

SUPPORT VECTOR MACHINES

Today

- Reading
 - AIMA 18.9

- Goals
 - Finish Backpropagation
 - Introduce support vector machines (SVMs)

Backpropagation

1. Begin with randomly initialized weights
2. Apply the neural network to each training example (each pass through examples is called an epoch)
3. If it misclassifies an example **modify the weights**
4. Continue until the neural network classifies all training examples correctly

(Derive gradient-descent update rule)

Backpropagation

```

function BACK-PROP-LEARNING(examples, network) returns a neural network
inputs: examples, a set of examples, each with input vector x and output vector y
         network, a multilayer network with L layers, weights  $w_{i,j}$ , activation function g
local variables:  $\Delta$ , a vector of errors, indexed by network node

repeat
  for each weight  $w_{i,j}$  in network do
     $w_{i,j} \leftarrow$  a small random number
  for each example (x, y) in examples do
    /* Propagate the inputs forward to compute the outputs */
    for each node i in the input layer do
       $a_i \leftarrow x_i$ 
    for  $\ell = 2$  to L do
      for each node j in layer  $\ell$  do
         $in_j \leftarrow \sum_i w_{i,j} a_i$ 
         $a_j \leftarrow g(in_j)$ 
    /* Propagate deltas backward from output layer to input layer */
    for each node j in the output layer do
       $\Delta[j] \leftarrow g'(in_j) \times (y_j - a_j)$ 
    for  $\ell = L - 1$  to 1 do
      for each node i in layer  $\ell$  do
         $\Delta[i] \leftarrow g'(in_i) \sum_j w_{i,j} \Delta[j]$ 
    /* Update every weight in network using deltas */
    for each weight  $w_{i,j}$  in network do
       $w_{i,j} \leftarrow w_{i,j} + \alpha \times a_i \times \Delta[j]$ 
  until some stopping criterion is satisfied
return network

```

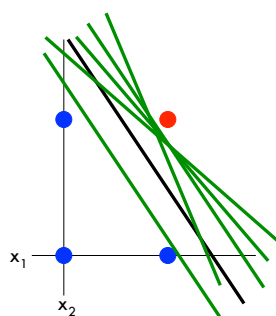
Figure 18.24 The back-propagation algorithm for learning in multilayer networks.

Support Vector Machines (SVMs)

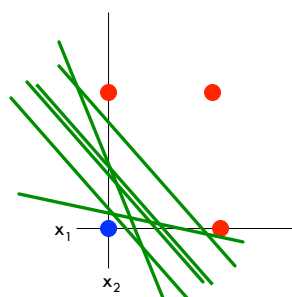
- SVMs are probably the most popular off-the-shelf classifier!
- Software Packages
 - LIBSVM (LIBLINEAR) – on the Resources page
 - SVM-Light

Linearly Separable

x_1	x_2	x_1 and x_2	
0	0	0	●
0	1	0	●
1	0	0	●
1	1	1	●

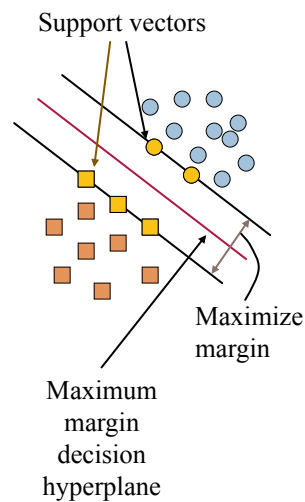


x_1	x_2	x_1 or x_2	
0	0	0	●
0	1	1	●
1	0	1	●
1	1	1	●



Support Vector Machines

- A **support vector machine** (SVM) is a linear classifier that finds the decision boundary btw. two classes that is *maximally far from any point in the training set*
- The **margin** is the distance from the decision boundary to the closest data point
- The **support vectors** are a subset of the training examples that fully determine the decision boundary



Basic Linear Algebra Notes (on board)

- Length of a vector
- Unit vector
- Dot product
- Hyperplane
- Given this knowledge, how do we find the hyperplane with the maximum margin?

Solving the Optimization Problem

$$\min_{w,b} \frac{1}{2} \|w\|^2 \text{ such that } y^{(i)}(w^\top x^{(i)} + b) \geq 1 \quad \forall i$$

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

Solving the Optimization Problem

$$\min_{w,b} \frac{1}{2} \|w\|^2 \text{ such that } y^{(i)}(w^\top x^{(i)} + b) \geq 1 \quad \forall i$$

$$\max_{\alpha} \min_{w,b} \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y^{(i)}(w^\top x^{(i)} + b) - 1] \quad \left. \vphantom{\max_{\alpha}} \right\} \text{Dual}$$

Lagrange multipliers ↗

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)} x^{(j)}$$

subject to $\alpha_i \geq 0$ and $\sum_i \alpha_i y^{(i)} = 0$

Solving the Optimization Problem

- The solution has the form:

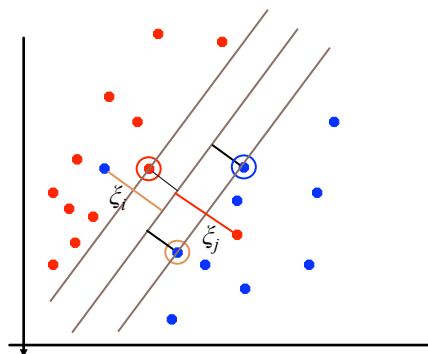
$$w = \sum_{i=1}^N \alpha_i y^{(i)} x^{(i)} \text{ and } b = y^{(i)} - w^T x^{(i)} \text{ for any } x^{(i)} \text{ s.t. } \alpha_i \neq 0$$

- Each non-zero alpha indicates corresponding x_i is a support vector
- The classifying function has the form: $g(x_i) = \text{sign}\left(\sum_i \alpha_i y^{(i)} x^{(i)} + b\right)$
- Relies on an inner product between the test point x and the support vectors x_i

Soft-margin Classification

If the training data is not linearly separable, *slack variables* ξ_i can be added to allow misclassification of difficult or noisy examples.

Still, try to minimize training set errors, and to place hyperplane “far” from each class (large margin)



How many support vectors?

- Determined by alphas in optimization
- Typically only a small proportion of the training data
- The number of support vectors determines the run time for prediction

How fast are SVMs?

Training

- Time for training is dominated by the time for solving the underlying quadratic programming problem
- Slower than Naïve Bayes
- Non-linear SVMs are worse

Testing (Prediction)

- Fast - as long as we don't have too many support vectors

Multi-label classification

- SVMs are inherently two-class classifiers
- Given C classes, common techniques are:
 - One-versus-all
 - Train C different SVMs where each SVM learns one class versus all the other classes
 - One-versus-one
 - Train $C(C-1)/2$ SVMs where each SVM learns to distinguish one class from another
- Multi-class SVMs
- Transductive SVMs

Linear SVMs Summary

- The classifier is a decision boundary (separating hyperplane)
- Most “important” training points are support vectors which define the hyperplane
- Quadratic optimization algorithms can identify which training points are support vectors (vectors with non-zero Lagrange multipliers)
- In the dual formation and in classifying an example, the training points appear only inside inner products