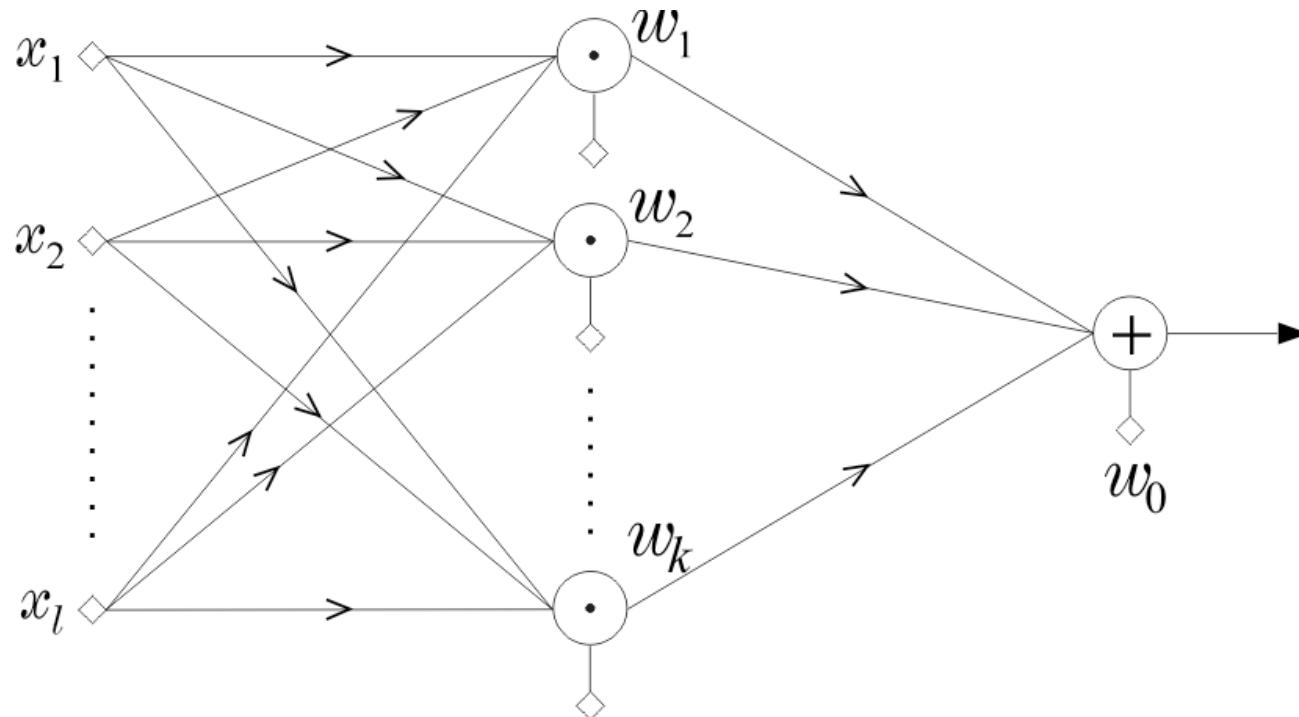
- For KDE we planted a kernel at each data point.
Thus there were N kernels.
- For RBF networks, we generally use far fewer kernels than the number of data points: $M \ll N$.
- This leads to greater efficiency and generalization.

RBF Networks

56

Probability & Bayesian Inference

- The Linear Classifier with nonlinear radial basis functions can be considered an artificial neural network where
 - The hidden nodes are nonlinear (e.g., Gaussian).
 - The output node is linear.



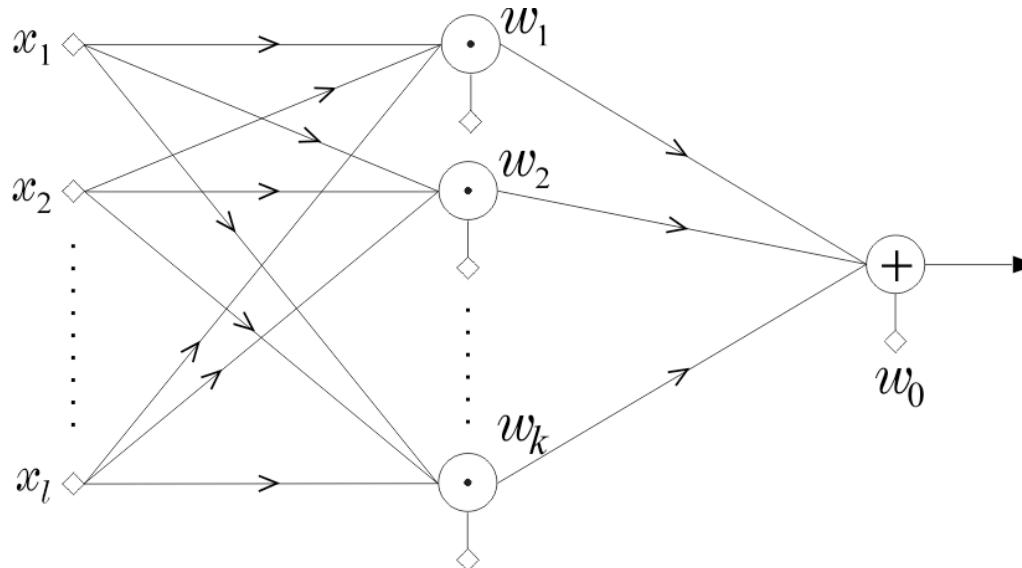
RBF Network for 2 Classes

RBF Networks vs Perceptrons

57

Probability & Bayesian Inference

- Recall that for a perceptron, the output of a hidden unit is invariant on a hyperplane.
- For an RBF, the output of a hidden unit is invariant on a circle centred on μ_i .
- Thus hidden units are global in a perceptron, but local in an RBF network.



RBF Network for 2 Classes

RBF Networks vs Perceptrons

58

Probability & Bayesian Inference

- This difference has consequences:
 - Multilayer perceptrons tend to learn slower than RBFs.
 - However, multilayer perceptrons tend to have better generalization properties, especially in regions of the input space where training data are sparse.
 - Typically, more neurons are needed for an RBF than for a multilayer perceptron to solve a given problem.

Parameters

59

Probability & Bayesian Inference

- There are two options for choosing the parameters (centres and scales) of the RBFs:

1. **Fixed.**

- For example, randomly select a subset of M of the input vectors and use these as centres. Use a common scale based upon your judgement.

2. **Learned.**

- Note that when the RBF parameters are fixed, the weights could be learned using linear classifier techniques (e.g., least squares).
- Thus the RBF parameters could be learned in an outer loop, by gradient descent.

Nonlinear Classification and Regression: Outline

60

Probability & Bayesian Inference

- Multi-Layer Perceptrons
 - The Back-Propagation Learning Algorithm
- Generalized Linear Models
 - Radial Basis Function Networks
 - **Sparse Kernel Machines**
 - **Nonlinear SVMs and the Kernel Trick**
 - Relevance Vector Machines

The Kernel Function

61

Probability & Bayesian Inference

- Recall that an SVM is the solution to the problem

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

subject to $0 \leq a_n \leq C$ and $\sum_{n=1}^N a_n t_n = 0$

- A new input \mathbf{x} is classified by computing

$$y(\mathbf{x}) = \sum_{n \in S} a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

- Where S is the set of support vectors.
- Here we introduced the kernel function $k(x, x')$, defined as

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

- This is more than a notational convenience!!

The Kernel Trick

62

Probability & Bayesian Inference

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

subject to $0 \leq a_n \leq C$ and $\sum_{n=1}^N a_n t_n = 0$

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

- Note that the basis functions and individual training vectors are no longer part of the objective function.
- Instead all we need is the kernel value (like a distance measure) for all pairs of training vectors.

The Kernel Function

63

Probability & Bayesian Inference

The kernel function $k(\mathbf{x}, \mathbf{x}')$ measures the 'similarity' of input vectors \mathbf{x} and \mathbf{x}' as an inner product in a feature space defined by the feature space mapping $\phi(\mathbf{x})$:

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

If $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ we say that the kernel is *stationary*

If $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$ we call it a *radial basis function*.

Constructing Kernels

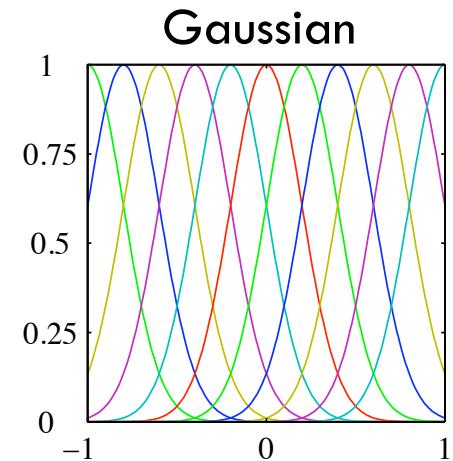
64

Probability & Bayesian Inference

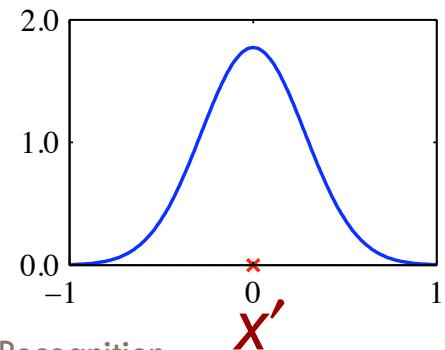
- We can construct a kernel by selecting a feature space mapping $\phi(x)$ and then defining

$$k(x, x') = \phi(x)^t \phi(x')$$

$$\phi(x')$$



$$k(x, x')$$



Constructing Kernels

65

Probability & Bayesian Inference

- Alternatively, we can construct the kernel function directly, ensuring that it corresponds to an inner product in some (possibly infinite-dimensional) feature space.

Constructing Kernels

66

Probability & Bayesian Inference

$$k(x) = \phi(x)^t \phi(x')$$

Example 1: $k(x, z) = x^t z$

Example 2: $k(x, z) = x^t z + c, c > 0$

Example 3: $k(x, z) = (x^t z)^2$

Kernel Properties

67

Probability & Bayesian Inference

- Kernels obey certain properties that make it easy to construct complex kernels from simpler ones.

Kernel Properties

68

Probability & Bayesian Inference

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ the following kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where $c > 0$, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\mathbf{x})$ is a mapping from $\mathbf{x} \rightarrow \mathbb{R}^M$, k_3 is a valid kernel on \mathbb{R}^M , A is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables such that $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ and k_a, k_b are valid kernels over their respective spaces.

Constructing Kernels

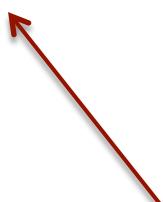
69

Probability & Bayesian Inference

□ Examples:

$$k(x, x') = (x^t x' + c)^M, c > 0 \quad (\text{Use 6.18})$$

$$k(x, x') = \exp\left(-\|x - x'\|^2 / 2\sigma^2\right) \quad (\text{Use 6.14 and 6.16.})$$

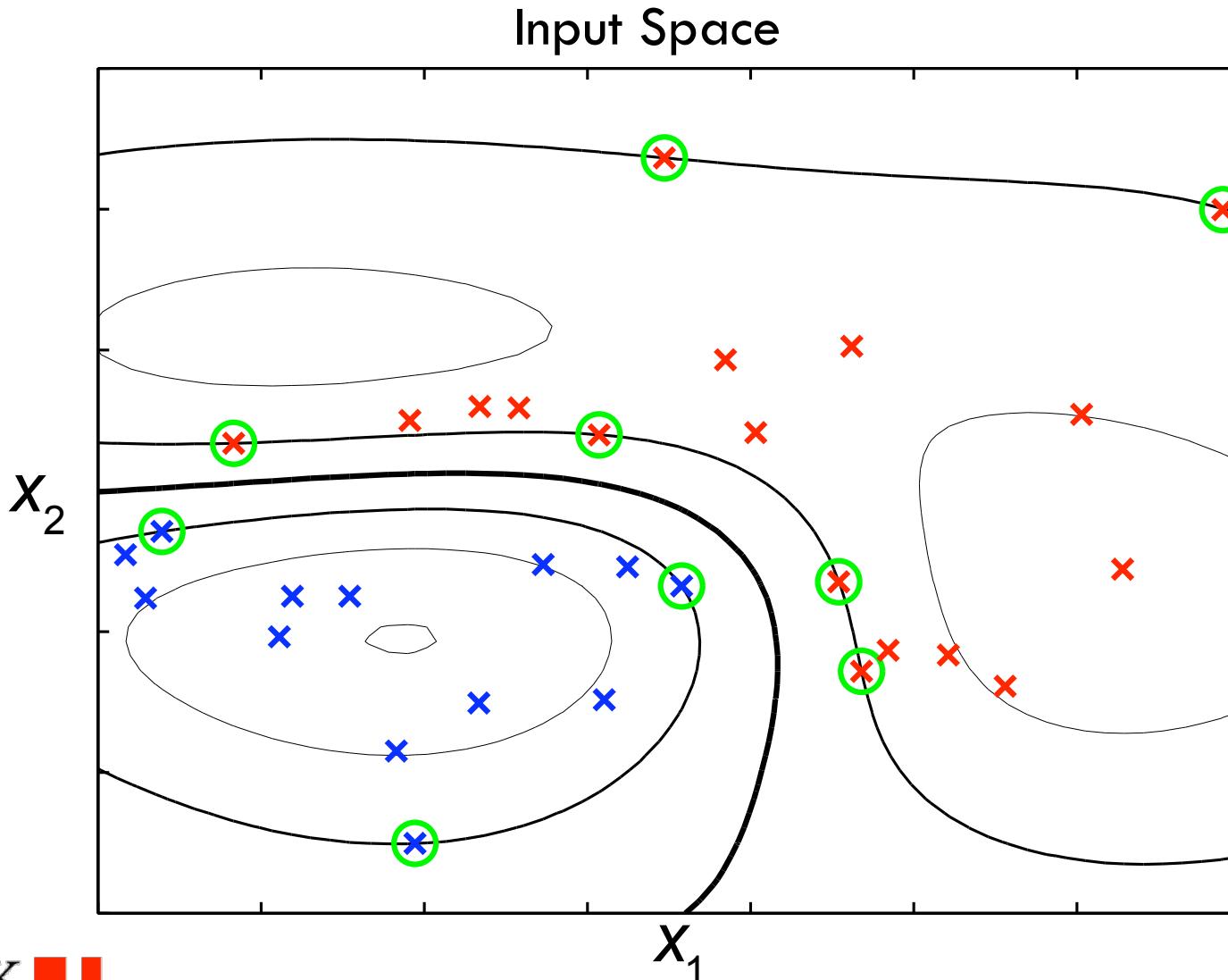


Corresponds to infinite-dimensional feature vector

Nonlinear SVM Example (Gaussian Kernel)

70

Probability & Bayesian Inference



SVMs for Regression

71

Probability & Bayesian Inference

In standard linear regression, we minimize

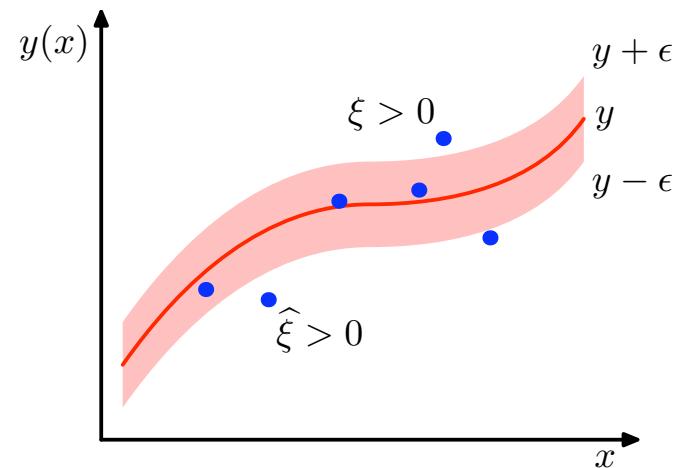
$$\frac{1}{2} \sum_{n=1}^N (y_n - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

This penalizes all deviations from the model.

To obtain sparse solutions, we replace the quadratic error function by an **ε -insensitive error function**, e.g.,

$$E_\varepsilon(y(\mathbf{x}) - t) = \begin{cases} 0, & \text{if } |y(\mathbf{x}) - t| < \varepsilon \\ |y(\mathbf{x}) - t| - \varepsilon, & \text{otherwise} \end{cases}$$

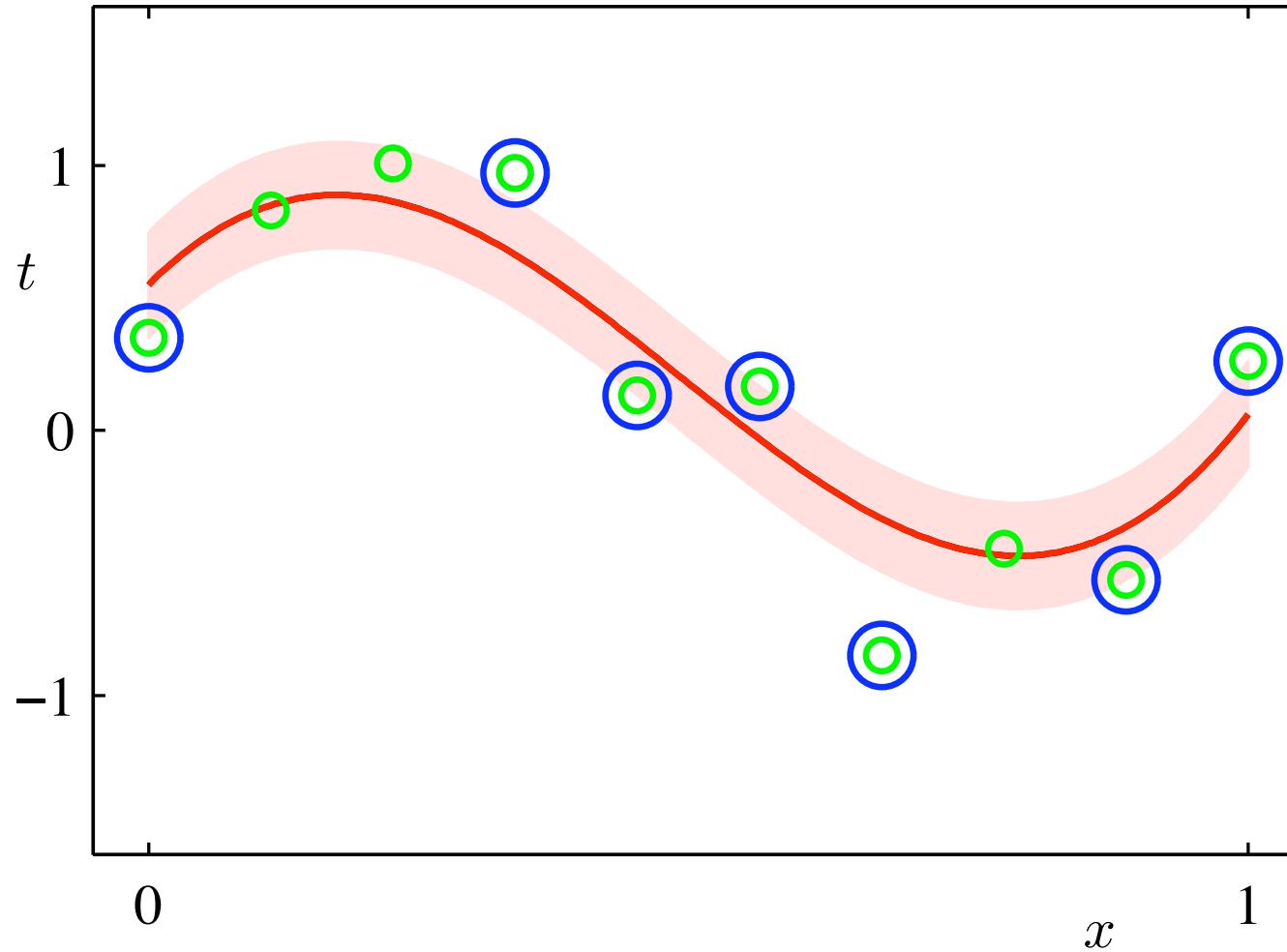
See text for details of solution.



Example

72

Probability & Bayesian Inference



Nonlinear Classification and Regression: Outline

73

Probability & Bayesian Inference

- Multi-Layer Perceptrons
 - The Back-Propagation Learning Algorithm
- Generalized Linear Models
 - Radial Basis Function Networks
 - Sparse Kernel Machines
 - Nonlinear SVMs and the Kernel Trick
 - **Relevance Vector Machines**

Relevance Vector Machines

74

Probability & Bayesian Inference

- Some drawbacks of SVMs:
 - Do not provide posterior probabilities.
 - Not easily generalized to $K > 2$ classes.
 - Parameters (C, ε) must be learned by cross-validation.
- The **Relevance Vector Machine** is a sparse Bayesian kernel technique that avoids these drawbacks.
- RVMs also typically lead to sparser models.

RVMs for Regression

75

Probability & Bayesian Inference

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}), \beta^{-1})$$

$$\text{where } y(\mathbf{x}) = \mathbf{w}^t \phi(\mathbf{x})$$

In an RVM, the basis functions $\phi(\mathbf{x})$ are kernels $k(\mathbf{x}, \mathbf{x}_n)$:

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

However, unlike in SVMs, the kernels need not be positive definite, and the \mathbf{x}_n need not be the training data points.

RVMs for Regression

76

Probability & Bayesian Inference

Likelihood:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta)$$

where the n^{th} row of \mathbf{X} is \mathbf{x}_n^t .

Prior:

$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M N(w_i | 0, \alpha_i^{-1})$$

- Note that each weight parameter has its own precision hyperparameter.

RVMs for Regression

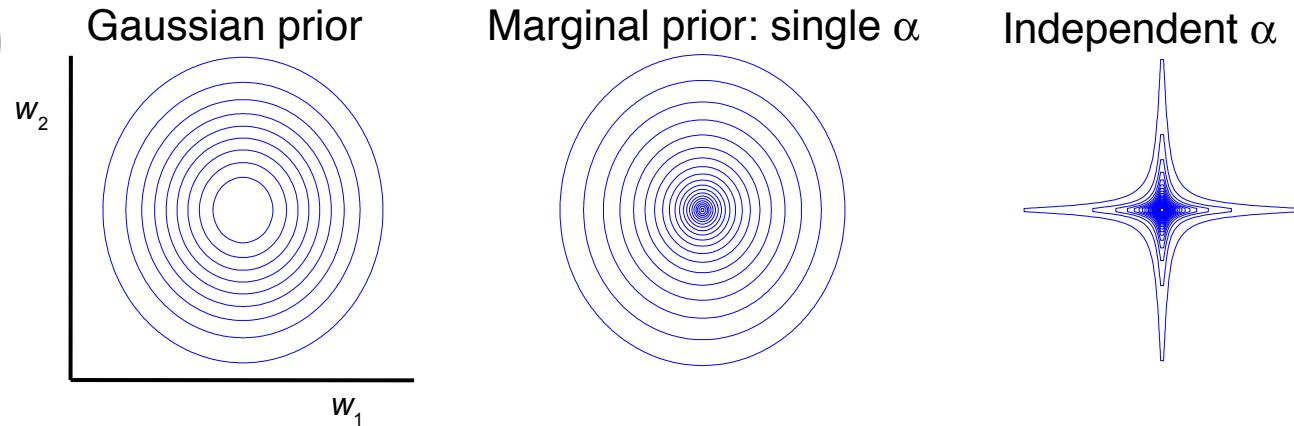
77

Probability & Bayesian Inference

$$p(w_i | \alpha_i) = N(w_i | 0, \alpha_i^{-1})$$

$$p(\alpha_i) = \text{Gam}(\alpha_i | a, b)$$

$$p(w_i) = \text{St}(w_i | 2a)$$



- The conjugate prior for the precision of a Gaussian is a gamma distribution.
- Integrating out the precision parameter leads to a Student's t distribution over w_i .
- Thus the distribution over \mathbf{w} is a product of Student's t distributions.
- As a result, maximizing the evidence will yield a sparse \mathbf{w} .
- Note that to achieve sparsity it is critical that each parameter w_i has a separate precision α_i .

RVMs for Regression

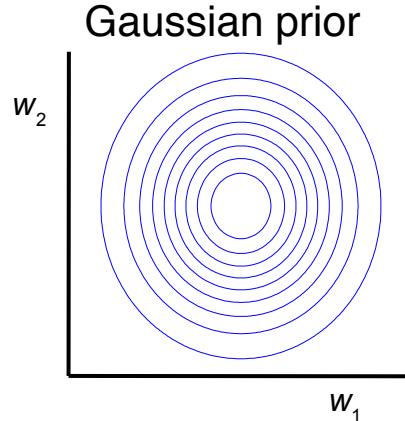
78

Probability & Bayesian Inference

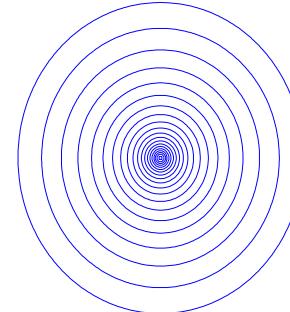
$$p(w_i | \alpha_i) = N(w_i | 0, \alpha_i^{-1})$$

$$p(\alpha_i) = \text{Gam}(\alpha_i | a, b)$$

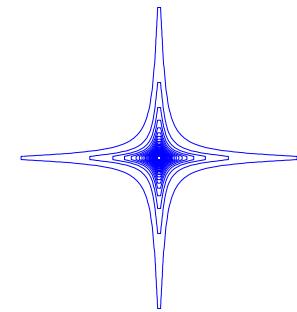
$$p(w_i) = \text{St}(w_i | 2a)$$



Marginal prior: single α



Independent α



$$\text{Gamma Distribution: } p(x | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}.$$

$$\text{Student's t Distribution: } p(x | v) = \frac{\Gamma\left(\frac{v+1}{2}\right)}{\sqrt{v\pi}\Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{x^2}{v}\right)^{-\frac{v+1}{2}}.$$

Also recall the rule for transforming densities:

If y is a monotonic function of x , then

$$p_Y(y) = p_X(x) \left| \frac{dx}{dy} \right|$$

Thus if we let $a \rightarrow 0, b \rightarrow 0$, then

$$p(\log \alpha_i) \rightarrow \text{uniform and } p(w_i) \rightarrow |w_i|^{-1}. \quad \text{Very sparse!!}$$

RVMs for Regression

79

Probability & Bayesian Inference

Likelihood:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta)$$

where the n^{th} row of \mathbf{X} is \mathbf{x}_n^t .

Prior:

$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M N(w_i | 0, \alpha_i^{-1})$$

- In practice, it is difficult to integrate α out exactly.
- Instead, we use an approximate maximum likelihood method, finding ML values for each α_i .
- When we maximize the evidence with respect to these hyperparameters, many will $\rightarrow \infty$.
- As a result, the corresponding weights will $\rightarrow 0$, yielding a sparse solution.

RVMs for Regression

80

Probability & Bayesian Inference

- Since both the likelihood and prior are normal, the posterior over \mathbf{w} will also be normal:

Posterior:

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) = N(\mathbf{w} | \mathbf{m}, \Sigma)$$

where

$$\mathbf{m} = \beta \Sigma \Phi^t \mathbf{t}$$

$$\Sigma = (\mathbf{A} + \beta \Phi^t \Phi)^{-1}$$

and

$$\Phi_{ni} = \phi_i(\mathbf{x}_n)$$

$$\mathbf{A} = \text{diag}(\alpha_i)$$

Note that when $\alpha_i \rightarrow \infty$, the i^{th} row and column of $\Sigma \rightarrow 0$, and

$$p(w_i | \mathbf{t}, \mathbf{X}, \alpha, \beta) = N(w_i | 0, 0)$$

RVMs for Regression

81

Probability & Bayesian Inference

- The values for α and β are determined using the evidence approximation, where we maximize

$$p(\mathbf{t} | \mathbf{X}, \alpha, \beta) = \int p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}$$

In general, this results in many of the precision parameters $\alpha_i \rightarrow \infty$, so that $w_i \rightarrow 0$.

Unfortunately, this is a non-convex problem.

Example

82

Probability & Bayesian Inference

