

Minimum Squared-Error Procedures: Pseudo-inverse (DHS 5.8)

- Again a 2-class problem
- Assume reflected prototypes
- Minimum mean square error (MSE) technique (\Rightarrow Least Mean Squares, LMS).
- Consider all samples

With all prototypes

$$\mathbf{Y}^T = [y_1^{(1)}, y_2^{(1)}, \dots, y_{n_1}^{(1)}, -y_1^{(2)}, -y_2^{(2)}, \dots, y_{n_2}^{(2)}]$$

$$= [\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_P], P=n_1+n_2$$

$$\mathbf{Y} = \begin{bmatrix} \tilde{y}_1^T \\ \tilde{y}_2^T \\ \cdot \\ \cdot \\ \cdot \\ \tilde{y}_P^T \end{bmatrix}$$

Previously find the solution of a set of linear inequalities

$\mathbf{Y}\mathbf{W} > 0$ then correct classification or $\mathbf{Y}\mathbf{W} > \mathbf{b} \cdot \mathbf{1}$ where \mathbf{b} is a safety margin and $\mathbf{1}$ is a vector of 1.

Now find the solution to a set of linear equalities,

$\mathbf{Y}\mathbf{W} = \mathbf{b}$ where \mathbf{b} is a target vector now.

That is to find \mathbf{W} given a suitable \mathbf{b}

New criterion function

$$J(\mathbf{W}) = \|\mathbf{Y}\mathbf{W} - \mathbf{b}\|^2$$

If \mathbf{Y} is nonsingular, $\mathbf{W} = \mathbf{Y}^{-1}\mathbf{b}$

But it is generally singular, that is \mathbf{Y} is rectangular (more rows than columns).

$\mathbf{Y}\mathbf{W} = \mathbf{b}$ is overdetermined (more equations than unknowns).

Pseudo-inverse (continues)

To minimize J ,

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 2\mathbf{Y}^T (\mathbf{Y}\mathbf{w} - \mathbf{b}) = \mathbf{0}$$

$$\mathbf{Y}^T \mathbf{Y} \mathbf{w} = \mathbf{Y}^T \mathbf{b}$$

$$\mathbf{w} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b}$$

$$\mathbf{w} = \mathbf{Y}^+ \mathbf{b}$$

$\mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$ is called pseudoinverse

If \mathbf{Y} is square and nonsingular, then pseudoinverse becomes regular inverse

Therefore MSE solution

$$\mathbf{w} = \mathbf{Y}^+ \mathbf{b}$$

$\underline{\mathbf{b}}$ = arbitrary, will get a solution whether data is separable or not, but no guarantee it is a good solution for separating prototypes.

If $\underline{\mathbf{b}}$ is carefully chosen, we may be able to get a good discriminant function for both separable and non-separable cases.

MSE solution depends on the target vector \mathbf{b}

Different choices for \mathbf{b} give the solution different properties

Example 1: Constructing a linear classifier by matrix pseudoinverse

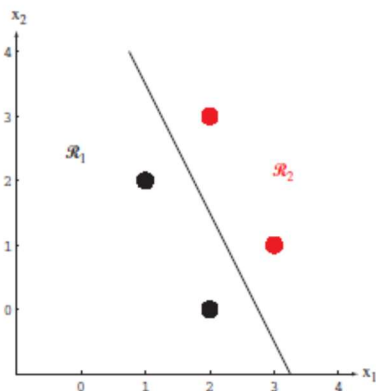
Suppose we have the following two-dimensional points for two categories: ω_1 : $(1, 2)^t$ and $(2, 0)^t$, and ω_2 : $(3, 1)^t$ and $(2, 3)^t$, as shown in black and red, respectively, in the figure.

Our matrix \mathbf{Y} is therefore

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{pmatrix}$$

and after a few simple calculations we find that its pseudoinverse is

$$\mathbf{Y}^\dagger \equiv \lim_{\epsilon \rightarrow 0} (\mathbf{Y}^t \mathbf{Y} + \epsilon \mathbf{I})^{-1} \mathbf{Y}^t = \begin{pmatrix} 5/4 & 13/12 & 3/4 & 7/12 \\ -1/2 & -1/6 & -1/2 & -1/6 \\ 0 & -1/3 & 0 & -1/3 \end{pmatrix}$$



Four training points and the decision boundary $\mathbf{a}^t \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0$, where \mathbf{a} was found by means of a pseudoinverse technique.

We arbitrarily let all the margins be equal, i.e., $\mathbf{b} = (1, 1, 1)^t$. Our solution is $\mathbf{a} = \mathbf{Y}^\dagger \mathbf{b} = (11/3, -4/3, -2/3)^t$, and leads to the decision boundary shown in the figure. Other choices for \mathbf{b} would typically lead to different decision boundaries, of course.

Windrow-Hoff (DHS 5.8.4)

Use this cost function, $J(\underline{w}) = \|\underline{Y}\underline{w} - \underline{b}\|^2$

\underline{b} = target vector

Advantages over pseudoinverse:

- Pseudoinverse can be very large
- It could be singular
- It can have truncation problems, errors.

Here

- One-at-a-time update
- Feedback scheme to reduce truncation errors.

$$\nabla J = 2Y^T(\underline{Y}\underline{w} - \underline{b})$$

$$\text{Rule 1} \Rightarrow \underline{w}(k+1) = \underline{w}(k) + \alpha(k)Y^T(Y\underline{w}(k) - \underline{b}) \quad \text{for all samples}$$

or

Rule 2 $\Rightarrow \underline{w}(k+1) = \underline{w}(k) + \alpha(k)[\underline{b}(k) - \underline{w}^T(k)y_k]$ y_k considering the samples sequentially

$$\underline{w}(0) = \text{arbitrary}$$

The size of Y^TY is smaller than Y^+ , storage requirements are less.

Update for all prototypes (misclassified & correctly-classified)

Usually updates never cease

$\alpha(k) = \alpha(0)/k$ for convergence.

Requires a good \underline{b} .

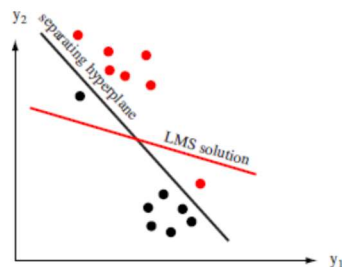


Figure 5.17: The LMS algorithm need not converge to a separating hyperplane, even if one exists. Since the LMS solution minimizes the sum of the squares of the distances of the training points to the hyperplane, for this exmple the plane is rotated clockwise compared to a separating hyperplane.

Ho-Kashyap Procedure (DHS 5.9.1)

Read Introduction in DHS 5.9.1. Check out the differences between the Perceptron and the MSE procedures in the case of linearly separable vs. nonseparable problems.

Task: Find \underline{w} and \underline{b} simultaneously

$$Y\underline{w} = \underline{b} > 0$$

$$J(\underline{w}, \underline{b}) = \|Y\underline{w} - \underline{b}\|^2$$

Minimize J w.r.t. \underline{w} and \underline{b} with constraint $\underline{b} > 0$

$$\nabla_{\underline{w}} J = 2 Y^T (Y\underline{w} - \underline{b})$$

$$\nabla_{\underline{b}} J = -2 (Y\underline{w} - \underline{b})$$

$$\underline{w} = Y^\dagger \underline{b} \Rightarrow \nabla_{\underline{w}} J = 0$$

Minimize J w.r.t. \underline{b} , with $\underline{w} = Y^\dagger \underline{b}$, subject to constraint $\underline{b} > 0$

Start with $\underline{b} > 0$

Only add positive elements when updating \underline{b}

Gradient descent:

$$\underline{b}(k+1) = \underline{b}(k) - 1/2 \alpha [\nabla_{\underline{b}} J - |\nabla_{\underline{b}} J|] \quad \alpha > 0$$

$$1/2 [a - |a|] = 0 \text{ if } a \geq 0$$

a if a < 0 to make it sure a positive update

$$|\underline{v}| \text{ means component-wise } |\cdot| \Rightarrow |\underline{v}| = [\dots, |v_i|, \dots]^T$$

Resulting Algorithm

$\underline{b}(0) > 0$ but otherwise arbitrary

$$\underline{w}(k) = Y^\dagger \underline{b}(k)$$

$$\text{Let } \underline{e}(k) = Y\underline{w}(k) - \underline{b}(k)$$

$$\underline{b}(k+1) = \underline{b}(k) + \alpha [\underline{e}(k) + |\underline{e}(k)|]$$

$$\alpha > 0$$

This is Ho-Kashyap Pseudoinverse.

Notes on Ho-Kashyap

1. Converges if samples are linearly separable (proved in DHS 5.9.2)
2. Generally required fewer steps to converge than Perceptron. However, each step requires more operations than Perceptron.
3. Update entire \underline{b} and \underline{w} , for both classes in each iteration
4. Nonseparability of data is indicated in the course of iterating. If $e(k) \leq 0$, not linearly separable.

Appropriate α

Option 1

$$0 < \alpha < 2$$

\Rightarrow converges fastest

$$\underline{w}(0) = (Y^T Y)^{-1} Y^T \underline{b}(0)$$

$$\text{and } \underline{b}(0) = \underline{1}$$

\Rightarrow solution $\underline{w}(k)$ is the best linear square fit for a given $\underline{b}(k)$

Option 2

$$0 < \alpha < \|Y^T Y\|^{-1}$$

$\|\cdot\|$ can be any of the following

$$\|A\| = \sum_{ij} |a_{ij}|$$

$$\|A\| = \max_i \sum_{j=1}^N |a_{ij}|$$

$$\|A\| = \text{tr}(AA^*)^{\frac{1}{2}} = \left[\sum_{ij} |a_{ij}|^2 \right]^{\frac{1}{2}}$$

This gives the simplest implementation but converges slower.

Ho-Kashyap Convergence (DHS 5.9.2)

If samples are linearly separable and if $0 < \alpha < 1$

- converges to solution in finite no of steps
- could add a halting condition for when prototypes are correctly classified

can show either

$e(k)=0$ within finite no of steps \rightarrow algorithm terminates with a solution vector

or $e(k) \rightarrow 0$ as $k \rightarrow \infty \Rightarrow Y_{\underline{w}(k)} > 0$ after finite no of steps

Same convergence properties for linearly separable prototypes

Different options on parameter α

Behavior of Ho-Kashyap Algorithm for Nonseparable Prototypes (DHS 5.9.3)

- If obtain an $\underline{e}(k)$ or converge to an $\underline{e}(k)$ such that $\underline{e}(k) \neq 0$ and no components of $\underline{e}(k)$ are positive, then the prototypes are not linearly separable.
- If the prototypes are not linearly separable, then either the algorithm will yield an $\underline{e}(k)$ such that $\underline{e}(k) \neq 0$ with no positive components, or will asymptotically approach it: $\underline{e}(k) \rightarrow \underline{e}(\infty) \neq 0$ with no components of $\underline{e}(\infty)$ being > 0

We have covered so far (see Table 5.1)

1. Fixed Increment in Perceptron
2. Variable Increment in Perceptron
3. Relaxation in Perceptron
4. Pseudo-Inverse
5. Windrow-Hoff
6. Ho-Kashyap

* Stochastic Approximation and Linear Programming (i.e., Simplex Algorithm) are not covered here.

Various Descent Algorithms

Table 5.1: Descent Procedures for Obtaining Linear Discriminant Functions

Name	Criterion	Algorithm	Conditions
Fixed Increment	$J_p = \sum_{\mathbf{a}^t \mathbf{y} \leq 0} (-\mathbf{a}^t \mathbf{y})$	$\mathbf{a}(k+1) = \mathbf{a}(k) + \mathbf{y}^k$ $(\mathbf{a}^t(k) \mathbf{y}^k \leq 0)$	—
Variable Increment	$J'_p = \sum_{\mathbf{a}^t \mathbf{y} \leq 0} -(\mathbf{a}^t \mathbf{y} - b)$	$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \mathbf{y}^k$ $(\mathbf{a}^t(k) \mathbf{y}^k \leq b)$	$\eta(k) \geq 0$ $\sum \eta(k) \rightarrow \infty$ $\frac{\sum \eta^2(k)}{(\sum \eta(k))^2} \rightarrow 0$
Relaxation	$J_r = \frac{1}{2} \sum_{\mathbf{a}^t \mathbf{y} \leq b} \frac{(\mathbf{a}^t \mathbf{y} - b)^2}{\ \mathbf{y}\ ^2}$	$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta \frac{b - \mathbf{a}^t(k) \mathbf{y}^k}{\ \mathbf{y}^k\ ^2} \mathbf{y}^k$ $(\mathbf{a}^t(k) \mathbf{y}^k \leq b)$	$0 < \eta < 2$
Widrow-Hoff (LMS)	$J_s = \sum_i (\mathbf{a}^t \mathbf{y}_i - b_i)^2$	$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)(b_k - \mathbf{a}^t(k) \mathbf{y}^k) \mathbf{y}^k$	$\eta(k) > 0$ $\eta(k) \rightarrow 0$
Stochastic Approx.	$J_m = \mathcal{E} [(\mathbf{a}^t \mathbf{y} - z)^2]$	$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)(z_k - \mathbf{a}^t(k) \mathbf{y}^k) \mathbf{y}^k$	$\sum \eta(k) \rightarrow \infty$ $\sum \eta^2(k) \rightarrow L < \infty$
		$\mathbf{a}(k+1) = \mathbf{a}(k) + \mathbf{R}(k)(z(k) - \mathbf{a}^t(k) \mathbf{y}^k) \mathbf{y}^k$	$\mathbf{R}^{-1}(k+1) = \mathbf{R}^{-1}(k) + \mathbf{y}_k \mathbf{y}_k^t$
Pseudo-inverse	$J_s = \ \mathbf{Y}\mathbf{a} - \mathbf{b}\ ^2$	$\mathbf{a} = \mathbf{Y}^\dagger \mathbf{b}$	—
Ho-Kashyap	$J_s = \ \mathbf{Y}\mathbf{a} - \mathbf{b}\ ^2$	$\mathbf{b}(k+1) = \mathbf{b}(k) + \eta(\mathbf{e}(k) + \mathbf{e}(k))$ $\mathbf{e}(k) = \mathbf{Y}\mathbf{a}(k) - \mathbf{b}(k)$ $\mathbf{a}(k) = \mathbf{Y}^\dagger \mathbf{b}(k)$	$0 < \eta < 1$ $\mathbf{b}(1) > 0$
		$\mathbf{b}(k+1) = \mathbf{b}(k) + \eta(\mathbf{e}(k) + (\mathbf{e}(k)))$ $\mathbf{a}(k+1) = \mathbf{a}(k) + \eta \mathbf{R} \mathbf{Y}^t \mathbf{e}(k) $	$\eta(k) = \frac{ \mathbf{e}(k) ^t \mathbf{Y} \mathbf{R} \mathbf{Y}^t \mathbf{e}(k) }{ \mathbf{e}(k) ^t \mathbf{Y} \mathbf{R} \mathbf{Y}^t \mathbf{Y} \mathbf{R} \mathbf{Y}^t \mathbf{e}(k) }$ is optimum; \mathbf{R} sym., pos. def.; $\mathbf{b}(1) > 0$
Linear Programming	$\tau = \max_{\mathbf{a}^t \mathbf{y}_i \leq b_i} [-(\mathbf{a}^t \mathbf{y}_i - b_i)]$	Simplex algorithm	$\mathbf{a}^t \mathbf{y}_i + \tau \geq b_i$ $b \geq 0$
	$J'_p = \sum_{i=1}^n \tau_i$ $= \sum_{\mathbf{a}^t \mathbf{y}_i \leq b_i} -(\mathbf{a}^t \mathbf{y}_i - b_i)$	Simplex algorithm	$\mathbf{a}^t \mathbf{y}_i + \tau \geq b_i$ $b \geq 0$

Support Vector Machines (or Maximum Margin Classifier) (DHS 5.11)

Concepts

- Recall linear machines with margins.
- SVMs are very much similar, but rely on preprocessing the data to represent patterns in a high dimension (much higher than original feature space)
- Typically a nonlinear mapping function (or a kernel function) $\phi(\cdot)$ is used. Thus transform a pattern x_k to $y_k = \phi(x_k)$.
- A linear discriminant can be expressed as $g(y_k) = w^T y_k$ in an augmented space.
- The goal of a SVM is to find a separating hyperplane with the largest margin.
- The support vectors are the training samples that define optimal separating hyperplane.
- The support vectors are the most difficult patterns to classify.
- See Fig. 5.19

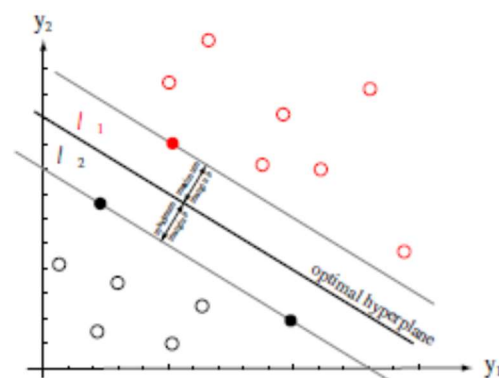


Figure 5.19: Training a Support Vector Machine consists of finding the optimal hyperplane, i.e., the one with the maximum distance from the nearest training patterns. The support vectors are those (nearest) patterns, a distance b from the hyperplane. The three support vectors are shown in solid dots.

Methods

- Modify the familiar Perceptron algorithm: train with the current worst-classified patterns. Of course finding the worst-classified patterns is difficult (computationally expensive)
- Training an SVM
 - Use the method of Lagrange Multipliers (not the focus of this class)
 - The cost function

$$L(w, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{k=1}^n \alpha_k [z_k w^T y_k - 1] \quad \text{with } z_k = \pm 1$$

Minimize L w.r.t. the weight vector w , and maximize it w.r.t. the multipliers $\alpha_k > 0$

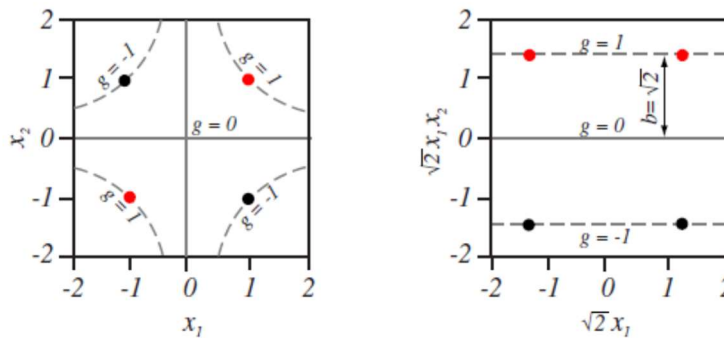
- This problem can be reformulated through the Kuhn-Tucker condition as

Maximizing $L(\alpha) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,j} \alpha_k \alpha_j z_k z_j y_j^T y_k$ with the constraints

$$\sum_{k=1}^n z_k \alpha_k = 0, \quad \alpha_k \geq 0, \quad k = 1, \dots, n$$

Example

- Example 2 (DHS p. 264)



The XOR problem in the original $x_1 - x_2$ feature space is shown at the left; the two red patterns are in category ω_1 and the two black ones in ω_2 . These four training patterns x are mapped to a six-dimensional space by $1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2$ and x_2^2 . In this space, the optimal hyperplane is found to be $g(x_1, x_2) = x_1x_2 = 0$ and the margin is $b = \sqrt{2}$. A two-dimensional projection of this space is shown at the right. The hyperplanes through the support vectors are $\sqrt{2}x_1x_2 = \pm 1$, and correspond to the hyperbolas $x_1x_2 = \pm 1$ in the original feature space, as shown.

- Try “svmtrain” under Bioinformatics Toolbox of Matlab