

# Linear Discriminant Functions: Gradient Descent and Perceptron Convergence

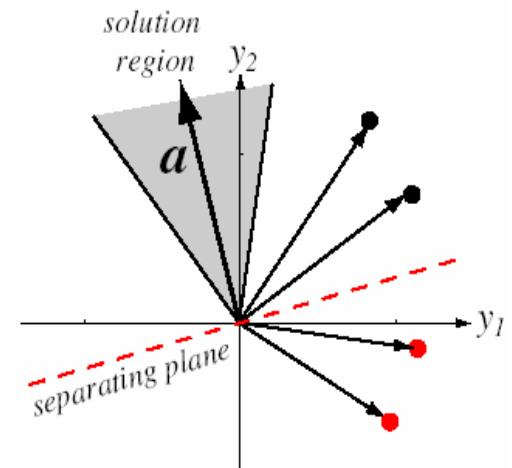
- The Two-Category Linearly Separable Case (5.4)
- Minimizing the Perceptron Criterion Function (5.5)

# Role of Linear Discriminant Functions

- A Discriminative Approach
  - as opposed to Generative approach of Parameter Estimation
- Leads to Perceptrons and Artificial Neural Networks
- Leads to Support Vector Machines

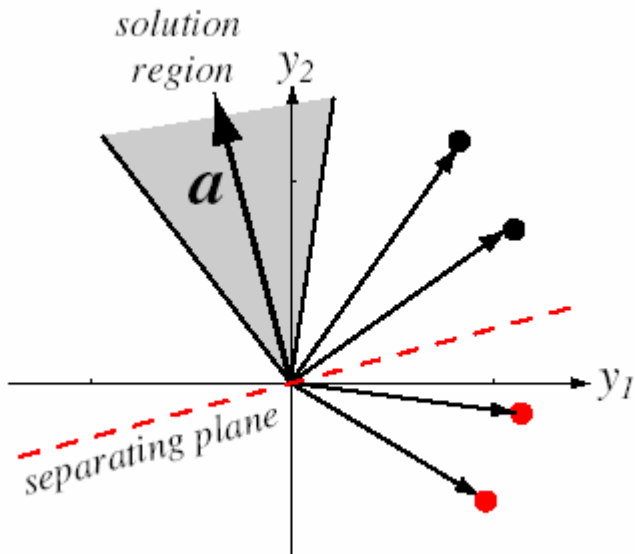
# Two-category Linearly Separable Case

- Set of  $n$  samples  $y_1, y_2, \dots, y_n$
- Some labelled  $\omega_1$  and some labelled  $\omega_2$
- Use samples to determine weight vector  $\mathbf{a}$  such that
$$\mathbf{a}^t \mathbf{y}_i > 0 \text{ implies class } \omega_1$$
$$\mathbf{a}^t \mathbf{y}_i < 0 \text{ implies class } \omega_2$$
- If such a weight vector exists then the samples are said to be linearly separable
- w.l.o.g. solution vector passes through origin (by mapping from x-space to y-space)



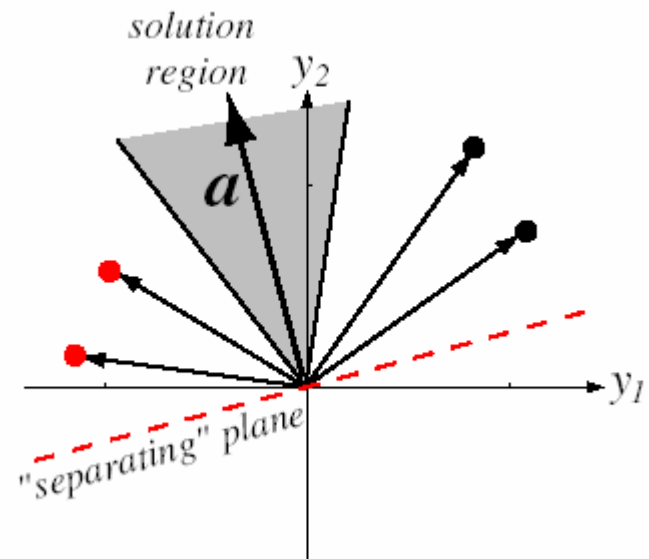
# Normalization

Raw Data



Solution vector separates red from black points

Normalized Data:  
red points changed in sign



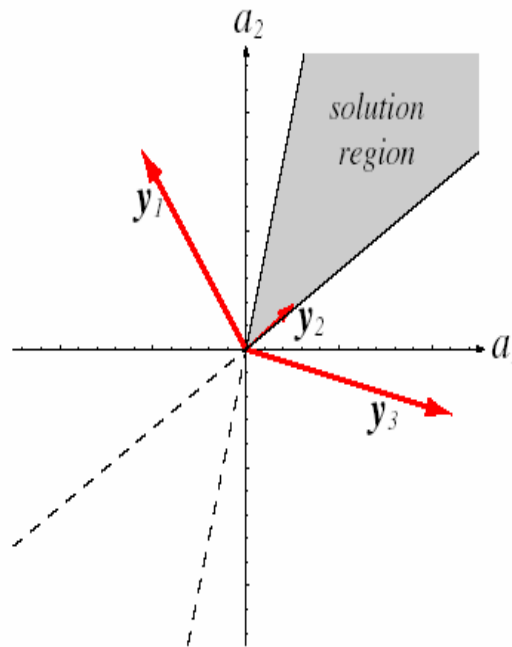
Solution vector places all vectors on same side of plane

Solution region is intersection of  $n$  half-spaces

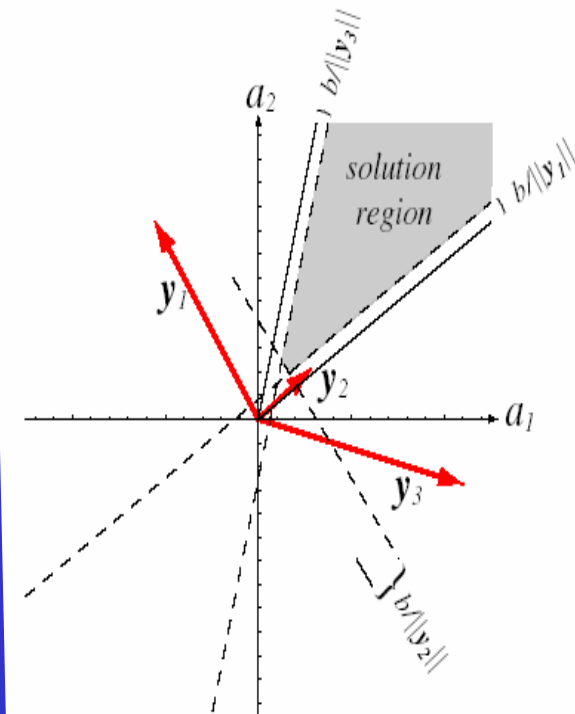
# Margin

- Solution vector is not unique
- Additional requirements to constrain solution vector. Two possible approaches:
  1. Seek unit length weight vector that maximizes minimum distance from samples to separating hyperplane
  2. Minimum length weight vector satisfying
$$a^t y_i \geq b > 0$$
where  $b$  is called the margin
- It is insulated from old boundaries by  $b / \|y_i\|$

No margin  
 $b = 0$



Margin  $b > 0$   
shrinks solution region  
by margins  $b / \|y_i\|$



# Gradient Descent Formulation

- To find solution to set of inequalities  $\mathbf{a}^t \mathbf{y}_i > 0$
- Define a criterion function  $J(\mathbf{a})$  that is minimized if  $\mathbf{a}$  is a solution vector
- We will define such a function  $J$  later

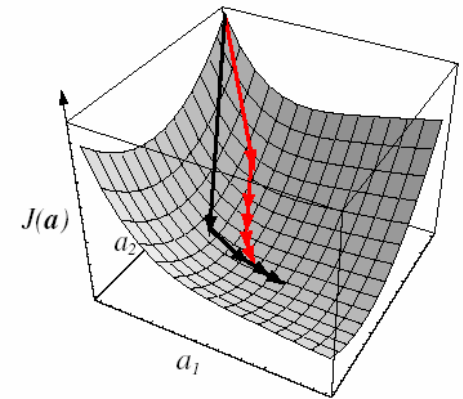
Start with an arbitrarily chosen weight  $\mathbf{a}(1)$  and

Compute gradient vector  $\nabla J(\mathbf{a}(1))$

Next value  $\mathbf{a}(2)$  obtained by moving some distance from  $\mathbf{a}(1)$  in the direction of the steepest descent, i.e., along the negative of the gradient. In general,

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k)),$$

Learning rate that  
Sets the step size



# Basic Gradient Descent Procedure

$$\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k)),$$

## ■ Algorithm 1. (Basic Gradient Descent)

```
1 begin initialize  $\mathbf{a}$ , threshold  $\theta$ ,  $\eta(\cdot)$ ,  $k \leftarrow 0$   
2       do  $k \leftarrow k + 1$   
3        $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$   
4       until  $|\eta(k) \nabla J(\mathbf{a})| < \theta$   
5   return  $\mathbf{a}$   
6 end
```

# Learning Rate

## Criterion for step size of Basic Gradient Descent

Criterion Function approximated by second-order expansion around  $\mathbf{a}(k)$ :

$$J(\mathbf{a}) \simeq J(\mathbf{a}(k)) + \nabla J^t (\mathbf{a} - \mathbf{a}(k)) + \frac{1}{2} (\mathbf{a} - \mathbf{a}(k))^t \mathbf{H} (\mathbf{a} - \mathbf{a}(k)),$$

substituting

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \eta(k) \nabla J(\mathbf{a}(k)),$$

$$J(\mathbf{a}(k+1)) \simeq J(\mathbf{a}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^t \mathbf{H} \nabla J.$$

Hessian matrix of  
Second order partial  
derivative

$$\partial^2 J / \partial a_i \partial a_j$$

Can be minimized by the choice:

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^t \mathbf{H} \nabla J},$$



# Criterion for step size of Basic Gradient Descent

$$J(\mathbf{a}) \simeq J(\mathbf{a}(k)) + \nabla J^t(\mathbf{a} - \mathbf{a}(k)) + \frac{1}{2}(\mathbf{a} - \mathbf{a}(k))^t \mathbf{H} (\mathbf{a} - \mathbf{a}(k)),$$

Hessian matrix of  
Second order partial  
derivative

$$\partial^2 J / \partial a_i \partial a_j$$

$$J(\mathbf{a}(k+1)) \simeq J(\mathbf{a}(k)) - \eta(k) \|\nabla J\|^2 + \frac{1}{2} \eta^2(k) \nabla J^t \mathbf{H} \nabla J.$$

Can be minimized by the choice:

$$\eta(k) = \frac{\|\nabla J\|^2}{\nabla J^t \mathbf{H} \nabla J},$$

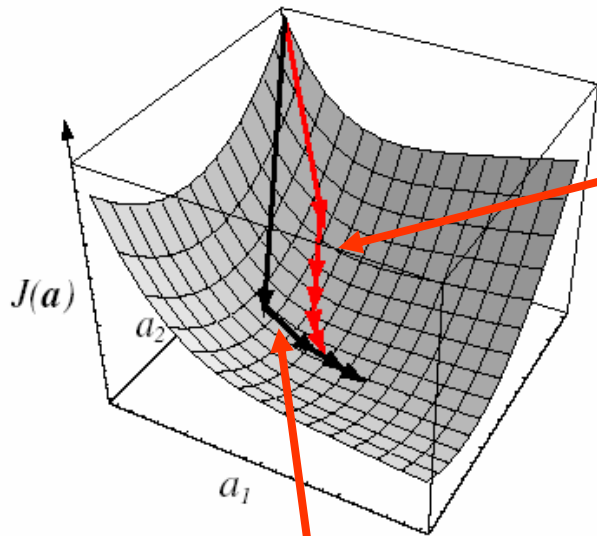
Alternative approach is Newton's algorithm

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \mathbf{H}^{-1} \nabla J,$$

## ■ Algorithm 2. (Newton Descent)

```
1 begin initialize a, threshold  $\theta$ 
2   do
3      $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \nabla J(\mathbf{a})$ 
4   until  $|\mathbf{H}^{-1} \nabla J(\mathbf{a})| < \theta$ 
5   return a
6 end
```

# Comparison of simple gradient descent and Newton's algorithms



Simple gradient descent method

Newton's second order method

Has greater improvement per step even when using optimal learning rates for both methods.  
However added computational burden of inverting the Hessian matrix.

# Perceptron Criterion Function

Obvious criterion function is no of samples misclassified. A poor choice since it is piecewise linear  
Alternative is the Perceptron criterion function

$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} (-\mathbf{a}^t \mathbf{y}),$$

Set of samples misclassified

If no samples are misclassified then criterion function evaluates to zero.  
Gradient function evaluates to:

$$\nabla J_p = \sum_{\mathbf{y} \in \mathcal{Y}} (-\mathbf{y}),$$

Update Rule becomes

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y},$$

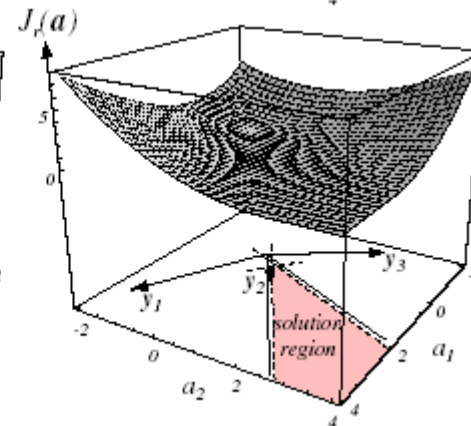
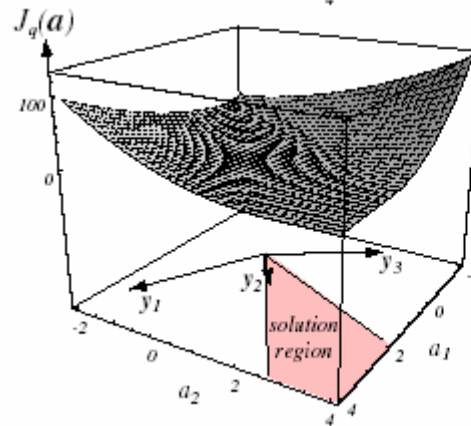
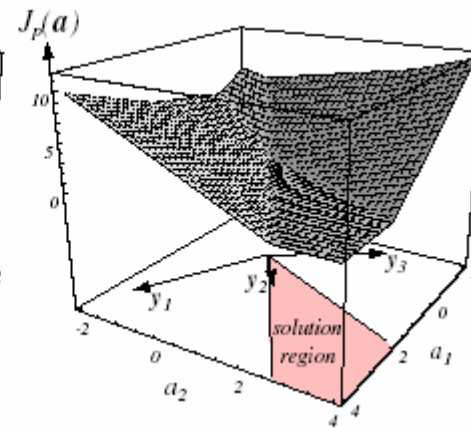
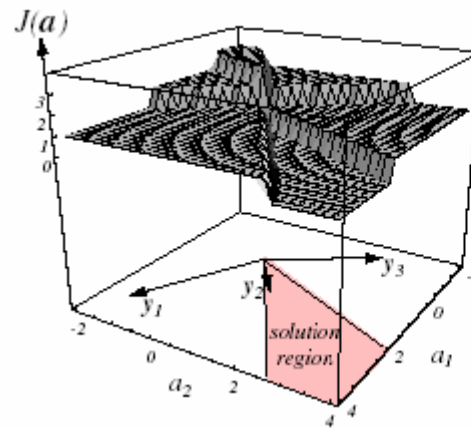
## ■ Algorithm 3. (Batch Perceptron)

```
1 begin initialize  $\mathbf{a}$ ,  $\eta(\cdot)$ , criterion  $\theta$ ,  $k \leftarrow 0$ 
2   do  $k \leftarrow k + 1$ 
3      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y}$ 
4   until  $|\eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y}| < \theta$ 
5   return  $\mathbf{a}$ 
6 end
```

# Comparison of Four Criterion functions

No of misclassified samples:  
Piecewise constant,  
unacceptable

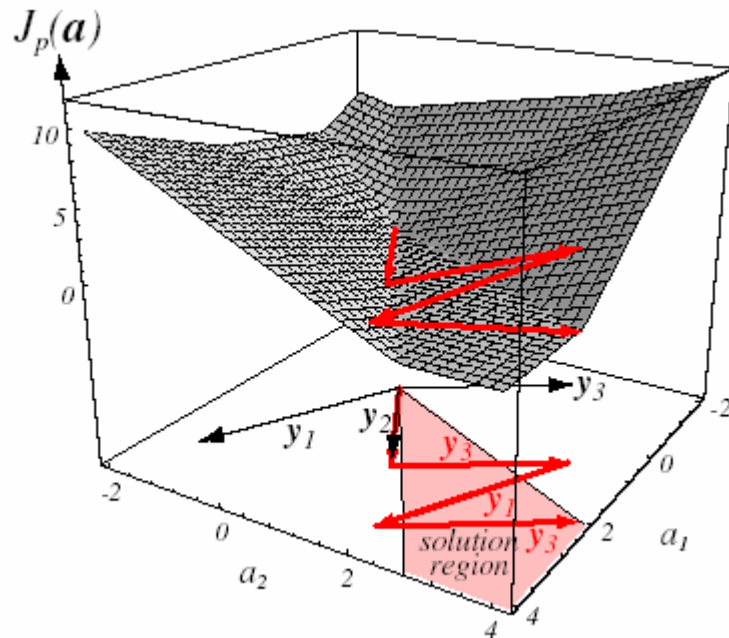
Perceptron criterion:  
Piecewise linear,  
acceptable for  
gradient descent



Squared error:  
Useful when patterns  
are not linearly separable

Squared Error  
with margin

# Perceptron Criterion as function of weights



Criterion function plotted as a function of weights  $a_1$  and  $a_2$

Starts at origin,

Sequence is  $y_2, y_3, y_1, y_3$

Second update by  $y_3$  takes solution further than first update by  $y_3$

# Convergence of Single-Sample Correction (simpler than batch)

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
 $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_1, \mathbf{y}_2, \dots$

$$\begin{aligned} \mathbf{a}(1) & \text{ arbitrary} \\ \mathbf{a}(k+1) &= \mathbf{a}(k) + \mathbf{y}^k \quad k \geq 1 \end{aligned}$$

## ■ Algorithm 4. (Fixed-Increment Single-Sample Perceptron)

```
1 begin initialize  $\mathbf{a}, k \leftarrow 0$ 
2   do  $k \leftarrow (k + 1) \bmod n$ 
3     if  $\mathbf{y}^k$  is misclassified by  $\mathbf{a}$  then  $\mathbf{a} \leftarrow \mathbf{a} + \mathbf{y}^k$ 
4   until all patterns properly classified
5   return  $\mathbf{a}$ 
6 end
```

# Perceptron Convergence Theorem

## Theorem (Perceptron Convergence)

If training samples are linearly separable, then the sequence of weight vectors given by Fixed Increment Error Correction Algorithm (Algorithm 4) will terminate at a solution vector

Let  $\hat{\mathbf{a}}$  be any solution vector, so that  $\hat{\mathbf{a}}^t \mathbf{y}_i$  is strictly positive for all  $i$ , and let  $\alpha$  be a positive scale factor. From Eq. 20 we have

$$\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}} = (\mathbf{a}(k) - \alpha \hat{\mathbf{a}}) + \mathbf{y}^k$$

and hence

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}(k) - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}(k) - \alpha \hat{\mathbf{a}})^t \mathbf{y}^k + \|\mathbf{y}^k\|^2.$$



# Proof for Single-Sample Correction

Because  $\mathbf{y}^k$  was misclassified,  $\mathbf{a}^t(k)\mathbf{y}^k \leq 0$  and thus

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(k) - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}^k + \|\mathbf{y}^k\|^2.$$

Because  $\hat{\mathbf{a}}^t \mathbf{y}^k$  is strictly positive, the second term will dominate the third if  $\alpha$  is sufficiently large. In particular, if we let  $\beta$  be the maximum length of a pattern vector,

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2, \quad (21)$$

and let  $\gamma$  be the smallest inner product of the solution vector with any pattern vector, that is,

$$\gamma = \min_i [\hat{\mathbf{a}}^t \mathbf{y}_i] > 0, \quad (22)$$



# Proof for Single-Sample Correction

then we have the inequality

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(k) - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha\gamma + \beta^2.$$

If we choose

$$\alpha = \frac{\beta^2}{\gamma}, \quad (23)$$

we obtain

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(k) - \alpha \hat{\mathbf{a}}\|^2 - \beta^2.$$

Thus, the squared distance from  $\mathbf{a}(k)$  to  $\alpha \hat{\mathbf{a}}$  is reduced by at least  $\beta^2$  at each correction, and after  $k$  corrections we obtain

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(1) - \alpha \hat{\mathbf{a}}\|^2 - k\beta^2. \quad (24)$$

Because this squared distance cannot become negative, it follows that the sequence of corrections must terminate after no more than  $k_0$  corrections, where

$$k_0 = \frac{\|\mathbf{a}(1) - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}. \quad (25)$$

# Bound on the number of corrections

$$k_0 = \frac{\alpha^2 \|\hat{\mathbf{a}}\|^2}{\beta^2} = \frac{\beta^2 \|\hat{\mathbf{a}}\|^2}{\gamma^2} = \frac{\max_i \|\mathbf{y}_i\|^2 \|\hat{\mathbf{a}}\|^2}{\min_i [\mathbf{y}_i^t \hat{\mathbf{a}}]^2}.$$

# Some Direct Generalizations

Correction whenever  $\mathbf{a}^t(k) \mathbf{y}^k$  fails to exceed a margin

$$\begin{array}{ll} \mathbf{a}(1) & \text{arbitrary} \\ \mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \mathbf{y}^k & k \geq 1, \end{array}$$

## ■ Algorithm 5. (Variable-Increment Perceptron with Margin)

```
1 begin initialize  $\mathbf{a}$ , threshold  $\theta$ , margin  $b$ ,  $\eta(\cdot)$ ,  $k \leftarrow 0$ 
2   do  $k \leftarrow (k+1) \bmod n$ 
3     if  $\mathbf{a}^t \mathbf{y}^k \leq b$  then  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \mathbf{y}^k$ 
4     until  $\mathbf{a}^t \mathbf{y}^k > b$  for all  $k$ 
5   return  $\mathbf{a}$ 
6 end
```

It can be shown that if the samples are linearly separable and if

$$\eta(k) \geq 0, \quad (28)$$

$$\lim_{m \rightarrow \infty} \sum_{k=1}^m \eta(k) = \infty \quad (29)$$

and

$$\lim_{m \rightarrow \infty} \frac{\sum_{k=1}^m \eta^2(k)}{(\sum_{k=1}^m \eta(k))^2} = 0, \quad (30)$$

then  $\mathbf{a}(k)$  converges to a solution vector  $\mathbf{a}$  satisfying  $\mathbf{a}^t \mathbf{y}_i > b$  for all  $i$  (Problem 19).

# Some Direct Generalizations: original gradient descent algorithm for $J_p$

$$\begin{aligned} \mathbf{a}(1) & \quad \text{arbitrary} \\ \mathbf{a}(k+1) &= \mathbf{a}(k) + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y}, \end{aligned}$$

## ■ Algorithm 6. (Batch Variable Increment Perceptron)

```
1 begin initialize  $\mathbf{a}$ ,  $\eta(\cdot)$ ,  $k \leftarrow 0$ 
2   do  $k \leftarrow (k + 1) \bmod n$ 
3      $\mathcal{Y}_k = \{\}$ 
4      $j = 0$ 
5     do  $j \leftarrow j + 1$ 
6       if  $\mathbf{y}_j$  is misclassified then Append  $\mathbf{y}_j$  to  $\mathcal{Y}_k$ 
7       until  $j = n$ 
8        $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \mathbf{y}$ 
9       until  $\mathcal{Y}_k = \{\}$ 
10    return  $\mathbf{a}$ 
11 end
```

# Some Direct Generalizations: Winnow Algorithm

Weights  $a^+$  and  $a^-$  associated with each of the categories to be learnt

## ■ Algorithm 7. (Balanced Winnow)

```
1 begin initialize  $a^+, a^-, \eta(\cdot), k \leftarrow 0, \alpha > 1$ 
2       if  $\text{Sgn}[a^{+t}y_k - a^{-t}y_k] \neq z_k$  (pattern misclassified)
3       then if  $z_k = +1$  then  $a_i^+ \leftarrow \alpha^{+y_i} a_i^+; a_i^- \leftarrow \alpha^{-y_i} a_i^-$  for all  $i$ 
4
5       if  $z_k = -1$  then  $a_i^+ \leftarrow \alpha^{-y_i} a_i^+; a_i^- \leftarrow \alpha^{+y_i} a_i^-$  for all  $i$ 
6       return  $a^+, a^-$ 
7 end
```

Advantages:

convergence is faster than in a Perceptron because of proper setting of learning rate

Each constituent value does not overshoot its final value

Benefit is pronounced when there are a large number of irrelevant or redundant features