

# Machine Learning: Models and Applications

## Lecture 7

Lecturer: Xinchao Wang

# Outline

- Linear Regression
  - Regression technique
- Logistic Regression (called regression, but in fact, classification)
  - Classification technique
- Gradient Descent/Ascent
  - Descent  $\rightarrow$  min.
  - Ascent  $\rightarrow$  max.
  - Optimization technique
- Perceptron
  - Classification technique

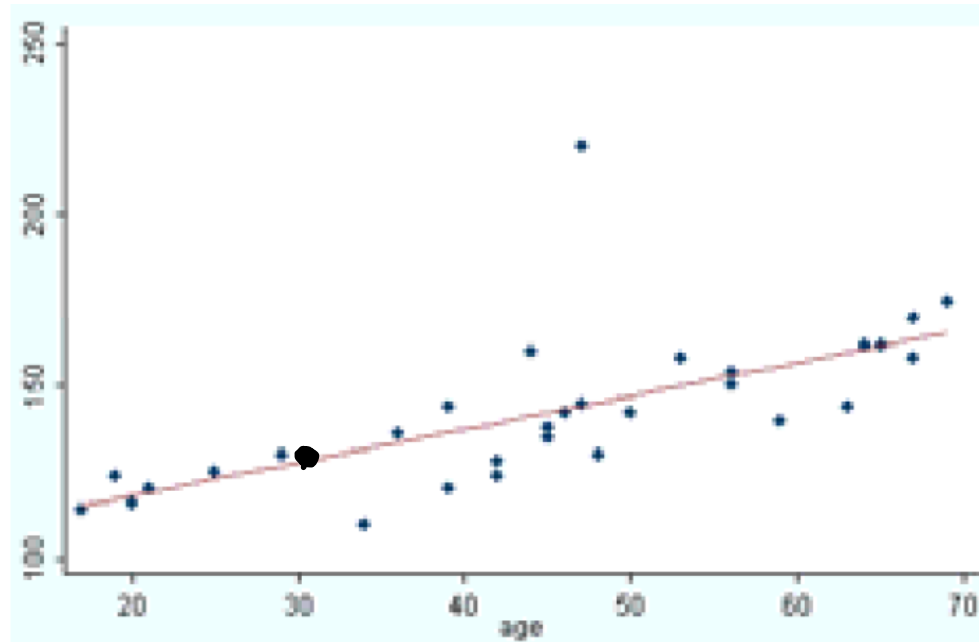
# Linear Regression

# Simple Linear Regression

- How does a single variable of interest relate to another (single) variable?
  - $Y$  = outcome variable (response, dependent...)
  - $X$  = explanatory variable (predictor, feature, independent...)
- Data:  $n$  pairs of continuous observations  $(X_1, Y_1) \dots (X_n, Y_n)$

# Example

- How does systolic blood pressure (SBP) relate to age?

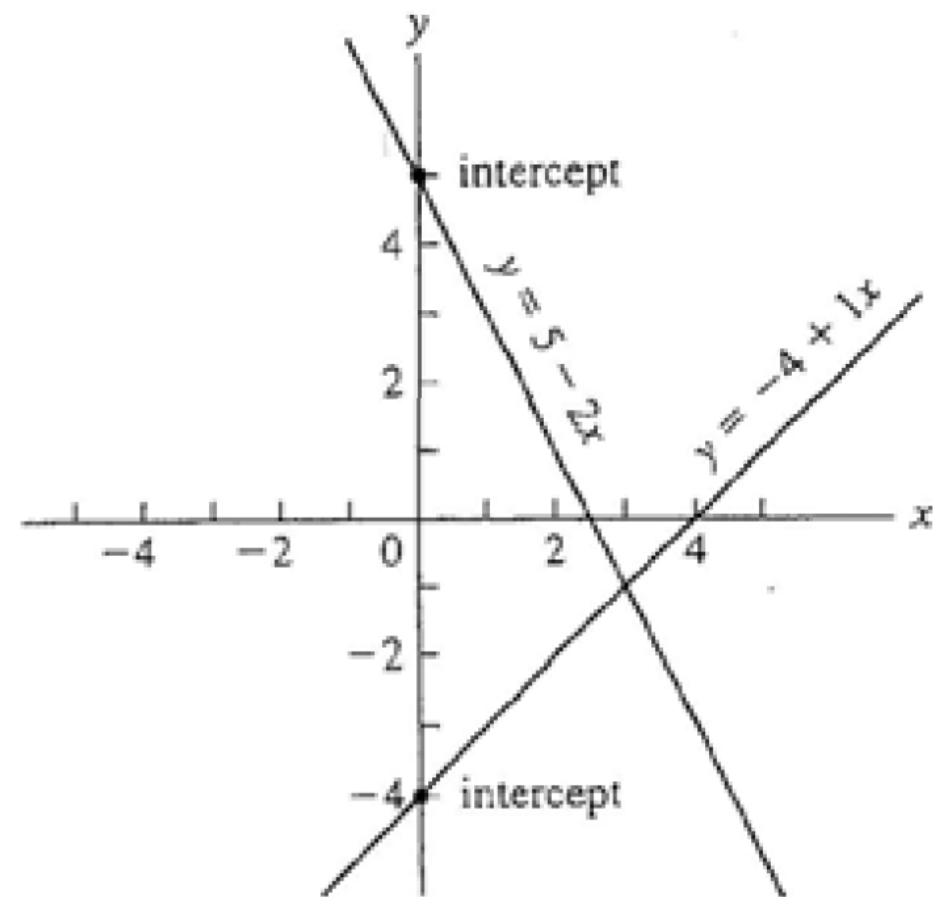


- Graph suggests that Y relates to X in an approximately linear way

# Regression: Step by Step

1. Assume a linear model:  $Y = \beta_0 + \beta_1 X$
2. Find the line which “best” fits the data, i.e. estimate parameters  $\beta_0$  and  $\beta_1$
3. Does variation in  $X$  help describe variation in  $Y$  ?
4. Check assumptions of model
5. Draw inferences and make predictions

# Straight-line Plots



# Assumptions of Linear Regression

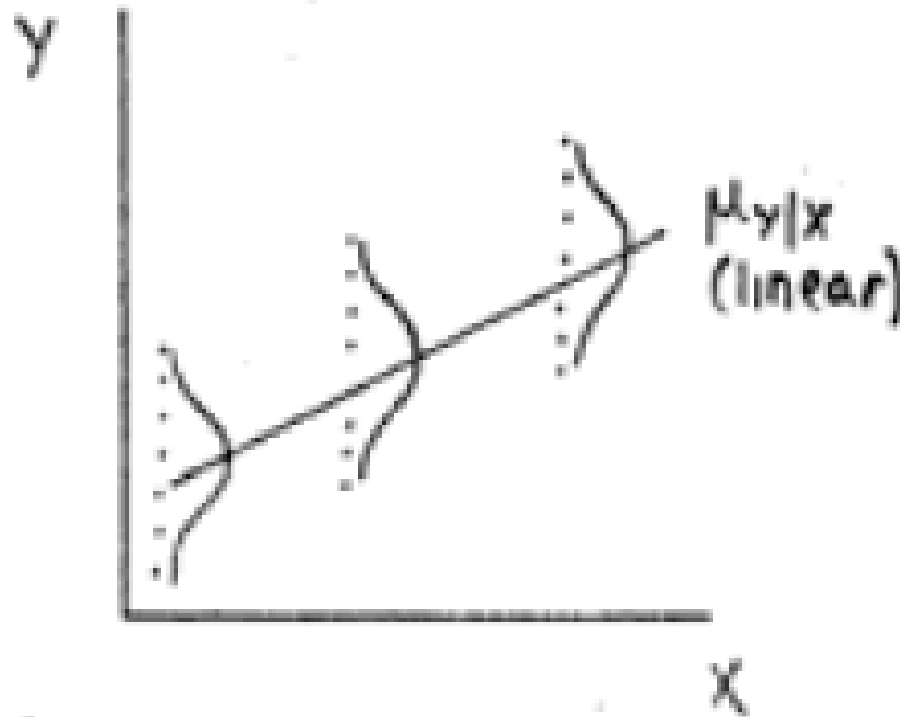
- Five basic assumptions
  1. Existence: for each fixed value of  $X$ ,  $Y$  is a random variable with finite mean and variance
  2. Independence: the set of  $Y_i$  are independent random variables given  $X_i$



# Assumptions of Linear Regression

3. Linearity: the mean value of  $Y$  is a linear function of  $X$

$$\mu_{Y|X} = E[Y|X] = \beta_0 + \beta_1 X$$



# Assumptions of Linear Regression

- 4. Homoscedasticity: the variance of Y is the same for any X
- 5. Normality: For each fixed value of X, Y has a normal distribution (by assumption 4,  $\sigma^2$  does not depend on X)

$$Y \sim N(\mu_{Y|X}, \sigma^2)$$

# Formulation

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{independent}$$

$$\Rightarrow Y_i \text{ are } N(\beta_0 + \beta_1 X_i, \sigma^2) \text{ given } X_i$$

$$\Rightarrow E(Y_i | X_i) = \beta_0 + \beta_1 X_i$$

$$\text{Var}(Y_i | X_i) = \sigma^2 = \text{variability of } Y_i \text{ about } \mu_{Y|X_i}$$

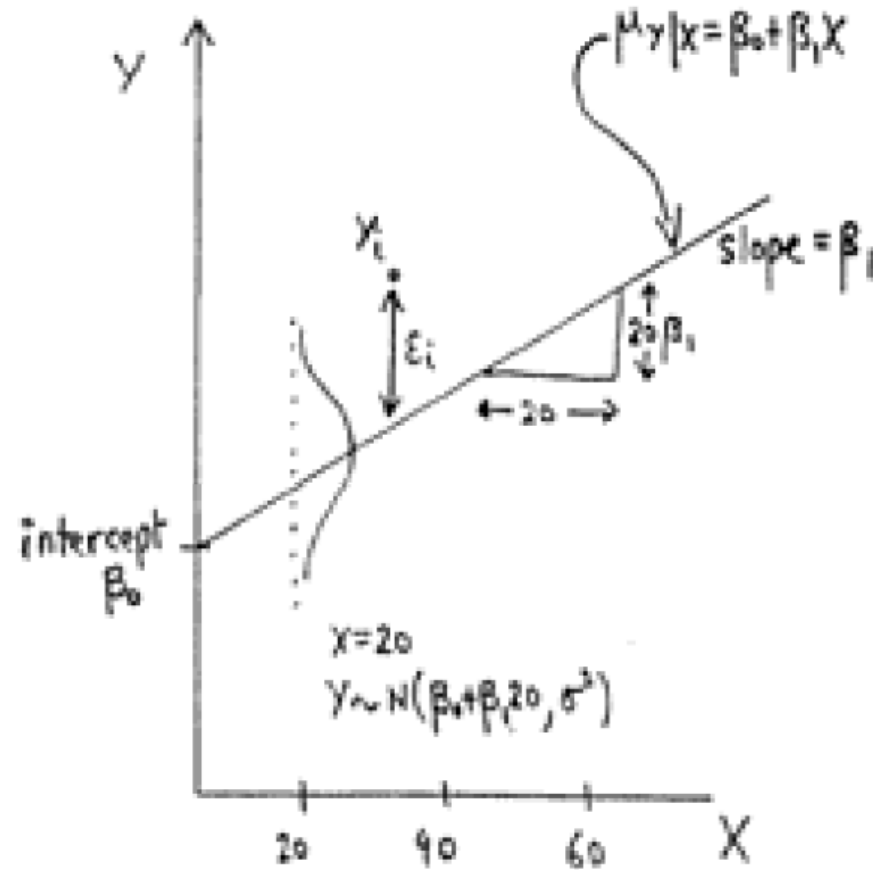
- $Y_i$  are linear function of  $X_i$  plus some random error

$$\varepsilon_i = \text{error} \left( = Y_i - \mu_{Y|X_i} \right)$$

$$\hat{\varepsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \text{ is called the "residual" for } Y_i$$

$$\hat{\varepsilon}_i = Y_i - \hat{Y}_i$$

# Linear Regression



# Estimating $\beta_0$ and $\beta_1$

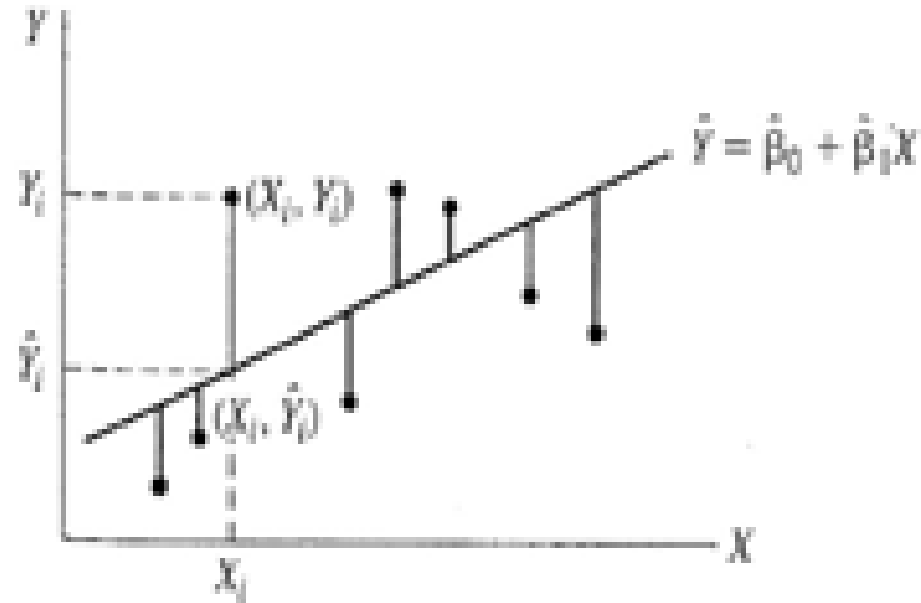
- Find “best” line
- Criterion for “best”: estimate  $\beta_0$  and  $\beta_1$  to minimize:

$$\begin{aligned} & \sum_{i=1}^n \left( Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \right)^2 \\ &= \sum_{i=1}^n \left( Y_i - \hat{Y}_i \right)^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 \end{aligned}$$

- This is the residual sum of squares, sum of squares due to error, or sum of squares about regression line
- Least Squares estimator

# Rationale for LS Estimates

- $\hat{\varepsilon}^2$  measures the “deviance” of  $Y_i$  from the estimated model
- The “best” model is the one from which the data deviate the least



# Least Squares Estimators

- Taking derivatives with respect to  $\beta$ , we obtain

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

# Example: SBP/age data

$$\bar{X} = 45.13 \quad \bar{Y} = 142.53$$

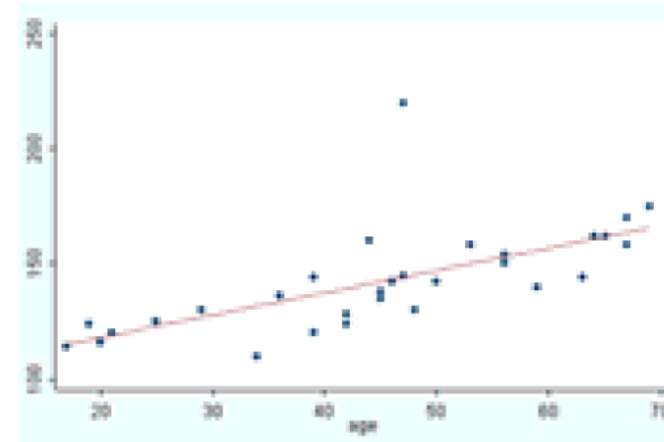
$$S_x^2 = 233.91 \quad S_y^2 = 509.91 \quad S_{xy} = 227.10$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_x^2} = \frac{227.10}{233.91} = 0.97$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 142.53 - 0.97(45.13) = 98.71$$

$$\hat{y} = 98.71 + .97x$$

.97mm Hg  $\uparrow$  for every  
1 yr  $\uparrow$  in age





# Using the Model

- Using the parameter estimates, our best guess for any Y given X is

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

- Hence at  $\bar{X}$

$$\hat{\beta}_0 + \hat{\beta}_1 \bar{X} = (\bar{Y} - \hat{\beta}_1 \bar{X}) + \hat{\beta}_1 \bar{X} = \bar{Y}$$

- Also  $\hat{Y} = \bar{Y} + \hat{\beta}_1 (X - \bar{X})$

# Regression and Correlation Coefficient

Regression  
Coefficient

Correlation  
Coefficient

$$\hat{\beta}_1 = \frac{S_{xy}}{S_x^2} \text{ and } r_{xy} = \frac{S_{xy}}{S_x S_y}$$

$$\Rightarrow r_{xy} = \hat{\beta}_1 \cdot \frac{S_x}{S_y}$$

# Example

Suppose we have:

$x$	$y$
1	2
2	4
-1	0
3	4
4	6

$$n = 5$$

$$\bar{x} = \frac{1 + 2 - 1 + 3 + 4}{5} = 1.8$$

$$\bar{y} = \frac{2 + 4 + 0 + 4 + 6}{5} = 3.2$$

Calculate  $r$ :

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1-1.8 = -0.8	2-3.2 = -1.2	(-0.8)(-1.2) = 0.96
2-1.8 = 0.2	4-3.2 = 0.8	(0.2)(0.8) = 0.16
-1-1.8 = -2.8	0-3.2 = -3.2	(-2.8)(-3.2) = 8.96
3-1.8 = 1.2	4-3.2 = 0.8	(1.2)(0.8) = 0.96
4-1.8 = 2.2	6-3.2 = 2.8	(2.2)(2.8) = 6.16
0	0	17.2

# Example

$$\sum (x_i - \bar{x})^2 = (-0.8)^2 + (0.2)^2 + (-2.8)^2 + (1.2)^2 + (2.2)^2 = 14.8$$

$$\sum (y_i - \bar{y})^2 = (-1.2)^2 + (0.8)^2 + (-3.2)^2 + (0.8)^2 + (2.8)^2 = 20.8$$

$$s_x = \sqrt{\frac{14.8}{4}} = 1.92$$

$$s_y = \sqrt{\frac{20.8}{4}} = 2.28$$

$$r = \frac{17.2}{\sqrt{(14.8)(20.8)}} = 0.98$$

$$\hat{\beta}_1 = r \cdot \frac{s_y}{s_x} = (0.98) \left( \frac{2.28}{1.92} \right) = 1.16$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 3.2 - (1.16)(1.8) = 1.11$$

Estimated regression line is  $\hat{y}_i = 1.11 + 1.16x_i$

# Example

Fitted values are the  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

$$\hat{y}_1 = 1.11 + 1.16(1) = 2.27 \quad \hat{y}_4 = 1.11 + 1.16(3) = 4.59$$

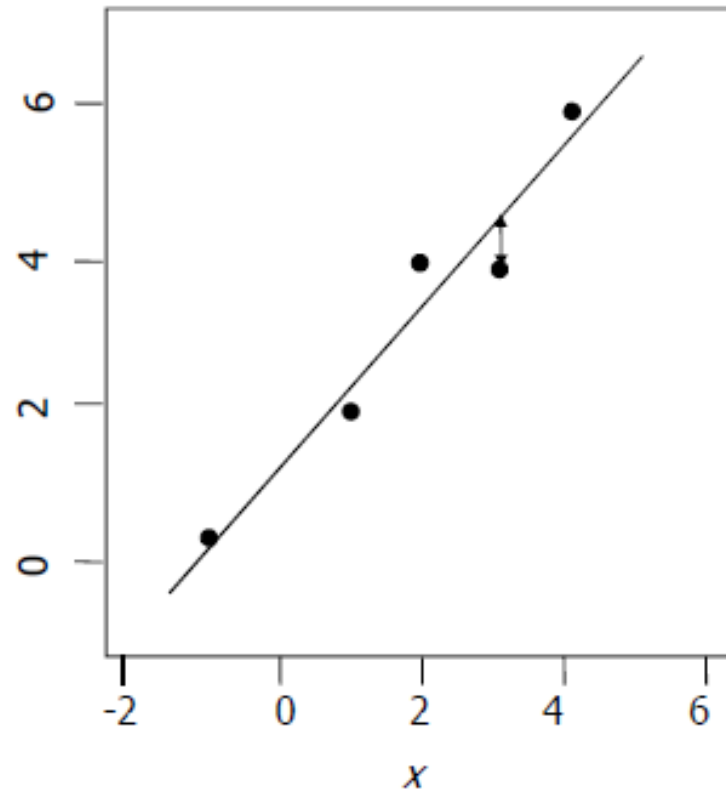
$$\hat{y}_2 = 1.11 + 1.16(2) = 3.43 \quad \hat{y}_5 = 1.11 + 1.16(4) = 5.75$$

$$\hat{y}_3 = 1.11 + 1.16(-1) = -0.05$$

Residuals =

$y_i - \hat{y}_i = \hat{\varepsilon}_i$	$(y_i - \hat{y}_i)^2$
2-2.27 = -0.27	$(-0.27)^2 = 0.073$
4-3.43 = 0.57	$(0.57)^2 = 0.345$
0-(-0.05) = 0.05	$(0.05)^2 = 0.0025$
4-4.59 = -0.59	$(-0.59)^2 = 0.348$
6-5.75 = 0.25	$(0.25)^2 = 0.0625$
$\sum \hat{\varepsilon}_i = \sum (y_i - \hat{y}_i) = 0$	$\sum (y_i - \hat{y}_i)^2 = 0.811$

# Example



$$\hat{y} = 1.11 + 1.16x$$

At  $x = 3$ :

$$y_i - \hat{y}_i = 4 - 4.59 = -.59$$

# Logistic Regression

# Logistic Regression

- Learn  $P(Y|X)$  directly
- 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 ( $\pi$ )
    - $Y$  is 1 with probability  $\pi$



# 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

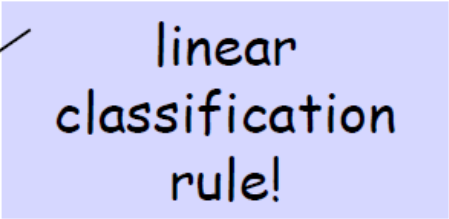
$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

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

implies

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

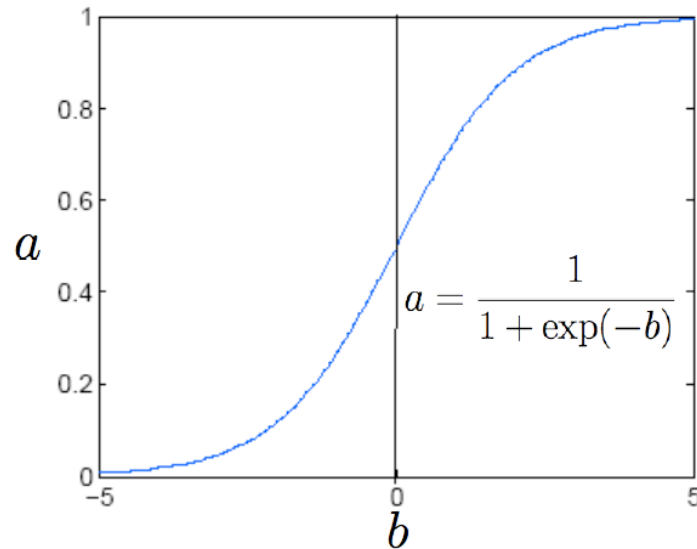


linear  
classification  
rule!

# 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



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

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

# 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: MCLE

- 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 (MCLE)

# 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

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

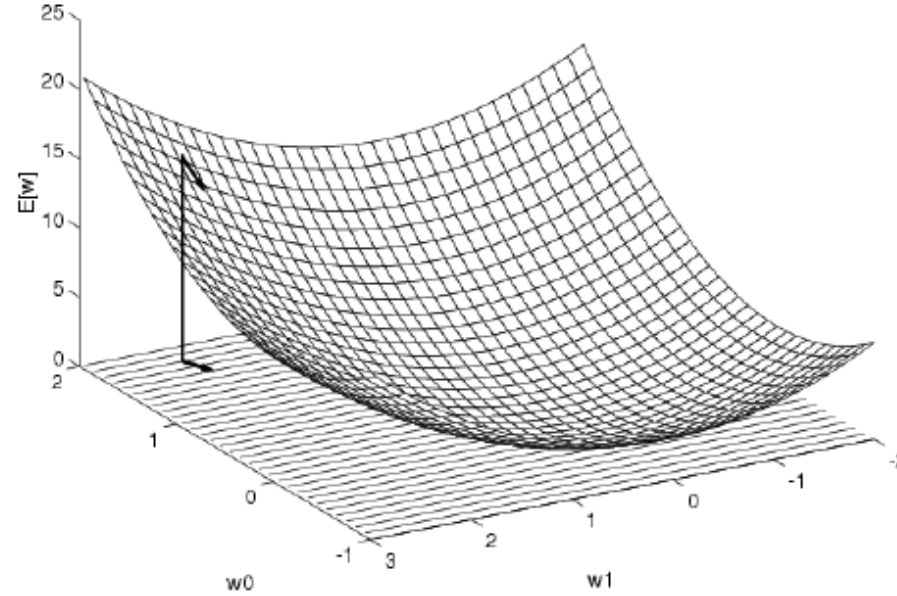
$$P(Y = 1|X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

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

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 ( $\pi$ )
- 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:  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$
- Can rewrite it:

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

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

*old problem*

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

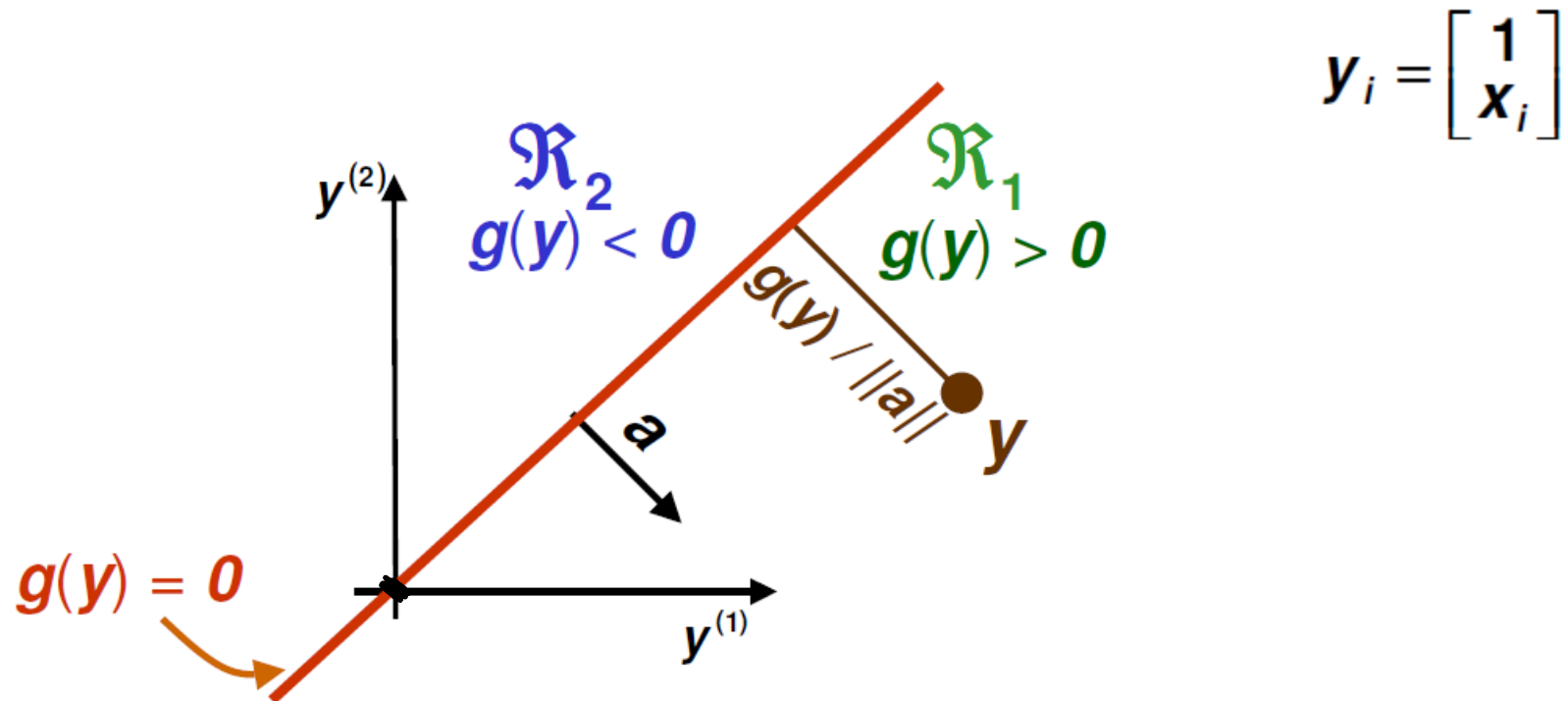
$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

*new problem*

$$g(\mathbf{y}) = \mathbf{a}^t \mathbf{y}$$

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

- Feature augmentation is done for simpler notation
- From now on, always assume that we have augmented feature vectors
  - Given samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  convert them to augmented samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by adding a new dimension of value 1





# Training Error

- For the rest of this part, assume we have 2 classes
  - Samples:  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , some in class 1, some in class 2
- Use samples to determine weights  $\mathbf{a}$  in the discriminant function  $\mathbf{g}(\mathbf{y}) = \mathbf{a}^t \mathbf{y}$
- What should the criterion for determining  $\mathbf{a}$  be?
- For now, suppose we want to minimize the training error (the number of misclassified samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$ )

- Recall that:
$$\begin{aligned} g(\mathbf{y}_i) > 0 &\Rightarrow \mathbf{y}_i \text{ classified as } \mathbf{c}_1 \\ g(\mathbf{y}_i) < 0 &\Rightarrow \mathbf{y}_i \text{ classified as } \mathbf{c}_2 \end{aligned}$$

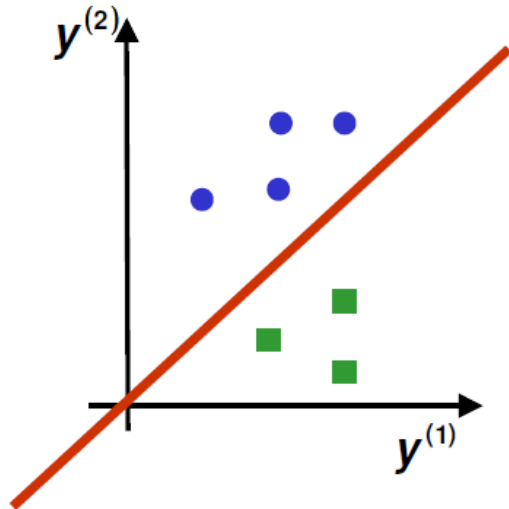
- Thus training error is 0 if
$$\begin{cases} g(\mathbf{y}_i) > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ g(\mathbf{y}_i) < 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$$

# “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

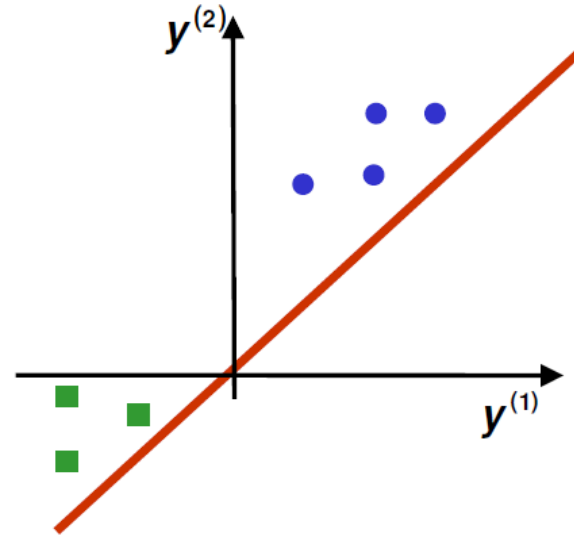
# Normalization

*before normalization*



- Seek a hyperplane that separates patterns from different categories

*after “normalization”*

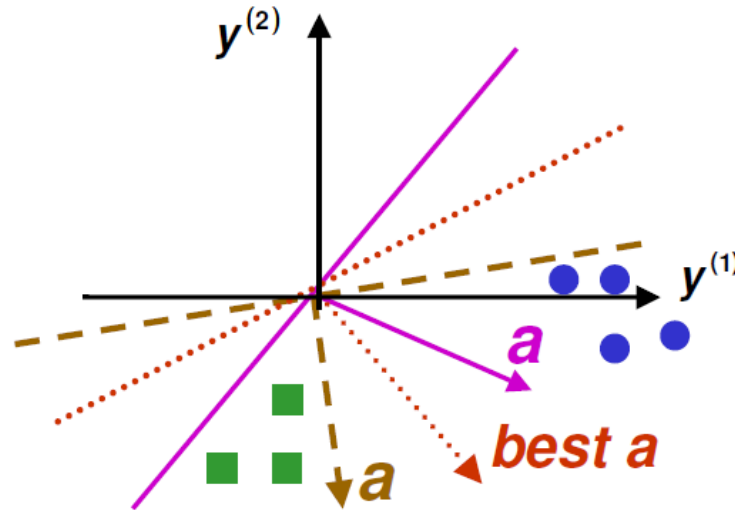


- Seek hyperplane that puts *normalized* patterns on the same (positive) side

# Solution Region

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

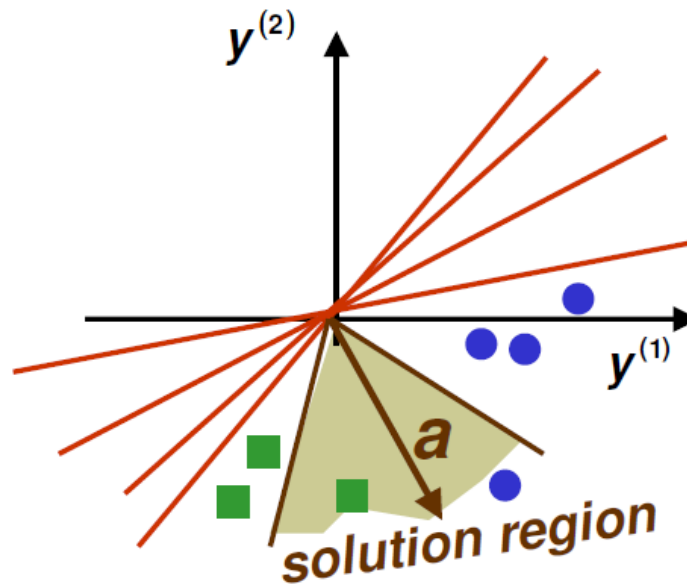
$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^d \mathbf{a}_k \mathbf{y}_i^{(k)} > 0$$



- In general, there can be many solutions

# Solution Region

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



# Optimization

- Need to minimize a function of many variables

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

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

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

*gradient*

# Optimization

- However solving analytically is not always easy

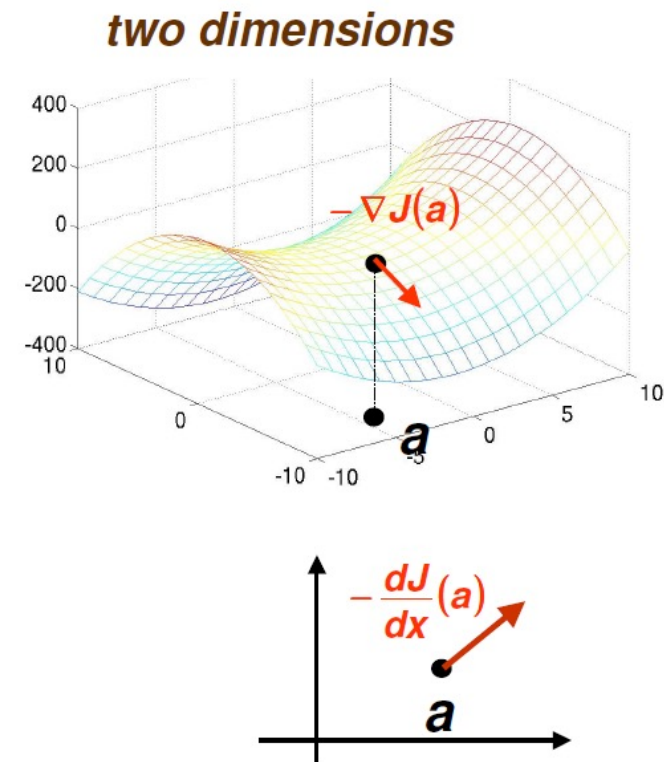
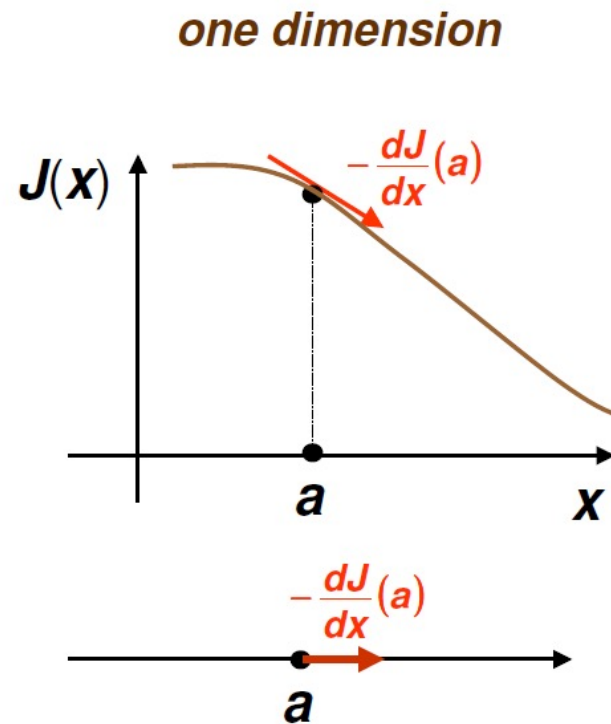
- For example:

$$\begin{cases} \sin(x_1^2 + x_2^3) + e^{x_4^2} = 0 \\ \cos(x_1^2 + x_2^3) + \log(x_3^3)^{x_4^2} = 0 \end{cases}$$

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

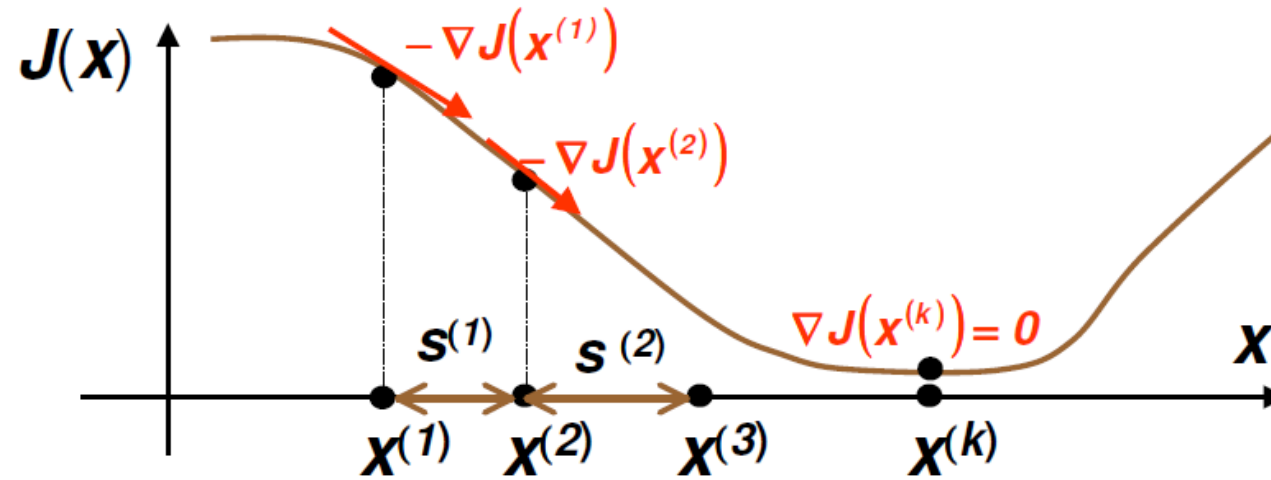
# Gradient Descent

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





# Gradient Descent



## Gradient Descent for minimizing any function $J(\mathbf{x})$

- Set  $k = 1$  and  $\mathbf{x}^{(1)}$  to some initial guess for the weight vector
- While  $\eta^{(k)} \|\nabla J(\mathbf{x}^{(k)})\| > \varepsilon$ 
  - Choose learning rate  $\eta^{(k)}$

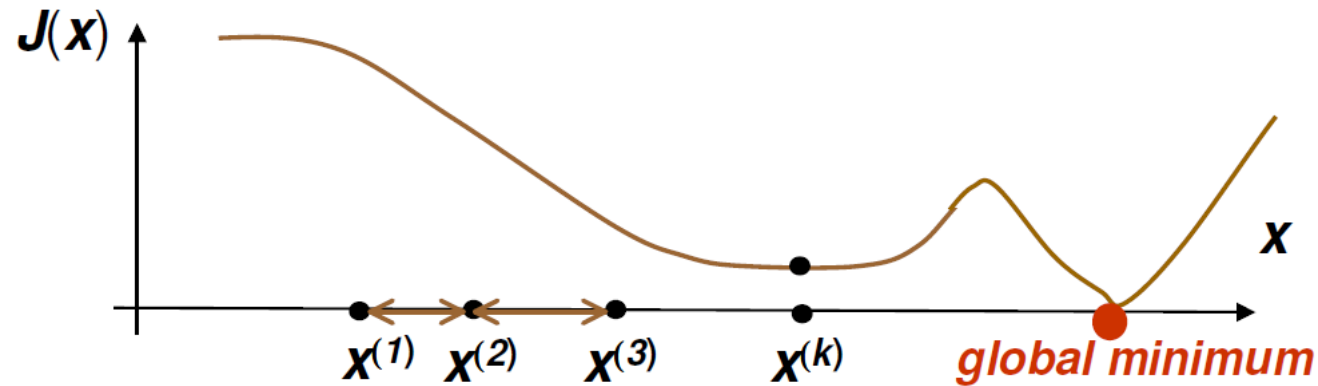
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta^{(k)} \nabla J(\mathbf{x})$$

$$k = k + 1$$

(update rule)

# Gradient Descent

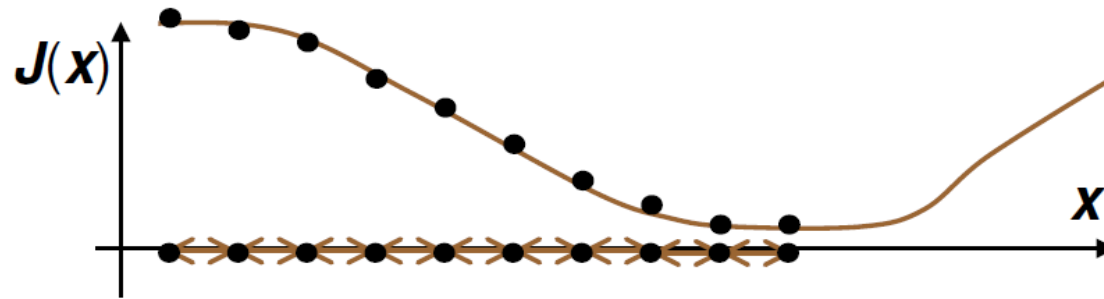
- Gradient descent is guaranteed to only find **local minima**



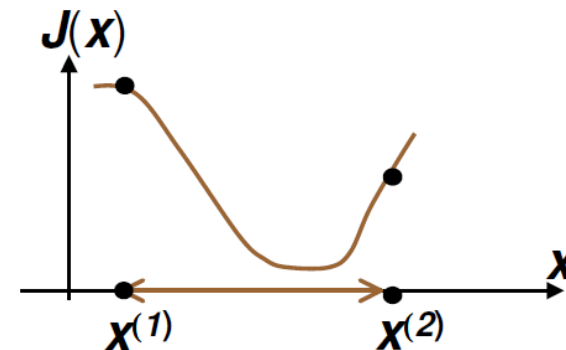
- Nevertheless gradient descent is very popular because it is simple and applicable to any function

# Gradient Descent

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



- If  $\eta$  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  $\mathbf{y}_1, \dots, \mathbf{y}_n$

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

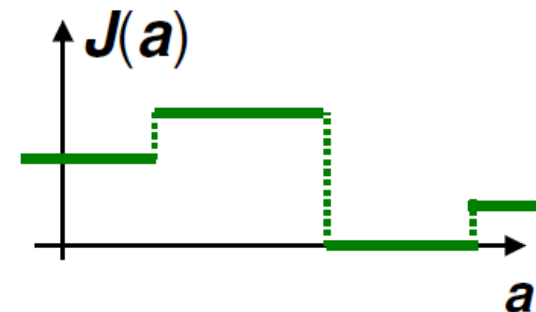
- Need criterion function  $J(\mathbf{a})$  which is minimized when  $\mathbf{a}$  is a solution vector
- Let  $\mathbf{Y}_M$  be the set of examples misclassified by  $\mathbf{a}$

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

- First natural choice: number of misclassified examples

$$J(\mathbf{a}) = |\mathbf{Y}_M(\mathbf{a})|$$

- Piecewise constant, gradient descent is useless

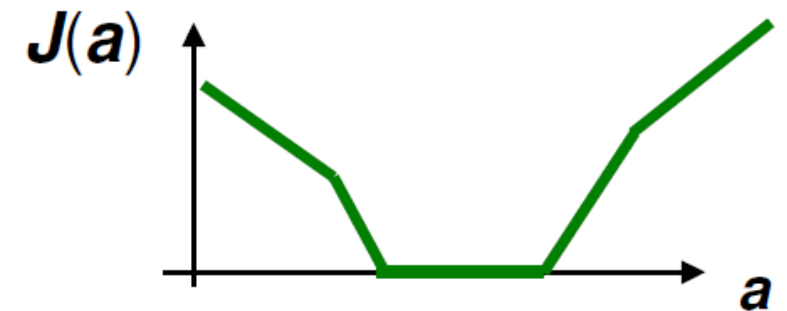
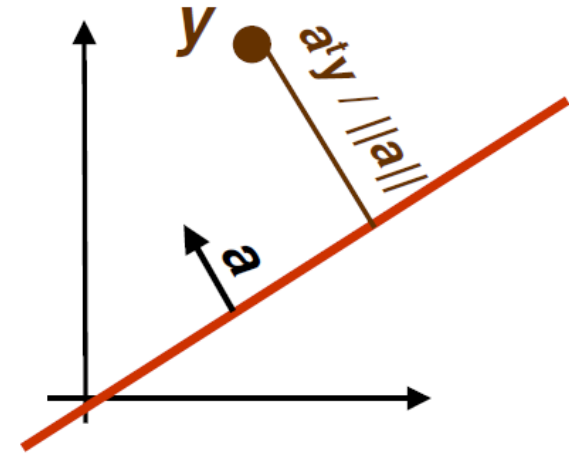


# Perceptron

# 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

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

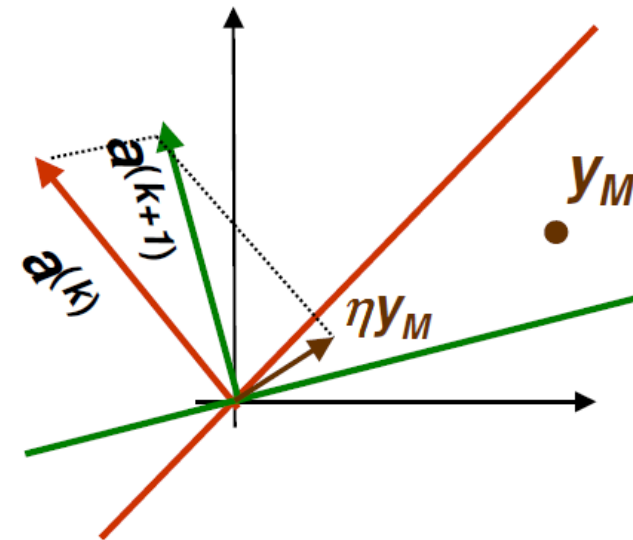
- Gradient of  $J_p(\mathbf{a})$  is:  $\nabla J_p(\mathbf{a}) = \sum_{y \in Y_M} (-y)$ 
  - $Y_M$  are samples misclassified by  $\mathbf{a}^{(k)}$
  - It is not possible to solve  $\nabla J_p(\mathbf{a}) = \mathbf{0}$  analytically because of  $Y_M$
- Update rule for gradient descent:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta^{(k)} \nabla J(\mathbf{x})$
- 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  $\mathbf{J}_p(\mathbf{a})$  is:

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

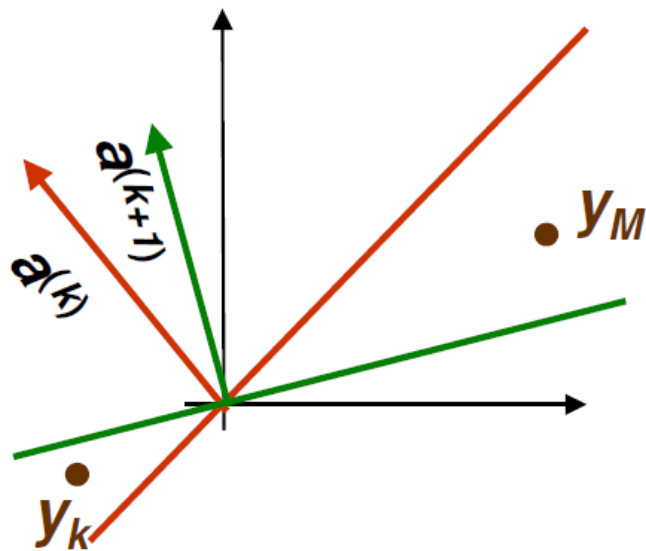
- Note that  $\mathbf{y}_M$  is one sample misclassified by  $\mathbf{a}^{(k)}$
- Must have a consistent way of visiting samples
- Geometric Interpretation:
  - $\mathbf{y}_M$  misclassified by  $\mathbf{a}^{(k)}$   $(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$
  - $\mathbf{y}_M$  is on the wrong side of decision hyperplane
  - Adding  $\eta \mathbf{y}_M$  to  $\mathbf{a}$  moves the new decision hyperplane in the right direction with respect to  $\mathbf{y}_M$



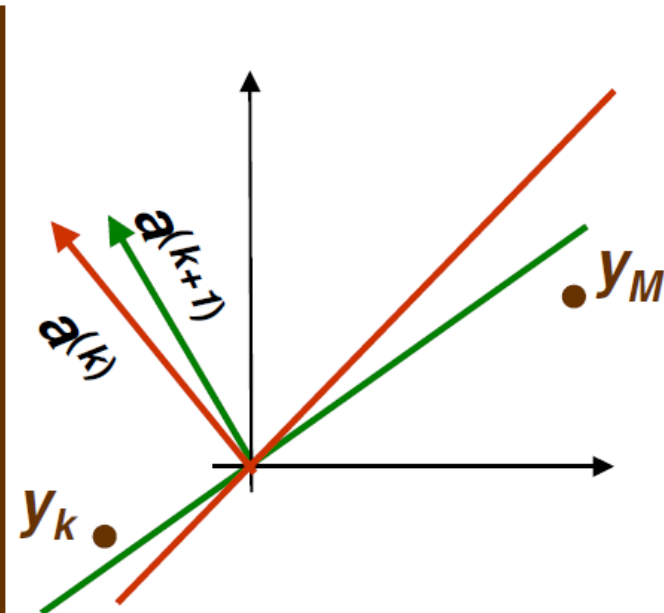


# Perceptron Single Sample Rule

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



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



$\eta$  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	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

- **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
- Seek ***a*** such that:

$$\begin{aligned}
 \mathbf{y}_i &\rightarrow -\mathbf{y}_i & \forall \mathbf{y}_i \in \mathbf{c}_2 \\
 \mathbf{a}^t \mathbf{y}_i &> 0 & \forall \mathbf{y}_i
 \end{aligned}$$

	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  $\eta$  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><math>\mathbf{a}^t \mathbf{y}</math></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>	<i><math>\mathbf{a}^t \mathbf{y}</math></i>	<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>	<i><math>\mathbf{a}^t \mathbf{y}</math></i>	<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>	<i>a<sup>t</sup>y</i>	<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:

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

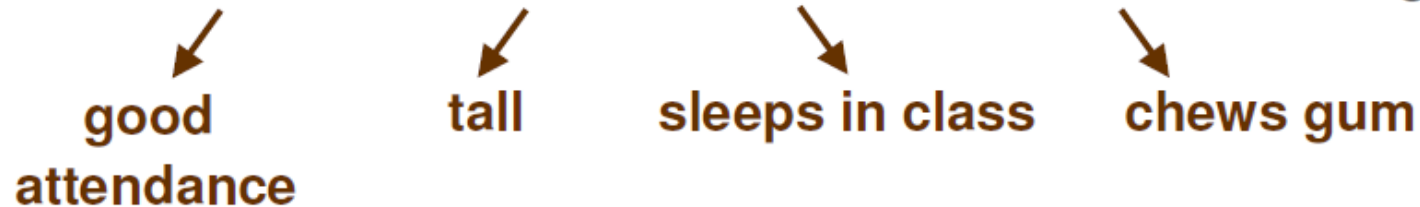
- Converting back to the original features  $\mathbf{x}$ :

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

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

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

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


  
 good attendance      tall      sleeps in class      chews gum

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

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

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

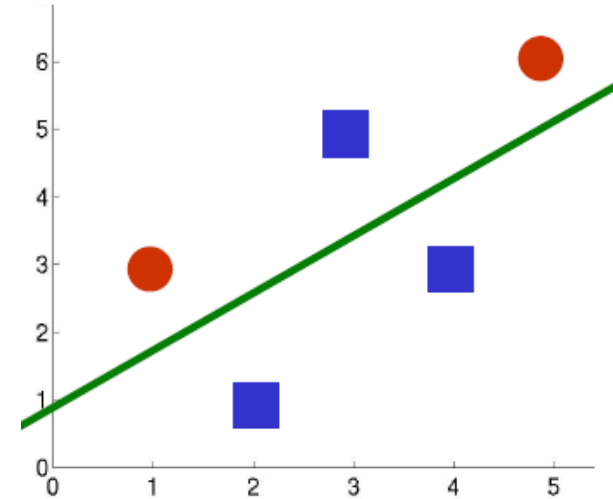
$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} > 1 \Rightarrow \text{grade A}$$

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

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

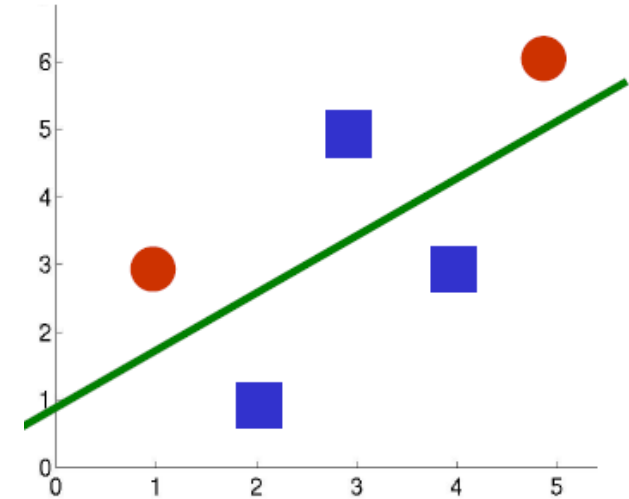
# LDF: Non-separable Example

- 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

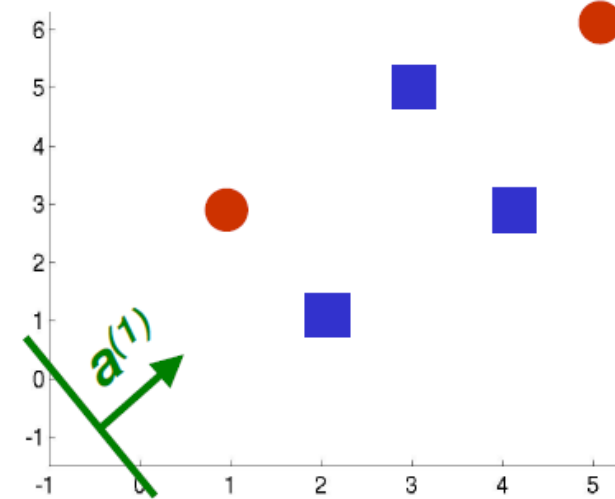
- Obtain  $y_1, y_2, y_3, y_4$  by adding extra feature and “normalizing”



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

# 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  $\eta = 1$



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

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

- $y_1^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 2 \ 1]^t > 0 \quad \checkmark$
- $y_2^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 4 \ 3]^t > 0 \quad \checkmark$
- $y_3^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 3 \ 5]^t > 0 \quad \checkmark$

# LDF: Non-separable Example

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

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

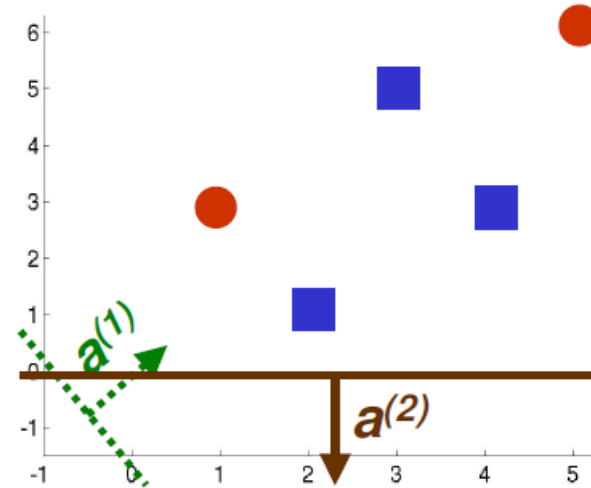
- $\mathbf{y}_4^t \mathbf{a}^{(1)} = [1 \ 1 \ 1]^* [-1 \ -1 \ -3]^t = -5 < 0$

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

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

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

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



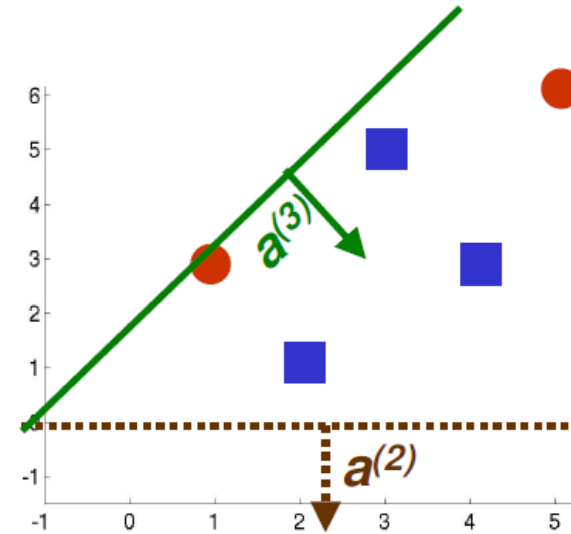
# LDF: Non-separable Example

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

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

- $\mathbf{y}_2^t \mathbf{a}^{(3)} = [1 \ 4 \ 3]^* [1 \ 2 \ -1]^t = 6 > 0 \quad \checkmark$
- $\mathbf{y}_3^t \mathbf{a}^{(3)} = [1 \ 3 \ 5]^* [1 \ 2 \ -1]^t > 0 \quad \checkmark$
- $\mathbf{y}_4^t \mathbf{a}^{(3)} = [-1 \ -1 \ -3]^* [1 \ 2 \ -1]^t = 0$

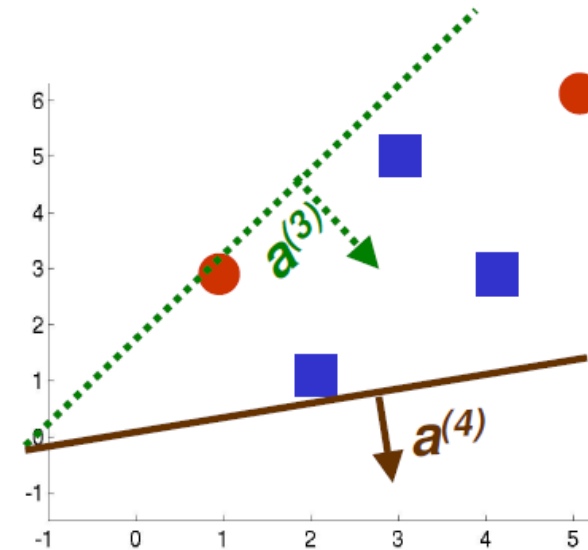
$$\mathbf{a}^{(4)} = \mathbf{a}^{(3)} + \mathbf{y}_M = [1 \ 2 \ -1] + [-1 \ -1 \ -3] = [0 \ 1 \ -4]$$



# LDF: Non-separable Example

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

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



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

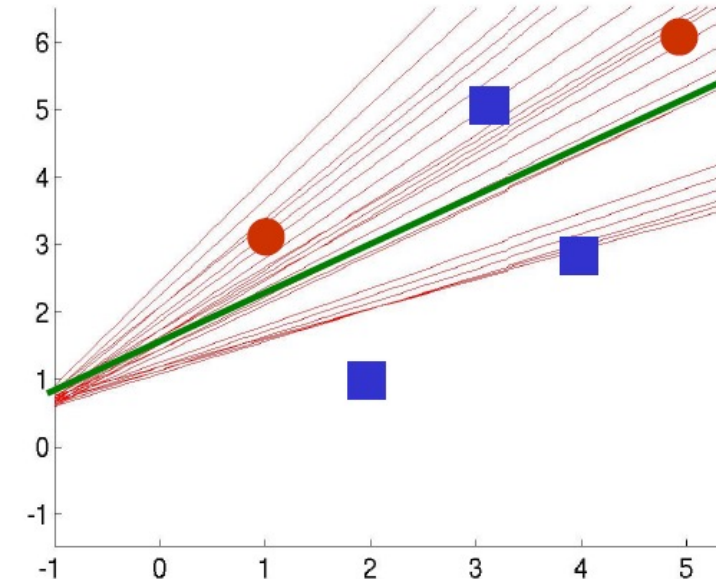


# LDF: Non-separable Example

- 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