

Statistical Methods in Artificial Intelligence

CSE471 - Monsoon 2016 : Lecture 05



Avinash Sharma
CVIT, IIIT Hyderabad

Lecture 05: Plan

- Recap
- Minimum Squared Error Procedures
- The Widrow-Hoff /LMS Procedure
- The Ho-Kashyap Procedure

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

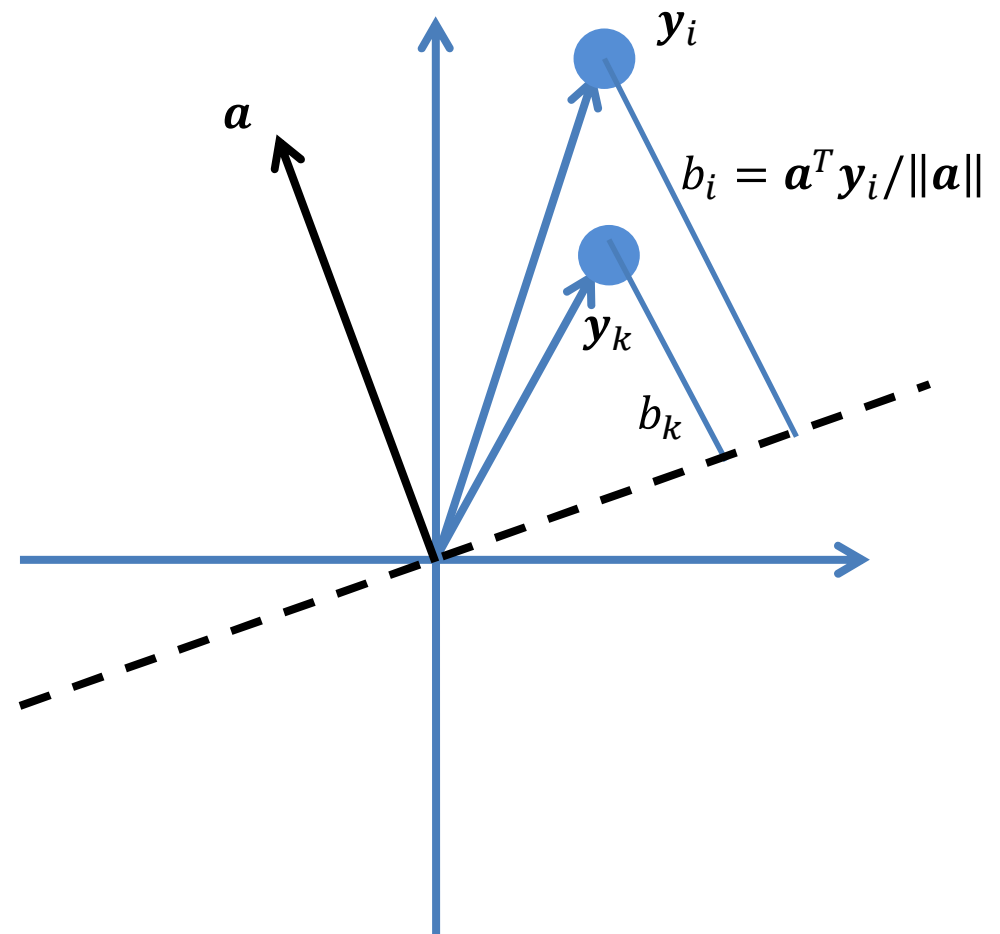
Minimum Squared Error Procedures

- MSE consider all samples instead of just misclassified ones.

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



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



Minimum Squared Error Procedures

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

Iterative Gradient Descent for MSE

- $(Y^T Y)$ is a square matrix and **often(?)** non-singular and hence invertible.
 - Singular if data points are highly correlated (rows Y are almost linear combination of each other).
- Computing inverse of $(Y^T Y)$ can be too expensive for high dimensional data (i.e., large matrix inversion).
- These issues can be avoided by adapting an iterative gradient descent solution.

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

$$\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(k) Y^T (Y \mathbf{a}(k) - \mathbf{b})$$

The Widrow-Hoff /LMS Procedure

- Reduce storage requirements significantly by considering single sample sequentially, i.e., $\hat{\mathbf{d}}$ -dimensional data point (\mathbf{y}_i) separately.
- Prof. Bernard Widrow and Dr. Ted Hoff proved the convergence of individual “Least Mean Square (LMS) Error” minimization in the stochastic sense.

1. Initialize $\mathbf{a}, \mathbf{b}, \theta$ (*threshold*), $\eta(\cdot), k = 0, i = 0$
2. do $i = \text{mod}(i + 1, n)$
3. $\mathbf{a}(k + 1) = \mathbf{a}(k) - \eta(i) \mathbf{y}_i (\mathbf{y}_i^T \mathbf{a}(k) - b_i)$
4. until $|\eta(k) \mathbf{y}_i (\mathbf{y}_i^T \mathbf{a}(k) - b_i)| < \theta$
5. return \mathbf{a}

*Use Annealing
for learning rate
 $\eta(k) = \eta(1)/k$

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}

The Ho-Kashyap Procedure

- In the MSE procedure, for linearly separable case one can find a separable hyperplane (\mathbf{a}) iff $\mathbf{b} > \mathbf{0}$, which however is chosen arbitrarily.
- Here, the idea is to find both \mathbf{a} and \mathbf{b} , simultaneously using the modified gradient descent procedure that minimizes $J(\mathbf{a}, \mathbf{b})$.
 - ❖ $J(\mathbf{a}, \mathbf{b}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$
 - ❖ Fix \mathbf{b} and minimize $J(\mathbf{a}, \mathbf{b})$ w.r.t. to \mathbf{a}
 - ❖ Fix \mathbf{a} and minimize $J(\mathbf{a}, \mathbf{b})$ w.r.t. to \mathbf{b}
- The partial derivatives of $J(\mathbf{a}, \mathbf{b})$ will be
 - ❖ $\nabla_{\mathbf{a}} J(\mathbf{a}, \mathbf{b}) = 2\mathbf{Y}^T(\mathbf{Y}\mathbf{a} - \mathbf{b})$
 - ❖ $\nabla_{\mathbf{b}} J(\mathbf{a}, \mathbf{b}) = -2(\mathbf{Y}\mathbf{a} - \mathbf{b})$ or $(\mathbf{Y}\mathbf{a} - \mathbf{b}) = -\frac{1}{2}\nabla_{\mathbf{b}} J(\mathbf{a}, \mathbf{b}) = \mathbf{e}$

The Ho-Kashyap Procedure

- The second step will look like

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

- However, during the iterative gradient descent we need to always ensure that $\mathbf{b} > \mathbf{0}$.
- This can be enforced by reducing the positive elements of $\nabla_{\mathbf{b}} J$ to zero.

$$\mathbf{b}(k + 1) = \mathbf{b}(k) - \eta(k) \frac{1}{2} (\nabla_{\mathbf{b}} J - |\nabla_{\mathbf{b}} J|) \text{ ** or}$$

$$\mathbf{b}(k + 1) = \mathbf{b}(k) + \eta(k) \left(-\frac{1}{2} \nabla_{\mathbf{b}} J + \frac{1}{2} |\nabla_{\mathbf{b}} J| \right) \text{ or}$$

** $|\mathbf{v}|$ is a vector with all components having absolute values of vector \mathbf{v}

$$\mathbf{b}(k + 1) = \mathbf{b}(k) + \eta(k)(\mathbf{e}_k + |\mathbf{e}_k|)$$

The Ho-Kashyap Procedure

1. Initialize $\mathbf{a}, \mathbf{b}, \eta(\cdot) < 1, k = 0, \text{threshold}(b_{min}, k_{max})$
2. do $k = k + 1$
3. $\mathbf{e}_k = (\mathbf{Y}\mathbf{a}(k) - \mathbf{b}(k))$
4. $\mathbf{b}(k + 1) = \mathbf{b}(k) + \eta(k)(\mathbf{e}_k + |\mathbf{e}_k|)$
5. $\mathbf{a}(k + 1) = \mathbf{Y}^\dagger \mathbf{b}(k + 1)$
6. if $|\mathbf{e}_k| \leq b_{min}$ then return \mathbf{a}, \mathbf{b} and exit
7. until $k < k_{max}$
8. Print "NO SOLUTION"

The Ho-Kashyap Procedure

- Since \mathbf{a} is determined by \mathbf{b} only (step 5), the procedure can be interpreted as the one which sequentially produce margin vectors.
- Not doing steepest descent anymore, but we are still doing descent and ensure that \mathbf{b} is positive.
- In case when none of the element of \mathbf{e}_k vector is positive then
 - Either they all are (close to) zero which means we have solution
 - Or there is no solution
- In Linearly Separable case: solution will be found in finite step when we reach to $\mathbf{e}_k = \mathbf{0}$.
- In the Non Separable case: \mathbf{e}_k will have only negative components.
 - No bound on number of iterations needed to prove the non-separability.