

INTRODUCTION AND APPLICATIONS OF MACHINE LEARNING

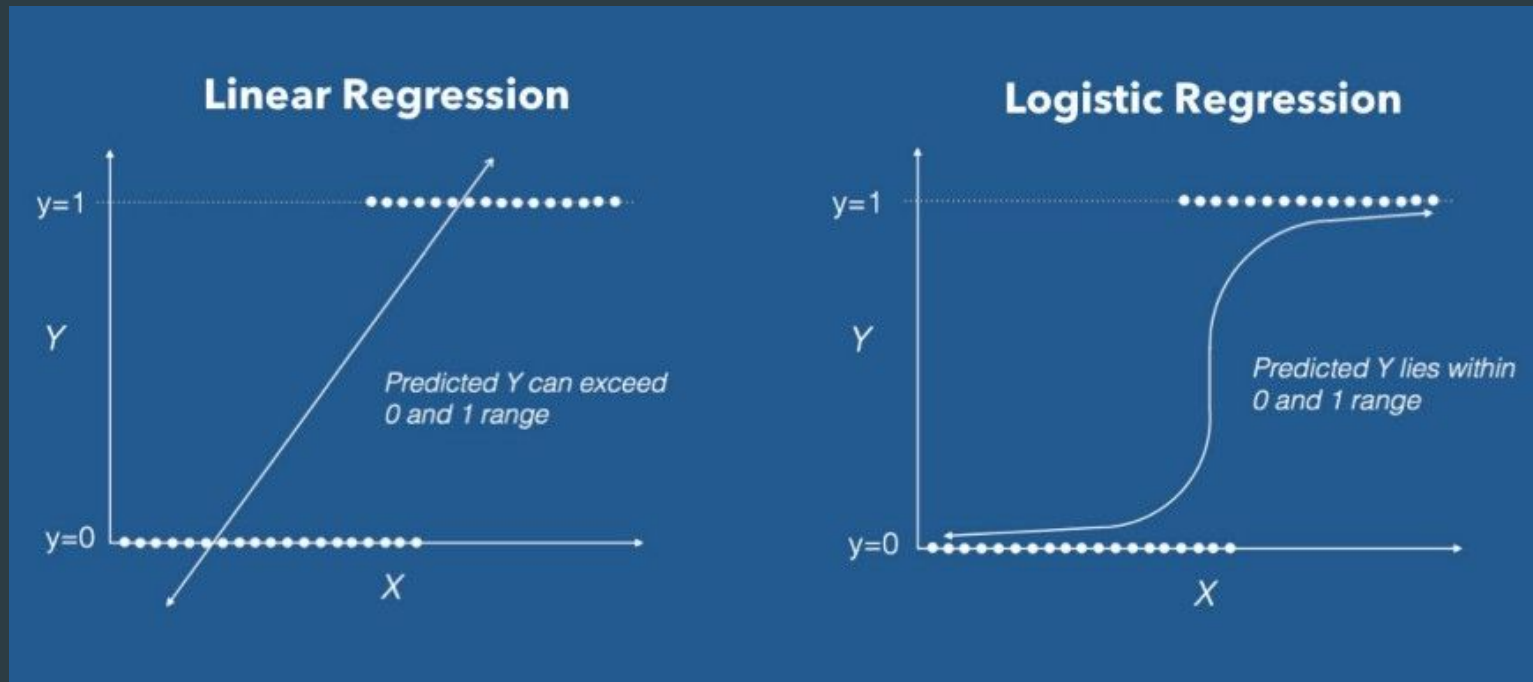


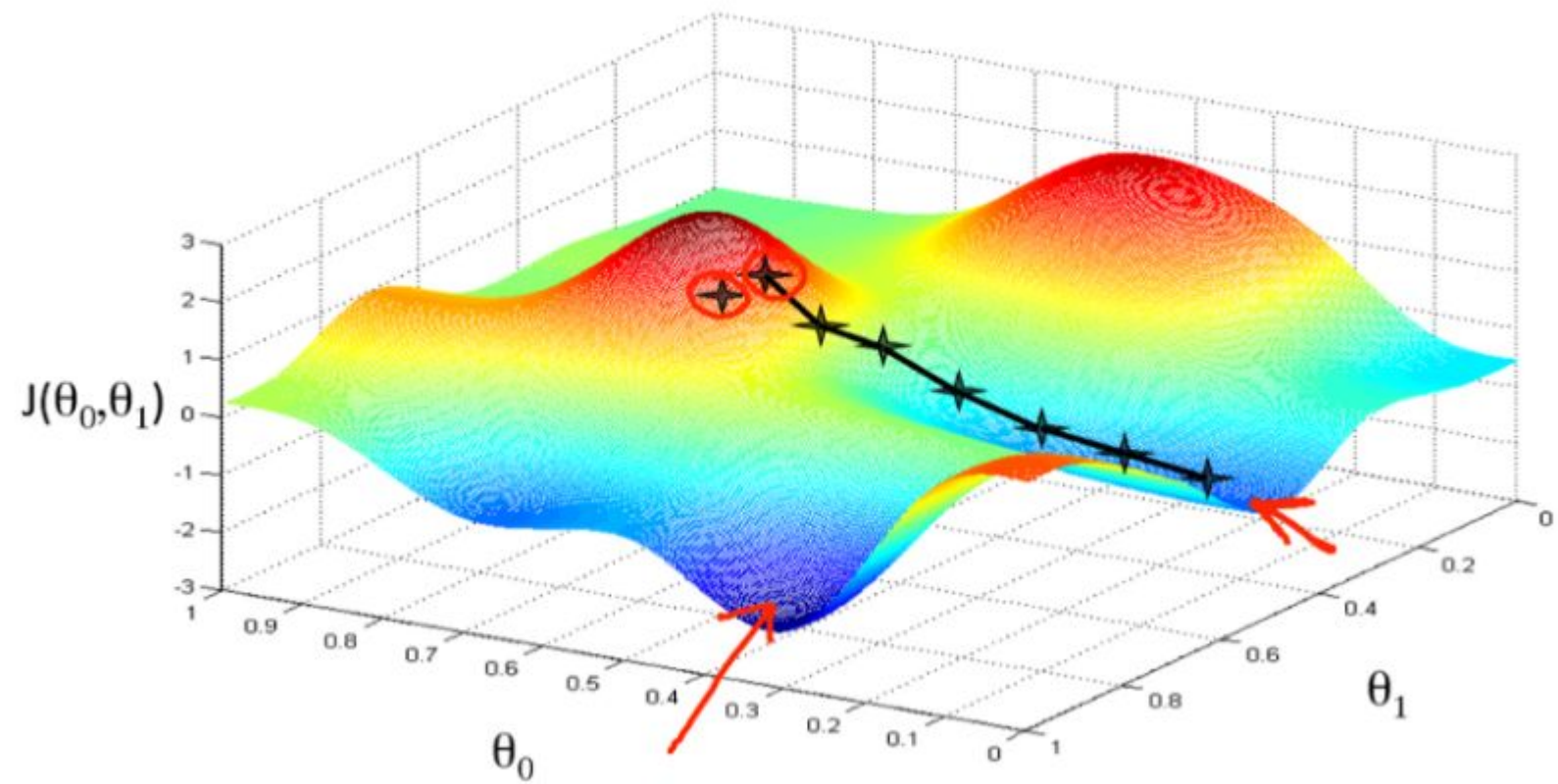
Roadmap

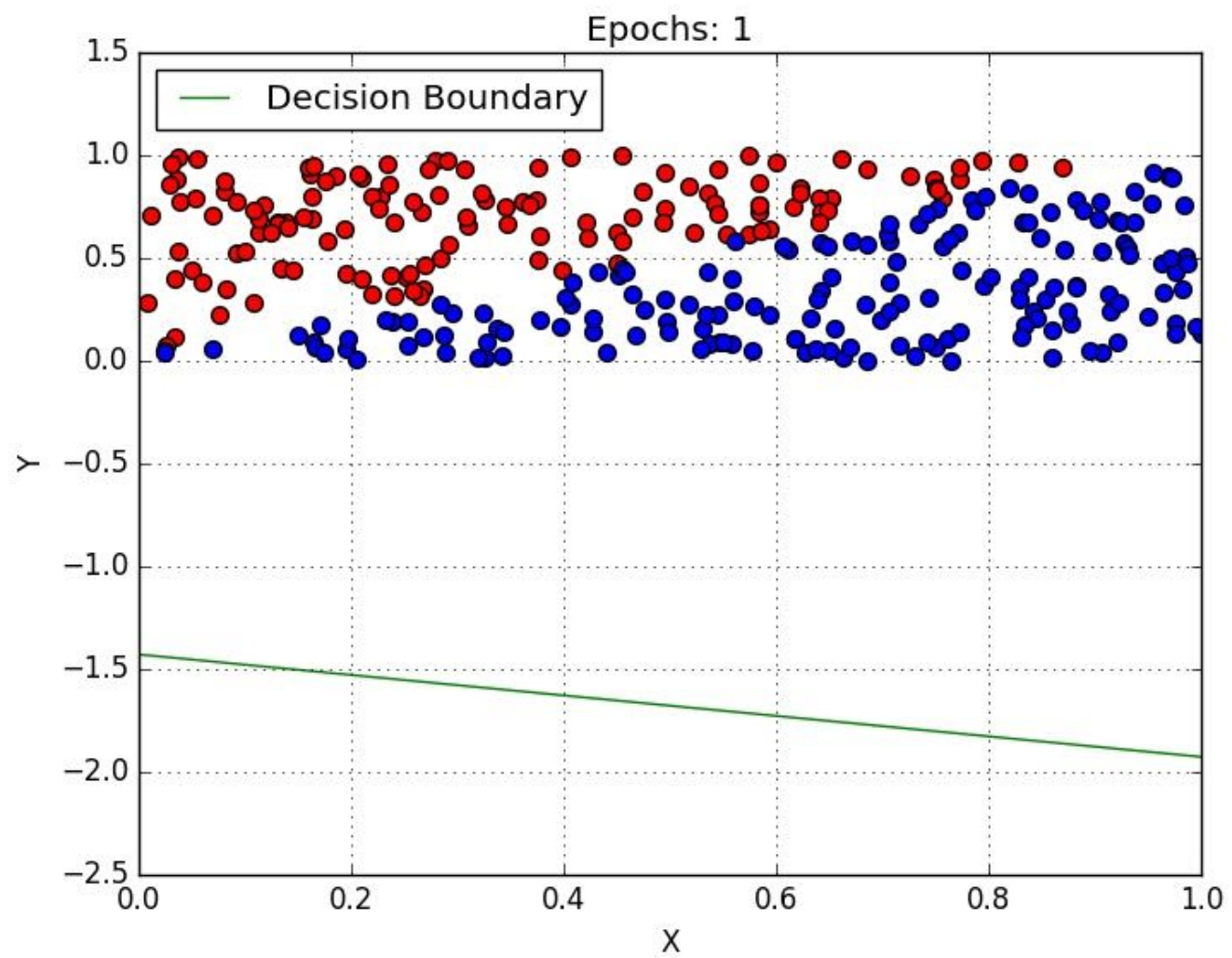
- Logistic regression
 - Hypothesis function
 - Loss function
 - Gradient descent update

Logistic regression

- Logistic regression is a classification algorithm used to assign observations to a discrete set of classes. Some of the examples of classification problems are Email spam or not spam, Online transactions Fraud or not Fraud, Tumor Malignant or Benign. Logistic regression transforms its output using the logistic sigmoid function to return a probability value.

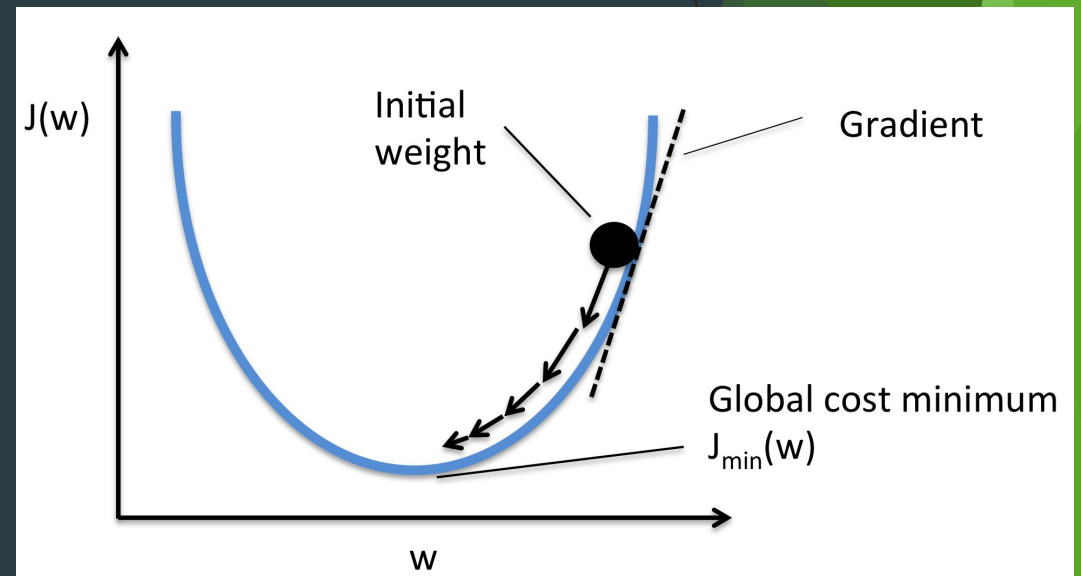




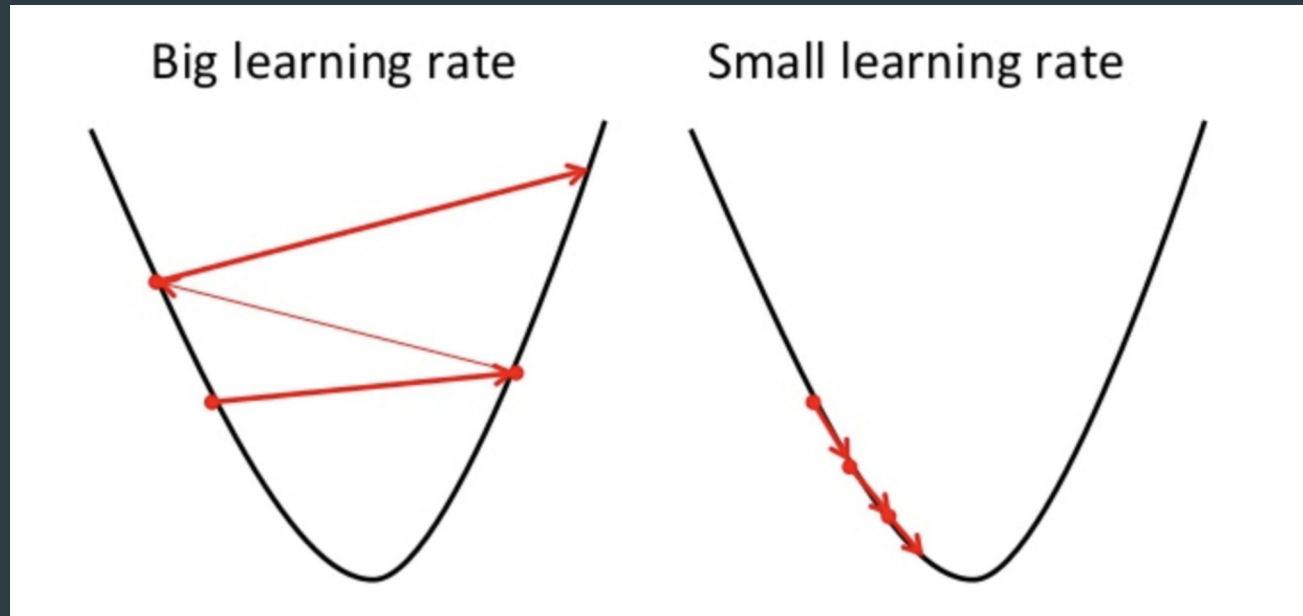


Gradient Descent

- Optimization method to find minima in a function (here reduce the cost function MSE)
- It is an iterative process
- Start at any point and move towards the minima
- This depends on
 - Step size (η)
 - direction(determined by the negative of the gradient)



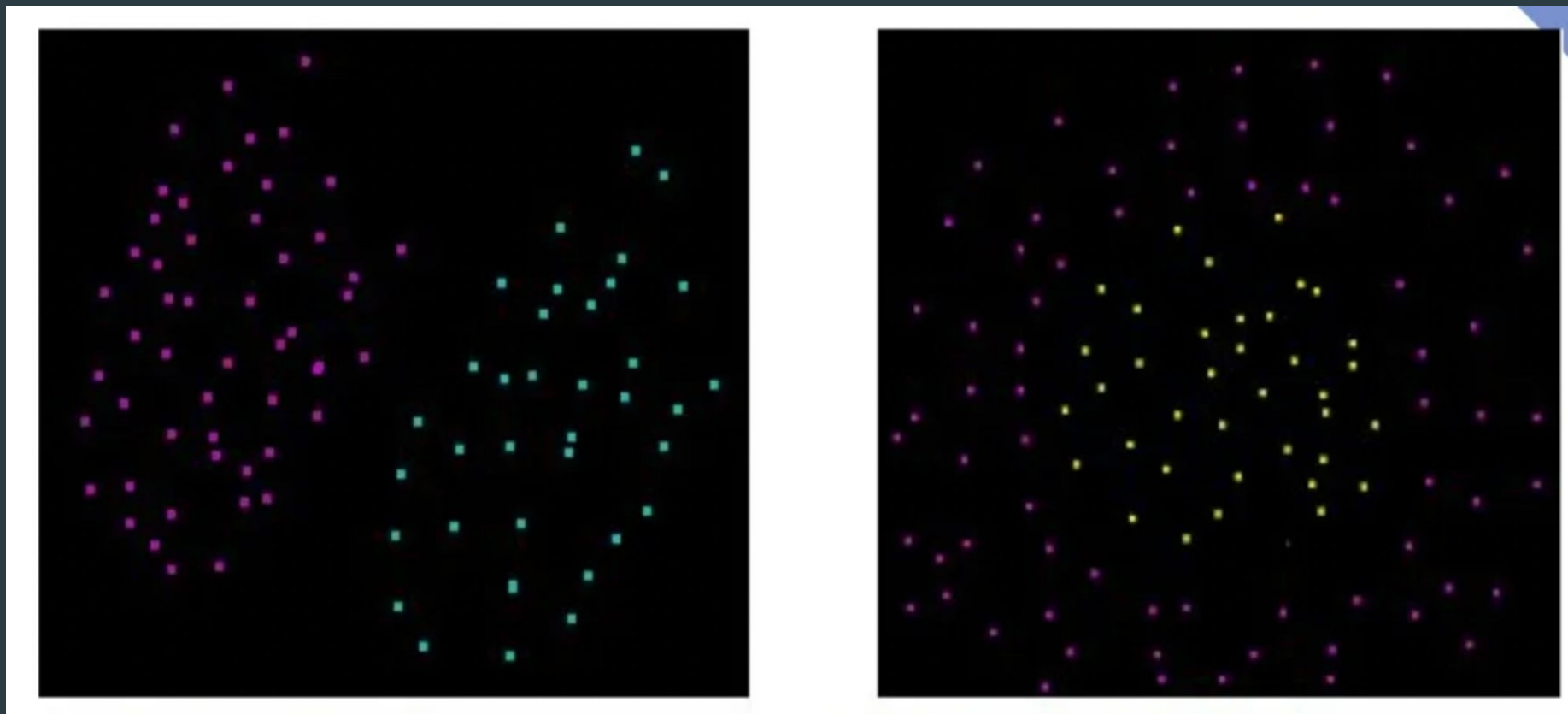
Gradient update



Support Vector Machines

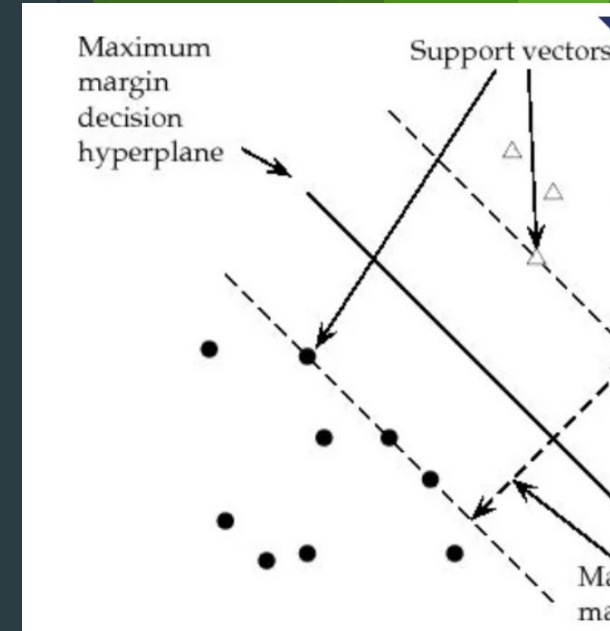
- ▶ A support vector machine is a discriminative classifier which intakes training data, the algorithm learns an optimal hyperplane which categorizes new examples
- ▶ supervised learning algorithm

Linearly-separable vs non



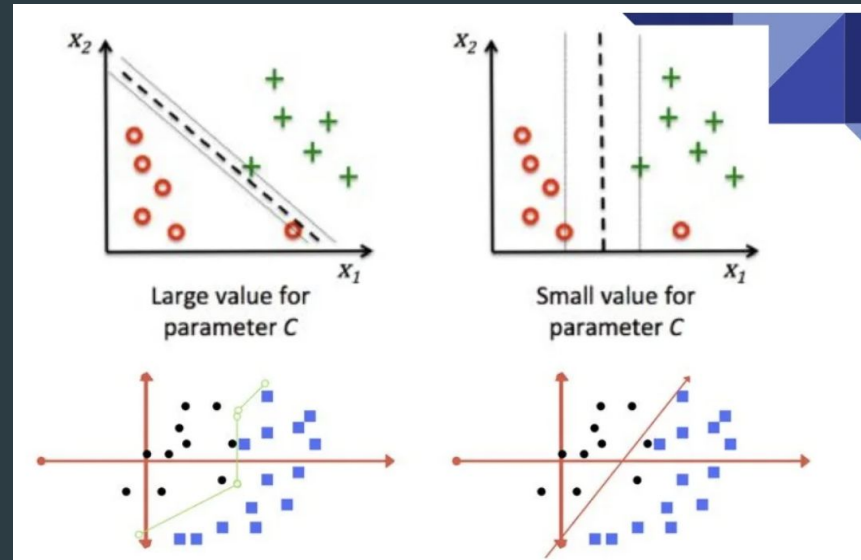
Margin

- ▶ Margin is the perpendicular distance between the closest data points and the hyperplane
- ▶ The best optimized line(hyperplane) with maximum margin is termed as Marginal Maximal hyperplane
- ▶ The closest points where the margin distance is calculated is known as the support vectors



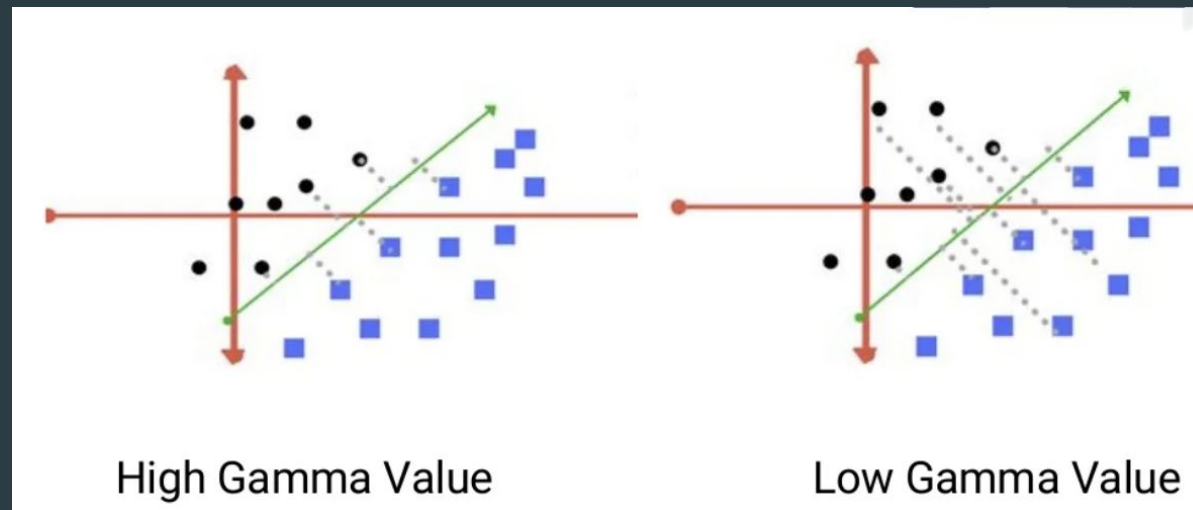
Regularization(c)

- ▶ Optimizes SVM classifier to avoid misclassifying
- ▶ $C \rightarrow$ large and margin of hyperplane \rightarrow small
- ▶ $C \rightarrow$ small and margin of hyperplane \rightarrow large



Gamma

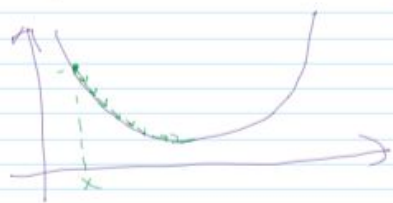
- ▶ Defines how far influences the calculation of plausible line of separation
- ▶ Low gamma -> points far from plausible line are considered for calculation
- ▶ High gamma -> points closer to plausible line are considered for calculation



Kernels

- ▶ Mathematical functions for transforming data
- ▶ Different SVM algorithms use different type of kernel functions
- ▶ Linear kernel
- ▶ Radial basis function
- ▶ Sigmoid
- ▶ Polynomial

gradient descent algorithm



$$y = f(x)$$

$$f(x) = (x-5)^2$$

$$x = x - \eta \left[\frac{\partial f(x)}{\partial x} \right]_0$$

(gradient learning rate)
gradient/slope

$$J_0 = \frac{1}{n} \sum_{i=1}^m [y^{(i)} - y_{\text{actual}}]^2$$

update θ using J_0

$$J_\theta = \frac{1}{n} \sum [\theta_0 + \theta_1 x]$$

$$\theta = \theta - \eta \nabla J(\theta)$$

$$\Rightarrow \theta_0 = \theta_0 - \eta \frac{\partial J(\theta)}{\partial \theta_0} \quad \theta_1 = \theta_1 - \eta \frac{\partial J(\theta)}{\partial \theta_1}$$

$$J(\theta) = \frac{1}{n} \sum_{i=0}^m [\theta_0 + \theta_1 x^{(i)} - y^{(i)}]_{\text{actual}}^2$$

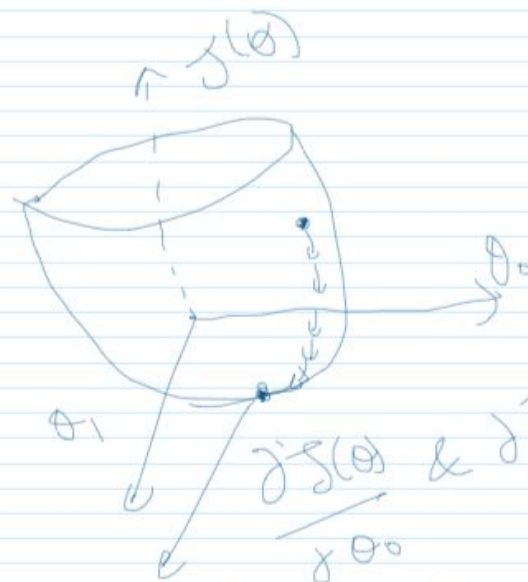
$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{n} \sum_{i=0}^m [\theta_0 + \theta_1 x^{(i)} - y^{(i)}]_{\text{actual}} \cdot 1$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{n} \sum_{i=0}^m [\theta_0 + \theta_1 x^{(i)} - y^{(i)}]_{\text{actual}} \cdot x^{(i)}$$

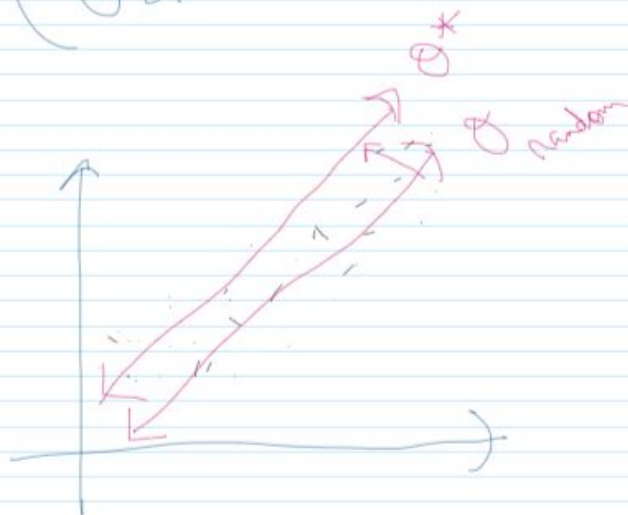
finally,

$$\theta_0 = \theta_0 - \eta \frac{1}{n} \sum_{i=1}^m [\hat{y}^{(i)} - y_{\text{actual}}^{(i)}]$$

$$\theta_1 = \theta_1 - \eta \frac{1}{n} \sum_{i=1}^m [\hat{y}^{(i)} - y_{\text{actual}}^{(i)}] x^{(i)}$$



$$(\theta_0, \theta_1) \rightarrow \theta^*$$





$$X = \begin{bmatrix} x_0^1 & x_1^1 & \dots & x_n^1 \\ \vdots & \vdots & \ddots & \vdots \\ x_0^i & x_1^i & \dots & x_n^i \\ \vdots & \vdots & \ddots & \vdots \\ x_0^m & x_1^m & \dots & x_n^m \end{bmatrix} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

x_0^i is the bias term (always 1).
 x_1, x_2, x_3, \dots are the features.
 (m rows, n features)

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$= \theta_0 + \sum_{i=1}^n \theta_i x_i$$

$$= \theta_0 x_0$$

$$= \sum_{i=0}^n \theta_i x_i$$

$x_0 = 1$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad X = \begin{bmatrix} x_0 & x_1 & x_2 & \dots & x_n \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$n+1$ parameters: $\theta_0, \dots, \theta_n$
 $m \times (n+1)$

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

$$J(\theta) = \frac{1}{m} \sum (y - y^{(i)})^2$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{\partial}{\partial \theta_0} (h_{\theta}(x) - y)$$

$$= 2 (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_0} (h_{\theta}(x))$$

$$= (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_0} \left[\sum_{j=0}^n \theta_j x_j \right]$$

$$= (h_{\theta}(x) - y) \cdot x_0$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \sum_{i=1}^m (y - y^{(i)}) \cdot x_0$$

gradient w.r.t

$j=0$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \sum_{i=1}^m (y - y^{(i)}) \cdot 1$$

θ_0

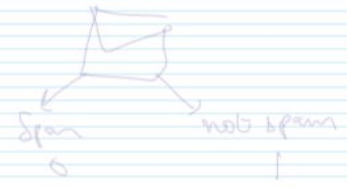
$j=1$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \sum_{i=1}^m (y - y^{(i)}) \cdot x_i$$

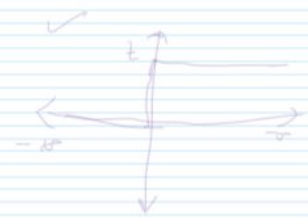
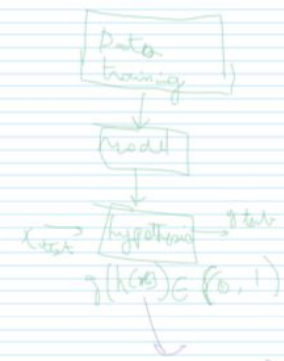
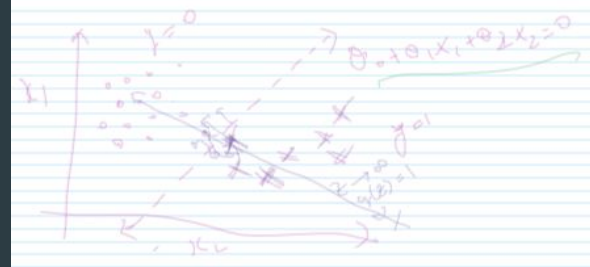
Gradient update

$$\theta_j = \theta_j - \left(\sum_{i=1}^m (y^{(i)} - y^{(i)}) x_j \right)$$

logistic regression
 → supervised algorithm (X, y)
 → Classification algorithm
 binary
 DT, SVMs



$$x \in \mathbb{R}^2 \rightarrow y \in \{0, 1\}$$



$$g(z) = 0 \quad z \leq 0$$

$$g(z) = 1 \quad z > 0$$

$g(z) \geq 0.5 \Rightarrow \text{class 1}$
 $g(z) < 0.5 \Rightarrow \text{class 0}$

$$g(z) = \frac{1}{1 + e^{-z}}$$

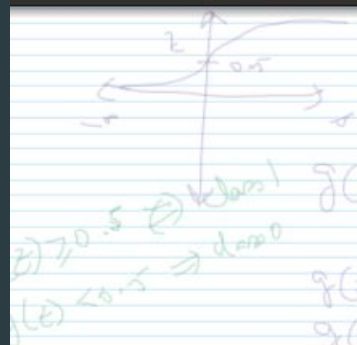
$$g(z) \rightarrow 1 \quad z \rightarrow \infty$$

$$g(z) \rightarrow 0 \quad z \rightarrow -\infty$$

$$h_0(z) = \sum_{i=0}^n \theta_i z_i$$

$$g(h_0(x)) = \frac{1}{1 + e^{-\sum_{i=0}^n \theta_i x_i}}$$

at=1 deg=0



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g(z) \rightarrow 1 \quad z \rightarrow \infty$$

$$g(z) \rightarrow 0 \quad z \rightarrow -\infty$$

$$h_0(x) = \sum_{i=0}^n \theta_i x_i$$

$$g(h_0(x)) = \frac{1}{1 + e^{-\sum_{i=0}^n \theta_i x_i}}$$

① (0, 1)

② Assign diff. values for diff. points
→ classifying

③ loss func

$$\text{Avg. sq. loss} = \frac{1}{m} \sum (\hat{y}^{(i)} - g(h_0(x^{(i)})))^2$$

④ Minimize the error

$$\theta = \theta - \eta \cdot \frac{\partial L}{\partial \theta}$$

Log loss Binary Cross entropy $\hat{y} = h_0(x)$

$$\text{loss} = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_0(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_0(x^{(i)})) \right]$$

$$y = \begin{bmatrix} 1 & 0 & 1 & 1 & 0 \\ 0.5 & 0.6 & 0.2 & 0.1 & 0.3 \end{bmatrix}$$

$$y^{(i)} = 0 \quad 0 \quad 0 \quad 0 \quad 0$$

$$(-1 \log 1 + 0 \log 0)$$

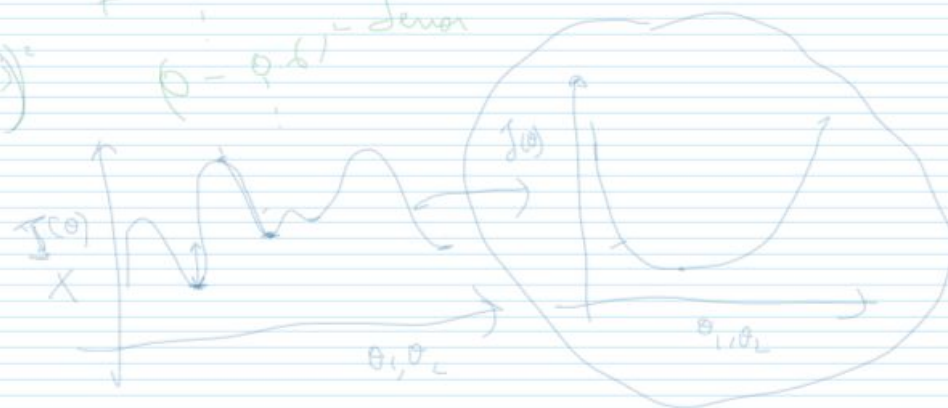
$$(1 - 0.8)^2$$

$$(1 - 0.2)^2$$

$$+ \dots$$

$$(0.6)^2$$

total error



$$y = \begin{pmatrix} 1 & 0 & 1 & 1 & 0 \\ 0.5 & 0.6 & 0.2 & 0.1 & 0.3 \end{pmatrix}$$

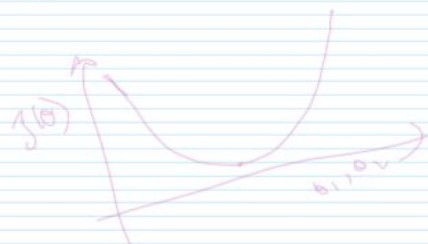
$$y^{(i)} = 1 \ 1 \ 1 \ 1 \ 1$$

$$\text{loss} = -\frac{1}{n} \sum_{i=1}^m \log(h_\theta(x^{(i)}))$$



$$y^{(i)} = 0 \ 0 \ 0 \ 0 \ 0$$

$$\text{loss} = -\log(1 - h_\theta(x^{(i)}))$$



$$J(\theta) = -\frac{1}{n} \sum_{i=1}^m y^{(i)} (\log(h_\theta(x^{(i)}))) + (1 - y^{(i)}) (1 - \log(h_\theta(x^{(i)})))$$

$$h_\theta(x) = g(\sum x_i \theta_i)$$

$$\Rightarrow \theta = \theta - \eta \cdot \frac{\partial J(\theta)}{\partial \theta}$$

THANK YOU!

