

CS 559: Machine Learning Fundamentals and Applications 9th Set of Notes

Instructor: Philippos Mordohai

Webpage: www.cs.stevens.edu/~mordohai

E-mail: Philippos.Mordohai@stevens.edu

Office: Lieb 215

Overview

- Logistic Regression
 - Notes by T. Mitchell
 - Barber Ch. 17
 - HTF Ch. 4
- Linear Discriminant Functions (Slides based on Olga Veksler's)
 - Optimization with gradient descent
 - Perceptron Criterion Function
 - Batch perceptron rule
 - Single sample perceptron rule
 - Minimum Squared Error (MSE) rule

Overview (cont.)

- Support Vector Machines (SVM)
 - Introduction
 - Linear Discriminant
 - Linearly Separable Case
 - Linearly Non Separable Case
 - Kernel Trick
 - Non Linear Discriminant
 - Multi-class SVMs
- See HTF Ch. 12

Logistic Regression

- Idea: generative models compute $P(Y|X)$ by learning $P(Y)$ and $P(X|Y)$
- Why not learn $P(Y|X)$ directly?

Logistic Regression

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
 - assume all X_i are conditionally independent given Y
 - model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i^2)$
 - model $P(Y)$ as Bernoulli (π)
 - Y is 1 with probability π

Derivation of P(Y|X)

$$P(Y = 1|X) = \frac{P(Y = 1)P(X|Y = 1)}{P(Y = 1)P(X|Y = 1) + P(Y = 0)P(X|Y = 0)}$$

$$= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}}$$

$$= \frac{1}{1 + \exp(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})}$$

$$= \frac{1}{1 + \exp((\ln \frac{1-\pi}{\pi}) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)})}$$

$$P(x | y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{-\frac{(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

$$\sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)$$

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Very Convenient

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

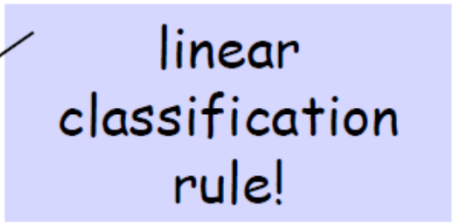
implies

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$\frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i)$$

linear
classification
rule!



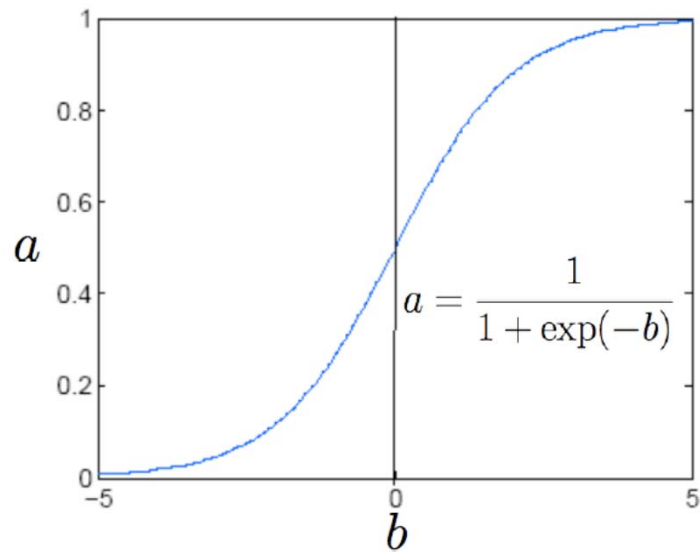
implies

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

Very Convenient

- Posteriors sum to 1 and remain in $[0, 1]$
- Logit: $l = \text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \alpha + \beta x$
 - l is linear in x
- Probability: $p = \frac{e^l}{1+e^l}$

Logistic Function



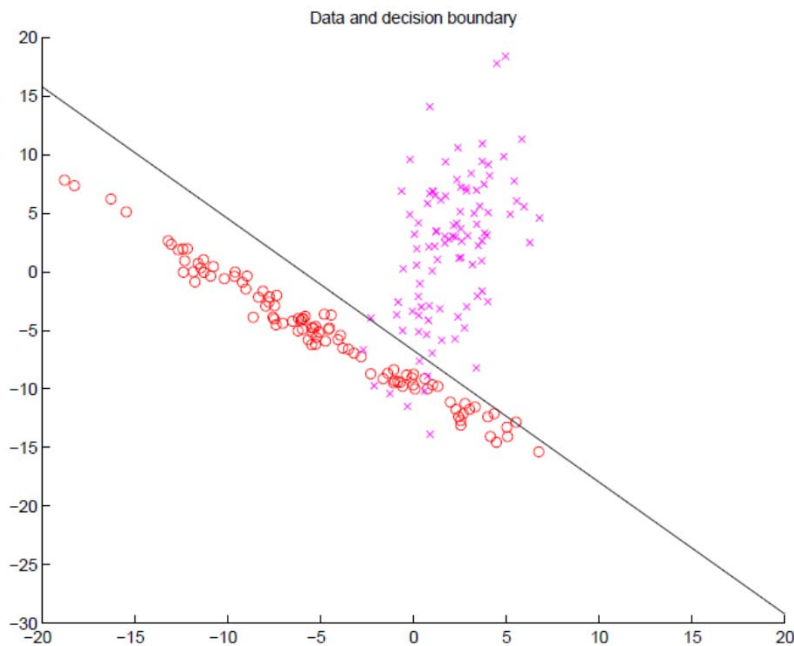
- $p = 0, \quad l = \log\left(\frac{p}{1-p}\right) = -\infty$
- $p = \frac{1}{2}, \quad l = \log\left(\frac{p}{1-p}\right) = 0$
- $p = 1, \quad l = \log\left(\frac{p}{1-p}\right) = +\infty$

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Decision Boundary

- How to make decisions given

$$P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$



Logistic Regression More Generally

- Logistic regression when Y is not boolean (but still discrete)
 - $y \in \{y_1 \dots y_R\}$: learn $R-1$ sets of weights

- for $k < R$
$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^n w_{ki} X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

- for $k = R$
$$P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

Training Logistic Regression: MLE

- We have L training examples:

$$\{\langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle\}$$

- Maximum likelihood estimate for parameters W

$$\begin{aligned} W_{MLE} &= \arg \max_W P(\langle X^1, Y^1 \rangle \dots \langle X^L, Y^L \rangle | W) \\ &= \arg \max_W \prod_l P(\langle X^l, Y^l \rangle | W) \end{aligned}$$

- Maximum conditional likelihood estimate

Training Logistic Regression: MCLE

- Choose parameters $\langle w_0 \dots w_n \rangle$ to maximize conditional likelihood of training data

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

- Training data $D = \{ \langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle \}$
- Data likelihood = $\prod_l P(X^l, Y^l | W)$
- Data conditional likelihood = $\prod_l P(Y^l | X^l, W)$

$$W_{MCLE} = \arg \max_W \prod_l P(Y^l | W, X^l)$$

Conditional Log Likelihood

$$l(W) \equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W)$$

$$P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\begin{aligned} l(W) &= \sum_l Y^l \ln P(Y^l = 1 | X^l, W) + (1 - Y^l) \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l \ln \frac{P(Y^l = 1 | X^l, W)}{P(Y^l = 0 | X^l, W)} + \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

Maximizing Conditional Log Likelihood

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

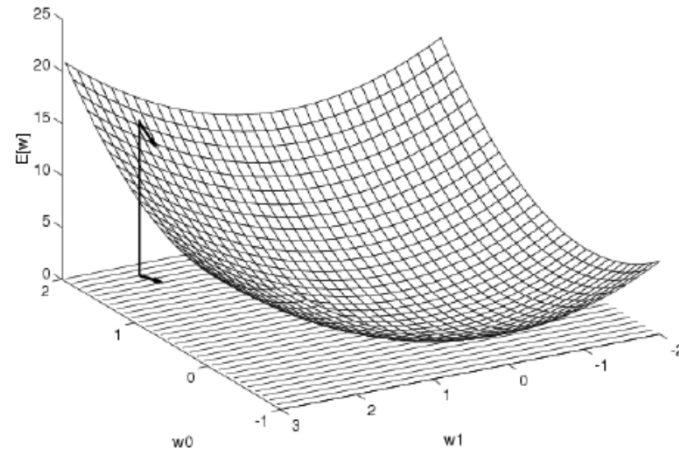
$$P(Y = 1|X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i w_i X_i^l)) \end{aligned}$$

Good news: $l(W)$ is concave function of W

Bad news: no closed-form solution to maximize $l(W)$

Gradient Descent



Gradient

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Maximize Conditional Log Likelihood: Gradient Ascent

$$\begin{aligned}l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l))\end{aligned}$$

$$\frac{\partial l(W)}{\partial w_i} = \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

Maximize Conditional Log Likelihood: Gradient Ascent

$$\begin{aligned}l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l))\end{aligned}$$

$$\frac{\partial l(W)}{\partial w_i} = \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

Gradient ascent algorithm: iterate until change $< \varepsilon$

For all i , repeat

$$w_i \leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

Logistic Regression: Summary

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
 - assume all X_i are conditionally independent given Y
 - model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i^2)$
 - model $P(Y)$ as Bernoulli (π)
- Then $P(Y|X)$ is of this form and we can directly estimate W

$$P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Linear Discriminant Functions

Augmented Feature Vector

- Linear discriminant function: $g(x) = w^t x + w_0$

- Can rewrite it:

$$g(x) = \underbrace{[w_0 \quad w^t]}_{\substack{\text{new weight} \\ \text{vector } a}} \underbrace{\begin{bmatrix} 1 \\ x \end{bmatrix}}_{\substack{\text{new feature} \\ \text{vector } y}} = a^t y = g(y)$$

- y is called the **augmented feature** vector
- Added a dummy dimension to get a completely equivalent new **homogeneous problem**

old problem

$$g(x) = w^t x + w_0$$

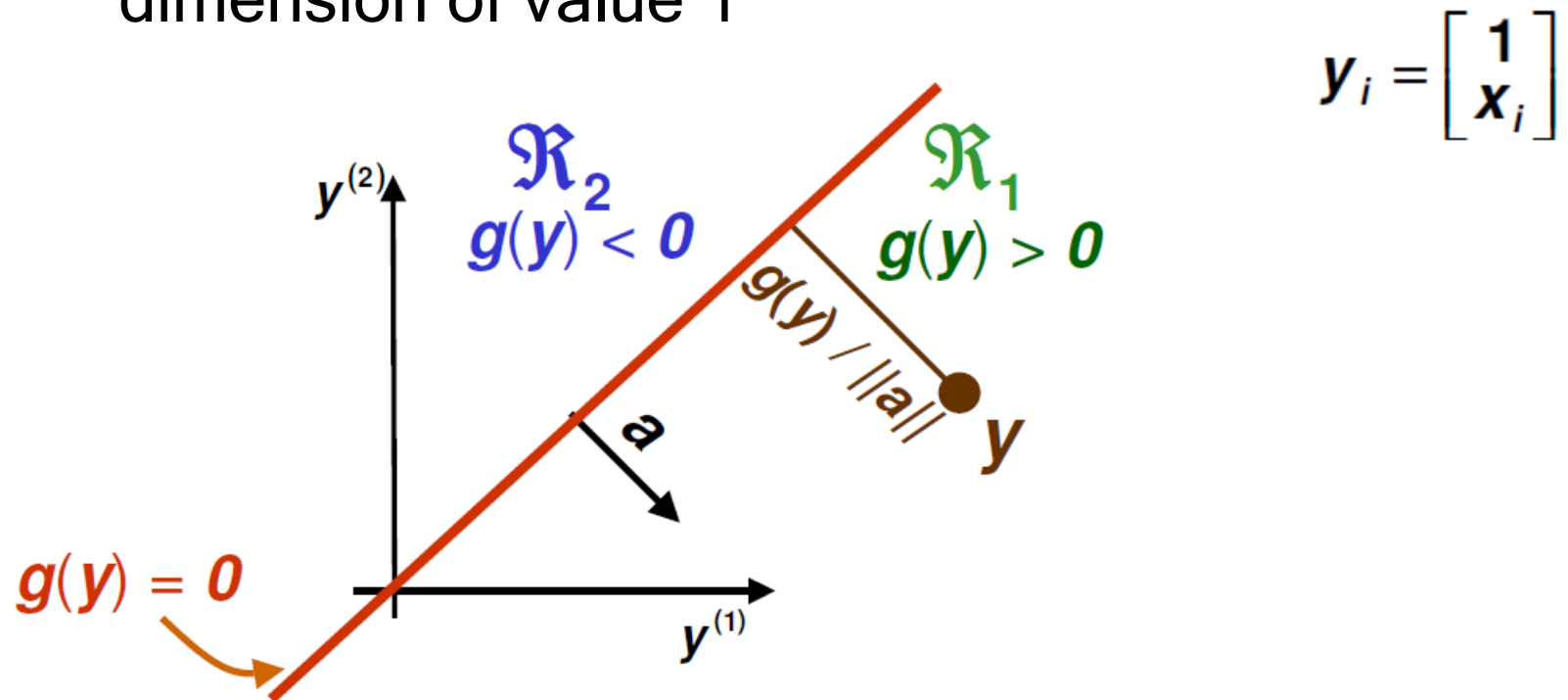
$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

new problem

$$g(y) = a^t y$$

$$\begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

- Feature augmentation is done for simpler notation
- From now on, always assume that we have augmented feature vectors
 - Given samples x_1, \dots, x_n convert them to augmented samples y_1, \dots, y_n by adding a new dimension of value 1



Training Error

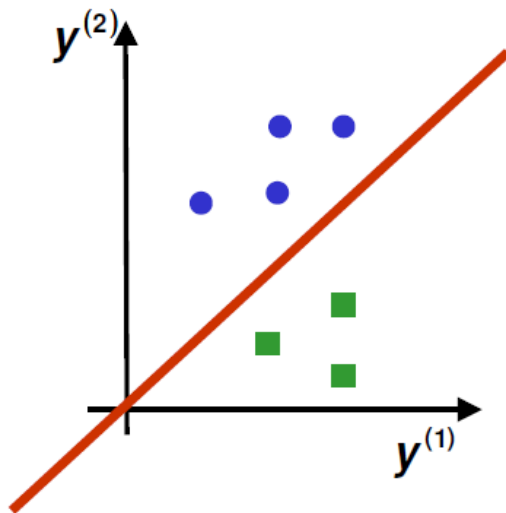
- For the rest of this part, assume we have 2 classes
 - Samples: y_1, \dots, y_n , some in class 1, some in class 2
- Use samples to determine weights a in the discriminant function $g(y) = a^t y$
- What should the criterion for determining a be?
- For now, suppose we want to minimize the training error (the number of misclassified samples y_1, \dots, y_n)
- Recall that: $g(y_i) > 0 \Rightarrow y_i$ *classified as c_1*
 $g(y_i) < 0 \Rightarrow y_i$ *classified as c_2*
- Thus training error is 0 if
$$\begin{cases} g(y_i) > 0 & \forall y_i \in c_1 \\ g(y_i) < 0 & \forall y_i \in c_2 \end{cases}$$

“Normalization”

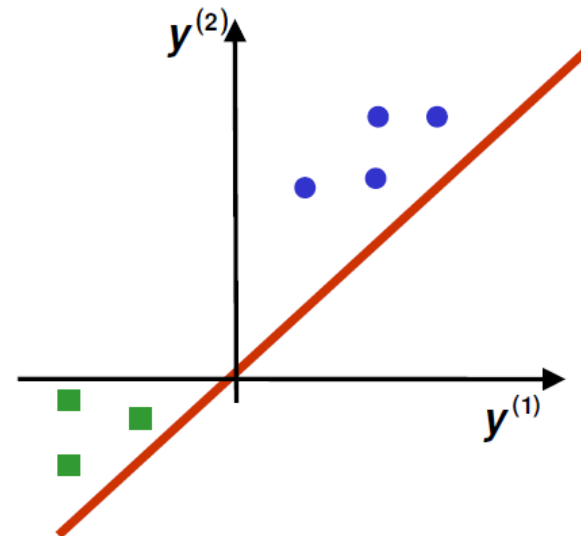
- Thus training error is 0 if: $\begin{cases} \mathbf{a}^t \mathbf{y}_i > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ \mathbf{a}^t \mathbf{y}_i < 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$
- Equivalently, training error is 0 if: $\begin{cases} \mathbf{a}^t \mathbf{y}_i > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ \mathbf{a}^t (-\mathbf{y}_i) > 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$
- This suggests “normalization” (a.k.a. reflection):
 1. Replace all examples from class 2 by:
$$\mathbf{y}_i \rightarrow -\mathbf{y}_i \quad \forall \mathbf{y}_i \in \mathbf{c}_2$$
 2. Seek weight vector \mathbf{a} such that
$$\mathbf{a}^t \mathbf{y}_i > 0 \quad \forall \mathbf{y}_i$$
 - If such \mathbf{a} exists, it is called a separating or solution vector
 - Original samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ can indeed be separated by a line

Normalization

before normalization



after "normalization"

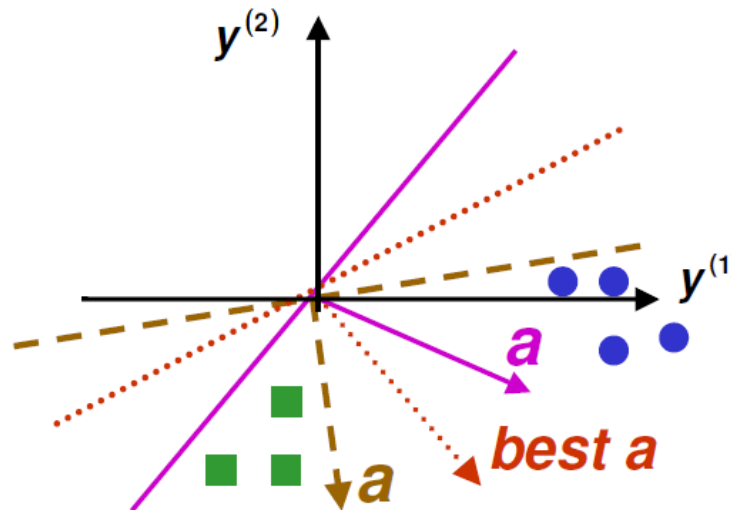


- Seek a hyperplane that separates patterns from different categories
- Seek hyperplane that puts *normalized* patterns on the same(positive) side

Solution Region

- Find weight vector \mathbf{a} such that for all samples:

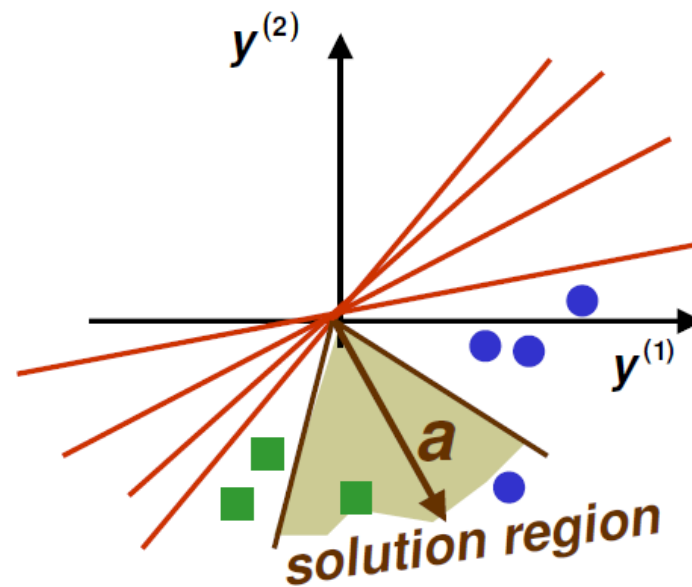
$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^d a_k y_i^{(k)} > 0$$



- In general, there can be many solutions

Solution Region

- **Solution region** for \mathbf{a} : set of all possible solutions defined in terms of normal \mathbf{a} to the separating hyperplane



Optimization

- Need to minimize a function of many variables

$$J(\mathbf{x}) = J(x_1, \dots, x_d)$$

- We know how to minimize $J(\mathbf{x})$
 - Take partial derivatives and set them to zero

$$\begin{bmatrix} \frac{\partial}{\partial x_1} J(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_d} J(\mathbf{x}) \end{bmatrix} = \nabla J(\mathbf{x}) = 0$$

gradient

Optimization

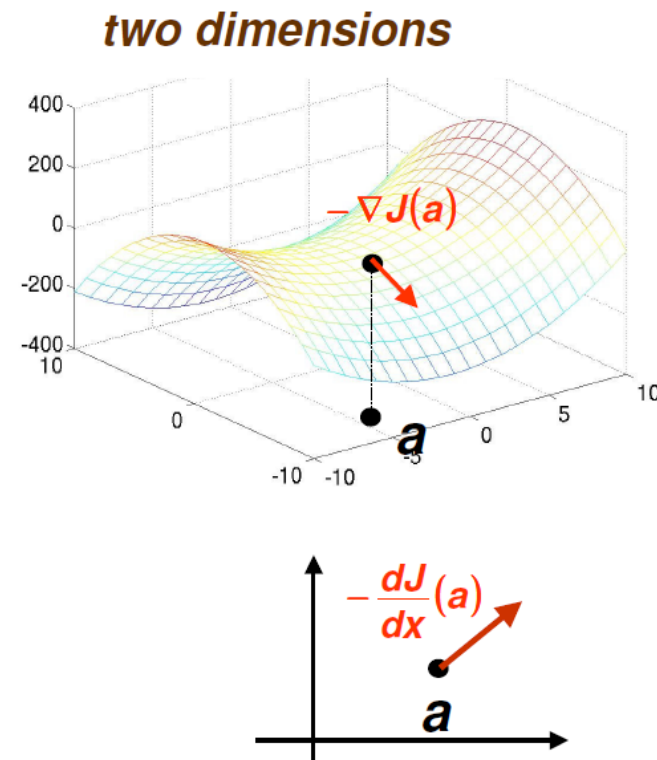
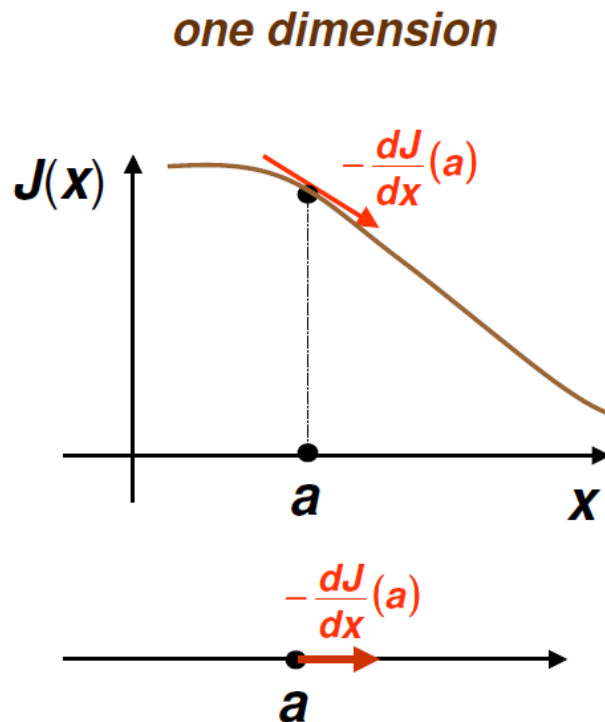
- However solving analytically is not always easy
 - For example:

$$\begin{cases} \sin(x_1^2 + x_2^3) + e^{x_4^2} = 0 \\ \cos(x_1^2 + x_2^3) + \log(x_5^3)^{x_4^2} = 0 \end{cases}$$

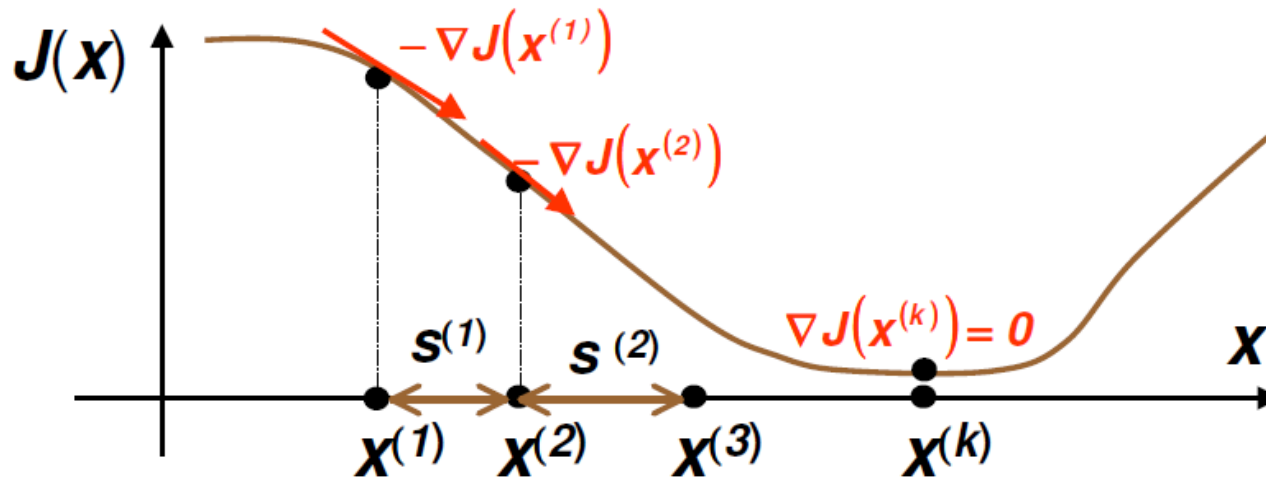
- Sometimes it is not even possible to write down an analytical expression for the derivative (example later today)

Gradient Descent

- Gradient $\nabla J(x)$ points in direction of steepest increase of $J(x)$, and $-\nabla J(x)$ in direction of steepest decrease



Gradient Descent



Gradient Descent for minimizing any function $J(x)$

– Set $k = 1$ and $x^{(1)}$ to some initial guess for the weight vector

– While $\eta^{(k)} |\nabla J(x^{(k)})| > \epsilon$

• Choose learning rate $\eta^{(k)}$

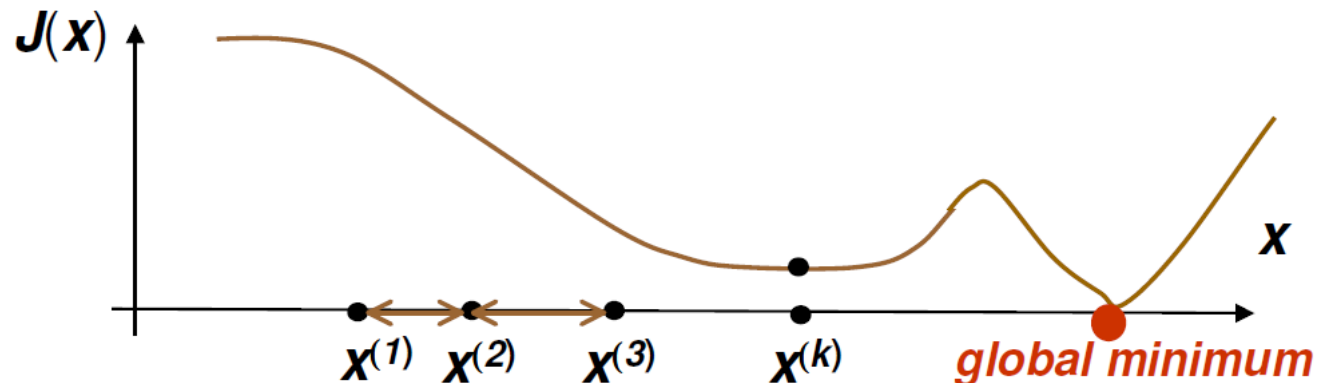
$$x^{(k+1)} = x^{(k)} - \eta^{(k)} \nabla J(x)$$

(update rule)

$$k = k + 1$$

Gradient Descent

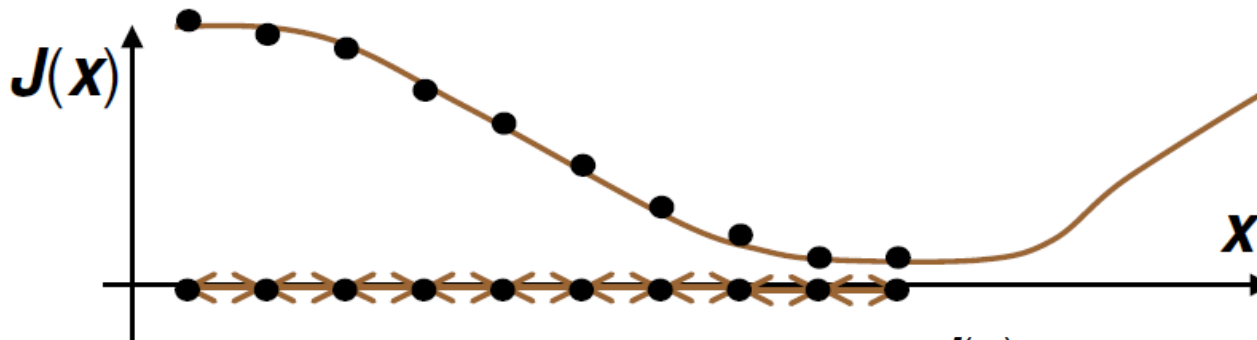
- Gradient descent is guaranteed to only find **local minima**



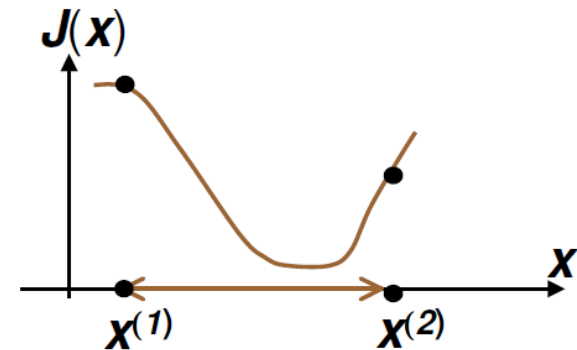
- Nevertheless gradient descent is very popular because it is simple and applicable to any function

Gradient Descent

- Main issue: how to set parameter η (learning rate)
 - If η is too small, too many iterations



- If η is too large may overshoot the minimum and possibly never find it



LDF Criterion Function

- Find weight vector \mathbf{a} such that for all samples $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^d \mathbf{a}_k \mathbf{y}_i^{(k)} > 0$$

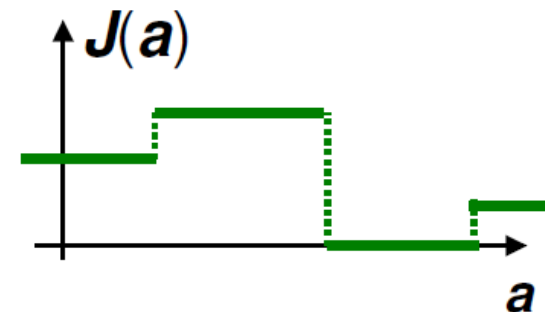
- Need criterion function $J(\mathbf{a})$ which is minimized when \mathbf{a} is a solution vector
- Let Y_M be the set of examples misclassified by \mathbf{a}

$$Y_M(\mathbf{a}) = \{ \mathbf{y}_i \text{ s.t. } \mathbf{a}^t \mathbf{y}_i < 0 \}$$

- First natural choice: number of misclassified examples

$$J(\mathbf{a}) = |Y_M(\mathbf{a})|$$

- Piecewise constant, gradient descent is useless

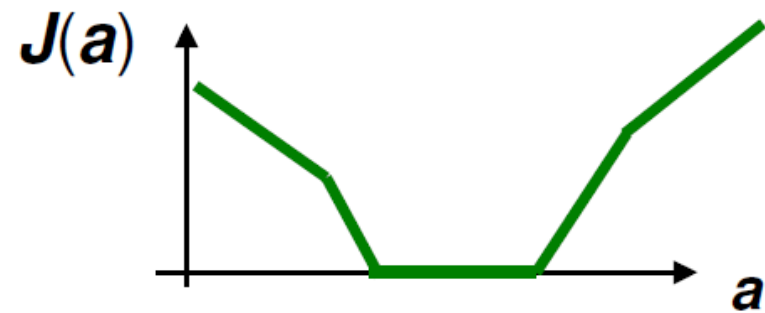
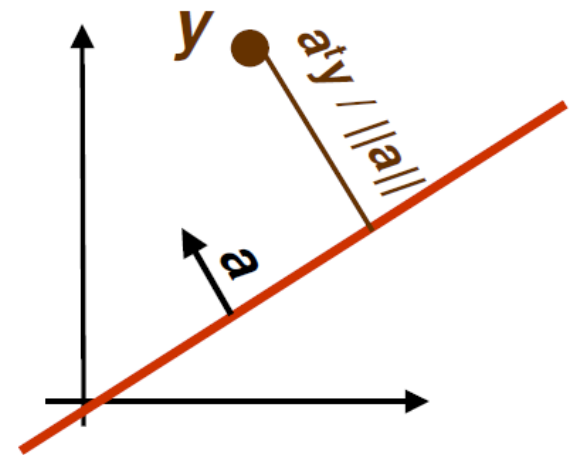


Perceptron

Perceptron Criterion Function

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

- If \mathbf{y} is misclassified, $\mathbf{a}^t \mathbf{y} < 0$
- Thus $J_p(\mathbf{a}) > 0$
- $J_p(\mathbf{a})$ is $\|\mathbf{a}\|$ times the sum of distances of misclassified examples to decision boundary
- $J_p(\mathbf{a})$ is piecewise linear and thus suitable for gradient descent



Perceptron Batch Rule

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

- Gradient of $J_p(\mathbf{a})$ is: $\nabla J_p(\mathbf{a}) = \sum_{y \in Y_M} (-\mathbf{y})$
 - Y_M are samples misclassified by $\mathbf{a}^{(k)}$
 - It is not possible to solve $\nabla J_p(\mathbf{a}) = \mathbf{0}$ analytically because of Y_M
- Update rule for gradient descent: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta^{(k)} \nabla J(\mathbf{x})$
- Thus the *gradient decent batch update rule* for $J_p(\mathbf{a})$ is:
$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \sum_{y \in Y_M} \mathbf{y}$$
- It is called batch rule because it is based on all misclassified examples

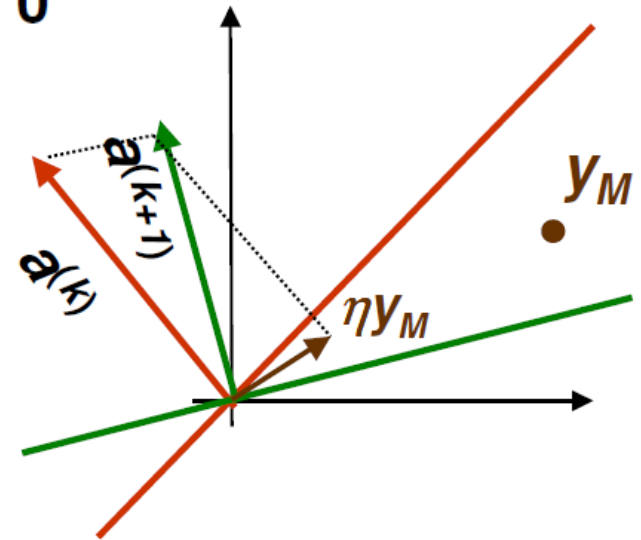
Perceptron Single Sample Rule

- The *gradient decent single sample rule* for $J_p(\mathbf{a})$ is:
$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$

- Note that \mathbf{y}_M is one sample misclassified by $\mathbf{a}^{(k)}$
- Must have a consistent way of visiting samples

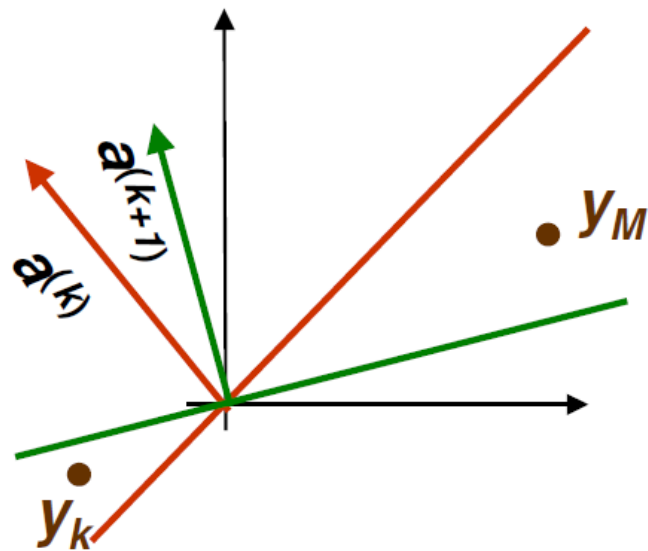
- Geometric Interpretation:

- \mathbf{y}_M misclassified by $\mathbf{a}^{(k)}$ $(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$
- \mathbf{y}_M is on the wrong side of decision hyperplane
- Adding $\eta \mathbf{y}_M$ to \mathbf{a} moves the new decision hyperplane in the right direction with respect to \mathbf{y}_M

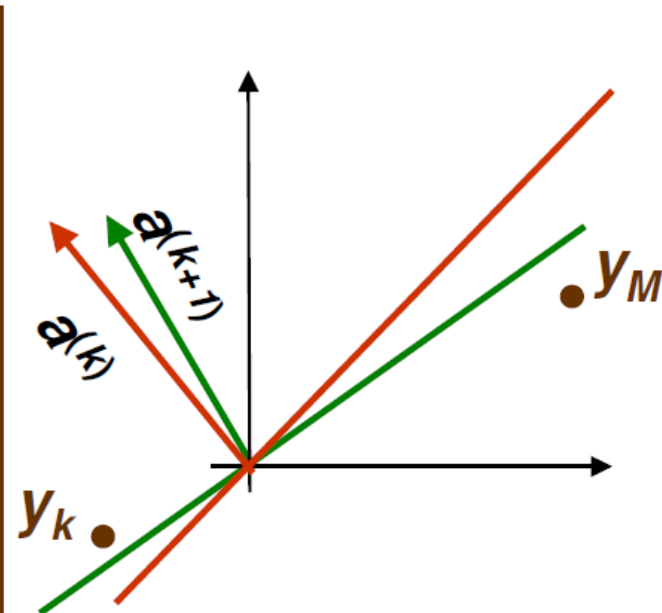


Perceptron Single Sample Rule

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$



η is too large, previously correctly classified sample \mathbf{y}_k is now misclassified



η is too small, \mathbf{y}_M is still misclassified

Perceptron Example

| | features | | | | grade |
|-------------|-------------------------|----------------|-------------------------|-------------------|----------|
| <i>name</i> | <i>good attendance?</i> | <i>tall?</i> | <i>sleeps in class?</i> | <i>chews gum?</i> | |
| Jane | <i>yes (1)</i> | <i>yes (1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>A</i> |
| Steve | <i>yes (1)</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>F</i> |
| Mary | <i>no (-1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>yes (1)</i> | <i>F</i> |
| Peter | <i>yes (1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>yes (1)</i> | <i>A</i> |

- Class 1: students who get A
- Class 2: students who get F

| | features | | | | | grade |
|-------------|--------------|-------------------------|----------------|-------------------------|-------------------|----------|
| <i>name</i> | <i>extra</i> | <i>good attendance?</i> | <i>tall?</i> | <i>sleeps in class?</i> | <i>chews gum?</i> | |
| Jane | <i>1</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>A</i> |
| Steve | <i>1</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>yes (1)</i> | <i>F</i> |
| Mary | <i>1</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>yes (1)</i> | <i>F</i> |
| Peter | <i>1</i> | <i>yes (1)</i> | <i>no (-1)</i> | <i>no (-1)</i> | <i>yes (1)</i> | <i>A</i> |

- **Augment samples** by adding an extra feature (dimension) equal to 1

| | features | | | | | grade |
|-------------|--------------|-------------------------|--------------|-------------------------|-------------------|-------|
| <i>name</i> | <i>extra</i> | <i>good attendance?</i> | <i>tall?</i> | <i>sleeps in class?</i> | <i>chews gum?</i> | |
| Jane | 1 | yes (1) | yes (1) | no (-1) | no (-1) | A |
| Steve | -1 | yes (-1) | yes (-1) | yes (-1) | yes (-1) | F |
| Mary | -1 | no (1) | no (1) | no (1) | yes (-1) | F |
| Peter | 1 | yes (1) | no (-1) | no (-1) | yes (1) | A |

- **Normalize:**

- Replace all examples from class 2 by their negative values

$$y_i \rightarrow -y_i \quad \forall y_i \in \mathbf{c}_2$$

- Seek \mathbf{a} such that:

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

| | features | | | | | grade |
|-------------|--------------|-------------------------|--------------|-------------------------|-------------------|-------|
| <i>name</i> | <i>extra</i> | <i>good attendance?</i> | <i>tall?</i> | <i>sleeps in class?</i> | <i>chews gum?</i> | |
| Jane | 1 | yes (1) | yes (1) | no (-1) | no (-1) | A |
| Steve | -1 | yes (-1) | yes (-1) | yes (-1) | yes (-1) | F |
| Mary | -1 | no (1) | no (1) | no (1) | yes (-1) | F |
| Peter | 1 | yes (1) | no (-1) | no (-1) | yes (1) | A |

- **Single Sample Rule**

- Sample is misclassified if $\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^4 \mathbf{a}_k \mathbf{y}_i^{(k)} < 0$
- Gradient descent single sample rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \sum_{\mathbf{y} \in Y_M} \mathbf{y}$
- Set η fixed learning rate to $\eta^{(k)} = 1$: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$

- Set equal initial weights
 $\mathbf{a}^{(1)} = [0.25, 0.25, 0.25, 0.25, 0.25]$
- Visit all samples sequentially, modifying the weights after each misclassified example

| <i>name</i> | <i>$\mathbf{a}^t \mathbf{y}$</i> | <i>misclassified?</i> |
|-------------|---|-----------------------|
| Jane | $0.25*1+0.25*1+0.25*1+0.25*(-1)+0.25*(-1) > 0$ | no |
| Steve | $0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1) < 0$ | yes |

- New weights

$$\begin{aligned}
 \mathbf{a}^{(2)} &= \mathbf{a}^{(1)} + \mathbf{y}_M = [0.25 \ 0.25 \ 0.25 \ 0.25 \ 0.25] + \\
 &\quad + [-1 \ -1 \ -1 \ -1 \ -1] = \\
 &= [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75]
 \end{aligned}$$

$$\mathbf{a}^{(2)} = [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75]$$

| <i>name</i> | $\mathbf{a}^t \mathbf{y}$ | <i>misclassified?</i> |
|-------------|---|-----------------------|
| Mary | $-0.75*(-1) - 0.75*1 - 0.75*1 - 0.75*1 - 0.75*(-1) < 0$ | yes |

- New weights

$$\begin{aligned} \mathbf{a}^{(3)} &= \mathbf{a}^{(2)} + \mathbf{y}_M = [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75] + \\ &\quad + [-1 \ 1 \ 1 \ 1 \ -1] = \\ &= [-1.75 \ 0.25 \ 0.25 \ 0.25 \ -1.75] \end{aligned}$$

$$\mathbf{a}^{(3)} = [-1.75 \quad 0.25 \quad 0.25 \quad 0.25 \quad -1.75]$$

| <i>name</i> | $\mathbf{a}^t \mathbf{y}$ | <i>misclassified?</i> |
|-------------|---|-----------------------|
| Peter | $-1.75 * 1 + 0.25 * 1 + 0.25 * (-1) + 0.25 * (-1) - 1.75 * 1 < 0$ | yes |

- New weights

$$\begin{aligned} \mathbf{a}^{(4)} &= \mathbf{a}^{(3)} + \mathbf{y}_M = [-1.75 \quad 0.25 \quad 0.25 \quad 0.25 \quad -1.75] + \\ &\quad + [1 \quad 1 \quad -1 \quad -1 \quad 1] = \\ &= [-0.75 \quad 1.25 \quad -0.75 \quad -0.75 \quad -0.75] \end{aligned}$$

$$\mathbf{a}^{(4)} = [-0.75 \quad 1.25 \quad -0.75 \quad -0.75 \quad -0.75]$$

| <i>name</i> | $\mathbf{a}^t \mathbf{y}$ | <i>misclassified?</i> |
|-------------|--|-----------------------|
| Jane | $-0.75 * 1 + 1.25 * 1 - 0.75 * 1 - 0.75 * (-1) - 0.75 * (-1) + 0$ | <i>no</i> |
| Steve | $-0.75 * (-1) + 1.25 * (-1) - 0.75 * (-1) - 0.75 * (-1) - 0.75 * (-1) > 0$ | <i>no</i> |
| Mary | $-0.75 * (-1) + 1.25 * 1 - 0.75 * 1 - 0.75 * 1 - 0.75 * (-1) > 0$ | <i>no</i> |
| Peter | $-0.75 * 1 + 1.25 * 1 - 0.75 * (-1) - 0.75 * (-1) - 0.75 * 1 > 0$ | <i>no</i> |

- Thus the discriminant function is:

$$g(\mathbf{y}) = -0.75 * y^{(0)} + 1.25 * y^{(1)} - 0.75 * y^{(2)} - 0.75 * y^{(3)} - 0.75 * y^{(4)}$$

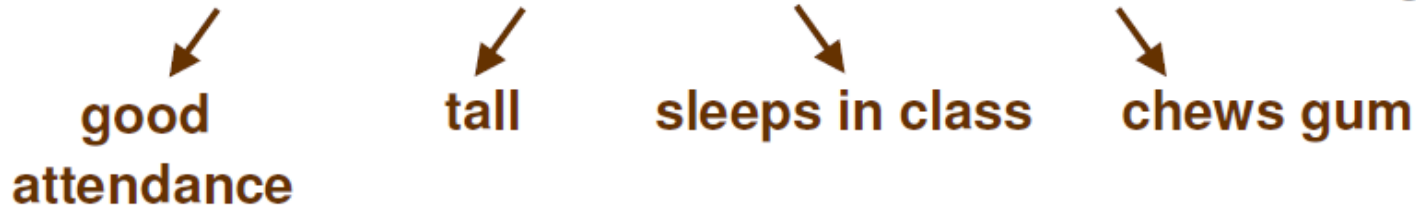
- Converting back to the original features \mathbf{x} :

$$g(\mathbf{x}) = 1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} - 0.75$$

- Converting back to the original features x :

$$1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} > 0.75 \Rightarrow \text{grade A}$$

$$1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} < 0.75 \Rightarrow \text{grade F}$$



- This is just one possible solution vector
- If we started with weights

$$a^{(1)} = [0, 0.5, 0.5, 0, 0],$$

- The solution would be $[-1, 1.5, -0.5, -1, -1]$

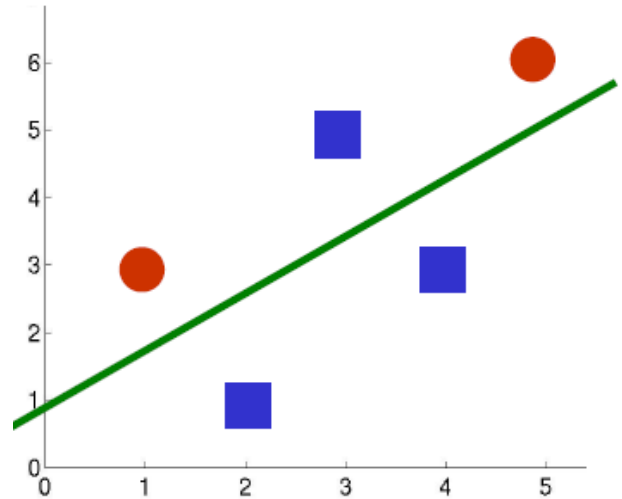
$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} > 1 \Rightarrow \text{grade A}$$

$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} < 1 \Rightarrow \text{grade F}$$

- In this solution, being tall is the least important feature

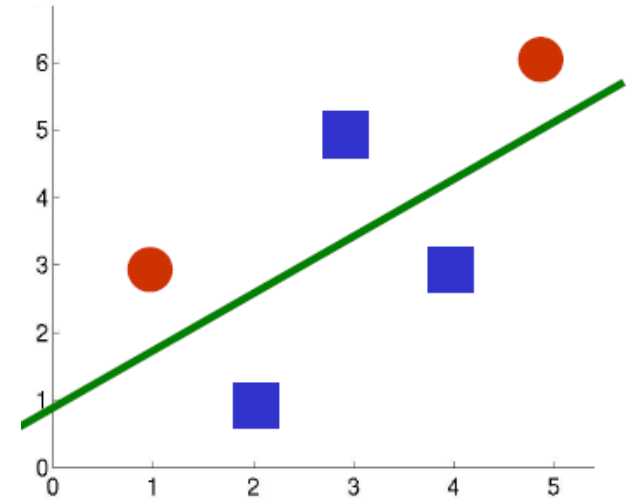
LDF: Non-separable Example

- Suppose we have 2 features and the samples are:
 - Class 1: $[2, 1]$, $[4, 3]$, $[3, 5]$
 - Class 2: $[1, 3]$ and $[5, 6]$
- These samples are not separable by a line
- Still would like to get approximate separation by a line
 - A good choice is shown in green
 - Some samples may be “noisy”, and we could accept them being misclassified



LDF: Non-separable Example

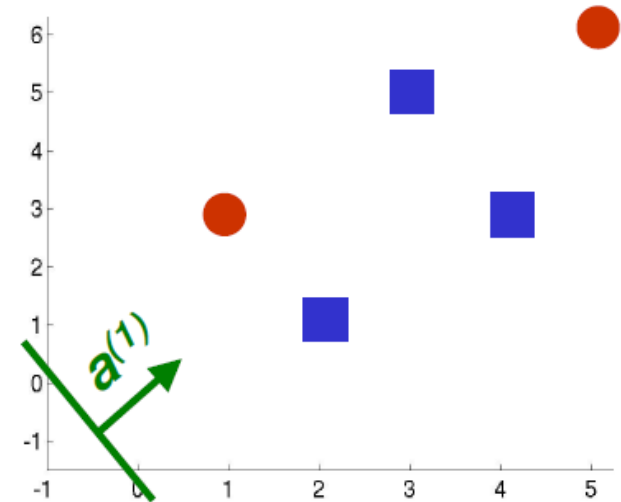
- Obtain y_1, y_2, y_3, y_4 by adding extra feature and “normalizing”



$$y_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad y_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad y_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

LDF: Non-separable Example

- Apply Perceptron single sample algorithm
- Initial equal weights $\mathbf{a}^{(1)} = [1 \ 1 \ 1]$
 - Line equation $x^{(1)} + x^{(2)} + 1 = 0$
- Fixed learning rate $\eta = 1$



$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

- $\mathbf{y}_1^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 2 \ 1]^t > 0 \quad \checkmark$
- $\mathbf{y}_2^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 4 \ 3]^t > 0 \quad \checkmark$
- $\mathbf{y}_3^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 3 \ 5]^t > 0 \quad \checkmark$

LDF: Non-separable Example

$$\mathbf{a}^{(1)} = [1 \ 1 \ 1] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

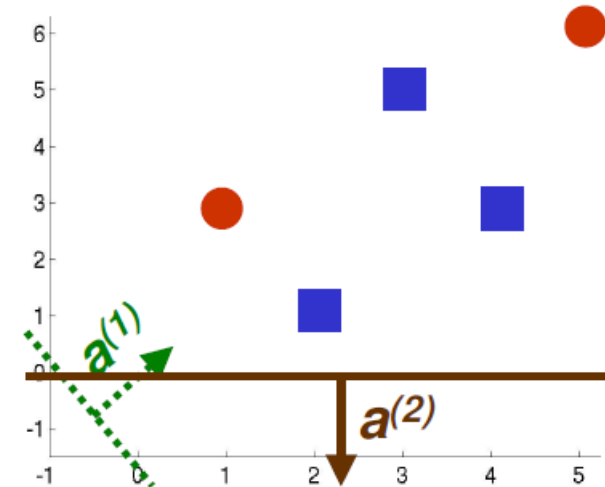
- $\mathbf{y}_4^t \mathbf{a}^{(1)} = [1 \ 1 \ 1]^* [-1 \ -1 \ -3]^t = -5 < 0$

$$\mathbf{a}^{(2)} = \mathbf{a}^{(1)} + \mathbf{y}_M = [1 \ 1 \ 1] + [-1 \ -1 \ -3] = [0 \ 0 \ -2]$$

- $\mathbf{y}_5^t \mathbf{a}^{(2)} = [0 \ 0 \ -2]^* [-1 \ -5 \ -6]^t = 12 > 0 \quad \checkmark$

- $\mathbf{y}_1^t \mathbf{a}^{(2)} = [0 \ 0 \ -2]^* [1 \ 2 \ 1]^t < 0$

$$\mathbf{a}^{(3)} = \mathbf{a}^{(2)} + \mathbf{y}_M = [0 \ 0 \ -2] + [1 \ 2 \ 1] = [1 \ 2 \ -1]$$



LDF: Non-separable Example

$$\mathbf{a}^{(3)} = [1 \ 2 \ -1] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

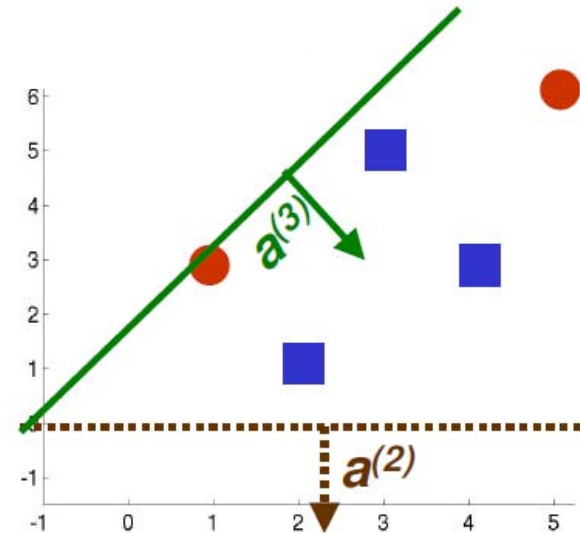
$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

- $\mathbf{y}_2^t \mathbf{a}^{(3)} = [1 \ 4 \ 3]^* [1 \ 2 \ -1]^t = 6 > 0 \quad \checkmark$

- $\mathbf{y}_3^t \mathbf{a}^{(3)} = [1 \ 3 \ 5]^* [1 \ 2 \ -1]^t > 0 \quad \checkmark$

- $\mathbf{y}_4^t \mathbf{a}^{(3)} = [-1 \ -1 \ -3]^* [1 \ 2 \ -1]^t = 0$

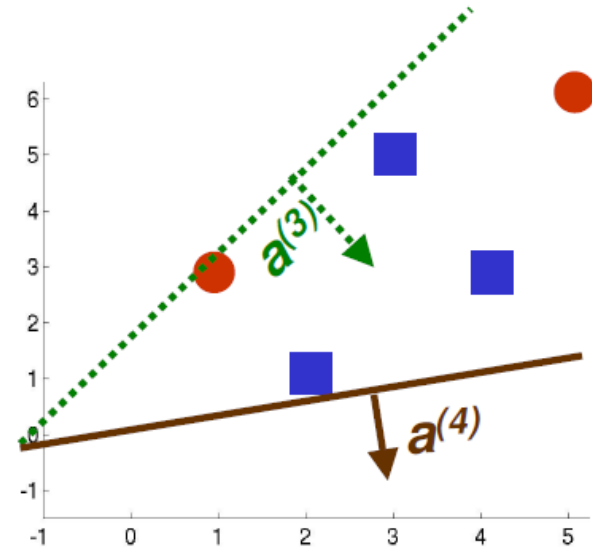
$$\mathbf{a}^{(4)} = \mathbf{a}^{(3)} + \mathbf{y}_M = [1 \ 2 \ -1] + [-1 \ -1 \ -3] = [0 \ 1 \ -4]$$



LDF: Non-separable Example

$$\mathbf{a}^{(4)} = [0 \ 1 \ -4] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$



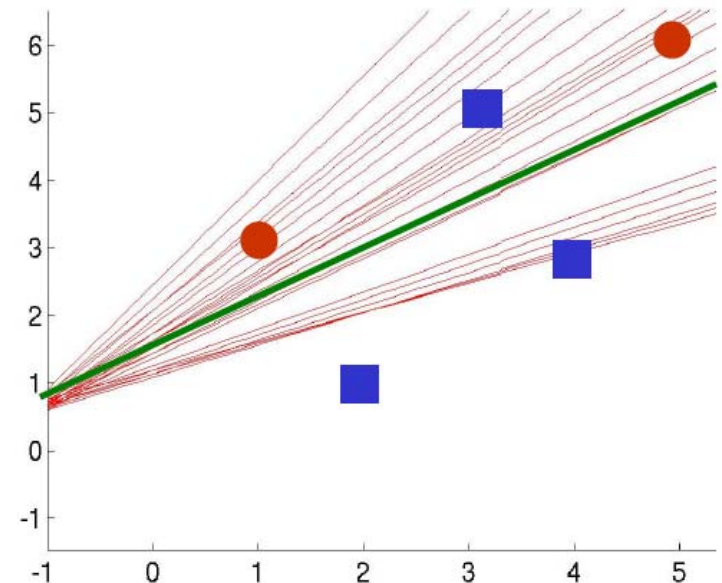
- $\mathbf{y}_5^t \mathbf{a}^{(4)} = [-1 \ -5 \ -6] \cdot [0 \ 1 \ -4] = 19 > 0$
- $\mathbf{y}_1^t \mathbf{a}^{(4)} = [1 \ 2 \ 1] \cdot [0 \ 1 \ -4] = -2 < 0$
-

LDF: Non-separable Example

- We can continue this forever
- There is no solution vector \mathbf{a} satisfying for all i

$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^5 \mathbf{a}_k \mathbf{y}_i^{(k)} > 0$$

- Need to stop but at a good point
- Solutions at iterations 900 through 915
 - Some are good and some are not
- How do we stop at a good solution?



Convergence of Perceptron Rules

- If classes are linearly separable and we use fixed learning rate, that is for $\eta^{(k)} = \text{const}$
- Then, *both the single sample and batch perceptron rules converge to a correct solution* (could be any \mathbf{a} in the solution space)
- If classes are not linearly separable:
 - The algorithm does not stop, it keeps looking for a solution which does not exist

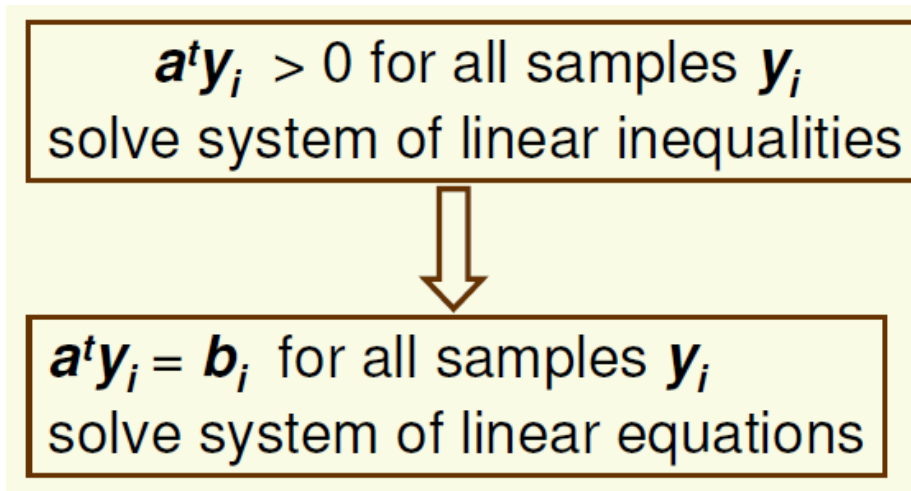
Convergence of Perceptron Rules

- If classes are not linearly separable:
 - By choosing appropriate learning rate, we can always ensure convergence: $\eta^{(k)} \rightarrow 0$ as $k \rightarrow \infty$
 - For example inverse linear learning rate: $\eta^{(k)} = \frac{\eta^{(1)}}{k}$
 - For inverse linear learning rate, convergence in the linearly separable case can also be proven
 - **No guarantee that we stopped at a good point**, but there are good reasons to choose inverse linear learning rate

Minimum Squared-Error Procedures

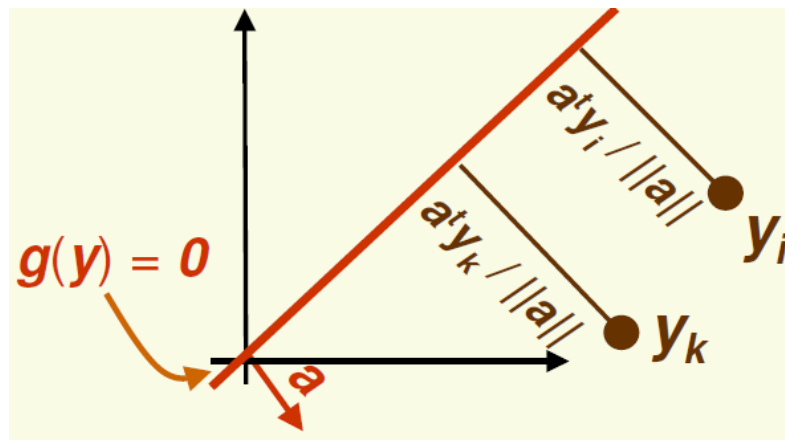
Minimum Squared-Error Procedures

- Idea: convert to easier and better understood problem



- MSE procedure
 - Choose positive constants $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n$
 - Try to find weight vector \mathbf{a} such that $\mathbf{a}^t \mathbf{y}_i = \mathbf{b}_i$ for all samples \mathbf{y}_i
 - If we can find such a vector, then \mathbf{a} is a solution because the \mathbf{b}_i 's are positive
 - Consider all the samples (not just the misclassified ones)

MSE Margins



- If $a^T y_i = b_i$, y_i must be at distance b_i from the separating hyperplane (normalized by $\|a\|$)
- Thus b_1, b_2, \dots, b_n give relative expected distances or “margins” of samples from the hyperplane
- Should make b_i small if sample i is expected to be near separating hyperplane, and large otherwise
- In the absence of any additional information, set $b_1 = b_2 = \dots = b_n = 1$

MSE Matrix Notation

- Need to solve n equations
- In matrix form $Y\mathbf{a}=\mathbf{b}$

$$\begin{cases} \mathbf{a}^t \mathbf{y}_1 = b_1 \\ \vdots \\ \mathbf{a}^t \mathbf{y}_n = b_n \end{cases}$$

$$\underbrace{\begin{bmatrix} \mathbf{y}_1^{(0)} & \mathbf{y}_1^{(1)} & \dots & \mathbf{y}_1^{(d)} \\ \mathbf{y}_2^{(0)} & \mathbf{y}_2^{(1)} & \dots & \mathbf{y}_2^{(d)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{y}_n^{(0)} & \mathbf{y}_n^{(1)} & \dots & \mathbf{y}_n^{(d)} \end{bmatrix}}_{\mathbf{Y}} \underbrace{\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix}}_{\mathbf{a}} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}}_{\mathbf{b}}$$

Exact Solution is Rare

- Need to solve a linear system $Y\mathbf{a} = \mathbf{b}$
 - Y is an $n \times (d + 1)$ matrix
- Exact solution only if Y is non-singular and square (the inverse Y^{-1} exists)
 - $\mathbf{a} = Y^{-1} \mathbf{b}$
 - (number of samples) = (number of features + 1)
 - Almost never happens in practice
 - Guaranteed to find the separating hyperplane

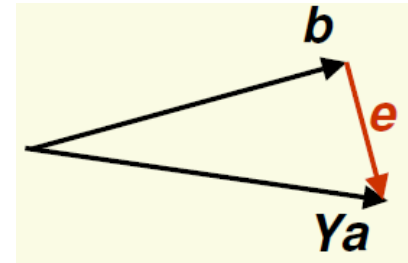
Approximate Solution

- Typically Y is overdetermined, that is it has more rows (examples) than columns (features)
 - If it has more features than examples, should reduce dimensionality
- Need $Ya = b$, but no exact solution exists for an overdetermined system of equations
 - More equations than unknowns
- Find an approximate solution
 - Note that approximate solution a does not necessarily give the separating hyperplane in the separable case
 - But the hyperplane corresponding to a may still be a good solution, especially if there is no separating hyperplane

MSE Criterion Function

- Minimum squared error approach: find \mathbf{a} which minimizes the length of the error vector \mathbf{e}

$$\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$$



- Thus minimize the **minimum squared error criterion** function:

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

- Unlike the perceptron criterion function, we can optimize the minimum squared error criterion function analytically by setting the gradient to 0

Computing the Gradient

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

$$\begin{aligned} \nabla J_s(\mathbf{a}) &= \begin{bmatrix} \frac{\partial J_s}{\partial a_0} \\ \vdots \\ \frac{\partial J_s}{\partial a_d} \end{bmatrix} = \frac{dJ_s}{d\mathbf{a}} = \sum_{i=1}^n \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i)^2 \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i) \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \mathbf{y}_i \\ &= 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{b}) \end{aligned}$$

Pseudo-Inverse Solution

$$\nabla J_s(\mathbf{a}) = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

- Setting the gradient to 0:

$$2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0 \Rightarrow \mathbf{Y}^t\mathbf{Y}\mathbf{a} = \mathbf{Y}^t\mathbf{b}$$

- The matrix $\mathbf{Y}^t\mathbf{Y}$ is square (it has $d + 1$ rows and columns) and it is often non-singular
- If $\mathbf{Y}^t\mathbf{Y}$ is non-singular, its inverse exists and we can solve for \mathbf{a} uniquely:

$$\mathbf{a} = \boxed{(\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t}\mathbf{b}$$

pseudo inverse of \mathbf{Y}
$$((\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t)\mathbf{Y} = (\mathbf{Y}^t\mathbf{Y})^{-1}(\mathbf{Y}^t\mathbf{Y}) = \mathbf{I}$$

MSE Procedures

- Only guaranteed separating hyperplane if $Y\mathbf{a} > 0$
 - That is if all elements of vector $Y\mathbf{a}$ are positive

$$Y\mathbf{a} = \begin{bmatrix} \mathbf{b}_1 + \varepsilon_1 \\ \vdots \\ \mathbf{b}_n + \varepsilon_n \end{bmatrix}$$

- where ε may be negative
- If $\varepsilon_1, \dots, \varepsilon_n$ are small relative to $\mathbf{b}_1, \dots, \mathbf{b}_n$, then each element of $Y\mathbf{a}$ is positive, and \mathbf{a} gives a separating hyperplane
 - If the approximation is not good, ε_i may be large and negative, for some i , thus $\mathbf{b}_i + \varepsilon_i$ will be negative and \mathbf{a} is not a separating hyperplane
- In linearly separable case, least squares solution \mathbf{a} does not necessarily give separating hyperplane

MSE Procedures

- We are free to choose \mathbf{b} . We may be tempted to make \mathbf{b} large as a way to ensure $\mathbf{Y}\mathbf{a} = \mathbf{b} > 0$
 - Does not work
 - Let β be a scalar, let's try $\beta\mathbf{b}$ instead of \mathbf{b}

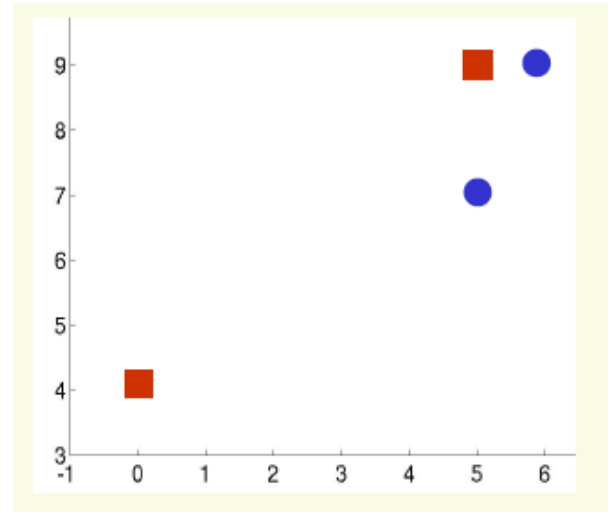
- If \mathbf{a}^* is a least squares solution to $\mathbf{Y}\mathbf{a} = \mathbf{b}$, then for any scalar β , the least squares solution to $\mathbf{Y}\mathbf{a} = \beta\mathbf{b}$ is $\beta\mathbf{a}^*$

$$\arg \min_a \|\mathbf{Y}\mathbf{a} - \beta\mathbf{b}\|^2 = \arg \min_a \beta^2 \|\mathbf{Y}(\mathbf{a} / \beta) - \mathbf{b}\|^2 = \beta\mathbf{a}^*$$

- Thus if the i^{th} element of $\mathbf{Y}\mathbf{a}$ is less than 0, that is $\mathbf{y}_i^t \mathbf{a} < 0$, then $\mathbf{y}_i^t (\beta\mathbf{a}) < 0$,
 - The relative difference between components of \mathbf{b} matters, but not the size of each individual component

LDF using MSE: Example 1

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 4)
- Add extra feature and “normalize”



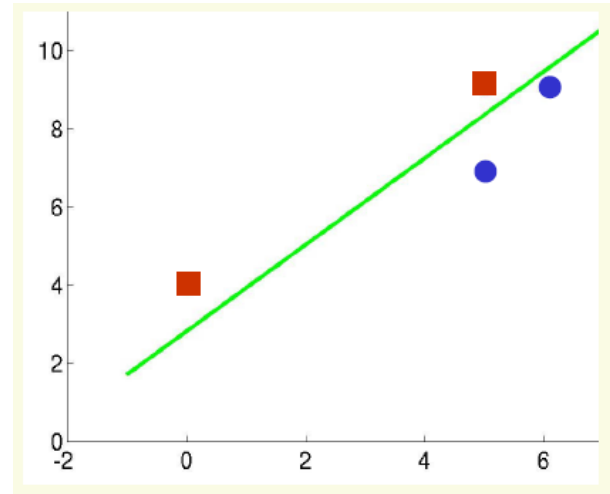
$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -4 \end{bmatrix}$$

$$Y = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -4 \end{bmatrix}$$

LDF using MSE: Example 1

- Choose $\mathbf{b}=[1 \ 1 \ 1 \ 1]^T$
- In Matlab, $\mathbf{a}=\mathbf{Y}\backslash\mathbf{b}$ solves the least squares problem

$$\mathbf{a} = \begin{bmatrix} 2.66 \\ 1.045 \\ -0.944 \end{bmatrix}$$

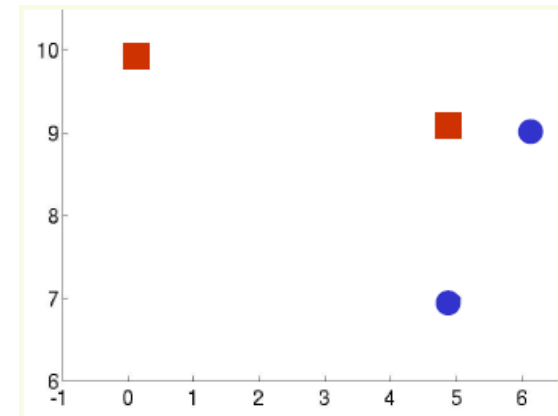


- Note \mathbf{a} is an approximation to $\mathbf{Y}\mathbf{a} = \mathbf{b}$, since no exact solution exists
- This solution gives a separating hyperplane since $\mathbf{Y}\mathbf{a} > 0$

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.44 \\ 1.28 \\ 0.61 \\ 1.11 \end{bmatrix}$$

LDF using MSE: Example 2

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 10)
- The last sample is very far compared to others from the separating hyperplane

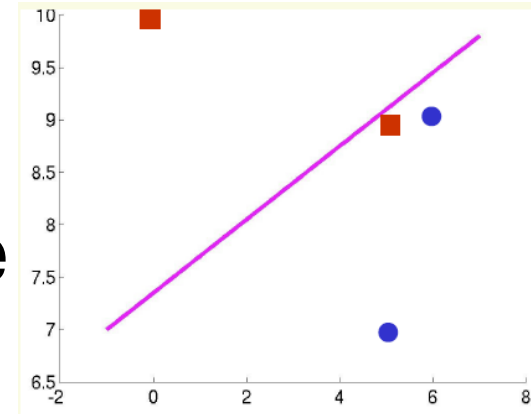


$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -10 \end{bmatrix}$$

$$Y = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -10 \end{bmatrix}$$

LDF using MSE: Example 2

- Choose $\mathbf{b}=[1 \ 1 \ 1 \ 1]^T$
- In Matlab, $\mathbf{a}=\mathbf{Y}\backslash\mathbf{b}$ solves the least squares problem



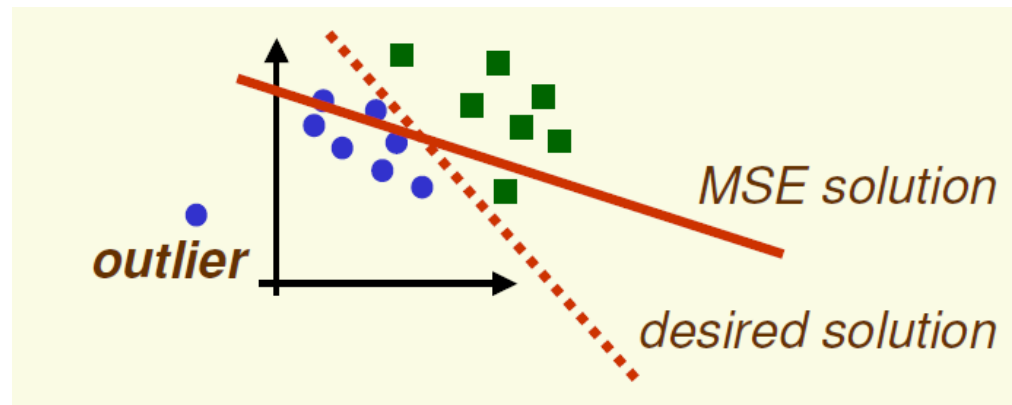
$$\mathbf{a} = \begin{bmatrix} 3.2 \\ 0.2 \\ -0.4 \end{bmatrix}$$

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.2 \\ 0.9 \\ -0.04 \\ 1.16 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

- This solution does not provide a separating hyperplane since $\mathbf{a}^T \mathbf{y}_3 < 0$

LDF using MSE: Example 2

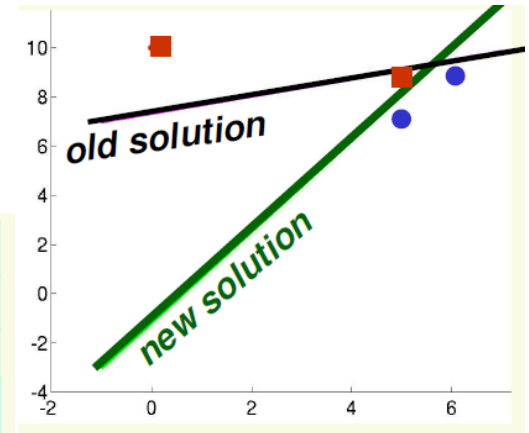
- MSE pays too much attention to isolated “noisy” examples
 - such examples are called outliers



- No problems with convergence
- Solution ranges from reasonable to good

LDF using MSE: Example 2

- We can see that the 4th point is vary far from separating hyperplane
 - In practice we don't know this
- A more appropriate \mathbf{b} could be $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$
- In Matlab, solve $\mathbf{a} = \mathbf{Y} \backslash \mathbf{b}$



$$\mathbf{a} = \begin{bmatrix} -1.1 \\ 1.7 \\ -0.9 \end{bmatrix}$$

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.9 \\ 1.0 \\ 0.8 \\ 10.0 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$$

- This solution gives the separating hyperplane since $\mathbf{Y}\mathbf{a} > 0$

Gradient Descent for MSE

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

- May wish to find MSE solution by gradient descent:
 1. Computing the inverse of $\mathbf{Y}^t\mathbf{Y}$ may be too costly
 2. $\mathbf{Y}^t\mathbf{Y}$ may be close to singular if samples are highly correlated (rows of \mathbf{Y} are almost linear combinations of each other) computing the inverse of $\mathbf{Y}^t\mathbf{Y}$ is not numerically stable
- As shown before, the gradient is:

$$\nabla J_s(\mathbf{a}) = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

Widrow-Hoff Procedure

$$\nabla J_s(\mathbf{a}) = 2Y^t(Y\mathbf{a} - \mathbf{b})$$

- Thus the update rule for gradient descent is:

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

- If $\eta^{(k)} = \eta^{(1)}/k$, then $\mathbf{a}^{(k)}$ converges to the MSE solution \mathbf{a} , that is $Y^t(Y\mathbf{a} - \mathbf{b}) = 0$
- The *Widrow-Hoff procedure* reduces storage requirements by considering single samples sequentially

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta^{(k)} \mathbf{y}_i (\mathbf{y}_i^t \mathbf{a}^{(k)} - b_i)$$

LDF Summary

- **Perceptron procedures**
 - Find a separating hyperplane in the linearly separable case,
 - Do not converge in the non-separable case
 - Can force convergence by using a decreasing learning rate, but are not guaranteed a reasonable stopping point
- **MSE procedures**
 - Converge in separable and not separable case
 - May not find separating hyperplane even if classes are linearly separable
 - Use pseudoinverse if Y^tY is not singular and not too large
 - Use gradient descent (Widrow-Hoff procedure) otherwise

Support Vector Machines

SVM Resources

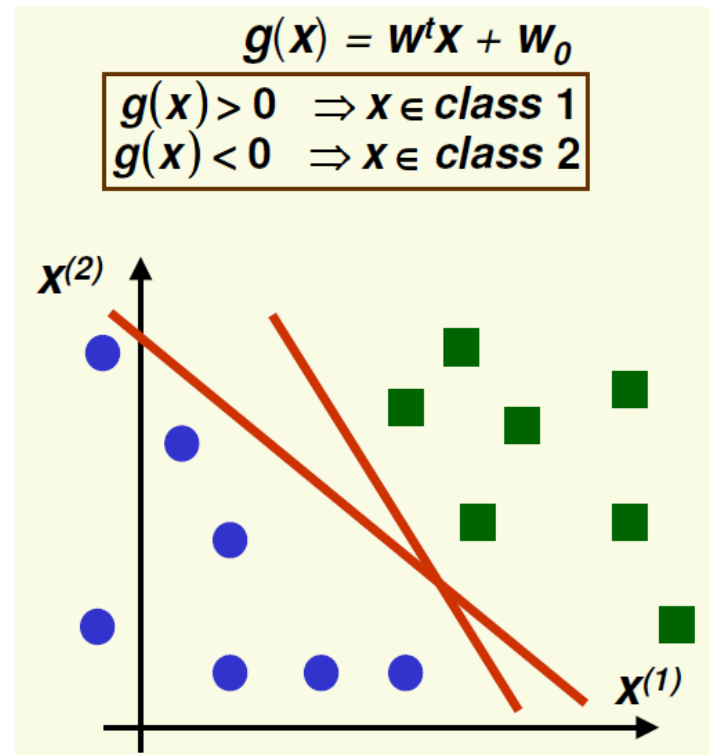
- Burges tutorial
 - <http://research.microsoft.com/en-us/um/people/cburges/papers/SVMTutorial.pdf>
- Shawe-Taylor and Christianini tutorial
 - <http://www.support-vector.net/icml-tutorial.pdf>
- Lib SVM
 - <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- LibLinear
 - <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>
- SVM Light
 - <http://svmlight.joachims.org/>
- Power Mean SVM (very fast for histogram features)
 - <https://sites.google.com/site/wujx2001/home/power-mean-svm>

SVMs

- One of the most important developments in pattern recognition in the last years
- Elegant theory
 - Has good generalization properties
- Have been applied to diverse problems very successfully

Linear Discriminant Functions

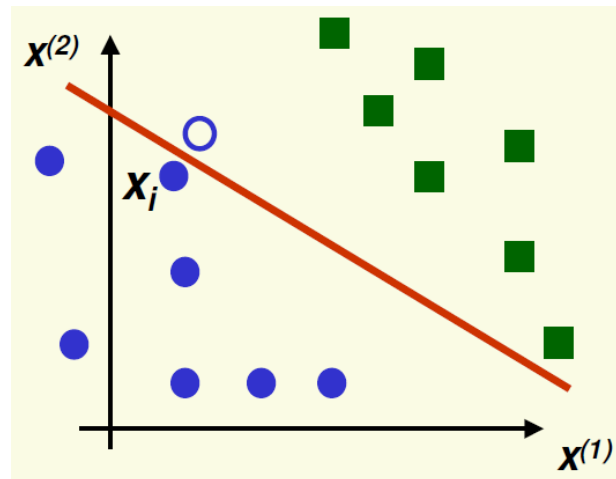
- A discriminant function is linear if it can be written as



- which separating hyperplane should we choose?

Linear Discriminant Functions

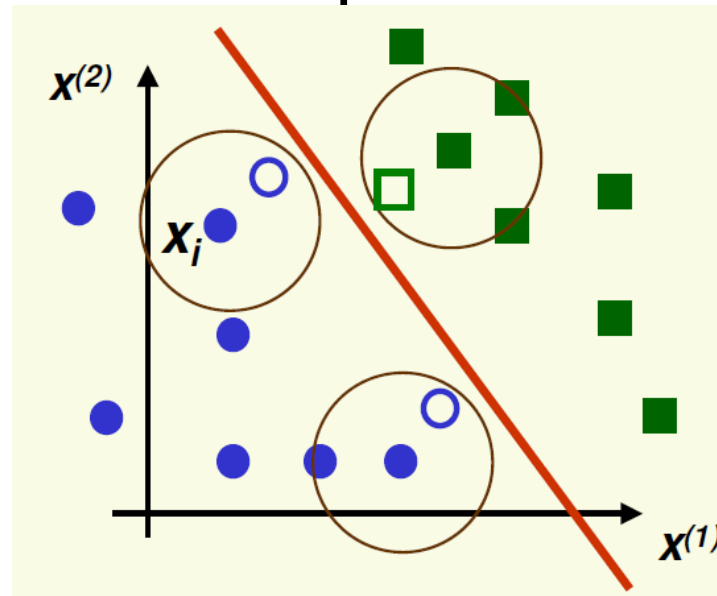
- Training data is just a subset of all possible data
 - Suppose hyperplane is close to sample x_i
 - If we see new sample close to x_i , it may be on the wrong side of the hyperplane



- Poor generalization (performance on unseen data)

Linear Discriminant Functions

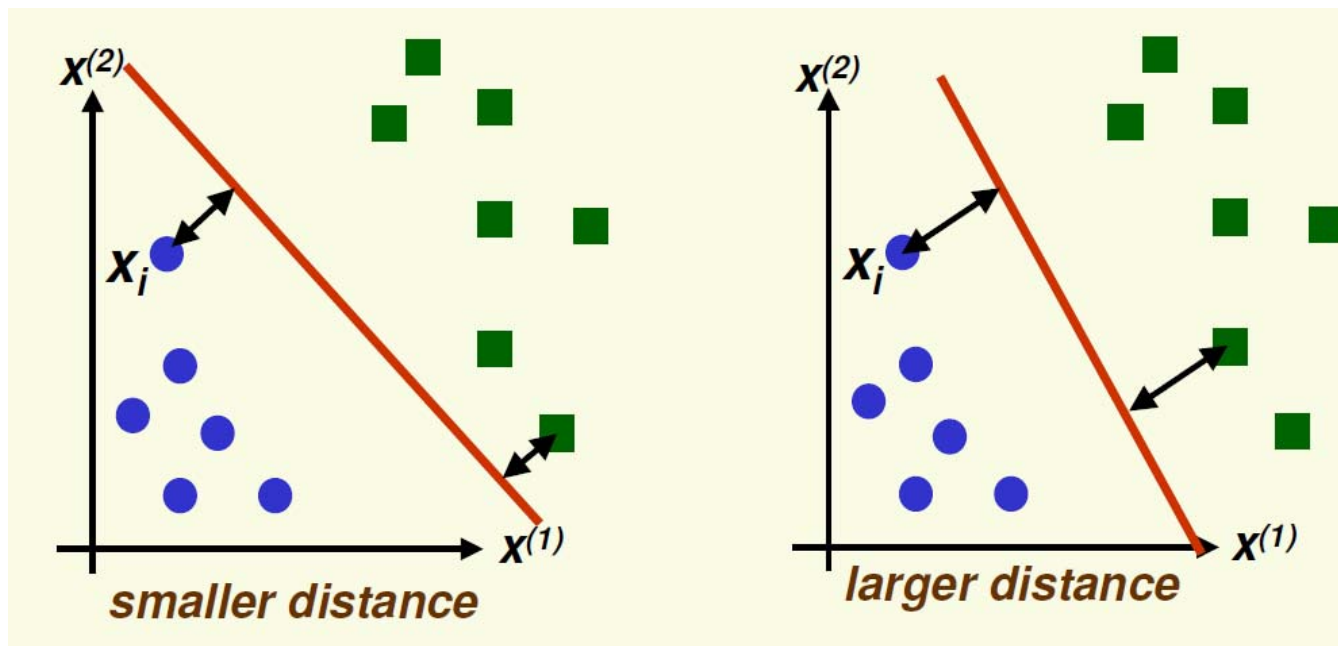
- Hyperplane as far as possible from any sample



- New samples close to the old samples will be classified correctly
- Good generalization

SVM

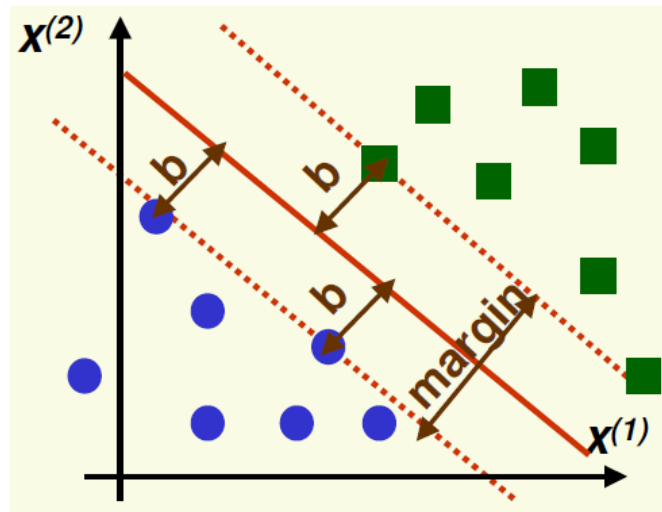
- Idea: maximize distance to the closest example



- For the optimal hyperplane
 - distance to the closest negative example = distance to the closest positive example

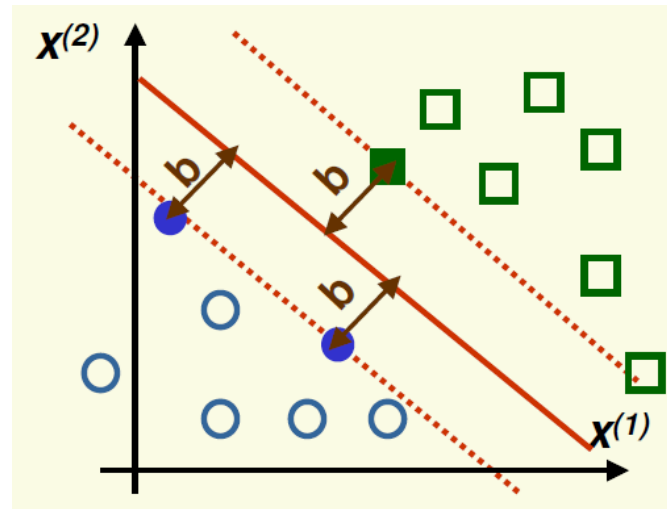
SVM: Linearly Separable Case

- SVM: maximize the margin



- The *margin* is twice the absolute value of distance b of the closest example to the separating hyperplane
- Better generalization (performance on test data)
 - in practice
 - and in theory

SVM: Linearly Separable Case



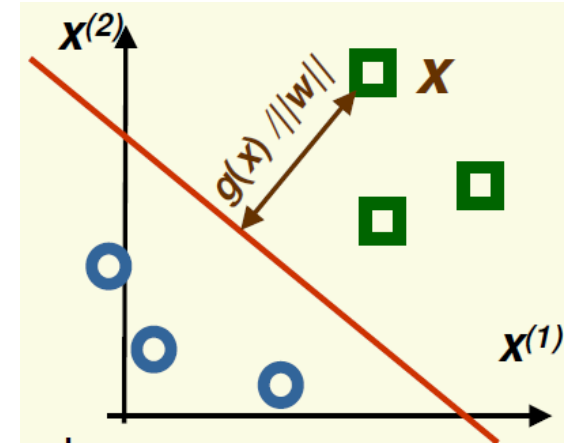
- ***Support vectors*** are the samples closest to the separating hyperplane
 - They are the most difficult patterns to classify
 - Recall perceptron update rule
- Optimal hyperplane is completely defined by support vectors
 - Of course, we do not know which samples are support vectors without finding the optimal hyperplane

SVM: Formula for the Margin

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

- Absolute distance between \mathbf{x} and the boundary $g(\mathbf{x}) = 0$

$$\frac{|\mathbf{w}^t \mathbf{x} + w_0|}{\|\mathbf{w}\|}$$



- Distance is unchanged for hyperplane

$$g_1(\mathbf{x}) = \alpha g(\mathbf{x}) \quad \frac{|\alpha \mathbf{w}^t \mathbf{x} + \alpha w_0|}{\|\alpha \mathbf{w}\|} = \frac{|\mathbf{w}^t \mathbf{x} + w_0|}{\|\mathbf{w}\|}$$

- Let \mathbf{x}_i be an example closest to the boundary (on the positive side). Set: $|\mathbf{w}^t \mathbf{x}_i + w_0| = 1$
- Now the largest margin hyperplane is unique

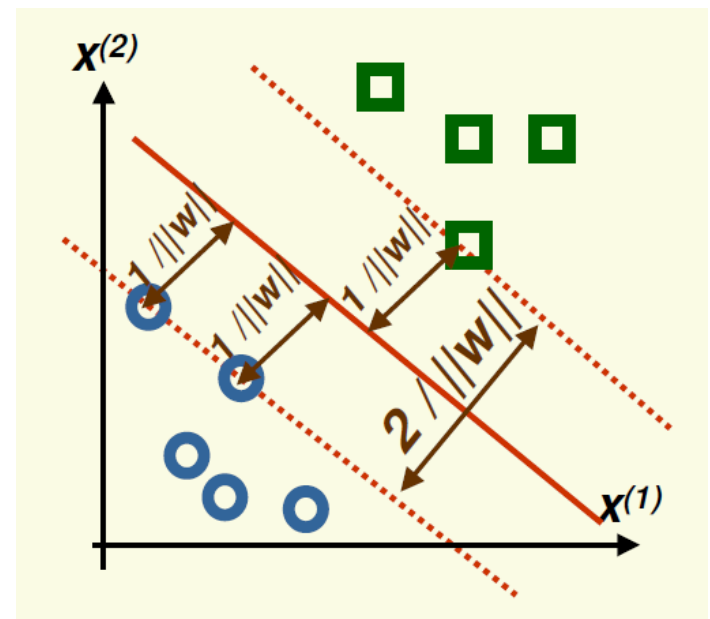
SVM: Formula for the Margin

- For uniqueness, set $|w^T x_i + w_0| = 1$ for any sample x_i closest to the boundary
- The distance from closest sample x_i to $g(x) = 0$ is

$$\frac{|w^T x_i + w_0|}{\|w\|} = \frac{1}{\|w\|}$$

- Thus the margin is

$$m = \frac{2}{\|w\|}$$



SVM: Optimal Hyperplane

- Maximize margin $m = \frac{2}{\|w\|}$
- Subject to constraints $\begin{cases} w^t x_i + w_0 \geq 1 & \text{if } x_i \text{ is positive example} \\ w^t x_i + w_0 \leq -1 & \text{if } x_i \text{ is negative example} \end{cases}$
- Let $\begin{cases} z_i = 1 & \text{if } x_i \text{ is positive example} \\ z_i = -1 & \text{if } x_i \text{ is negative example} \end{cases}$
- Can convert our problem to minimize
$$\begin{aligned} &\text{minimize } J(w) = \frac{1}{2} \|w\|^2 \\ &\text{constrained to } z_i (w^t x_i + w_0) \geq 1 \quad \forall i \end{aligned}$$
- $J(w)$ is a quadratic function, thus there is a single global minimum

SVM: Optimal Hyperplane

- Use Kuhn-Tucker theorem to convert our problem to:

$$\begin{aligned} &\text{maximize} && L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{x}_i^t \mathbf{x}_j \\ &\text{constrained to} && \alpha_i \geq 0 \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{aligned}$$

- $\mathbf{a} = \{\mathbf{a}_1, \dots, \mathbf{a}_n\}$ are new variables, one for each sample
- Optimized by quadratic programming

SVM: Optimal Hyperplane

- After finding the optimal $\mathbf{a} = \{\mathbf{a}_1, \dots, \mathbf{a}_n\}$
- Final discriminant function:

$$g(\mathbf{x}) = \left(\sum_{\mathbf{x}_i \in S} \alpha_i \mathbf{z}_i \mathbf{x}_i \right)^t \mathbf{x} + w_0$$

- where S is the set of support vectors

$$S = \{\mathbf{x}_i \mid \alpha_i \neq 0\}$$

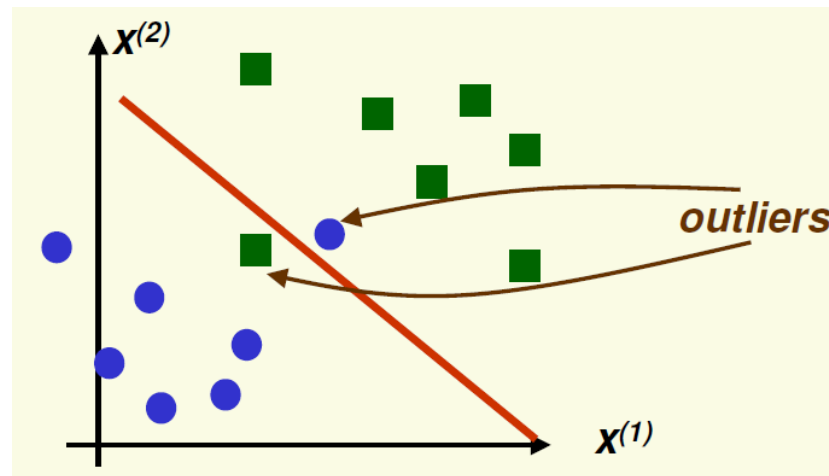
SVM: Optimal Hyperplane

$$\begin{aligned} \text{maximize} \quad & L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{x}_i^t \mathbf{x}_j \\ \text{constrained to} \quad & \alpha_i \geq 0 \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{aligned}$$

- $L_D(\mathbf{a})$ depends on the number of samples, not on dimension
 - samples appear only through the dot products $\mathbf{x}_j^t \mathbf{x}_i$
- This will become important when looking for a nonlinear discriminant function, as we will see soon

SVM: Non-Separable Case

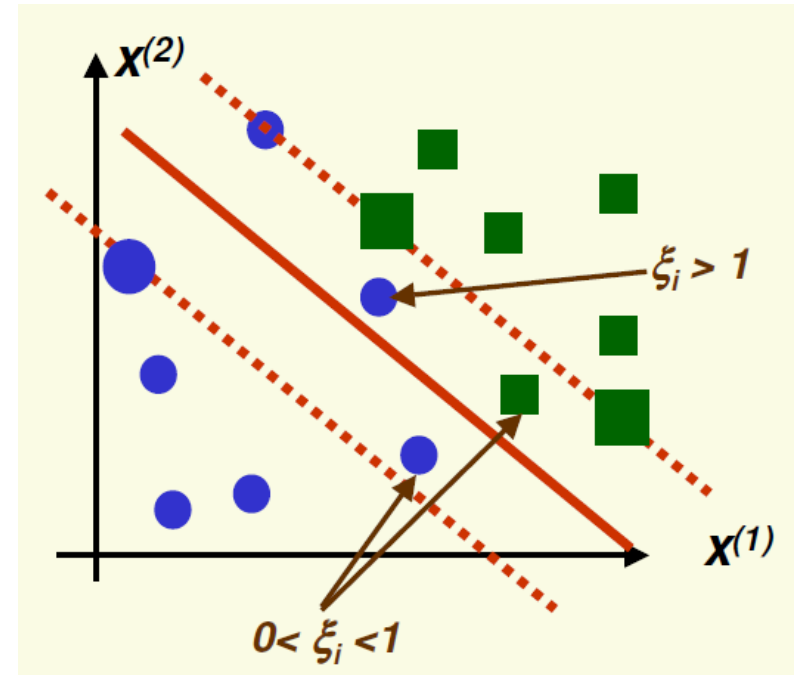
- Data are most likely to be not linearly separable, but linear classifier may still be appropriate



- Can apply SVM in non linearly separable case
- Data should be “almost” linearly separable for good performance

SVM: Non-Separable Case

- Use **slack variables** ξ_1, \dots, ξ_n (one for each sample)
- Change constraints from $z_i(w^t x_i + w_0) \geq 1 \quad \forall i$ to $z_i(w^t x_i + w_0) \geq 1 - \xi_i \quad \forall i$
- ξ_i is a measure of deviation from the ideal for x_i
 - $\xi_i > 1$: x_i is on the wrong side of the separating hyperplane
 - $0 < \xi_i < 1$: x_i is on the right side of separating hyperplane but within the region of maximum margin
 - $\xi_i < 0$: is the ideal case for x_i



SVM: Non-Separable Case

- We would like to minimize

$$J(w, \xi_1, \dots, \xi_n) = \frac{1}{2} \|w\|^2 + \beta \sum_{i=1}^n I(\xi_i > 0)$$

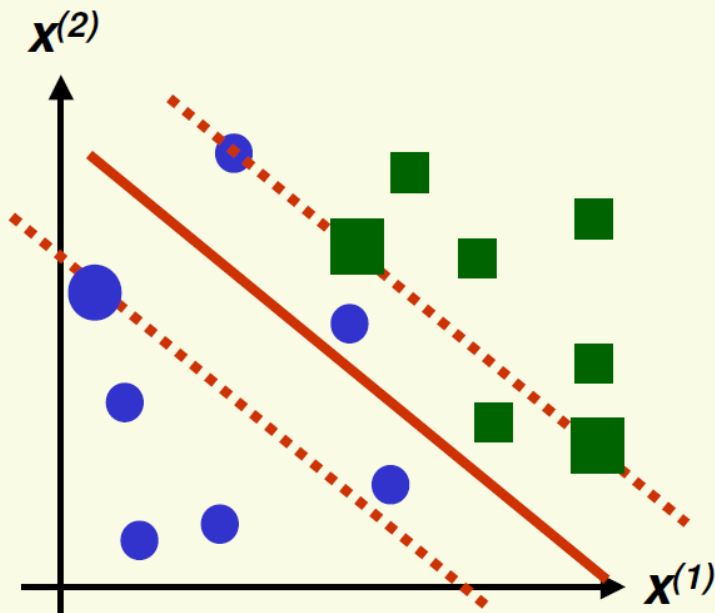
*# of samples
not in ideal location*

- where $I(\xi_i > 0) = \begin{cases} 1 & \text{if } \xi_i > 0 \\ 0 & \text{if } \xi_i \leq 0 \end{cases}$
- Constrained to $z_i(w^t x_i + w_0) \geq 1 - \xi_i$ and $\xi_i \geq 0 \quad \forall i$
- β is a constant that measures the relative weight of first and second term
 - If β is small, we allow a lot of samples to be in not ideal positions
 - If β is large, few samples can be in non-ideal positions

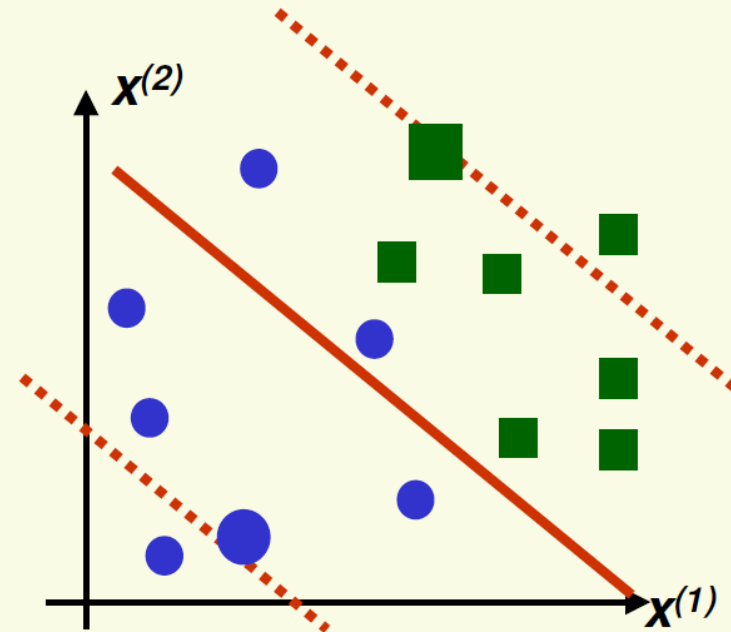
SVM: Non-Separable Case

$$J(w, \xi_1, \dots, \xi_n) = \frac{1}{2} \|w\|^2 + \beta \sum_{i=1}^n I(\xi_i > 0)$$

of examples not in ideal location



large β , few samples not in ideal position



small β , a lot of samples not in ideal position

SVM: Non-Separable Case

- Unfortunately this minimization problem is NP-hard due to the discontinuity of $I(\xi_i)$
- Instead, we minimize

$$J(\mathbf{w}, \xi_1, \dots, \xi_n) = \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^n \xi_i$$

a measure of # of misclassified examples

- Subject to
$$\begin{cases} \mathbf{z}_i (\mathbf{w}^t \mathbf{x}_i + w_0) \geq 1 - \xi_i & \forall i \\ \xi_i \geq 0 & \forall i \end{cases}$$

SVM: Non-Separable Case

- Use Kuhn-Tucker theorem to convert to:

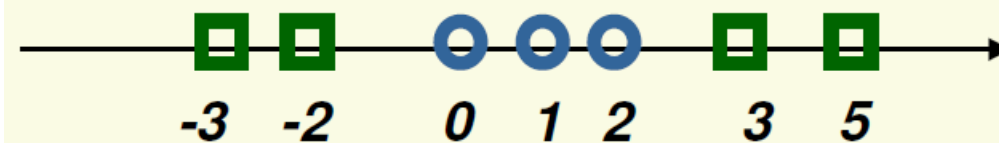
$$\begin{aligned} \text{maximize} \quad & L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^t x_j \\ \text{constrained to} \quad & 0 \leq \alpha_i \leq \beta \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{aligned}$$

- w is computed using: $w = \sum_{i=1}^n \alpha_i z_i x_i$

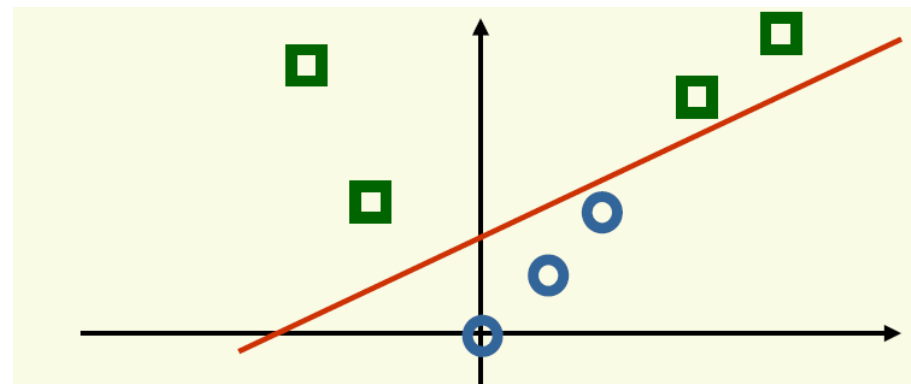
- Remember that $g(x) = \left(\sum_{x_i \in S} \alpha_i z_i x_i \right)^t x + w_0$

Nonlinear Mapping

- Cover's theorem: *"a pattern-classification problem cast in a high dimensional space non-linearly is more likely to be linearly separable than in a low-dimensional space"*
- One dimensional space, not linearly separable

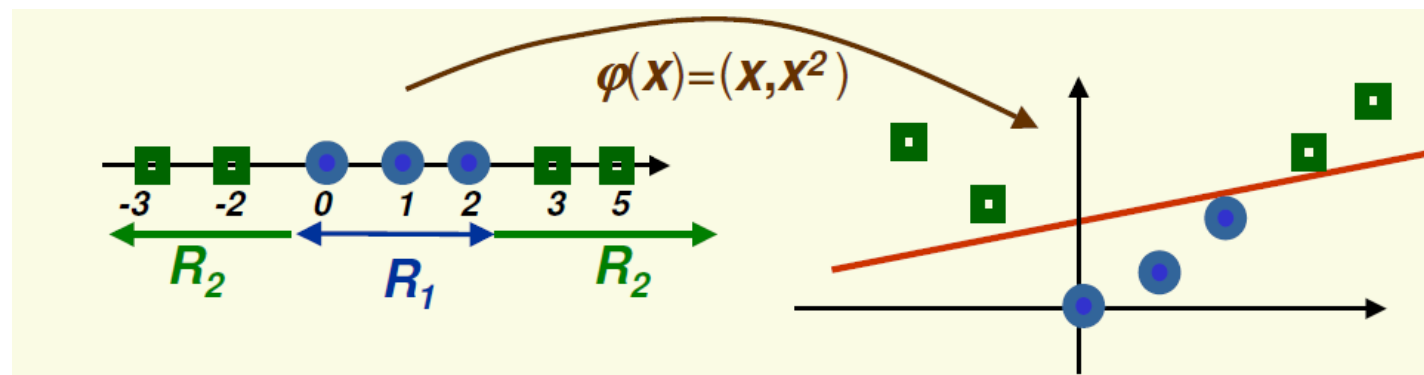


- Lift to two dimensional space with $\varphi(x)=(x,x^2)$



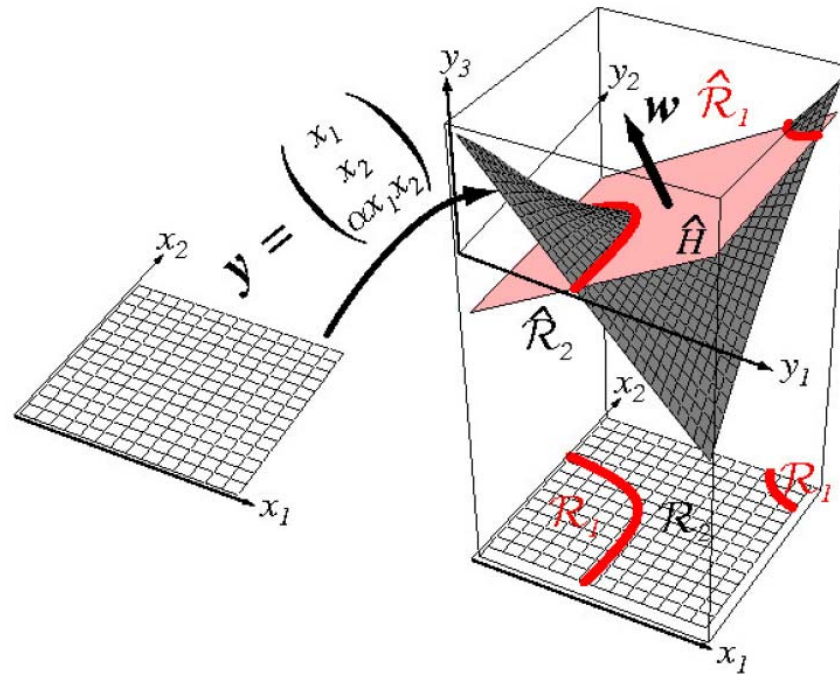
Nonlinear Mapping

- To solve a non linear classification problem with a linear classifier
 - Project data x to high dimension using function $\varphi(x)$
 - Find a linear discriminant function for transformed data $\varphi(x)$
 - Final nonlinear discriminant function is $g(x) = w^t\varphi(x) + w_0$



- In 2D, the discriminant function is linear $g\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix}\right) = [w_1 \ w_2] \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix} + w_0$
- In 1D, the discriminant function is not linear $g(x) = w_1x + w_2x^2 + w_0$

Nonlinear Mapping



- However, there always exists a mapping of N samples to an N -dimensional space in which the samples are separable by hyperplanes

Nonlinear SVM

- Can use any linear classifier after lifting data to a higher dimensional space. However we will have to deal with the curse of dimensionality
 - Poor generalization to test data
 - Computationally expensive
- SVM avoids the curse of dimensionality problems
 - Enforcing largest margin permits good generalization
 - It can be shown that generalization in SVM is a function of the margin, independent of the dimensionality
 - Computation in the higher dimensional case is performed only implicitly through the use of *kernel functions*

Kernels

- SVM optimization:

maximize
$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{x}_i^t \mathbf{x}_j$$

- Note this optimization depends on samples \mathbf{x}_i only through the dot product $\mathbf{x}_i^t \mathbf{x}_j$
- If we lift \mathbf{x}_i to high dimension using $\varphi(\mathbf{x})$, we need to compute high dimensional product $\varphi(\mathbf{x}_i)^t \varphi(\mathbf{x}_j)$

maximize
$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j \varphi(\mathbf{x}_i)^t \varphi(\mathbf{x}_j)$$

$K(\mathbf{x}_i, \mathbf{x}_j)$

- Idea: find kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$ s.t. $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^t \varphi(\mathbf{x}_j)$

Kernel Trick

- Then we only need to compute $K(\mathbf{x}_i, \mathbf{x}_j)$ instead of $\boldsymbol{\varphi}(\mathbf{x}_i)^t \boldsymbol{\varphi}(\mathbf{x}_j)$
- “**kernel trick**”: do not need to perform operations in high dimensional space explicitly

Kernel Example

- Suppose we have two features and $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^t \mathbf{y})^2$
- Which mapping $\boldsymbol{\varphi}(\mathbf{x})$ does this correspond to?

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^t \mathbf{y})^2 = \left(\begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} \end{bmatrix} \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{bmatrix} \right)^2 = (\mathbf{x}^{(1)} \mathbf{y}^{(1)} + \mathbf{x}^{(2)} \mathbf{y}^{(2)})^2 \\ &= (\mathbf{x}^{(1)} \mathbf{y}^{(1)})^2 + 2(\mathbf{x}^{(1)} \mathbf{y}^{(1)})(\mathbf{x}^{(2)} \mathbf{y}^{(2)}) + (\mathbf{x}^{(2)} \mathbf{y}^{(2)})^2 \\ &= \left[(\mathbf{x}^{(1)})^2 \quad \sqrt{2} \mathbf{x}^{(1)} \mathbf{x}^{(2)} \quad (\mathbf{x}^{(2)})^2 \right] \left[(\mathbf{y}^{(1)})^2 \quad \sqrt{2} \mathbf{y}^{(1)} \mathbf{y}^{(2)} \quad (\mathbf{y}^{(2)})^2 \right]^t \end{aligned}$$

$$\boldsymbol{\varphi}(\mathbf{x}) = \left[(\mathbf{x}^{(1)})^2 \quad \sqrt{2} \mathbf{x}^{(1)} \mathbf{x}^{(2)} \quad (\mathbf{x}^{(2)})^2 \right]$$

Choice of Kernel

- How to choose kernel function $K(x_i, x_j)$?
 - $K(x_i, x_j)$ should correspond to $\varphi(x_i)^t \varphi(x_j)$ in a higher dimensional space
 - Mercer's condition tells us which kernel function can be expressed as dot product of two vectors
 - If K and K' are kernels $aK+bK'$ is a kernel
- Intuitively: Kernel should measure the similarity between x_i and x_j
 - As inner product measures similarity of unit vectors
 - May be problem-specific

Choice of Kernel

- Some common choices:
 - Polynomial kernel

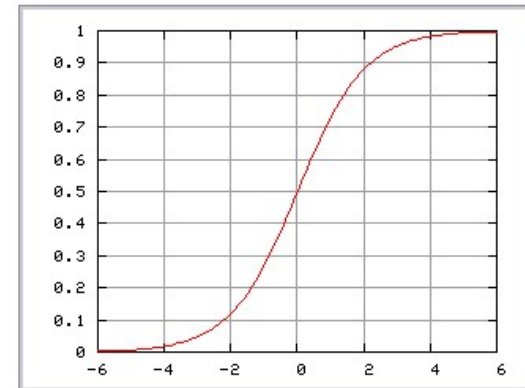
$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^t \mathbf{x}_j + 1)^p$$

- Gaussian radial Basis kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$$

- Hyperbolic tangent (sigmoid) kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(k \mathbf{x}_i^t \mathbf{x}_j + c)$$



- The mappings $\phi(\mathbf{x}_i)$ never have to be computed!!

Intersection Kernel

- Feature vectors are histograms

$$K(x_i, x_j) = \sum_{k=1}^n \min(x_{ik}, x_{jk})$$

- When $K(x_i, x_j)$ is small, x_i and x_j are dissimilar
- When $K(x_i, x_j)$ is large, x_i and x_j are similar
- The mapping $\phi(x)$ does not exist

More Additive Kernels

- χ^2 kernel $K_{\chi^2} = \sum_{k=1}^n \frac{2x_k y_k}{x_k + y_k}$
- Hellinger's kernel $K_H = \sum_{k=1}^n \sqrt{x_k y_k}$
- Designed for feature vectors that are histograms
 - Can be used for other feature vectors
- Offer very large speed-ups

The Kernel Matrix

- a.k.a the Gram matrix

$K =$

| | | | | |
|----------|----------|----------|-----|----------|
| $K(1,1)$ | $K(1,2)$ | $K(1,3)$ | ... | $K(1,m)$ |
| $K(2,1)$ | $K(2,2)$ | $K(2,3)$ | ... | $K(2,m)$ |
| | | | | |
| ... | ... | ... | ... | ... |
| $K(m,1)$ | $K(m,2)$ | $K(m,3)$ | ... | $K(m,m)$ |

- Contains all necessary information for the learning algorithm
- Fuses information about the data and the kernel (similarity measure)

Bad Kernels

- The kernel matrix is mostly diagonal
 - All points are orthogonal to each other
- Bad similarity measure
- Too many irrelevant features in high dimensional space

- We need problem-specific knowledge to choose appropriate kernel

Nonlinear SVM Step-by-Step

- Start with data x_1, \dots, x_n which live in feature space of dimension d
- Choose kernel $K(x_i, x_j)$ or function $\phi(x_i)$ which lifts sample x_i to a higher dimensional space
- Find the maximum margin linear discriminant function in the higher dimensional space by using quadratic programming package to solve:

$$\begin{aligned} &\text{maximize } L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j K(x_i, x_j) \\ &\text{constrained to } 0 \leq \alpha_i \leq \beta \quad \forall i \quad \text{and} \quad \sum_{i=1}^n \alpha_i z_i = 0 \end{aligned}$$

Nonlinear SVM Step-by-Step

- Weight vector w in the high dimensional space:

$$w = \sum_{x_i \in S} \alpha_i z_i \varphi(x_i)$$

– where S is the set of support vectors

- Linear discriminant function of maximum margin in the high dimensional space:

$$g(\varphi(x)) = w^t \varphi(x) = \left(\sum_{x_i \in S} \alpha_i z_i \varphi(x_i) \right)^t \varphi(x)$$

- Non linear discriminant function in the original space:

$$g(x) = \left(\sum_{x_i \in S} \alpha_i z_i \varphi(x_i) \right)^t \varphi(x) = \sum_{x_i \in S} \alpha_i z_i \varphi^t(x_i) \varphi(x) = \sum_{x_i \in S} \alpha_i z_i K(x_i, x)$$

- decide class 1 if $g(x) > 0$, otherwise decide class 2

Nonlinear SVM

- Nonlinear discriminant function

$$g(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathcal{S}} \alpha_i z_i K(\mathbf{x}_i, \mathbf{x})$$

$$g(\mathbf{x}) = \sum \left[\begin{array}{|l} \text{weight of support} \\ \text{vector } \mathbf{x}_i \end{array} \right] \left[\begin{array}{|l} \mp 1 \end{array} \right] \left[\begin{array}{|l} \text{"inverse distance"} \\ \text{from } \mathbf{x} \text{ to} \\ \text{support vector } \mathbf{x}_i \end{array} \right]$$

most important training samples, i.e. support vectors

$$K(\mathbf{x}_i, \mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}\|^2\right)$$

SVM Example: XOR Problem

- Class 1: $\mathbf{x}_1 = [1, -1]$, $\mathbf{x}_2 = [-1, 1]$
- Class 2: $\mathbf{x}_3 = [1, 1]$, $\mathbf{x}_4 = [-1, -1]$
- Use polynomial kernel of degree 2:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^t \mathbf{x}_j + 1)^2$$

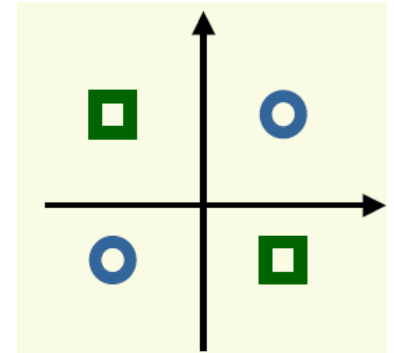
- This kernel corresponds to the mapping

$$\varphi(\mathbf{x}) = \left[1 \quad \sqrt{2}x^{(1)} \quad \sqrt{2}x^{(2)} \quad \sqrt{2}x^{(1)}x^{(2)} \quad (x^{(1)})^2 \quad (x^{(2)})^2 \right]^t$$

- Need to maximize

$$L_D(\alpha) = \sum_{i=1}^4 \alpha_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j z_i z_j (\mathbf{x}_i^t \mathbf{x}_j + 1)^2$$

constrained to $0 \leq \alpha_i \quad \forall i$ and $\alpha_1 + \alpha_2 - \alpha_3 - \alpha_4 = 0$



SVM Example: XOR Problem

- After some manipulation ...
- The solution is $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0.25$
 - satisfies the constraints

$$\forall i, 0 \leq \alpha_i \text{ and } \alpha_1 + \alpha_2 - \alpha_3 - \alpha_4 = 0$$

- All samples are support vectors

SVM Example: XOR Problem

$$\varphi(\mathbf{x}) = \left[1 \quad \sqrt{2}x^{(1)} \quad \sqrt{2}x^{(2)} \quad \sqrt{2}x^{(1)}x^{(2)} \quad (x^{(1)})^2 \quad (x^{(2)})^2 \right]^t$$

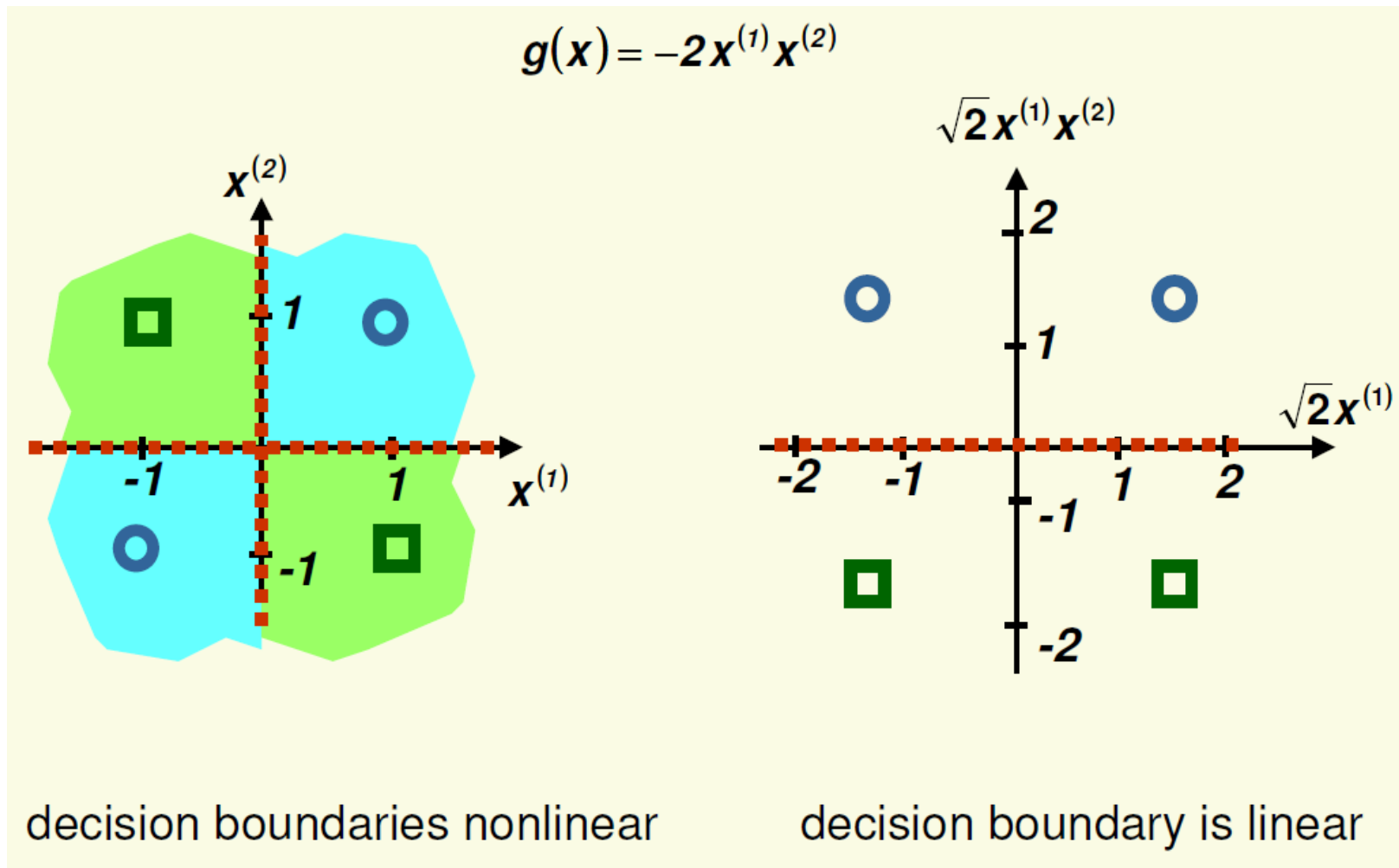
- The weight vector \mathbf{w} is:

$$\begin{aligned} \mathbf{w} &= \sum_{i=1}^4 \alpha_i \mathbf{z}_i \varphi(\mathbf{x}_i) = 0.25(\varphi(\mathbf{x}_1) + \varphi(\mathbf{x}_2) - \varphi(\mathbf{x}_3) - \varphi(\mathbf{x}_4)) \\ &= [0 \quad 0 \quad 0 \quad -\sqrt{2} \quad 0 \quad 0] \end{aligned}$$

- Thus the nonlinear discriminant function is:

$$g(\mathbf{x}) = \mathbf{w}\varphi(\mathbf{x}) = \sum_{i=1}^6 w_i \varphi_i(\mathbf{x}) = -\sqrt{2}(\sqrt{2}x^{(1)}x^{(2)}) = -2x^{(1)}x^{(2)}$$

SVM Example: XOR Problem



SVM Summary

- Advantages:
 - Based on very strong theory
 - Excellent generalization properties
 - Objective function has no local minima
 - Can be used to find non linear discriminant functions
 - Complexity of the classifier is characterized by the number of support vectors rather than the dimensionality of the transformed space
- Disadvantages:
 - Directly applicable to two-class problems
 - Quadratic programming is computationally expensive
 - Need to choose kernel

Multi-Class SVMs

- One against all
- Pairwise

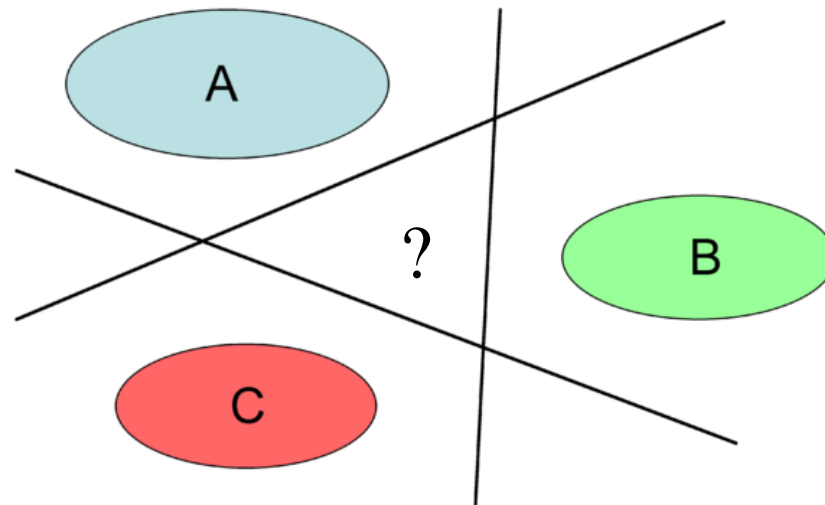
- These ideas apply to all binary classifiers when faced with multi-class problems

One-Against-All

- SVMs can only handle two-class outputs
- What can be done?
- Answer: learn N SVM's
 - SVM 1 learns “Output==1” vs “Output != 1”
 - SVM 2 learns “Output==2” vs “Output != 2”
 - ...
 - SVM N learns “Output==N” vs “Output != N”

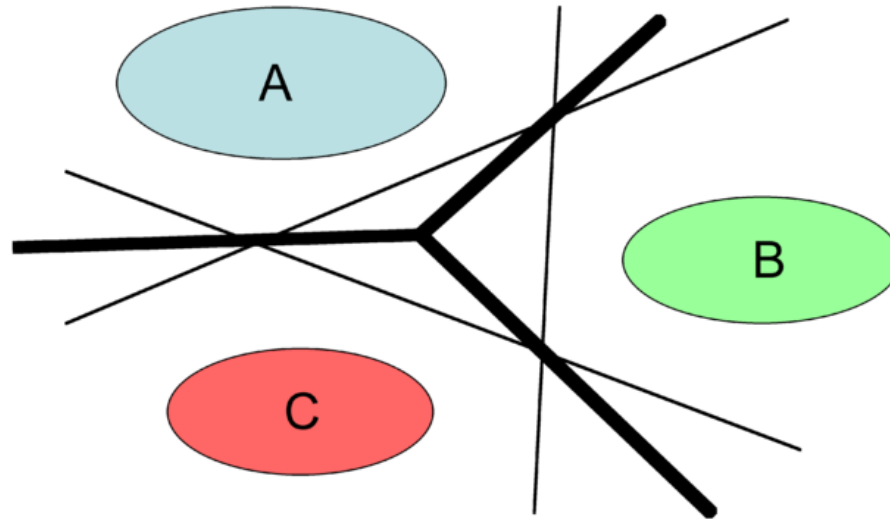
One-Against-All

- Original idea (Vapnik, 1995): classify x as ω_i if and only if the corresponding SVM accepts x and all other SVMs reject it



One-Against-All

- Modified idea (Vapnik, 1998): classify x according to the SVM that produces the highest value (use more than sign of decision function)



Pairwise SVMs

- Learn $N(N-1)/2$ SVM's
 - SVM 1 learns “Output==1” vs “Output == 2”
 - SVM 2 learns “Output==1” vs “Output == 3”
 - ...
 - SVM M learns “Output==N-1” vs “Output == N”

Pairwise SVMs

- To classify a new input, apply each SVM and choose the label that “wins” most often

