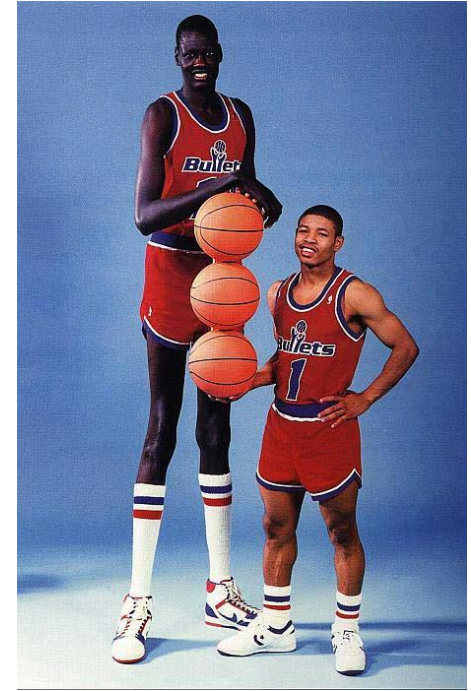


# Machine Learning

## 4771

Instructor: Itsik Pe'er

# Reminder: Parameter Estimation



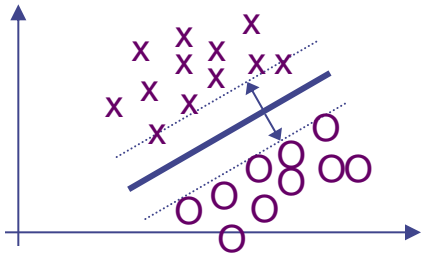
## A-Posteriori $\rightarrow$ Regularization

$$R_{regularized}(\theta) = \frac{1}{N} \sum_{i=1}^N \text{Loss}(y_i, f(x_i; \theta)) + \frac{\lambda}{2} \|\theta\|^2$$

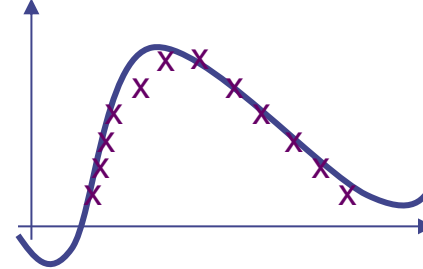
$$\theta^* = (X^T X + \lambda N I)^{-1} X^T y$$

# Regression

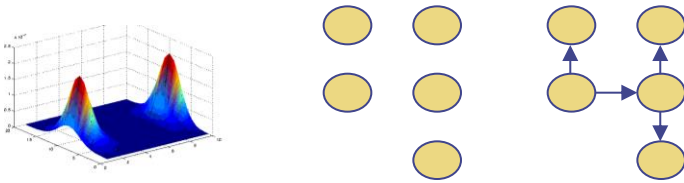
## Classification



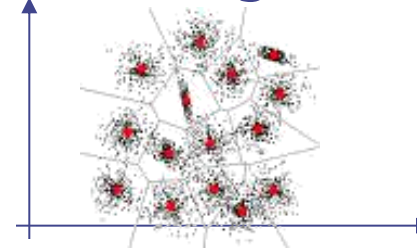
## Regression, $f(x)=y$



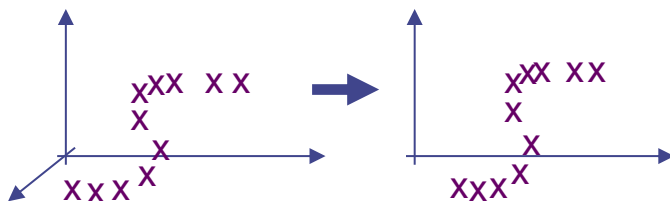
## Density/Structure Estimation



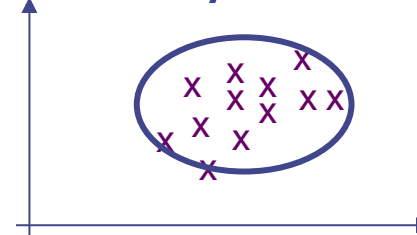
## Clustering



## Feature Selection



## Anomaly Detection

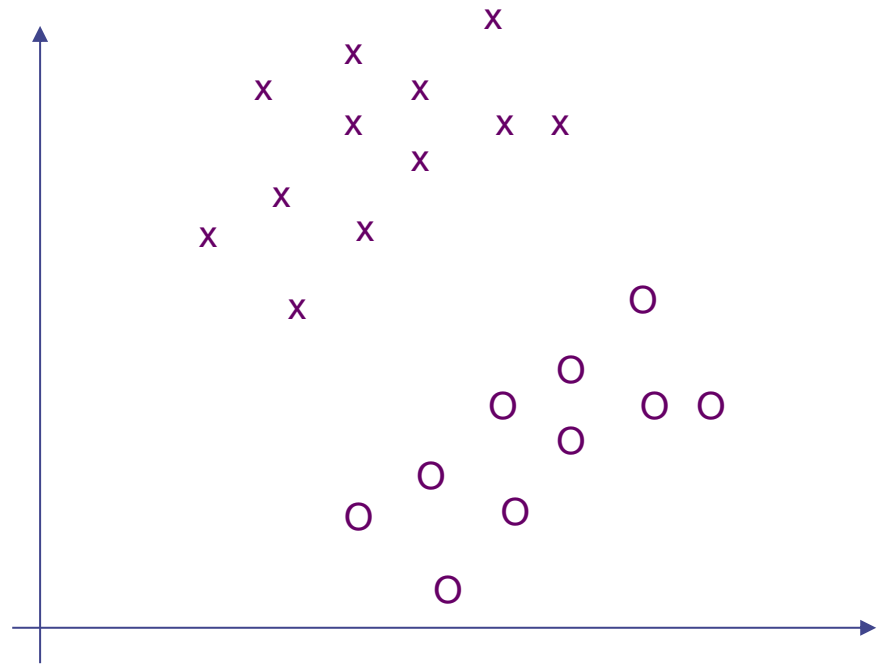


Supervised

Unsupervised

# Class 6

- **Classification**
- Logistic Regression
- Gradient Descent



# Classification Problems

- ◆ Determine student admission to Columbia based on GPA, prev. school rank, tests



# Classification Problems

- ◆ Determine student admission to Columbia based on GPA, prev. school rank, tests
- ◆ Decide malignant or benign tumors based on size, density, speed of growth

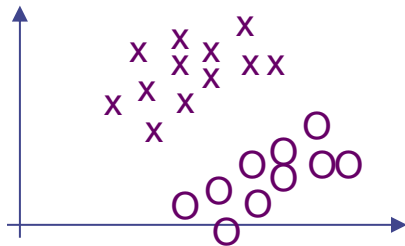


Kim et al. '07

# Formalizing Classification

- Classification is another important learning problem

Classification:  $X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_i \in \mathbf{R}^D, y_i \in \{0, 1\}$

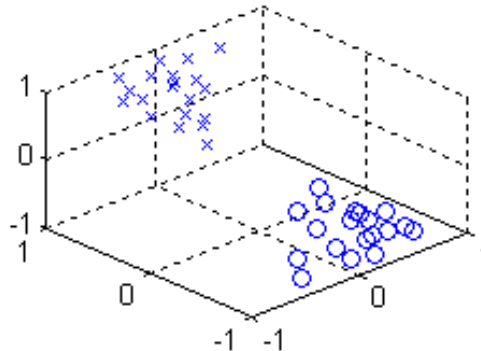
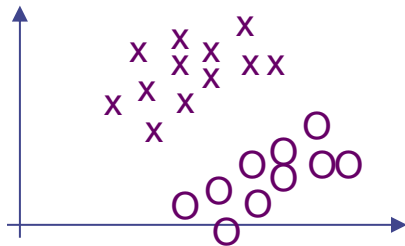


# Classification is like Regression

- Classification is another important learning problem

Classification:  $X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_i \in \mathbf{R}^D, y_i \in \{0, 1\}$

Regression:  $X = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}, \mathbf{x}_i \in \mathbf{R}^D, t_i \in \mathbf{R}$





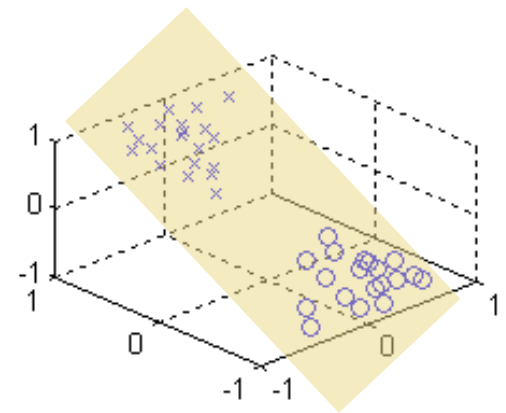
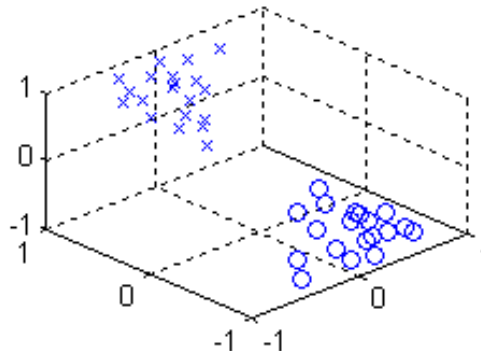
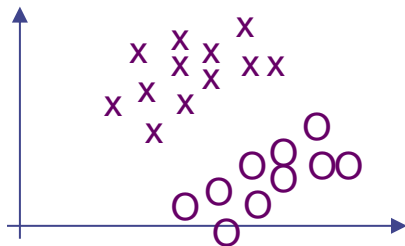
# Classification is like Regression

- Classification is another important learning problem

Classification:  $X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_i \in \mathbf{R}^D, y_i \in \{0, 1\}$

Regression:  $X = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}, \mathbf{x}_i \in \mathbf{R}^D, t_i \in \mathbf{R}$

- Should we solve this as a least squares regression problem?



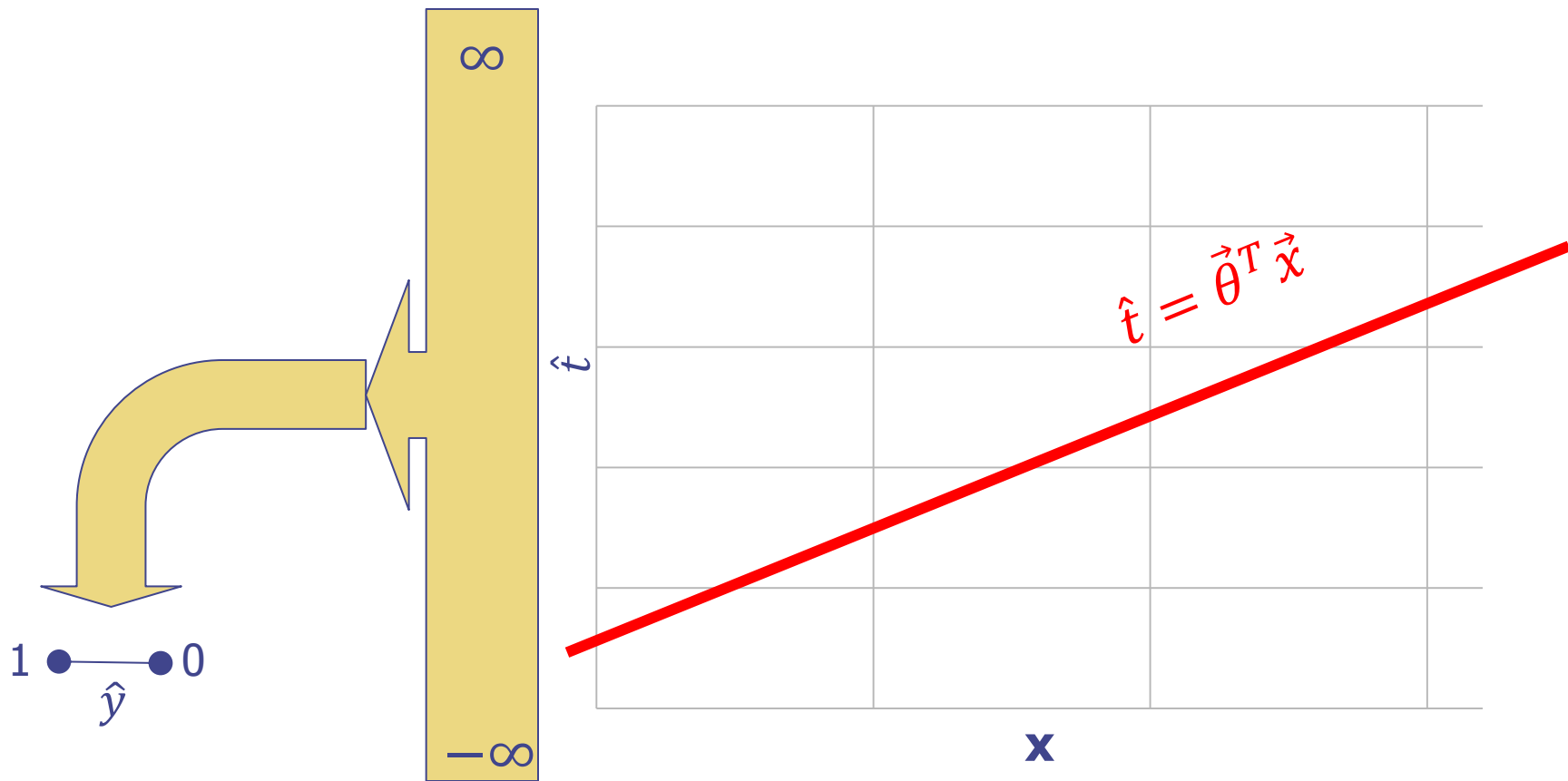
# Short hand for Linear Functions

- Hiding the intercept by notation

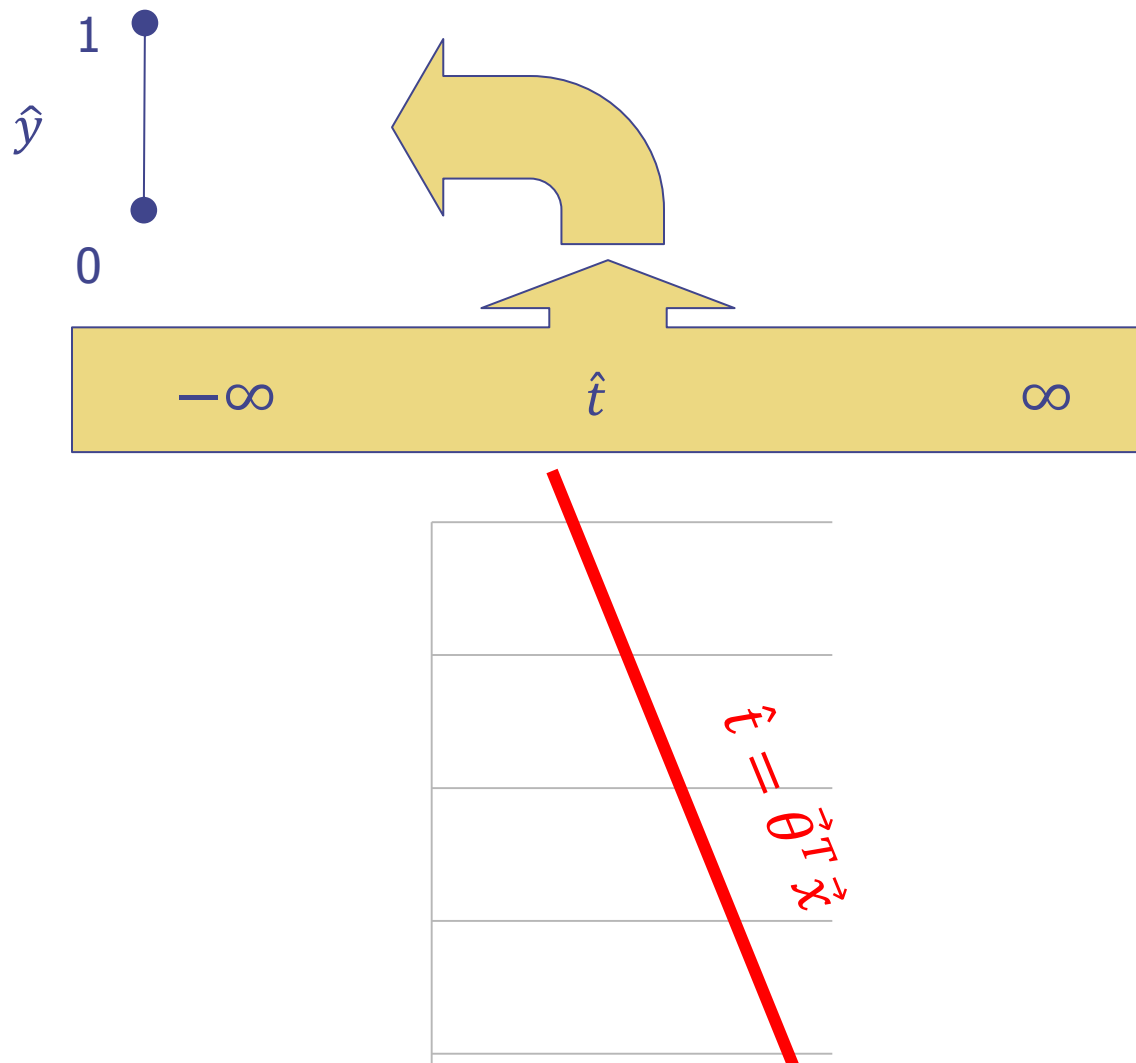
$$f(\mathbf{x}; \theta) = \theta^T \mathbf{x} + \theta_0$$

$$= \begin{bmatrix} \theta(1) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix} \begin{bmatrix} \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} + \theta_0 = \begin{bmatrix} \theta_0 \\ \theta(1) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix} \begin{bmatrix} 1 \\ \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} = \vec{\theta}^T \vec{\mathbf{x}}$$

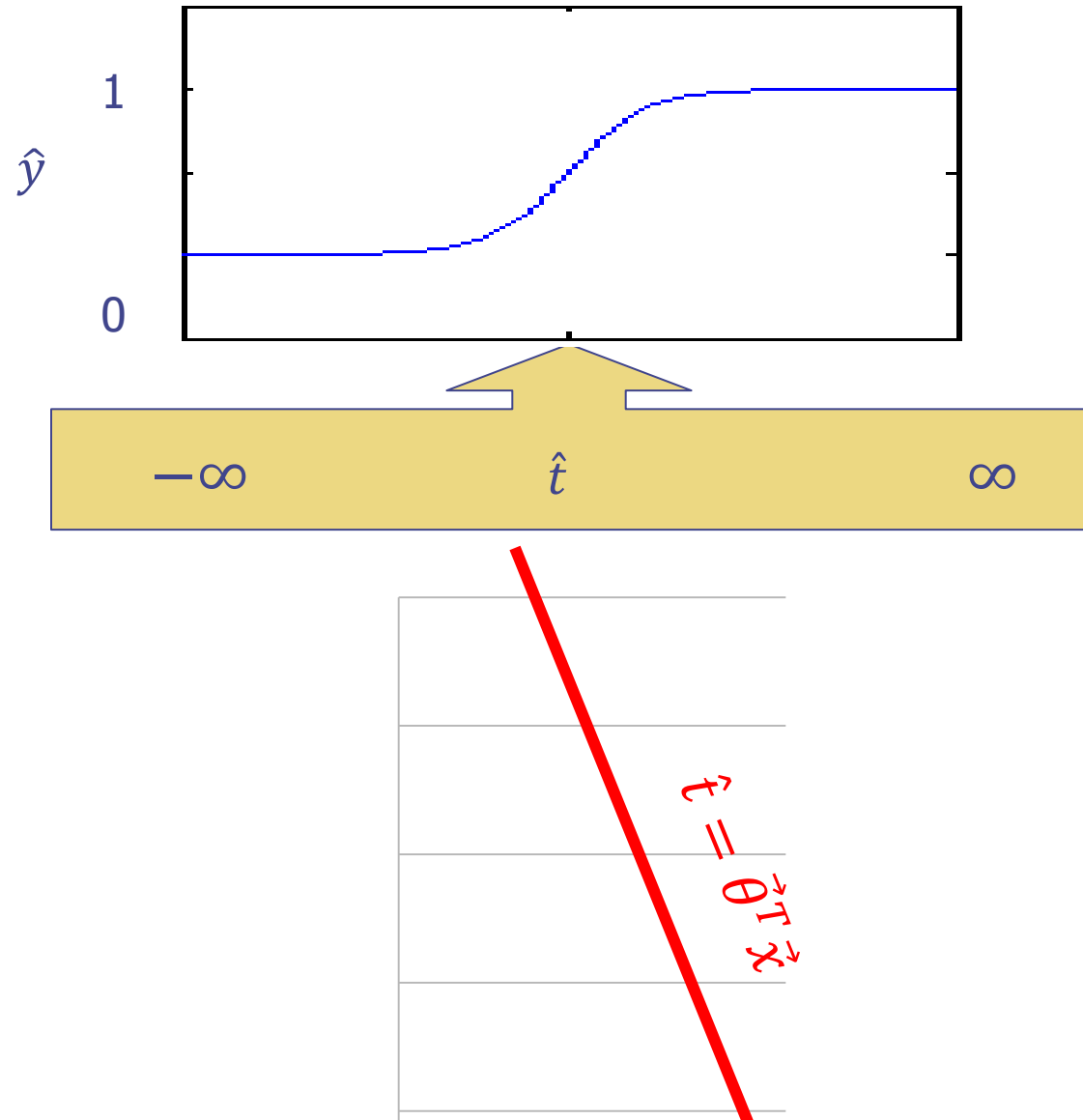
# Recalibrating Range



# Recalibrating Range



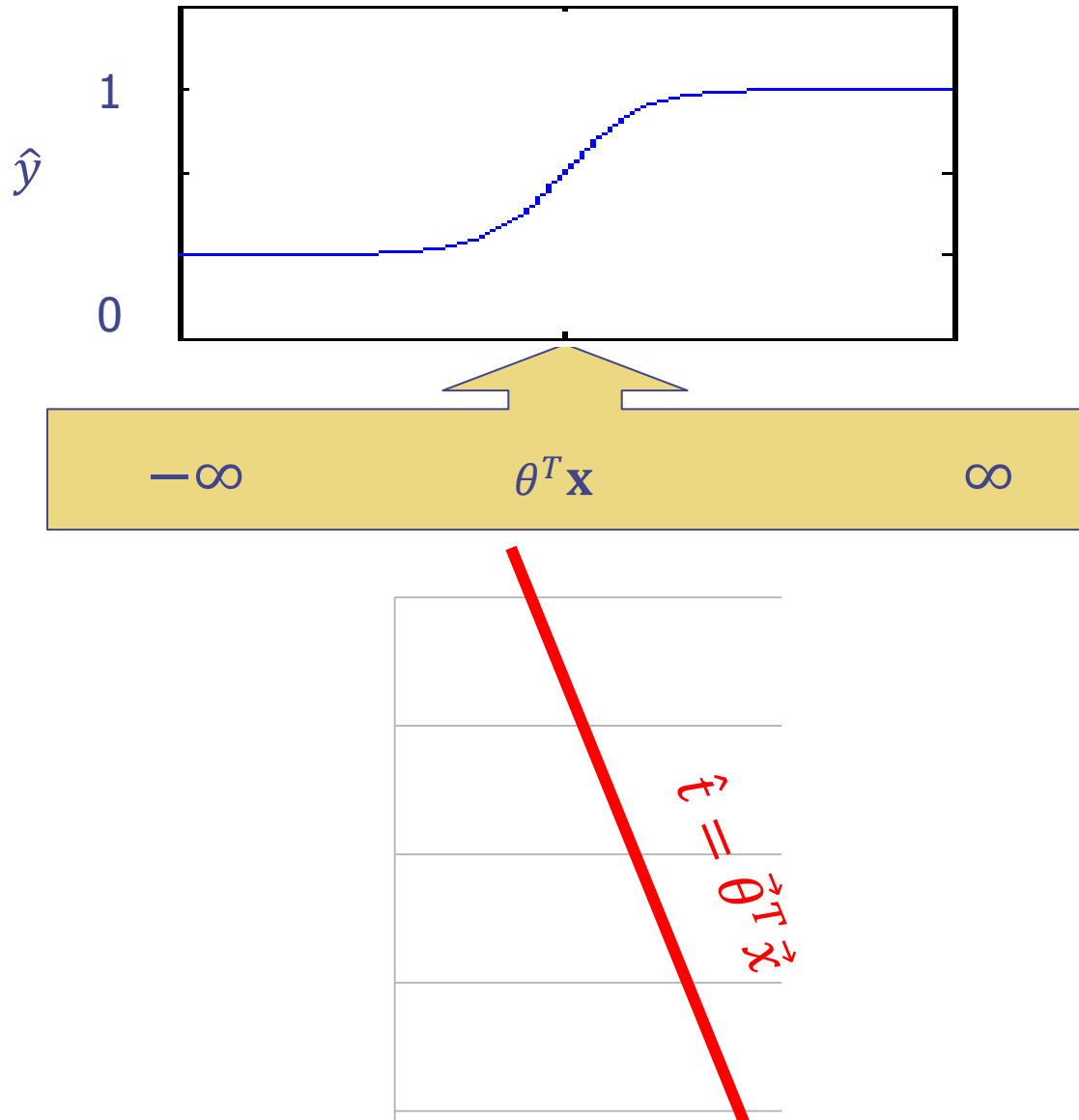
# Recalibrating Range



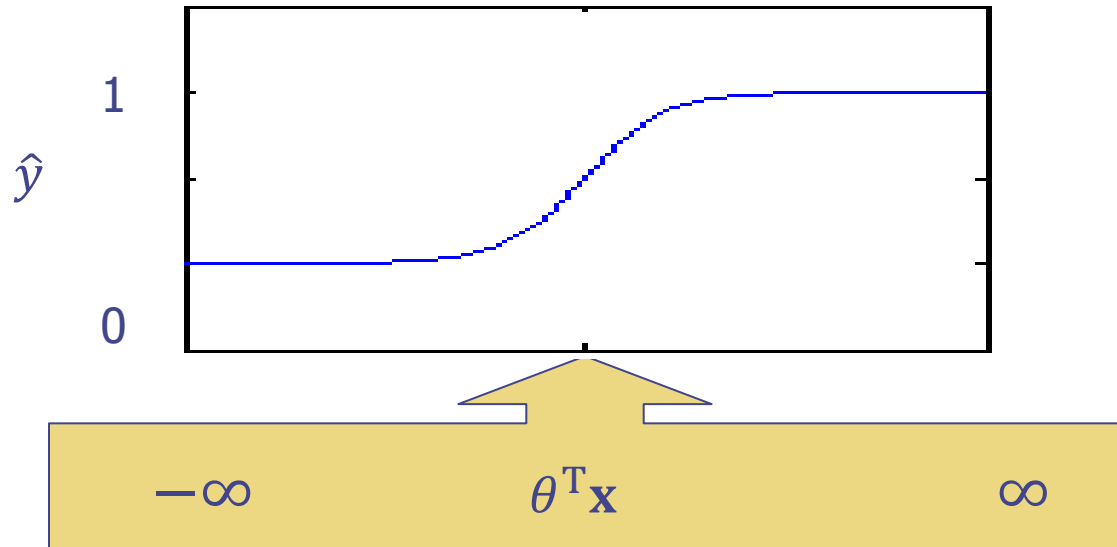
# Recalibrating Range

$$\lim_{\theta^T \mathbf{x} \rightarrow \infty} \hat{y} = 1$$

$$\lim_{\theta^T \mathbf{x} \rightarrow -\infty} \hat{y} = 0$$



# Recalibrating Range



$$\lim_{\theta^T \mathbf{x} \rightarrow \infty} \hat{y} = 1$$

$$\hat{y} \underset{\theta^T \mathbf{x} \rightarrow \infty}{\cong} 1 - \exp(-\theta^T \mathbf{x})$$

$$\lim_{\theta^T \mathbf{x} \rightarrow -\infty} \hat{y} = 0$$

$$\hat{y} \underset{\theta^T \mathbf{x} \rightarrow -\infty}{\cong} \exp(\theta^T \mathbf{x})$$



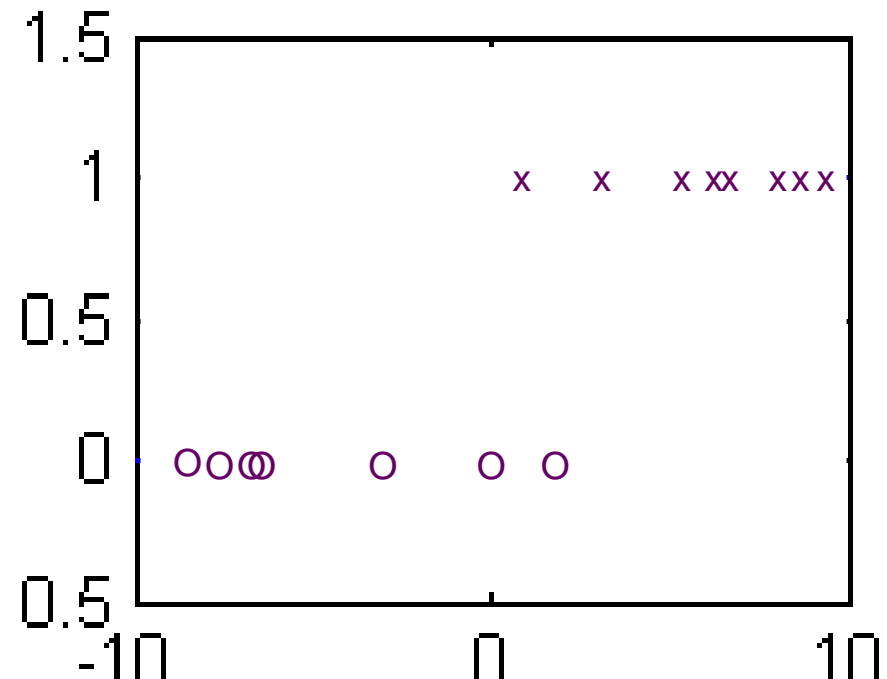
# Logistic Regression

- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_i \in \mathbf{R}^D, y_i \in \{0, 1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$





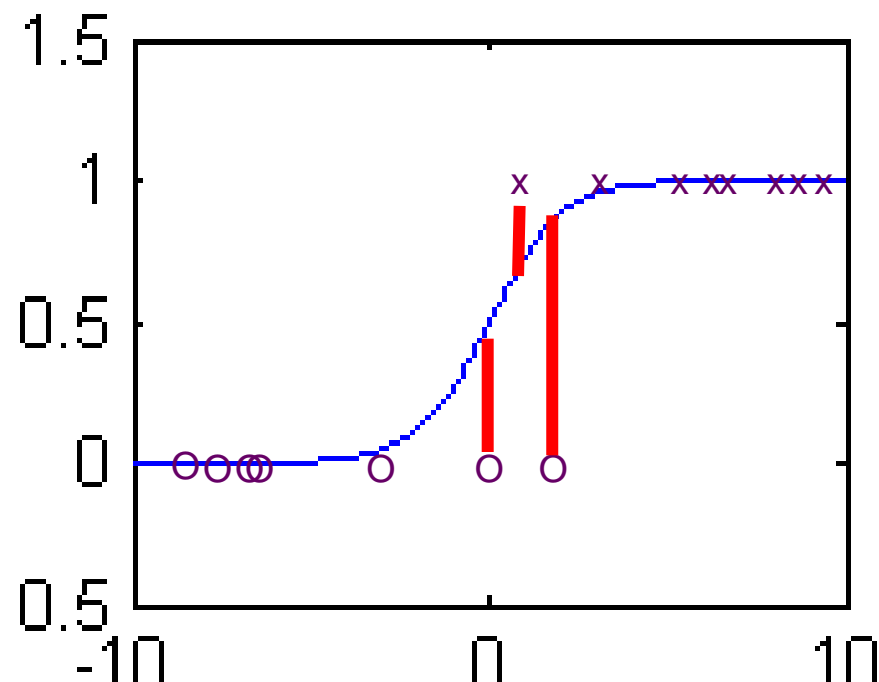
# Logistic Regression

- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0, 1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$



# Logistic Regression

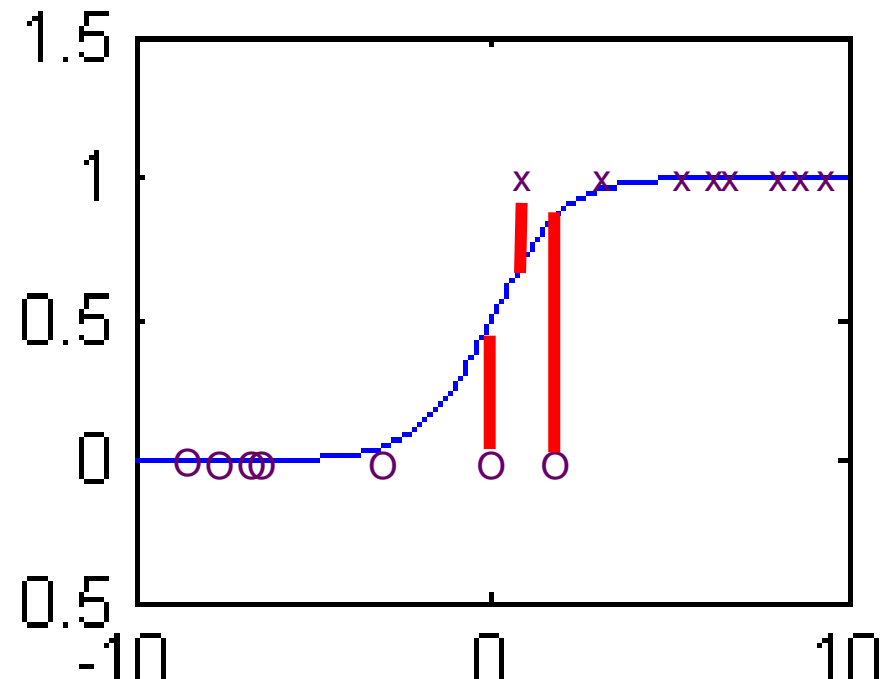
- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0,1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

- Assume  $\Pr(y = 1) = f(\mathbf{x}; \theta)$



# Logistic Regression

- Given a classification problem with binary outputs

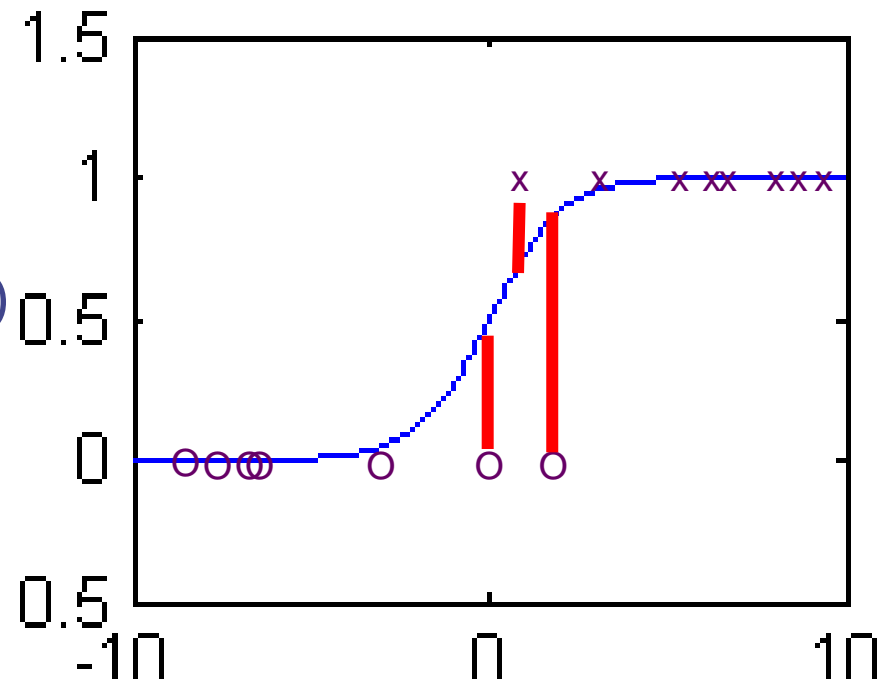
$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0, 1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

- Assume  $\Pr(y_i = 1) = f(\mathbf{x}_i; \theta)$

- $\Pr(y | f(\mathbf{x}; \theta)) =$   
 $= \prod_{y_i=0} (1 - f(\mathbf{x}_i; \theta)) \prod_{y_i=1} f(\mathbf{x}_i; \theta)$



# Logistic Regression

- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0,1\}$$

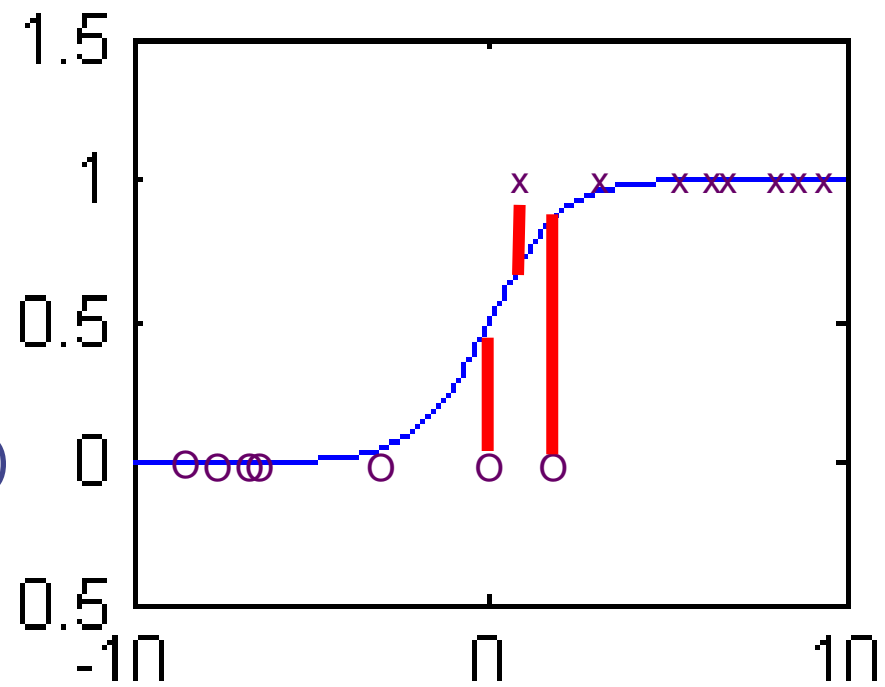
- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

- Assume  $\Pr(y_i = 1) = f(\mathbf{x}_i; \theta)$

$$\Pr(y|\mathbf{x}; \theta) = \prod_{y_i=0} (1 - f(\mathbf{x}_i; \theta)) \prod_{y_i=1} f(\mathbf{x}_i; \theta)$$

$$\log \Pr(y|\mathbf{x}; \theta) = \sum_i [y_i \log f(\mathbf{x}_i; \theta) + (1 - y_i) \log(1 - f(\mathbf{x}_i; \theta))]$$



# Logistic Regression

- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0, 1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

- Instead of squared loss, use **Logistic Loss** (i.e. negative binomial likelihood)

$$Loss_{log}(y, f(\mathbf{x}; \theta)) = (y - 1) \log(1 - f(\mathbf{x}; \theta)) - y \log(f(\mathbf{x}; \theta))$$

- The resulting method is called **Logistic Regression**.
- Empirical Risk:

# Logistic Regression

- Given a classification problem with binary outputs

$$X = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x} \in \mathbf{R}^D, y \in \{0, 1\}$$

- Use this function and output 1 if  $f(\mathbf{x}) > 0.5$  and 0 otherwise

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

- Instead of squared loss, use **Logistic Loss** (i.e. negative binomial likelihood)

$$Loss_{log}(y, f(\mathbf{x}; \theta)) = (y - 1) \log(1 - f(\mathbf{x}; \theta)) - y \log(f(\mathbf{x}; \theta))$$

- The resulting method is called **Logistic Regression**.
- Empirical Risk:

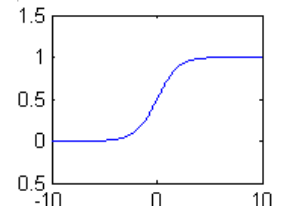
$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N [(y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))]$$

# Logistic Regression

- With empirical logistic risk has no closed form solution:

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$



# Logistic Regression

- With empirical logistic risk has no closed form solution:

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

$$\nabla_{\theta} R = \frac{1}{N} \sum_{i=1}^N \left( \frac{1 - y_i}{1 - f(\mathbf{x}_i; \theta)} - \frac{y_i}{f(\mathbf{x}_i; \theta)} \right) f'(\mathbf{x}_i; \theta) = 0 \quad \text{??????}$$

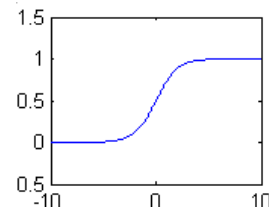
where

$$f(\mathbf{x}; \theta) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})} = g(\theta^T \mathbf{x})$$

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

$$g'(z) = g(z)(1 - g(z))$$

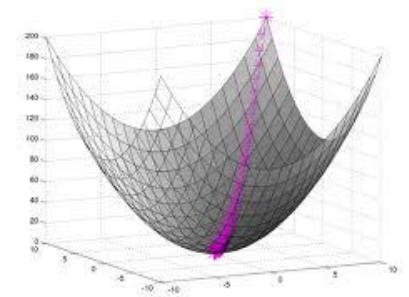
**Find best  $\theta$  numerically!**





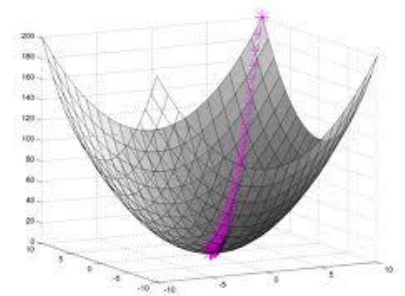
# Gradient Descent

- Useful when we can't get minimum solution in closed form
- Gradient points in direction of fastest increase
- Take step in the opposite direction!



# Gradient Descent

- Useful when we can't get minimum solution in closed form
- Gradient points in direction of fastest increase
- Take step in the opposite direction!



- Gradient Descent Algorithm

*choose scalar step size  $\eta$ , & tolerance  $\epsilon$*   
*initialize  $\theta^0 = \text{small random vector}$*

$$\theta^1 \leftarrow \theta^0 - \eta_0 \nabla_{\theta} R_{emp} |_{\theta^0} ; t \leftarrow 1$$

*while  $\|\theta^t - \theta^{t-1}\| \geq \epsilon$  {*

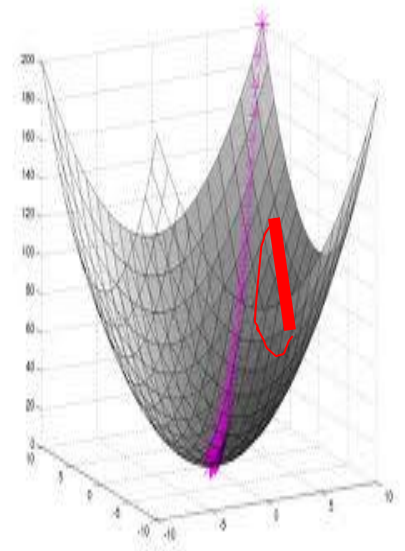
$$\theta^{t+1} \leftarrow \theta^t - \eta_t \nabla_{\theta} R_{emp} |_{\theta^t} ; t \leftarrow t + 1 \}$$

- For appropriate  $\{\eta_t\}$ , this will converge to local minimum

# Gradient Descent Convergence

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

is a convex function, so local minimum is global



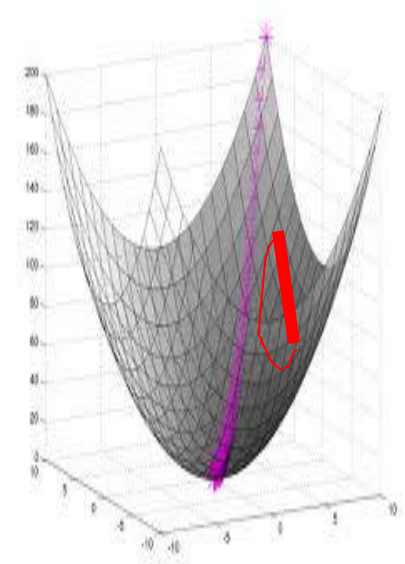
# Gradient Descent Convergence

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

is a convex function, so local minimum is global  
Proof:

Convex combination of  
convex functions

$-\log(1 - f(\mathbf{x}_i; \theta))$  and  $-\log(f(\mathbf{x}_i; \theta))$



# Convexity of Regression Loss

$-\log(f(\mathbf{x}_i; \theta))$  is a convex function

# Convexity of Regression Loss

$-\log(f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

$$\begin{aligned}
 \nabla_{\theta} [-\log(f(\mathbf{x}_i; \theta))] &= \\
 &= \nabla_{\theta} \left[ -\log \left( \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} \right) \right] \\
 &= \nabla_{\theta} [\log(1 + \exp(-\theta^T \mathbf{x}_i))] \\
 &= \frac{-\exp(-\theta^T \mathbf{x}_i) \mathbf{x}_i}{1 + \exp(-\theta^T \mathbf{x}_i)} \\
 &= \left( \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} - 1 \right) \mathbf{x}_i
 \end{aligned}$$

# Convexity of Regression Loss

$-\log(f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

$$\begin{aligned} \nabla_{\theta}^2 [-\log(f(\mathbf{x}_i; \theta))] \\ = \nabla_{\theta} \left[ \left( \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} - 1 \right) \mathbf{x}_i \right] = \end{aligned}$$

# Convexity of Regression Loss

$-\log(f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

$$\begin{aligned} & \nabla_{\theta}^2 [-\log(f(\mathbf{x}_i; \theta))] \\ &= \nabla_{\theta} \left[ \left( \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} - 1 \right) \mathbf{x}_i \right] \\ &= \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} \frac{\exp(-\theta^T \mathbf{x}_i)}{1 + \exp(-\theta^T \mathbf{x}_i)} \mathbf{x}_i \mathbf{x}_i^T \end{aligned}$$



# Convexity of Regression Loss

$-\log(f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

$$\begin{aligned} \mathbf{z}^T \nabla_{\theta}^2 [-\log(f(\mathbf{x}_i; \theta))] \mathbf{z} &= \\ \mathbf{z}^T \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} \frac{\exp(-\theta^T \mathbf{x}_i)}{1 + \exp(-\theta^T \mathbf{x}_i)} \mathbf{x}_i \mathbf{x}_i^T \mathbf{z} &= \\ \frac{1}{1 + \exp(-\theta^T \mathbf{x}_i)} \frac{\exp(-\theta^T \mathbf{x}_i)}{1 + \exp(-\theta^T \mathbf{x}_i)} (\mathbf{x}_i^T \mathbf{z})^2 \end{aligned}$$

# Convexity of Regression Loss

$-\log(1 - f(\mathbf{x}_i; \theta))$  is a convex function

# Convexity of Regression Loss

$-\log(1 - f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

$$\begin{aligned}\nabla_{\theta} [-\log(1 - f(\mathbf{x}_i; \theta))] &= \\ &= \nabla_{\theta} \left[ -\log \left( \frac{\exp(-\theta^T \mathbf{x}_i)}{1 + \exp(-\theta^T \mathbf{x}_i)} \right) \right] \\ &= \nabla_{\theta} [\theta^T \mathbf{x}_i + \log(1 + \exp(-\theta^T \mathbf{x}_i))] \\ &= \mathbf{x}_i + \nabla_{\theta} [\log(1 + \exp(-\theta^T \mathbf{x}_i))]\end{aligned}$$

# Convexity of Regression Loss

$-\log(1 - f(\mathbf{x}_i; \theta))$  is a convex function

Proof:

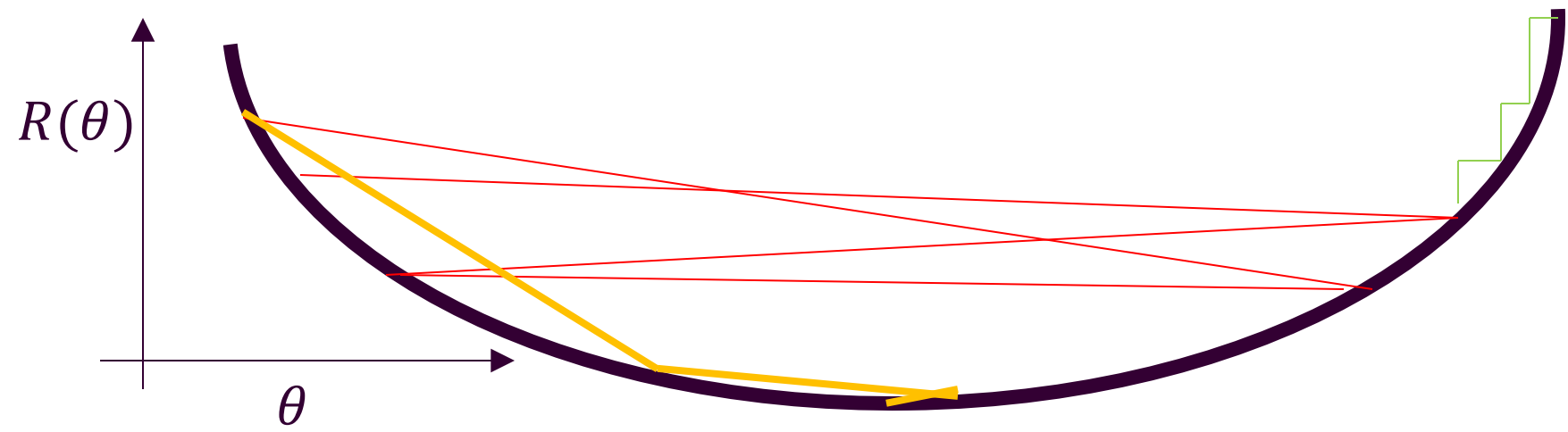
$$\begin{aligned}\nabla_{\theta}^2 [-\log(1 - f(\mathbf{x}_i; \theta))] &= \\ &= \nabla_{\theta}^2 [\log(1 + \exp(-\theta^T \mathbf{x}_i))] \end{aligned}$$

# How fast to descend?



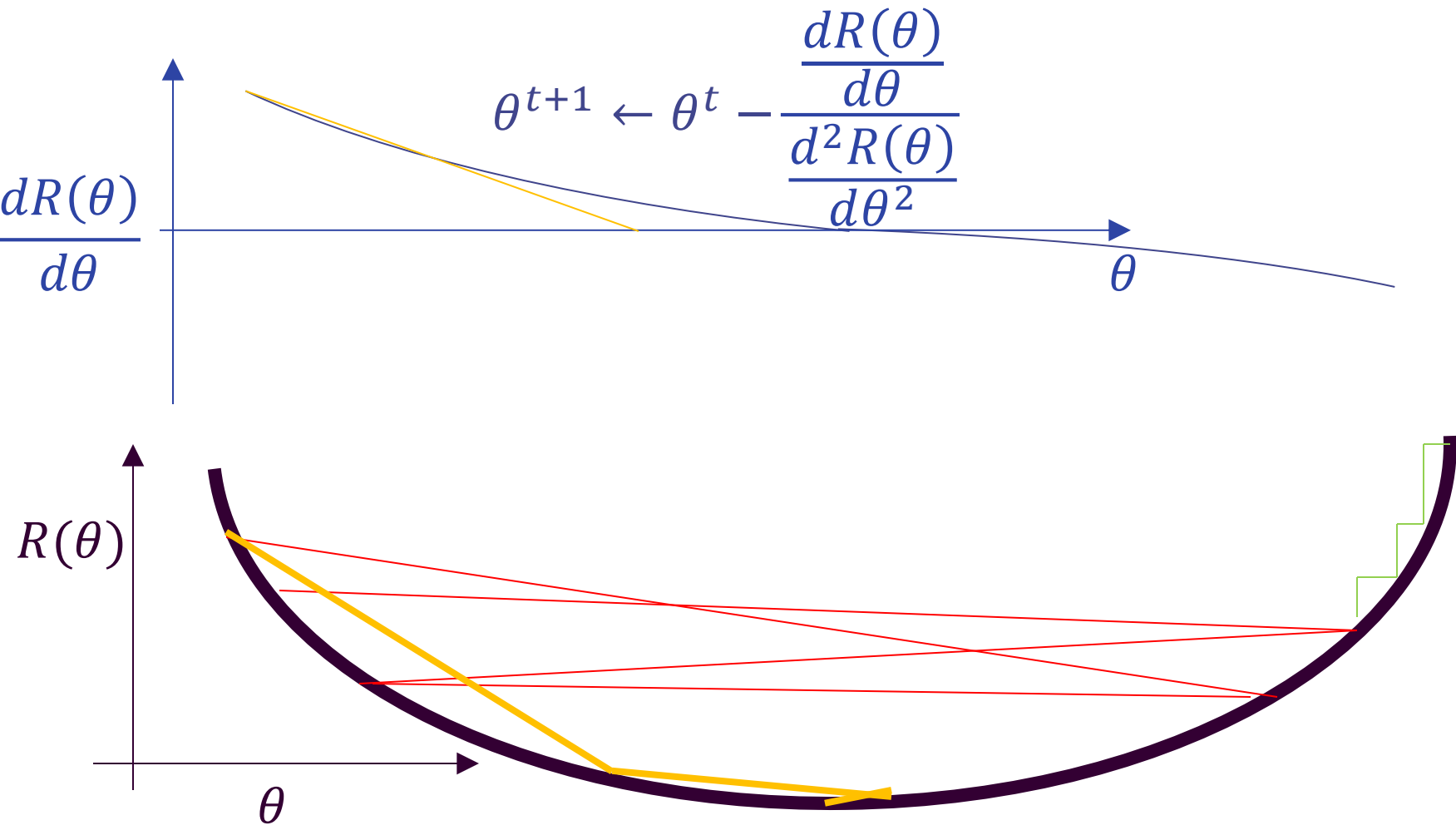
# How fast to descend?

◆ Not too **fast**, not too **slow**, just **right**



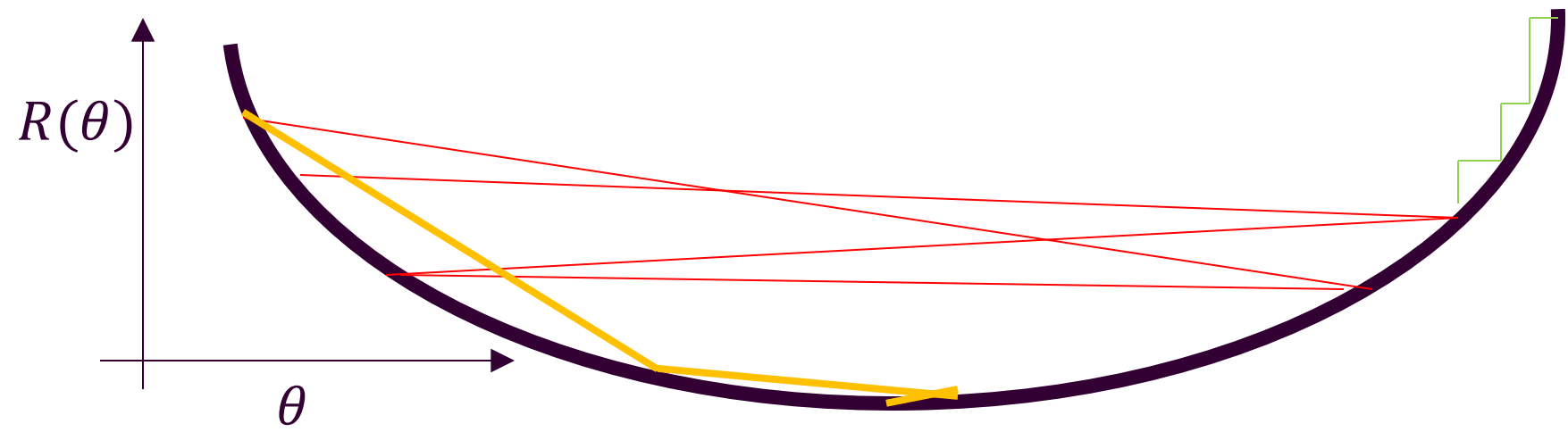
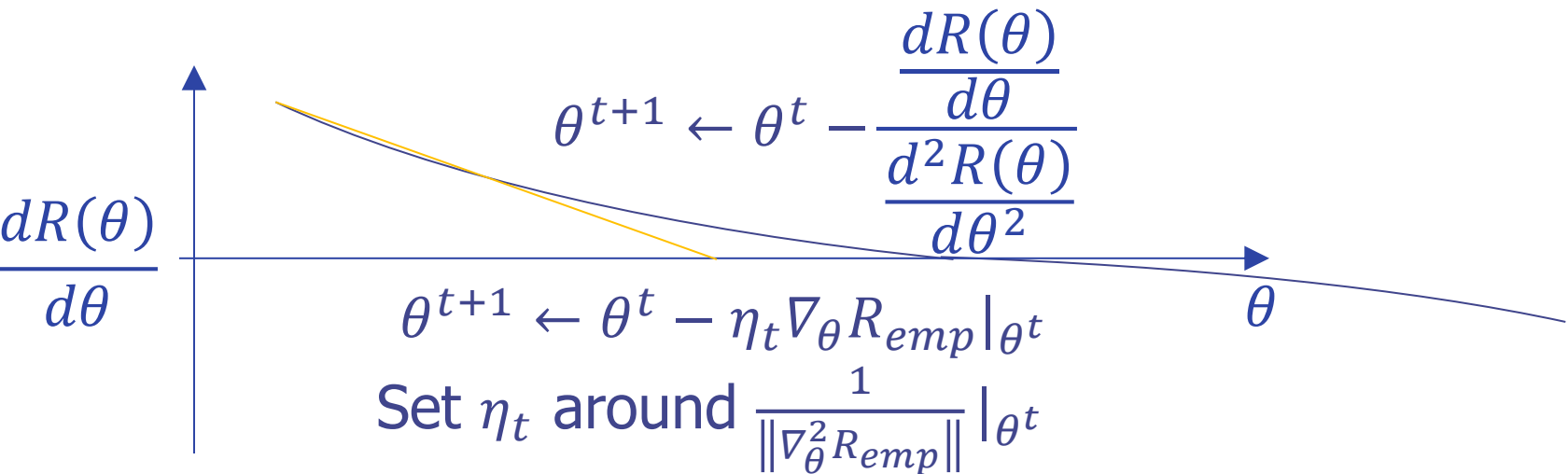
# How fast to descend?

◆ Newton's method for the derivative



# How fast to descend?

◆ Newton's method for the derivative





# Summary

- Classification
- Logistic Regression
- Gradient Descent