

Statistical Methods in Artificial Intelligence

CSE471 - Monsoon 2015 : Lecture 04



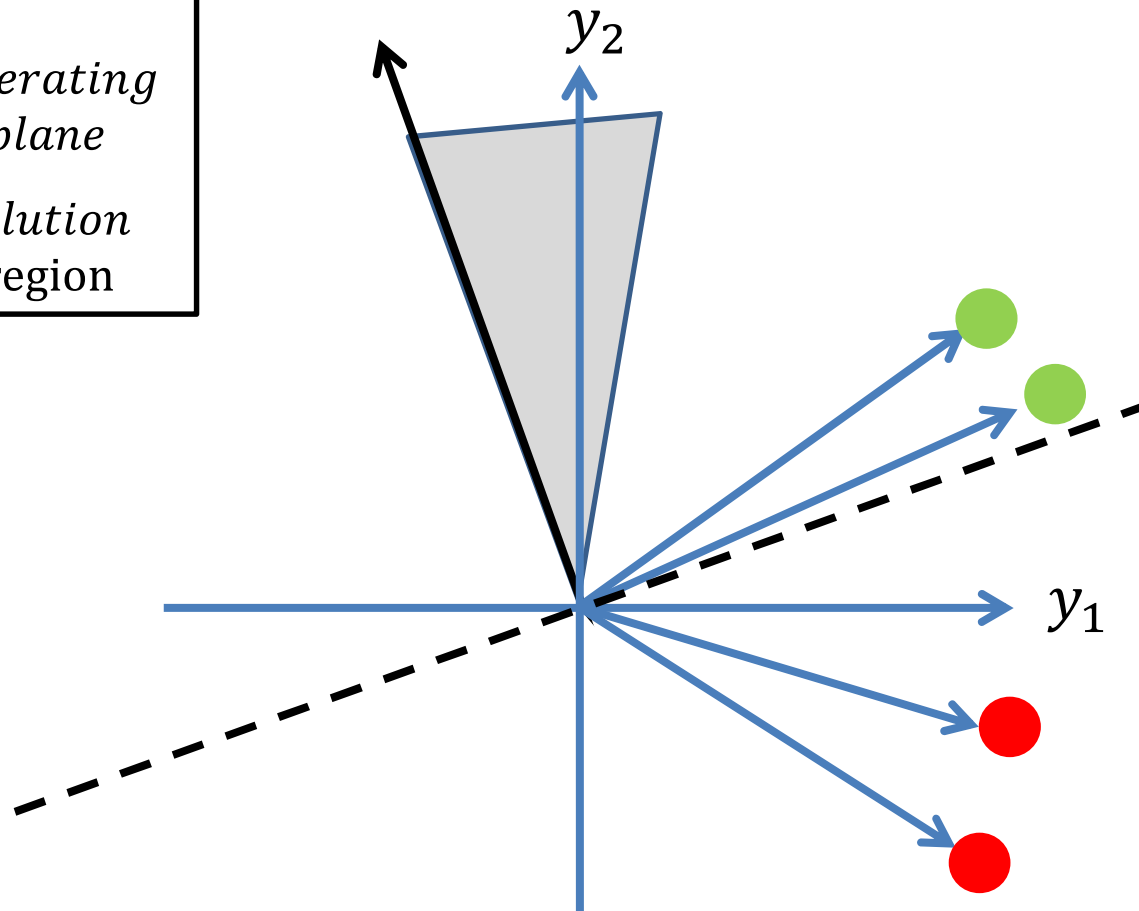
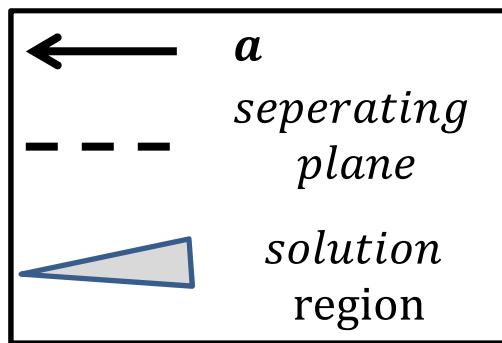
Avinash Sharma
CVIT, IIIT Hyderabad

Lecture 04: Plan

- Recap
- Learning LDF: Basic Gradient Descend
- Perceptron Criterion Function
- Batch Perceptron
- Single Sample Perceptron
- Relaxation Procedures
- Non Separable Behavior
- Minimum Squared Error Procedures
- LMS and Ho-Kashyap Procedures

Two-Category Linearly Separable Case

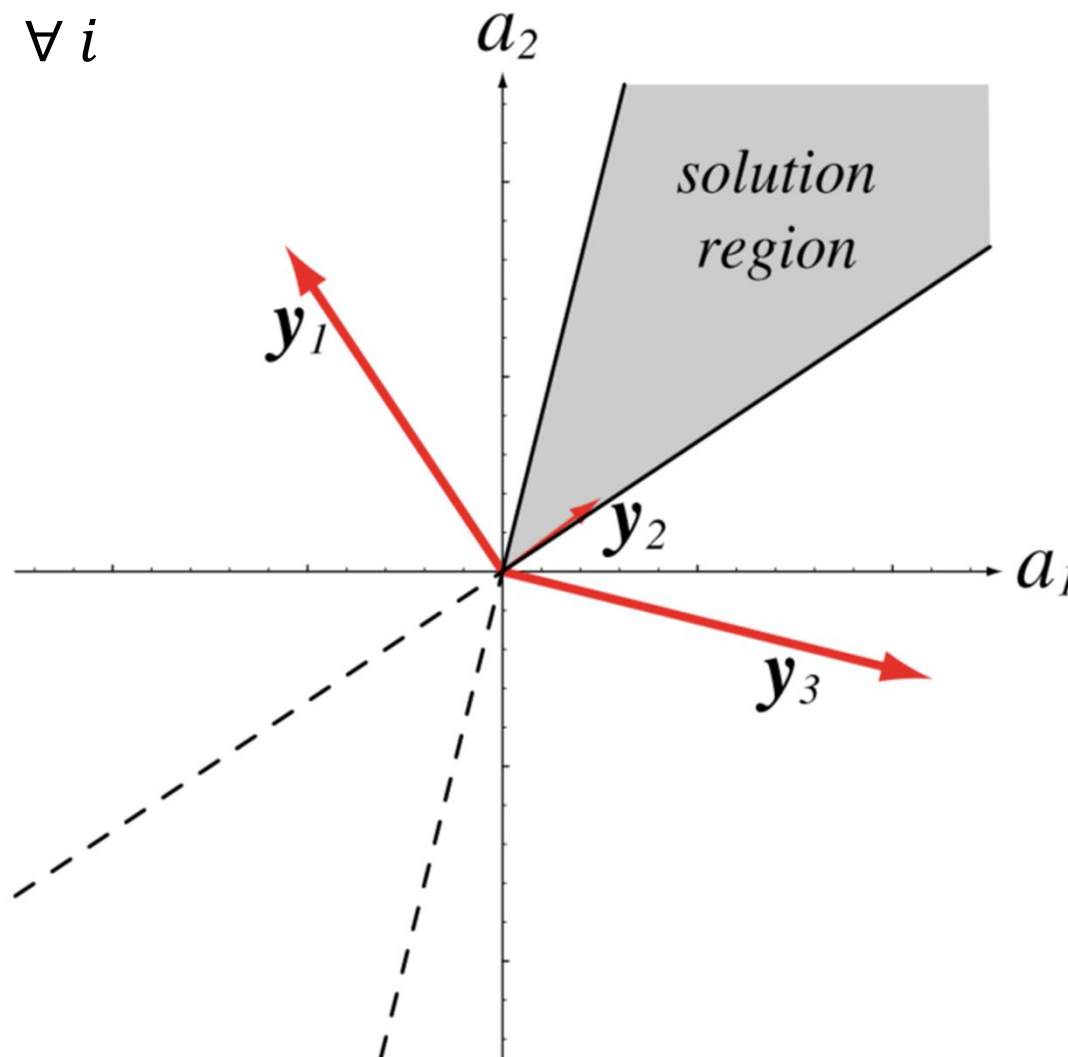
$$g(\mathbf{X}) = \mathbf{a}^T \mathbf{Y} = \sum_{i=1}^{\hat{d}} a_i y_i = \begin{cases} > 0 & (+ve) & \text{class A} \\ < 0 & (-ve) & \text{class B} \\ = 0 & \text{Decision Boundary} \end{cases}$$



Two-Category Linearly Separable Case

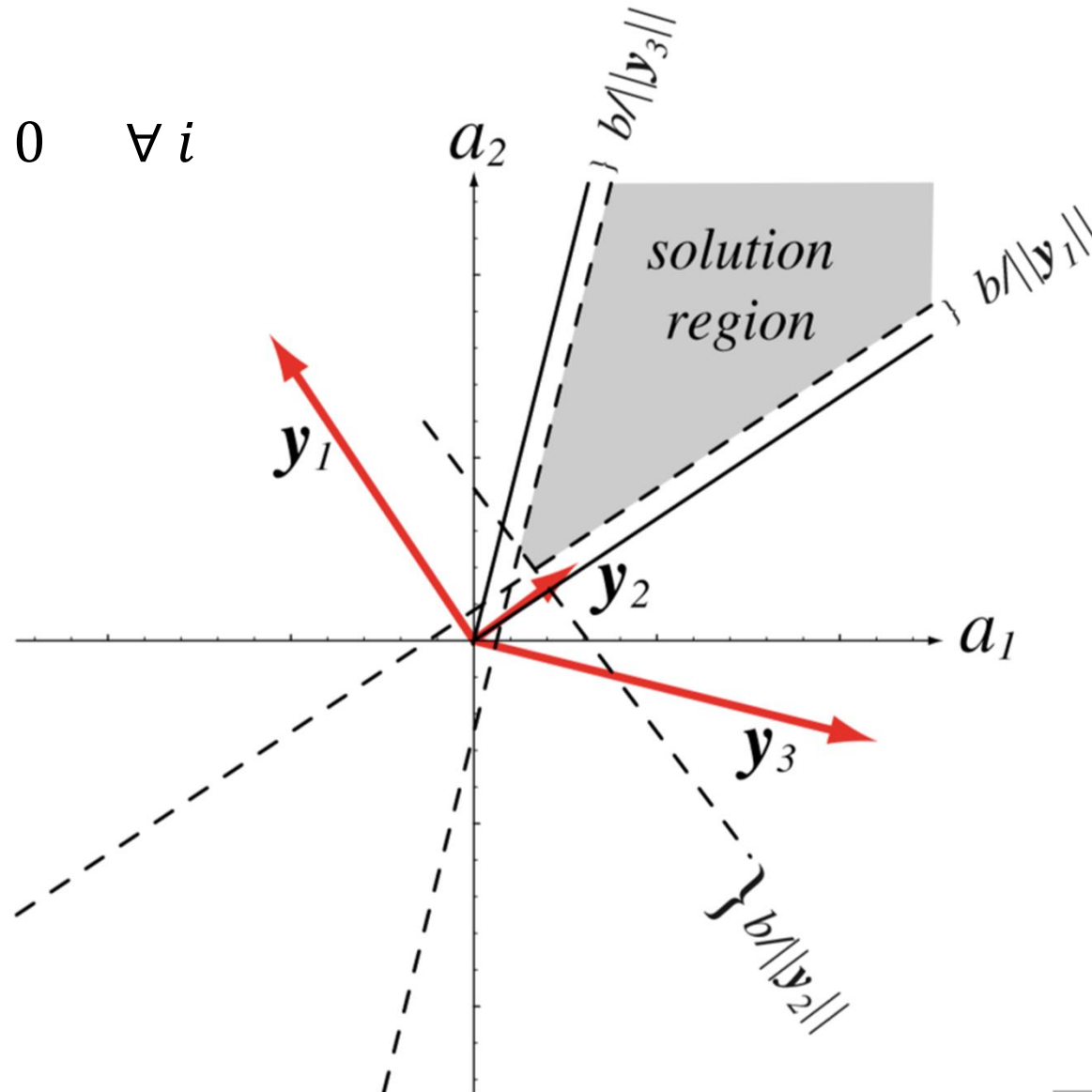
$$\mathbf{a}^T \mathbf{y}_i > 0 \quad \forall i$$

Data vector



Two-Category Linearly Separable Case

$$\mathbf{a}^T \mathbf{y}_i \geq b > 0 \quad \forall i$$



Learning LDF: Basic Gradient Descend

- Define a scalar function $J(\mathbf{a})$ which captures classification error for specific boundary plane described by parameter \mathbf{a}
- Minimize $J(\mathbf{a})$ using **gradient descent**.
 - Start with arbitrary value of $\mathbf{a}(1)$ for $k = 1$.
 - Iteratively refine estimate of \mathbf{a} :
$$\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(k)\nabla J(\mathbf{a}(k))$$
- η is the positive scale factor known as **learning rate**
 - A too small η makes the convergence very slow
 - A too large η can diverge due to overshooting correct.

Basic Gradient Descend Algorithm

1. Initialize \mathbf{a} , θ (threshold), $\eta(\cdot)$, $k = 0$

2. do $k = k + 1$

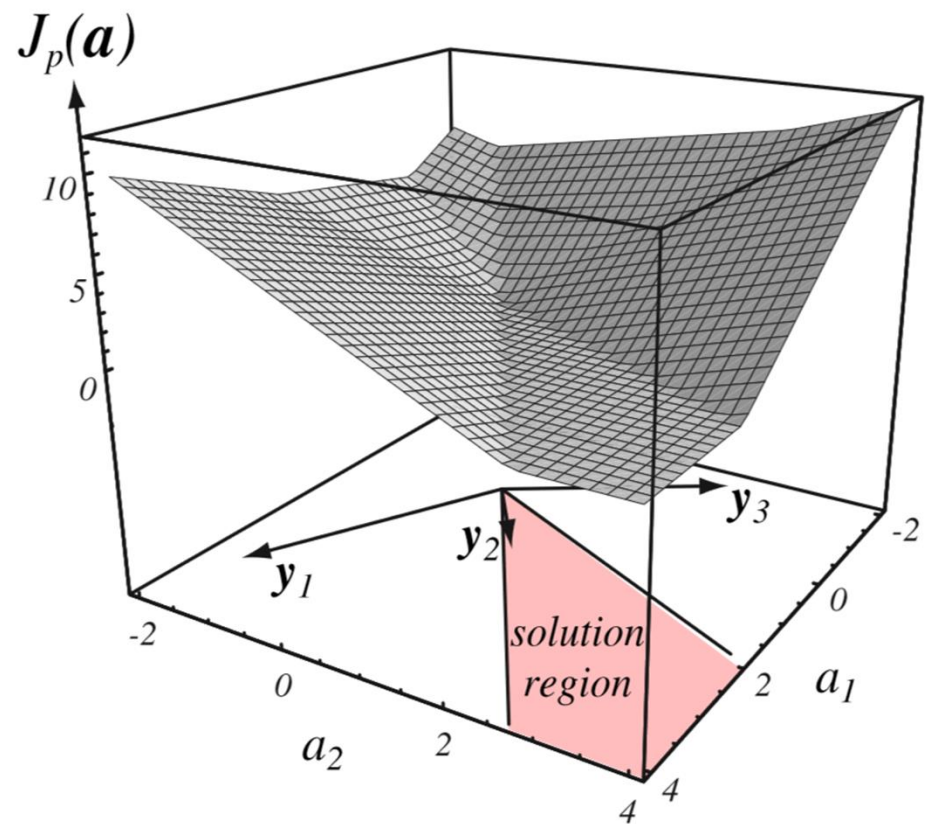
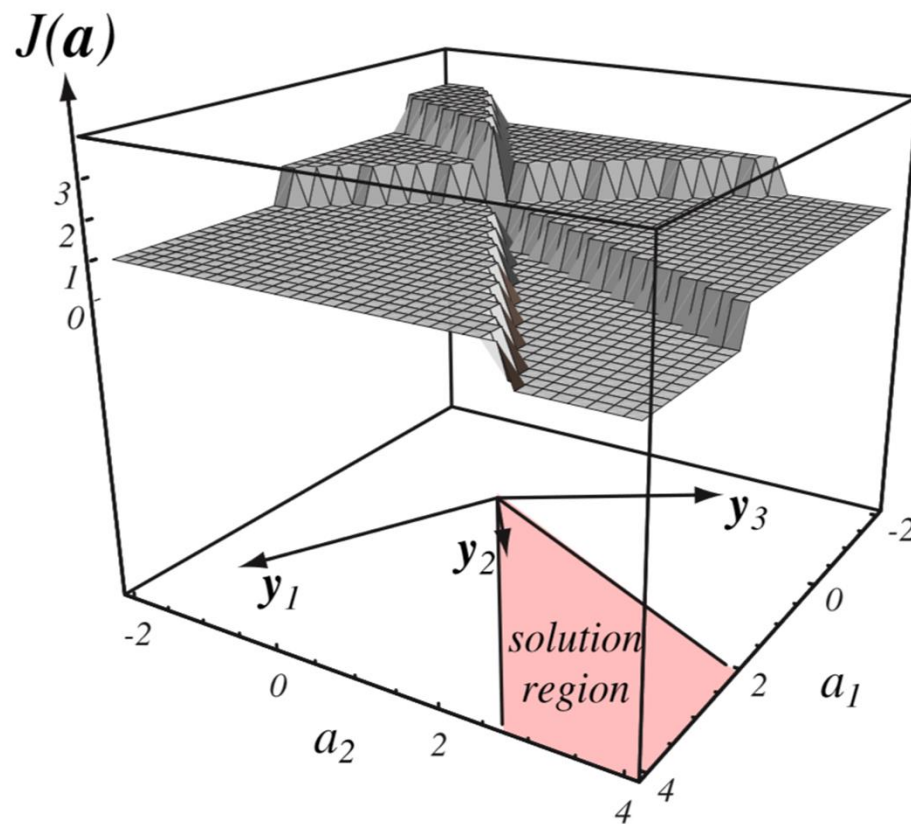
$$\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(k)\nabla J(\mathbf{a}(k))$$

3. untill $|\eta(k)\nabla J(\mathbf{a}(k))| < \theta$

4. return \mathbf{a}

Perceptron Criterion Function

- Discrete v/s continuous function



Perceptron Criterion Function

- Discrete v/s continuous function
- $J_p(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} (-\mathbf{a}^T \mathbf{y})$ where \mathcal{Y} is the set of misclassified samples
- $J_p(\mathbf{a})$ is proportional to the sum of the distances from all the miss classified samples to the decision boundary.

- Derivative of $J_p(\mathbf{a})$

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

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

Batch Perceptron

1. Initialize \mathbf{a} , θ (threshold), $\eta(\cdot)$, $k = 0$
2. do $k = k + 1$

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

3. untill $|\eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} (-\mathbf{y})| < \theta$
4. return \mathbf{a}

Single Sample Perceptron

1. Initialize \mathbf{a} , $k = 0$

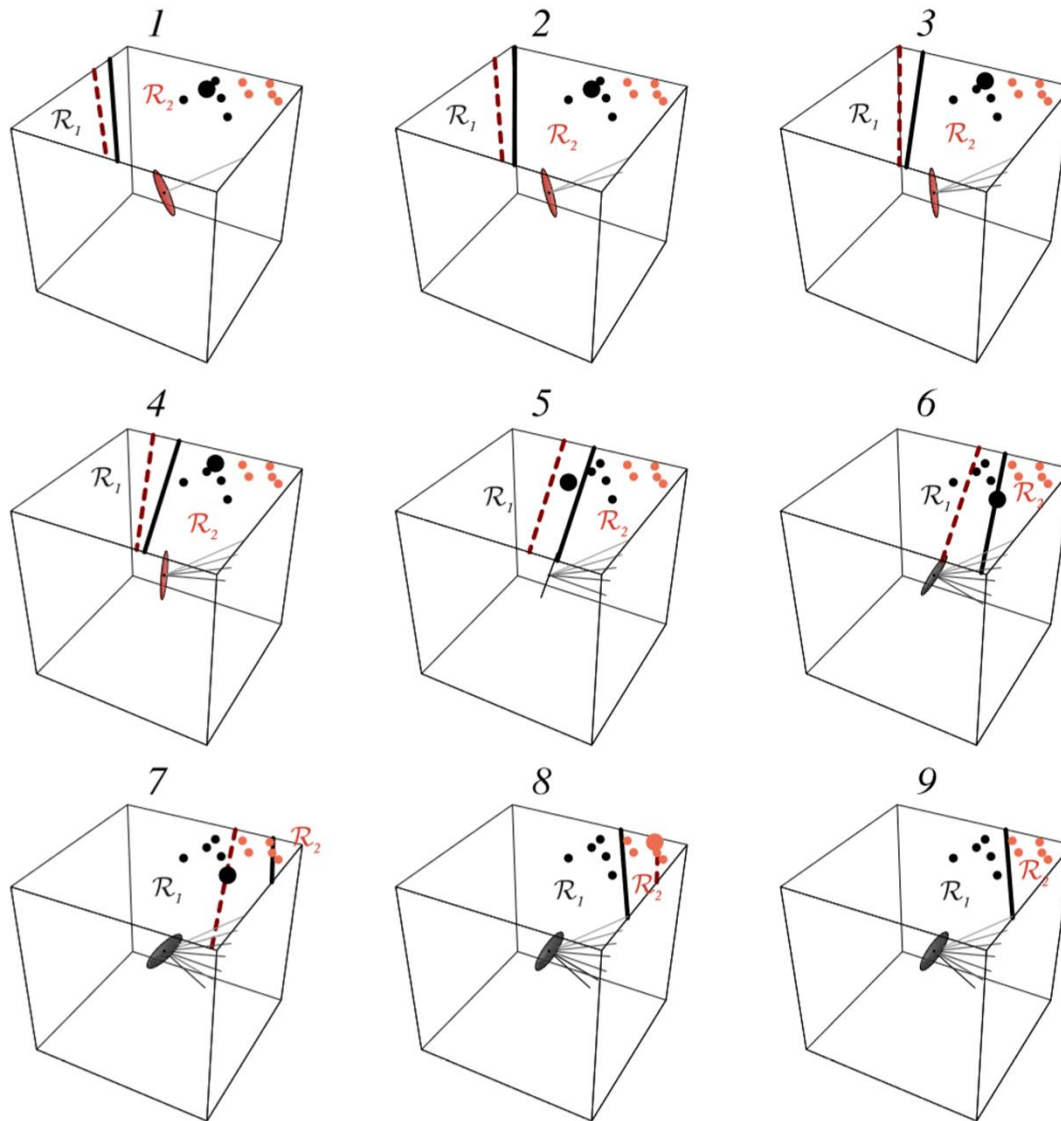
2. do $k = \text{mod}(k + 1, n)$

$$\mathbf{a} = \mathbf{a} + \mathbf{y}^k$$

3. untill all patterns are correctly classified

4. return \mathbf{a}

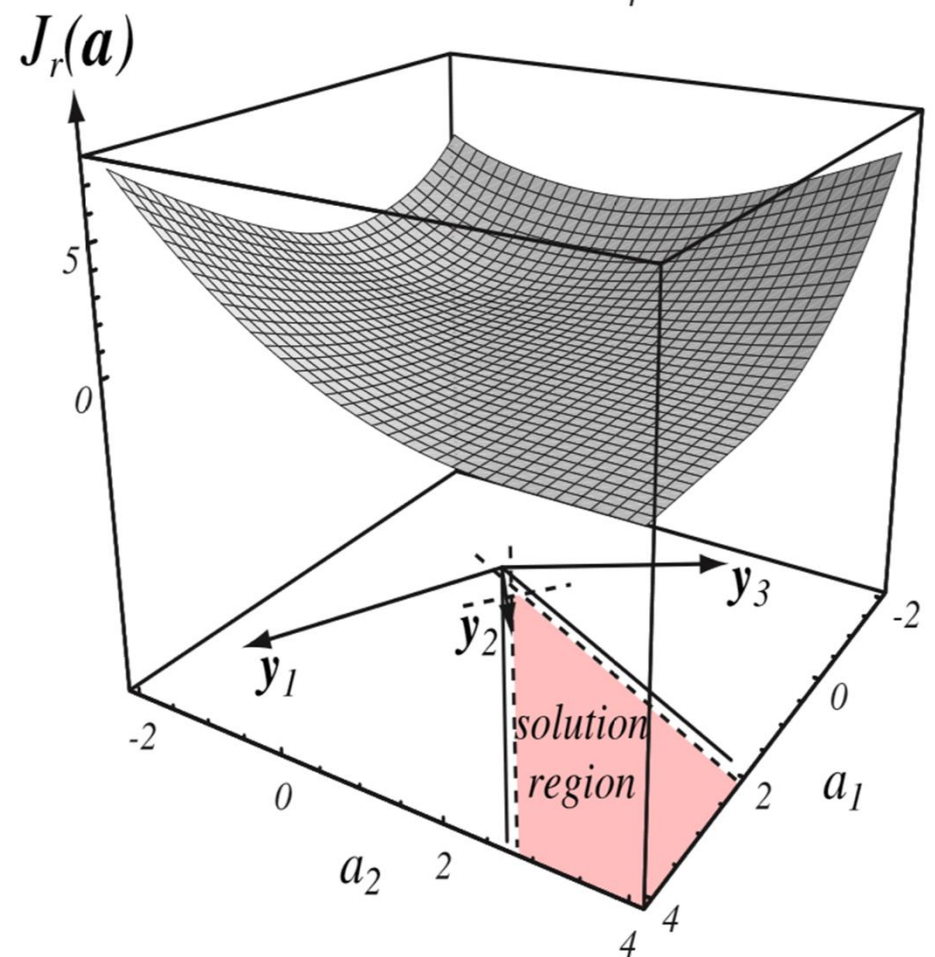
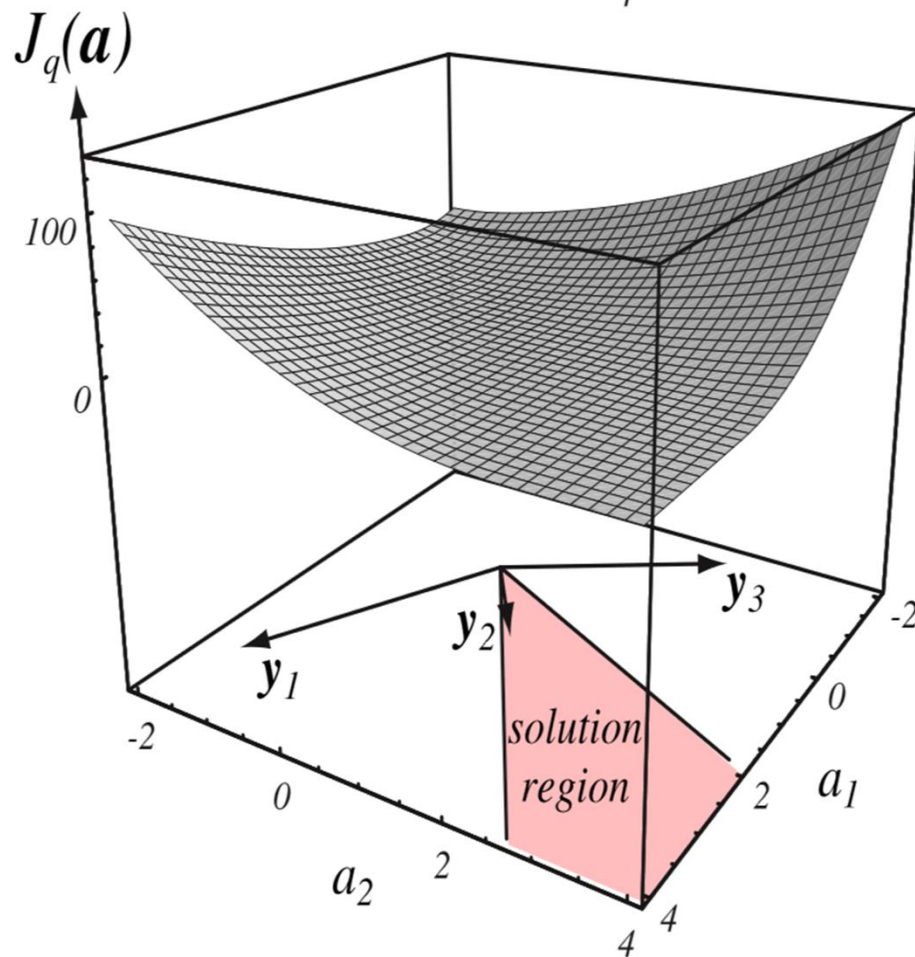
Single Sample Perceptron



Relaxation Procedures

- These are broader class of criterion functions and associated minimization methods.
- $J_q(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} (\mathbf{a}^T \mathbf{y})^2$
- Problems:
 - convergence to boundary
 - dominated by the longest sample vector
- $J_r(\mathbf{a}) = \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{Y}} \frac{(\mathbf{a}^T \mathbf{y} - b)^2}{\|\mathbf{y}\|^2}$ and $\nabla J_r(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{Y}} \frac{\mathbf{a}^T \mathbf{y} - b}{\|\mathbf{y}\|^2} \mathbf{y}$

Relaxation Procedures



Single Sample Relaxation with Margin

- Single Sample relaxation with margin

1. Initialize $\mathbf{a}, \eta(\cdot), k = 0$

2. do $k = \text{mod}(k + 1, n)$

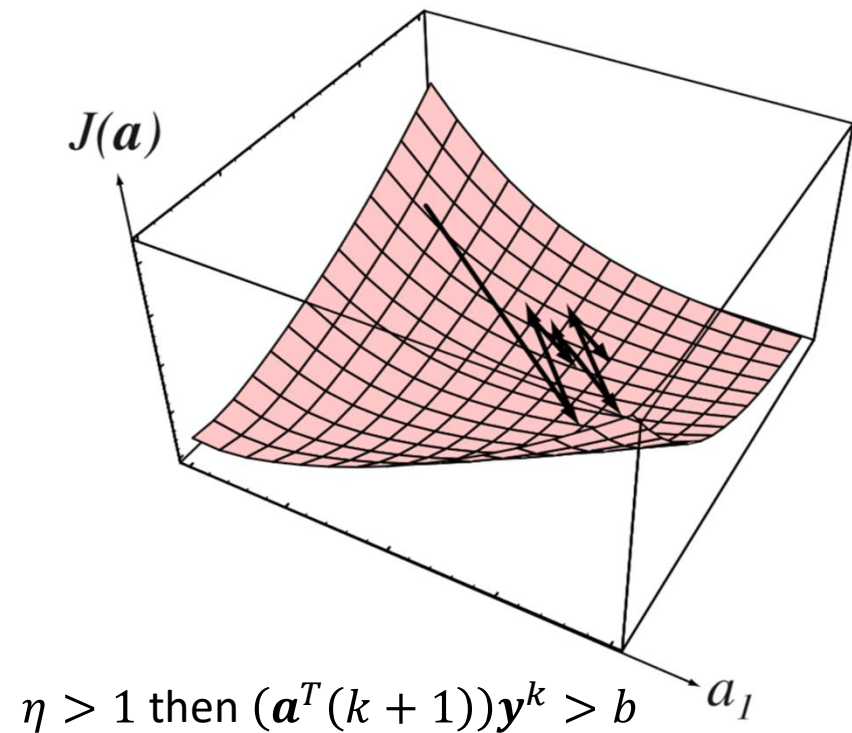
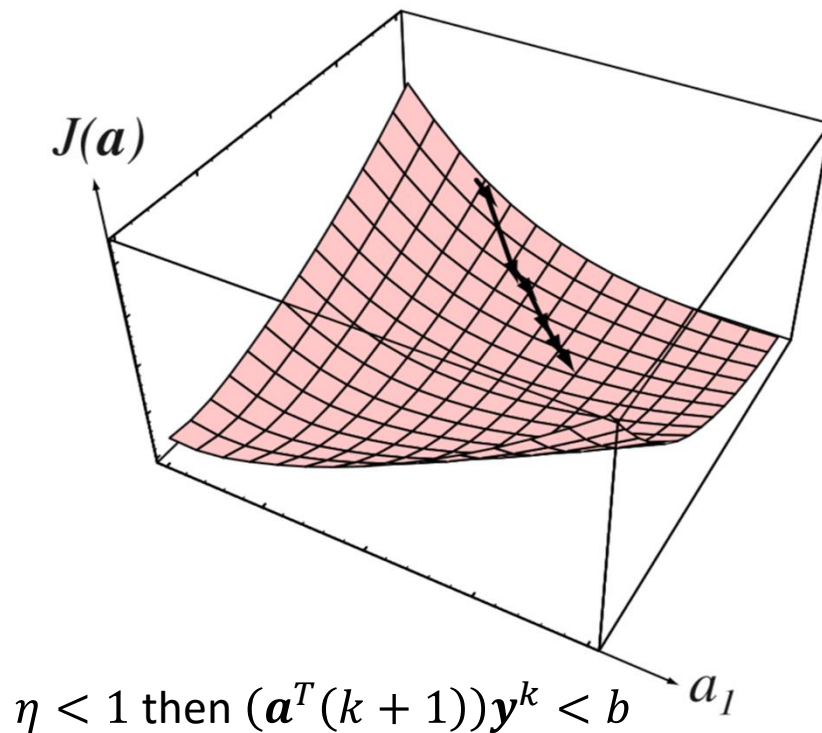
if $\mathbf{a}^T \mathbf{y}^k \leq b$ then $\mathbf{a} = \mathbf{a} + \eta(k) \frac{(b - \mathbf{a}^T \mathbf{y}^k)^2}{\|\mathbf{y}^k\|^2} \mathbf{y}^k$

3. until $\mathbf{a}^T \mathbf{y}^k > b$ for all \mathbf{y}^k

4. return \mathbf{a}

Over/Under-Relaxation

- $$r(k) = \frac{(b - \mathbf{a}^T \mathbf{y}^k)^2}{\|\mathbf{y}^k\|^2} \mathbf{y}^k$$



Non Separable Behavior

- Error Correcting Procedures
- Generalization to unseen test data not guaranteed
- Fails to handle non-separable case
- Many Heuristic exists to handle non-separable cases:
 - Forced termination of loop
 - Annealing of η with increasing k

Minimum Squared Error Procedures

- MSE consider all samples instead of just misclassified ones.
- Moved from problem of finding solution to a set of linear inequalities to a set of linear equations, i.e., $\mathbf{a}^T \mathbf{y}_i = b_i$ instead of $\mathbf{a}^T \mathbf{y}_i > 0$

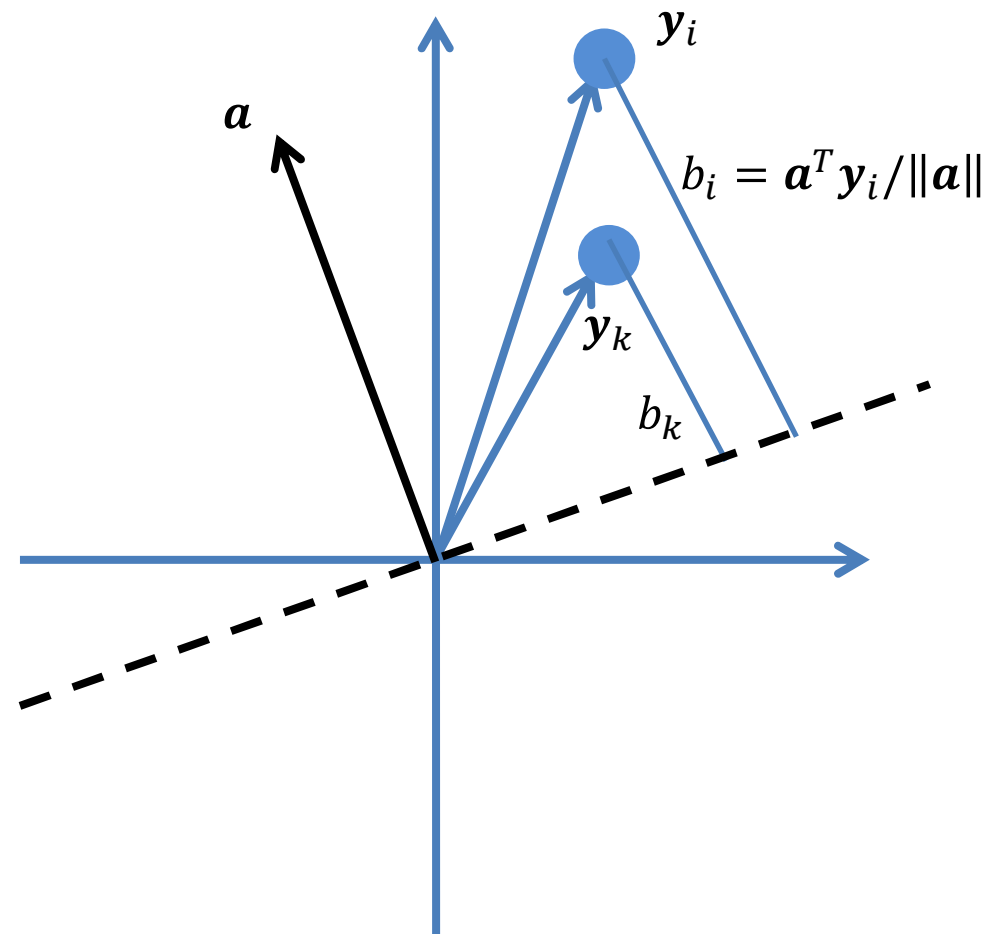
Minimum Squared Error Procedures

- MSE consider all samples instead of just misclassified ones.

$$\mathbf{a}^T \mathbf{y}_i > 0 \text{ for all samples } \mathbf{y}_i$$



$$\mathbf{a}^T \mathbf{y}_i = b_i \text{ for all samples } \mathbf{y}_i \\ \text{where } b_i > 0$$



Minimum Squared Error Procedures

- $\mathbf{Y} = [\mathbf{y}_1 \quad \cdots \quad \mathbf{y}_n]^T$ be the set of all data points where $\mathbf{y}_i = [y_{i0} \quad \cdots \quad y_{id}]^T \in \mathbb{R}^{\hat{d}=(d+1)}$
- Let $\mathbf{a} = [a_0 \quad \cdots \quad a_d]^T$ and $\mathbf{b} = [b_1 \quad \cdots \quad b_n]^T$
- $\mathbf{Y}\mathbf{a} = \mathbf{b}$ (over-determined problem as $n \gg \hat{d}$)
- $\mathbf{a} = \mathbf{Y}^{-1}\mathbf{b}$ not possible (Y is rectangular and possibly singular)
- No exact solution ! We look for approximate solution.

Minimum Squared Error Procedures

- $\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$ (Error definition)
- $J(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$
- $\nabla J(\mathbf{a}) = 2\mathbf{Y}^T(\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0$
- $\mathbf{Y}^T\mathbf{Y}\mathbf{a} = \mathbf{Y}^T\mathbf{b}$
- $\mathbf{a} = (\mathbf{Y}^T\mathbf{Y})^{-1}\mathbf{Y}^T\mathbf{b} = \mathbf{Y}^\dagger\mathbf{b}$

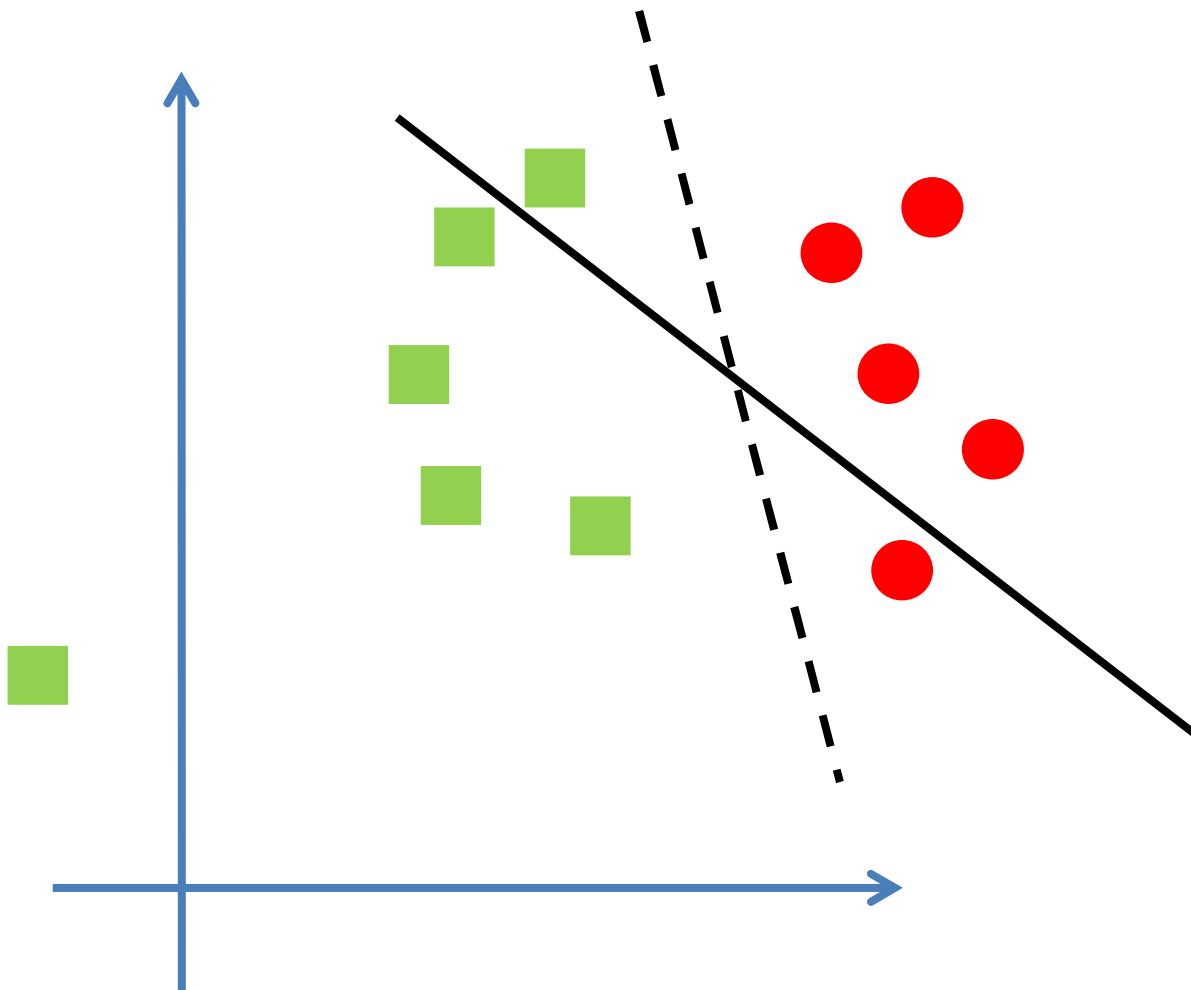
Minimum Squared Error Procedures

- $(Y^T Y)$ is a square matrix and often non-singular and hence invertible.
- There could be many solutions for weight vector \mathbf{a} based on choice of vector \mathbf{b} .
- A separating hyperplane is guaranteed if $(Y\mathbf{a} > \mathbf{0})$, i.e., $\forall i (\mathbf{a}^T \mathbf{y}_i) > 0$
- However, we have $Y\mathbf{a} \approx \mathbf{b}$, i.e., $Y\mathbf{a} = \mathbf{b} + \boldsymbol{\varepsilon}$.
- In practice, some entries of \mathbf{b} can be negative if $|b_i| < |\varepsilon_i|$ and $\varepsilon_i < 0$.

Minimum Squared Error Procedures

- Thus, even in linearly separable case, least square solution \mathbf{a} might not yield a separating hyperplane but a **reasonable** one.
- An arbitrary scaling of \mathbf{b} to overcome the $-\epsilon$ values is **not helpful** as it translates to scaling up the \mathbf{a} vector.
- However, relative difference in elements of \mathbf{b} can help in improving the classification, especially to handle the case of outlier data points.

Minimum Squared Error Procedures



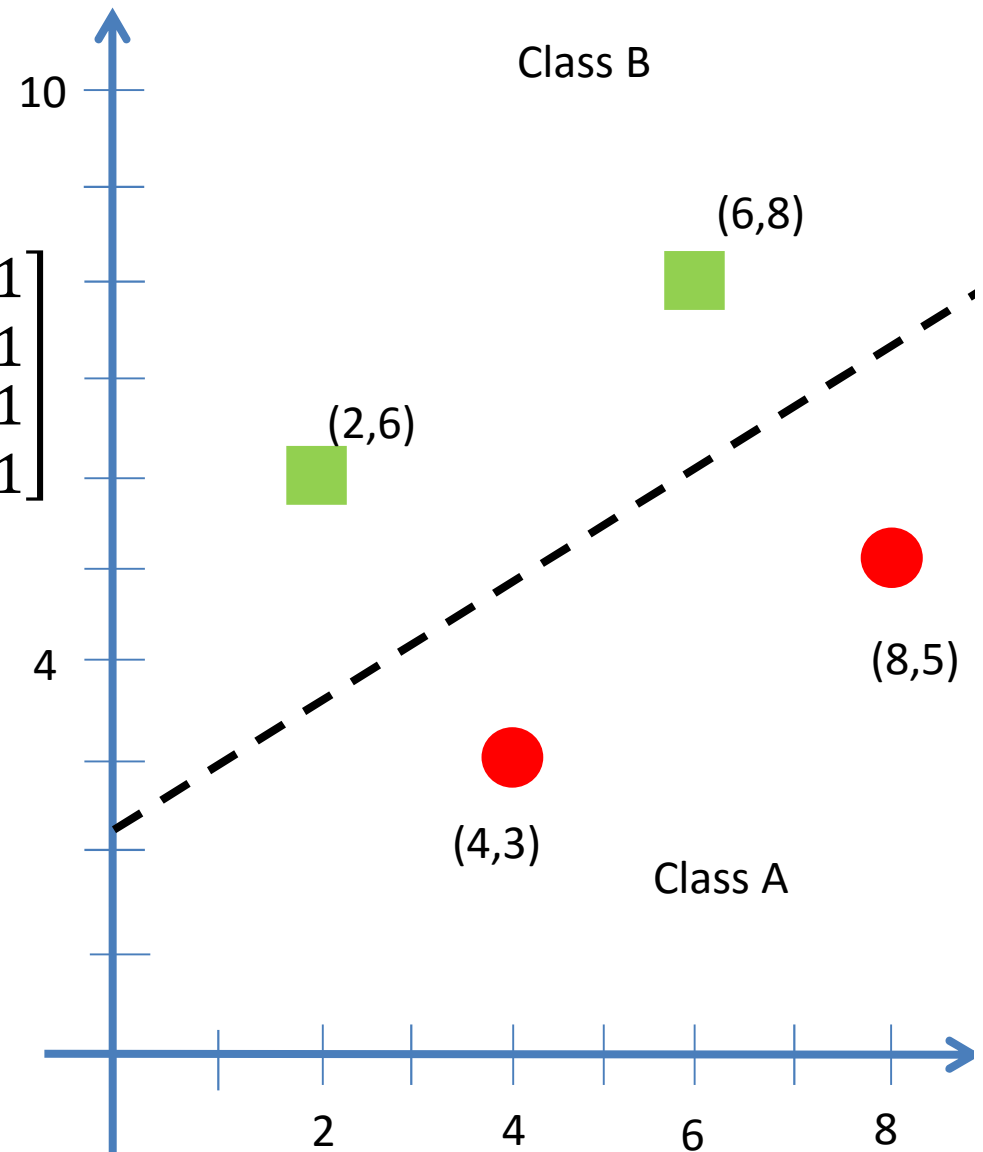
Example Walkthrough

- Class A: (8,5), (4,3)
- Class B: (2,6), (6,8)

- $Y^T = \begin{bmatrix} 1 & 1 & -1 & -1 \\ 8 & 4 & -2 & -6 \\ 5 & 3 & -6 & -8 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

- $\mathbf{a} = Y^\dagger \mathbf{b} = \begin{bmatrix} 1.5 \\ 0.25 \\ -0.5 \end{bmatrix}$

- $Y\mathbf{a} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$



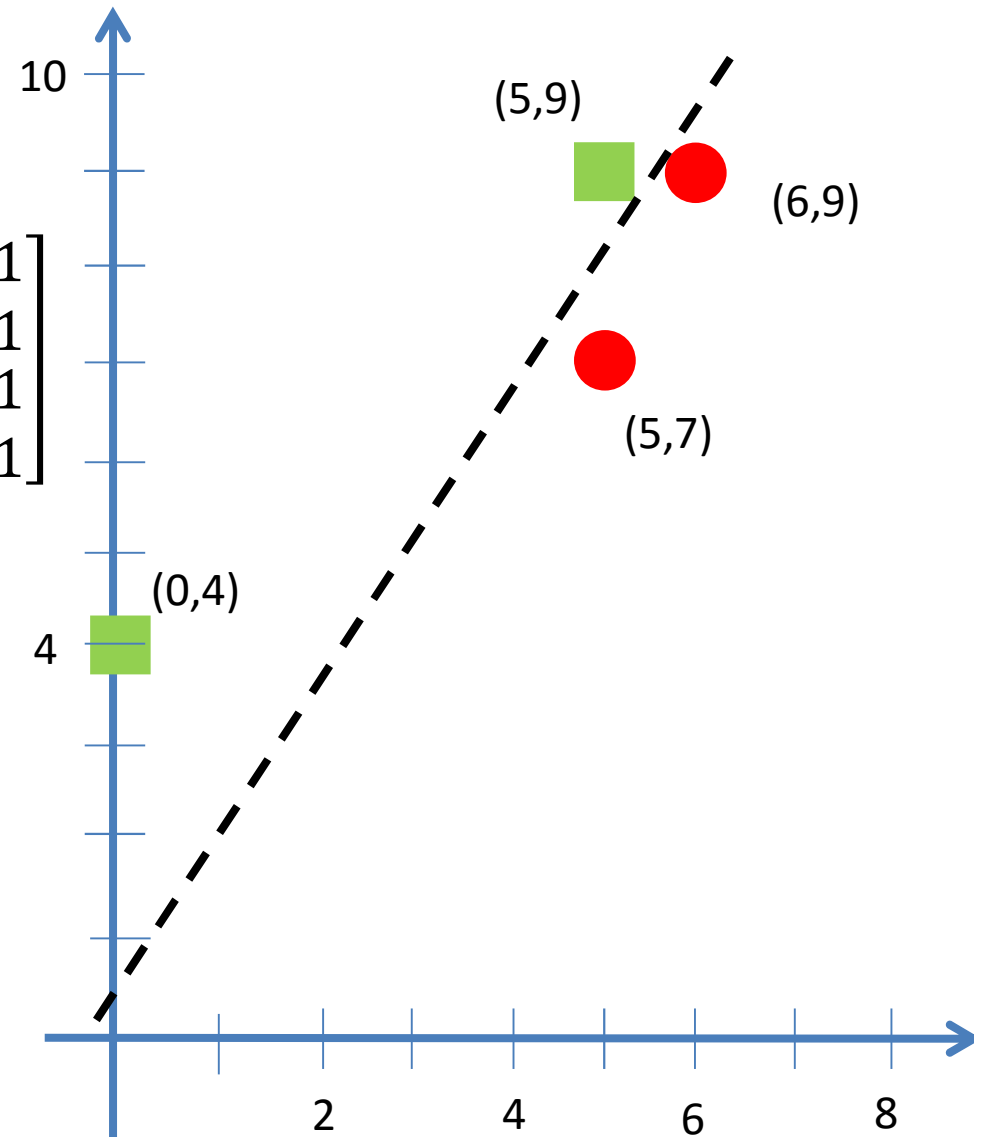
Example Walkthrough

- Class A: (6,9), (5,7)
- Class B: (5,9), (0,4)

- $Y^T = \begin{bmatrix} 1 & 1 & -1 & -1 \\ 6 & 5 & -5 & 0 \\ 9 & 7 & -9 & -4 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

- $\mathbf{a} = Y^\dagger \mathbf{b} = \begin{bmatrix} 2.66 \\ 1.04 \\ -0.94 \end{bmatrix}$

- $Y\mathbf{a} = \begin{bmatrix} 0.43 \\ 1.28 \\ 0.60 \\ 1.11 \end{bmatrix}$



Example Walkthrough

- Class A: (6,9), (5,7)
- Class B: (5,9),(0,10)

$$\bullet \quad Y^T = \begin{bmatrix} 1 & 1 & -1 & -1 \\ 6 & 5 & -5 & 0 \\ 9 & 7 & -9 & -10 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\bullet \quad \mathbf{a} = Y^\dagger \mathbf{b} = \begin{bmatrix} 3.21 \\ 0.15 \\ -0.43 \end{bmatrix}$$

$$\bullet \quad Y\mathbf{a} = \begin{bmatrix} 0.19 \\ 0.91 \\ -0.04 \\ 1.16 \end{bmatrix}$$

