

**tai911s**

José Quenum  
PhD Computer Science



linear models

# principle

**from a dataset  $(X_l, Y_l)$  predict  $Y_J$  corresponding to input  $X_J$**

**recap...**

- step 0: **get, clean** and **understand** the dataset

- step 1:
  - decide which **function** to **learn** (or approximate)...
  - **select** the useful **features**

- step 2:
  - choose a training **algorithm**

- step 3:
  - **train, test and evaluate** the model

# tasks

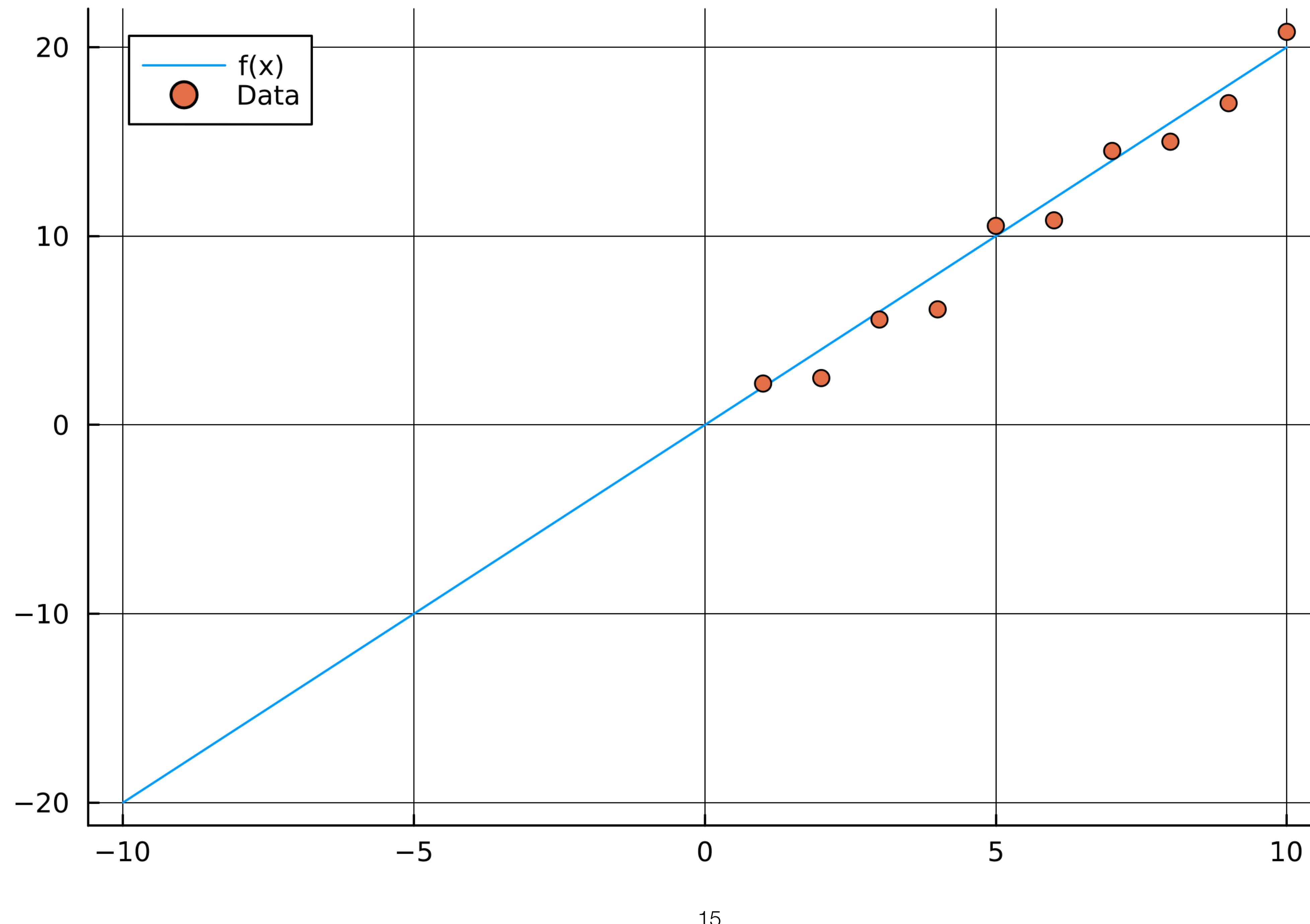
- **regression**

- **fit a continuous functional relation** between  $X_i$  and  $Y_i$
- e.g.,
  - weather forecasting, house price prediction

- **classification**
  - find out **decision boundaries** between the classes in a dataset
  - e.g.,
    - spam filter, hand-written character recognition



**linear regression**

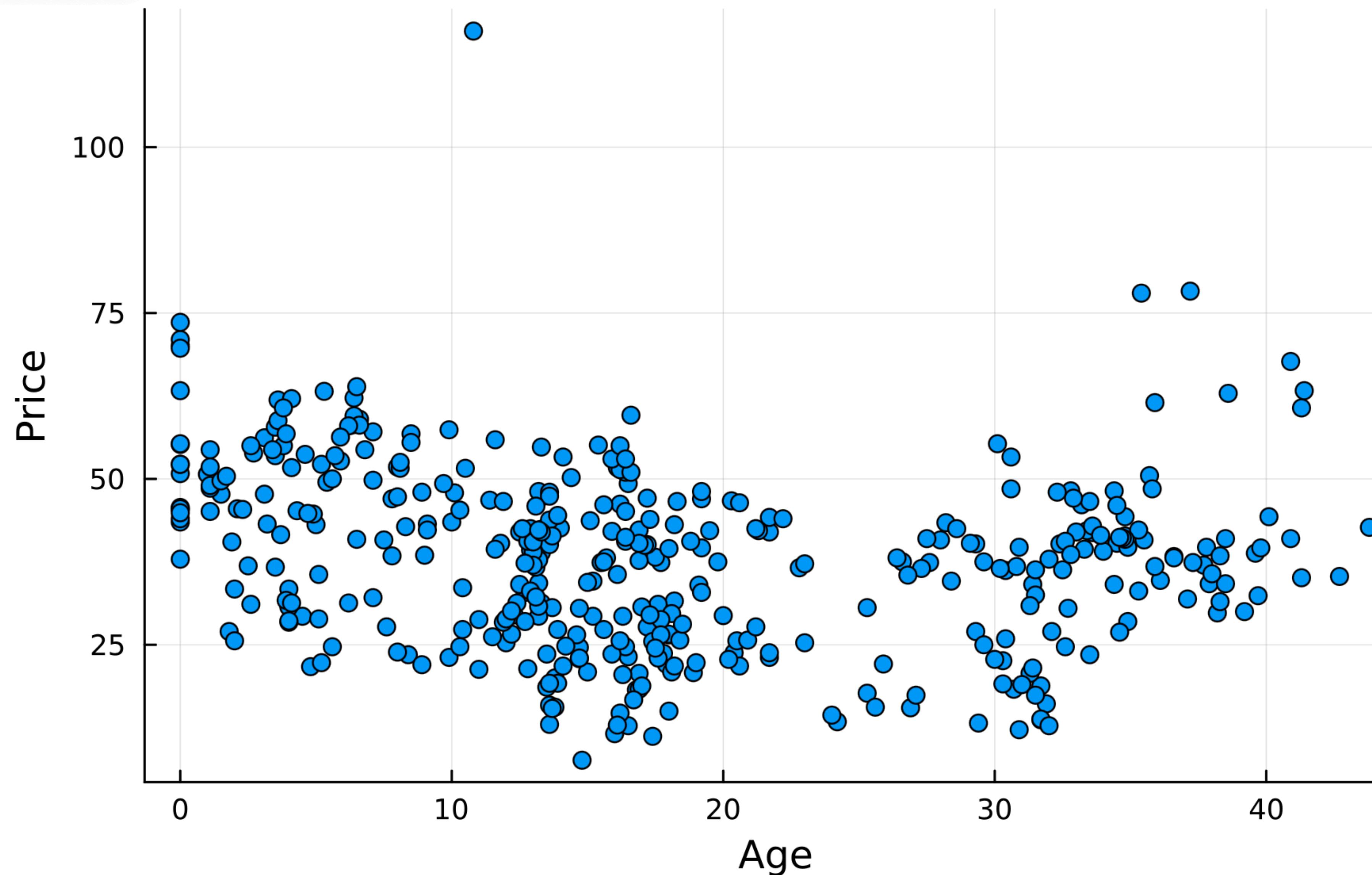


- univariate (simple) linear regression
- multi-variate (multiple) linear regression

transaction_date	house_age	distance_to_train_st	convenience_stores	latitude	longitude	house_price
2012.9.2	32	84.8788	10	24.983	121.54	37.9
2013.5.8	19.5	306.595	9	24.9803	121.54	42.2

Show/hide code

# Scatter Plot House Price vs House Age



$$h(X) = X \times \Theta^\top$$

$$h(X) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n$$

$$\Theta = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n]$$
$$X = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \cdots & x_1^n \\ 1 & x_2^1 & x_2^2 & \cdots & x_2^n \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_m^1 & x_m^2 & \cdots & x_m^n \end{bmatrix}$$

- residual difference between the hypothesis ( $h(X_i)$ ) and the actual value ( $Y_i$ )
- loss function
  - mean square error
  - mean absolute error

$$\mathcal{L}(h(X), Y) = \frac{1}{m} \sum_{i=1}^m (h(X^i) - Y^i)^2$$

$$\mathcal{L}(h(X), Y) = \frac{1}{m} \sum_{i=1}^m |h(X^i) - Y^i|$$

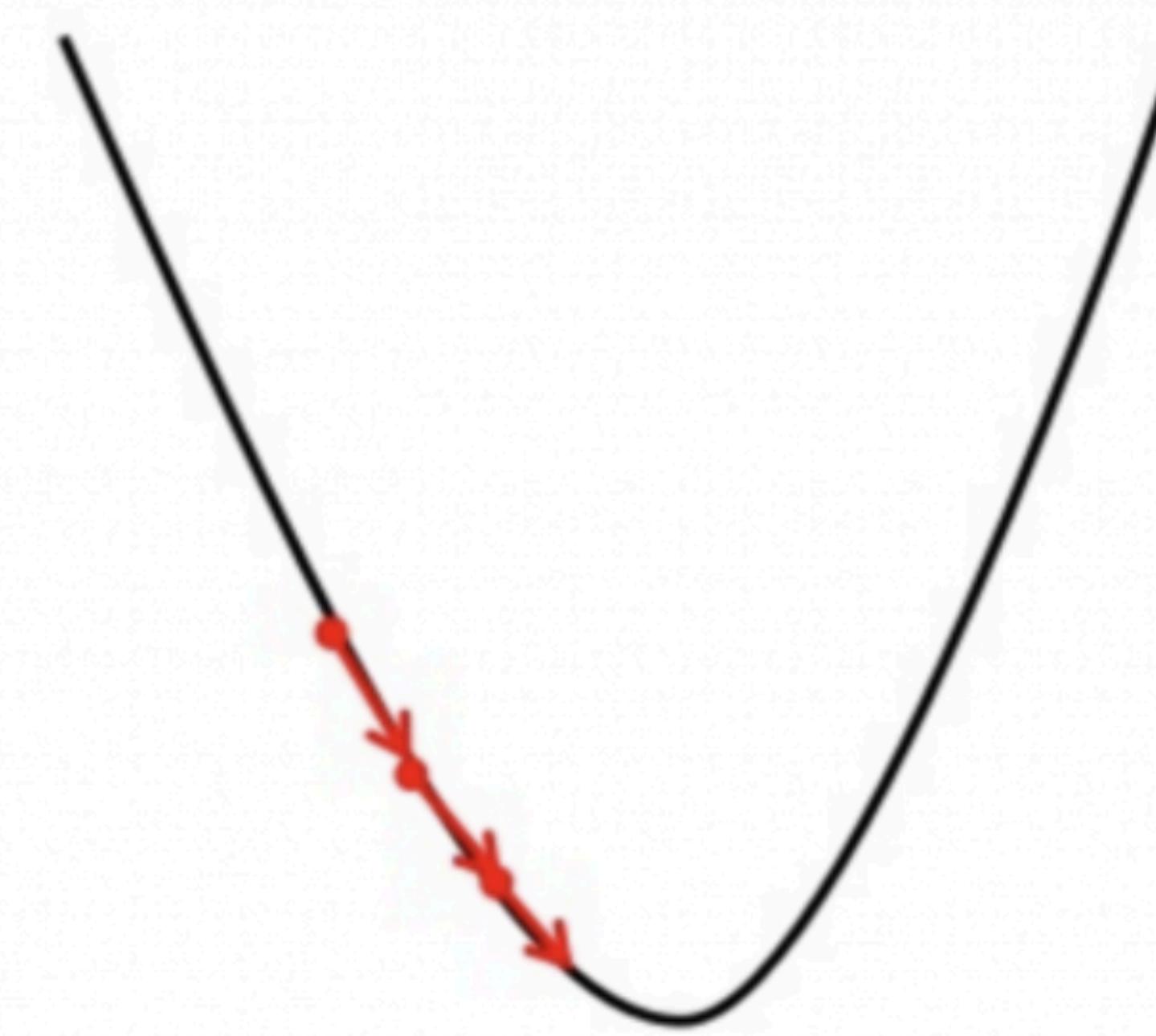
$$\operatorname{argmin} \mathcal{L}(h(X), Y) = \frac{1}{m} \sum_{i=1}^m (h(X^i) - Y^i)^2$$

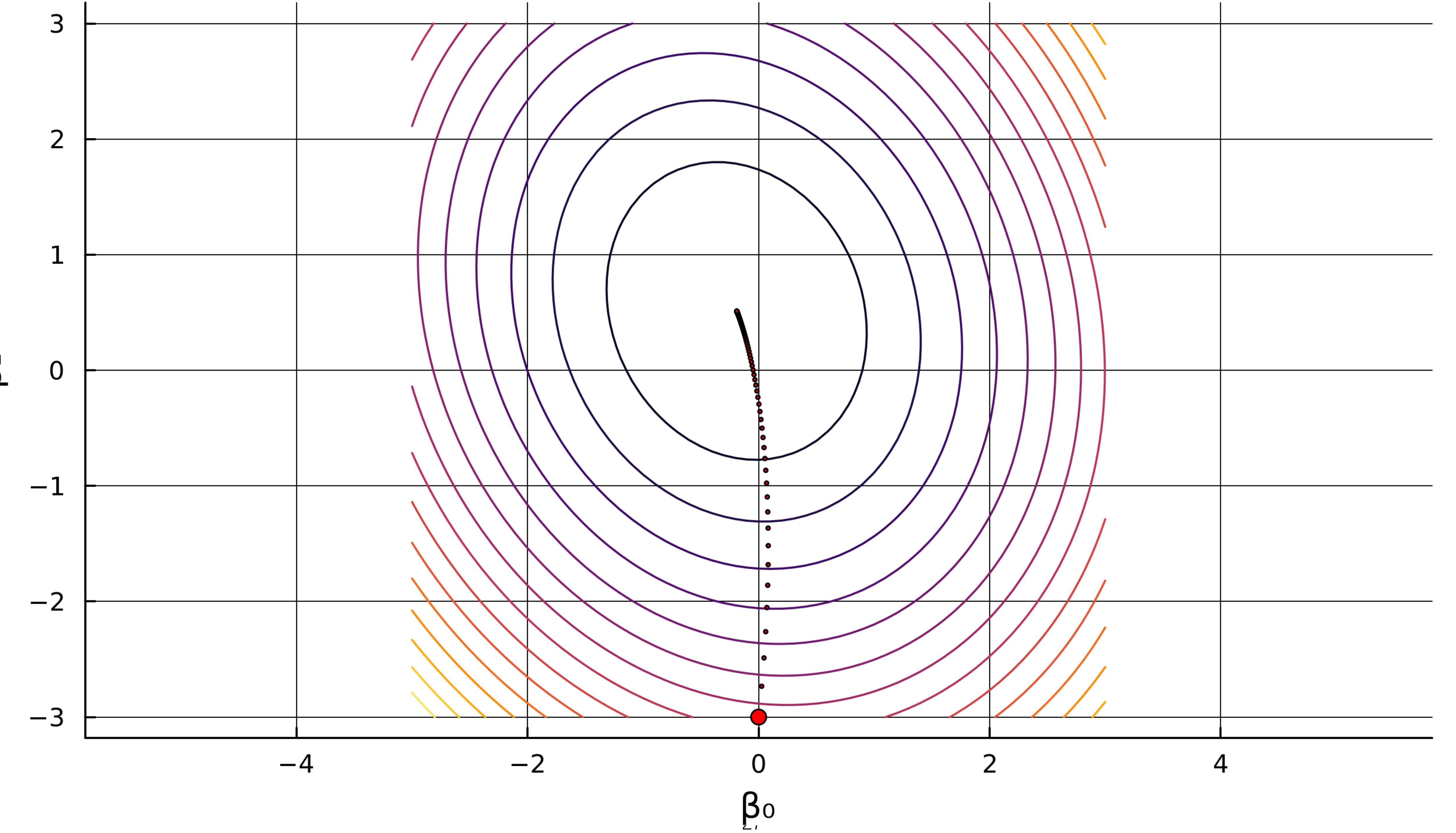
- vars in this optimisation fn?

$$\Theta = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n]$$

# optimisation method

- gradient descent
- stochastic gradient descent
- adam
- Etc.





$$\theta' = \theta - \alpha \nabla_{\theta} (\mathcal{L})$$

$$\theta'_j = \theta_j - \alpha \nabla_{\theta_j} (\mathcal{L})$$

$$\nabla_{\theta_j} (\mathcal{L}) = \frac{\partial \mathcal{L}_j}{\partial \theta_j}$$

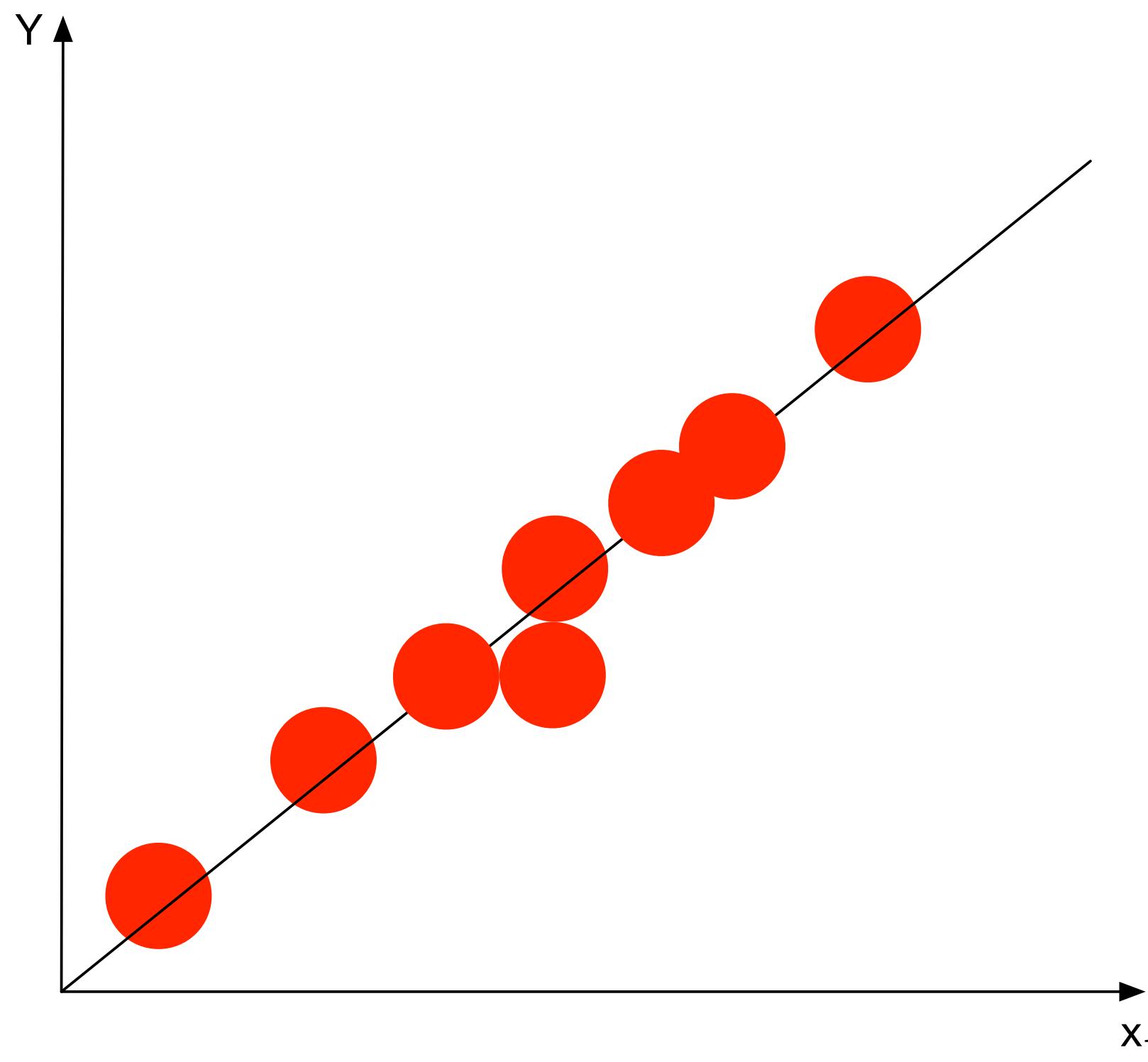
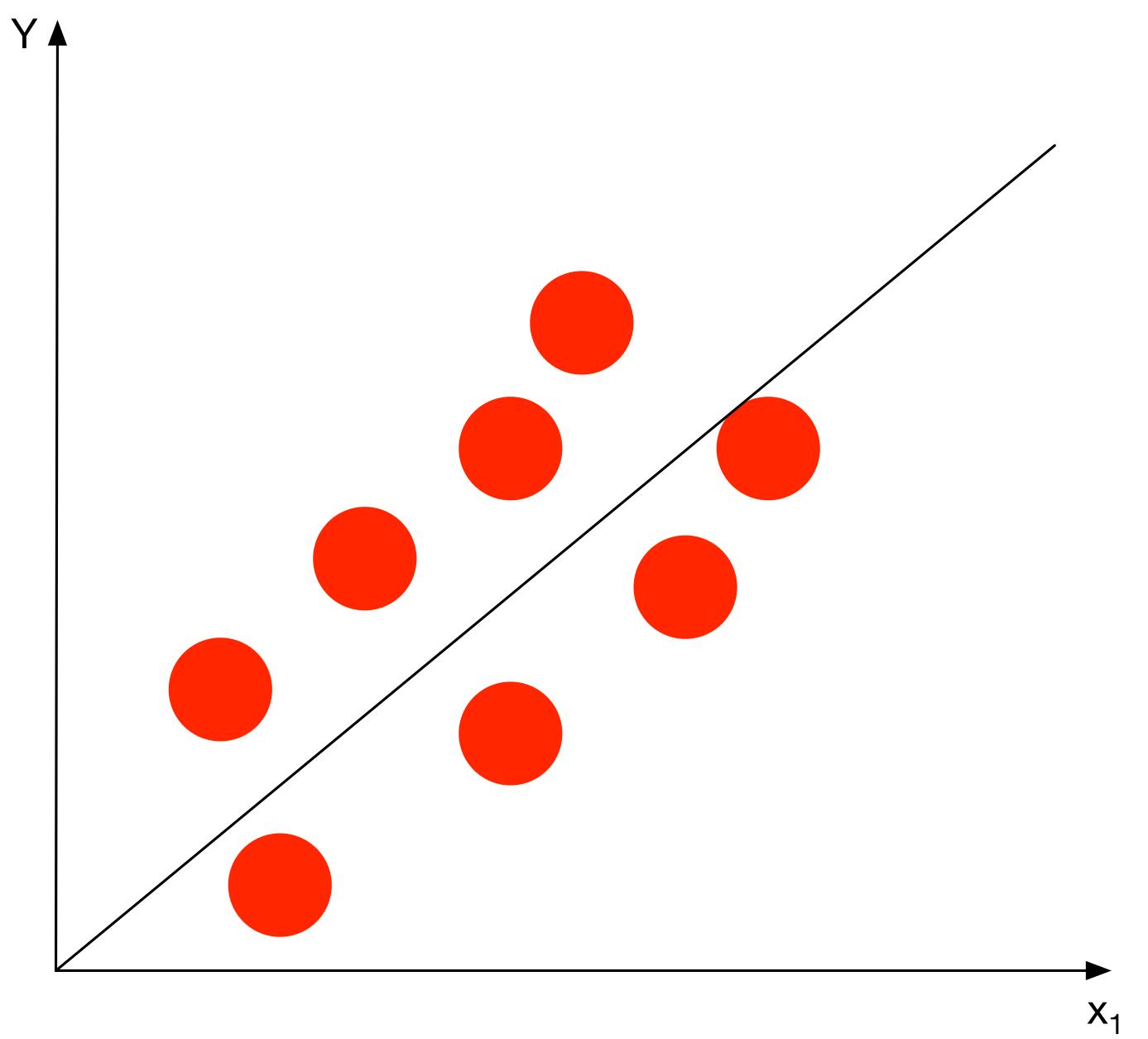
$$\theta_J = 0$$

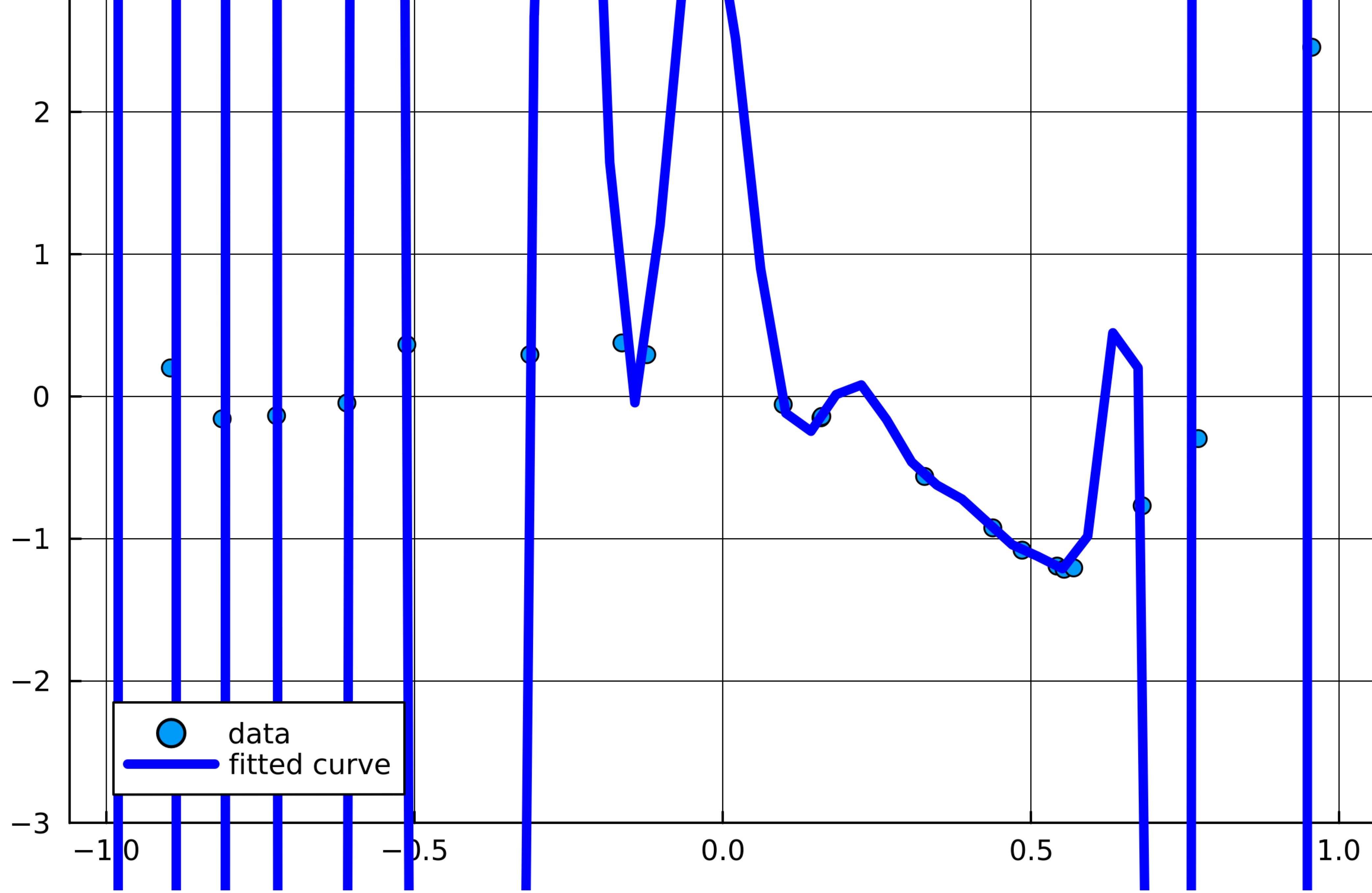
$$\theta'_J = \theta_J - \alpha \frac{1}{m} \sum_{i=1}^m (h(X^i) - Y^i) X_J^i$$

*for t = 1..T*

# but...

- overfitting
  - low bias and high variance
- underfitting
  - high bias and low variance





- regularisation
  - Lasso: absolute value of weight
  - Ridge: square weight

$$J(\Theta) = \frac{1}{2m} \sum_{i=1}^m (h(X^i) - Y^i)^2 + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2$$

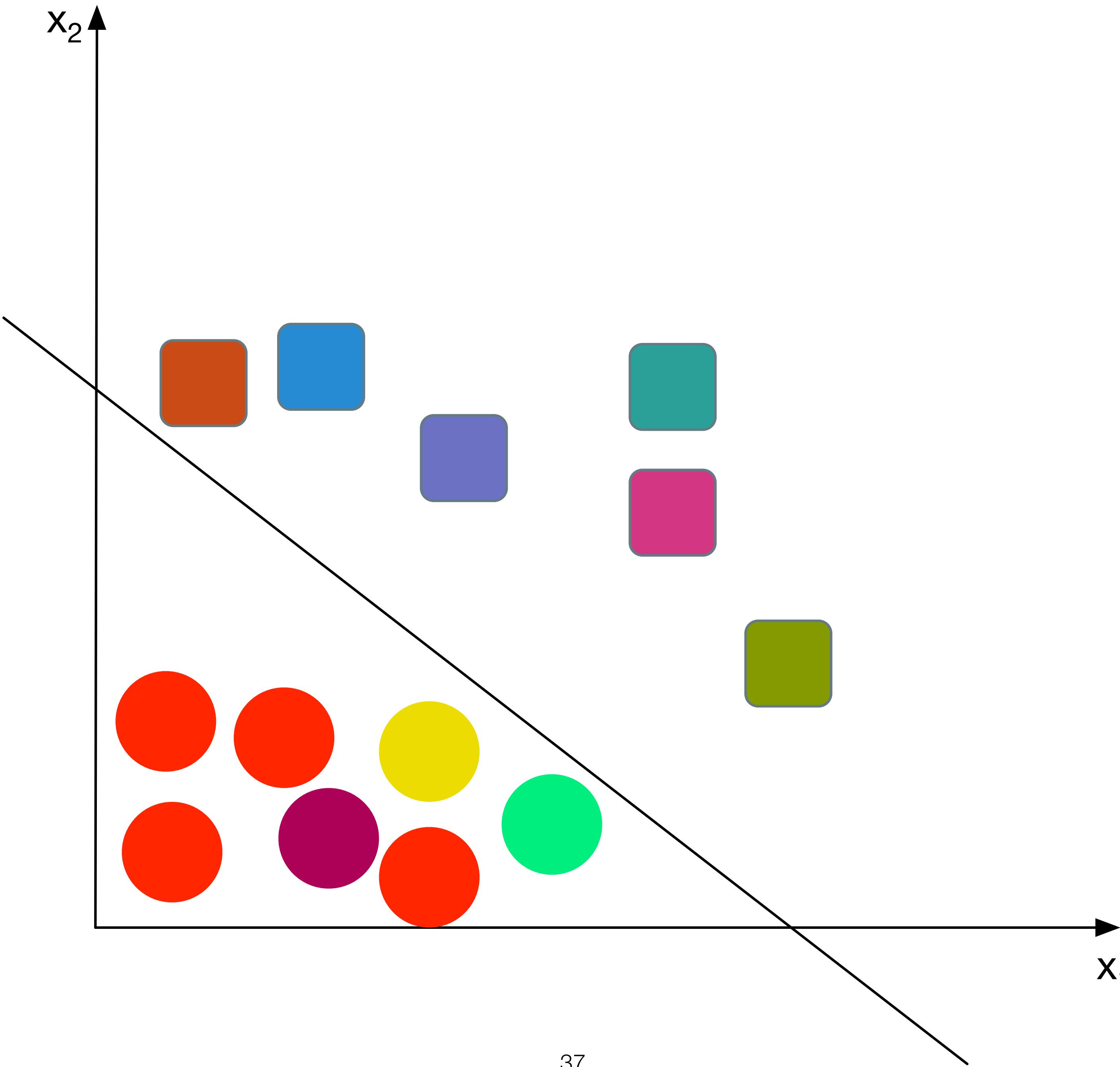
$$\theta'_j = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(X^i) - Y^i) X_j^i + \frac{\lambda}{n} \theta_j \right]$$

for  $j > 0$

$$\theta'_0 = \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h(X^i) - Y^i) X_0^i$$

- metrics
  - mean square error

# Logistic regression



- binomial classification
- multinomial classification

## **binomial** classification

$$y = \begin{cases} 0 \\ 1 \end{cases}$$

$$z = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n$$

$$z = X \times \Theta$$

$$\Theta = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n]$$

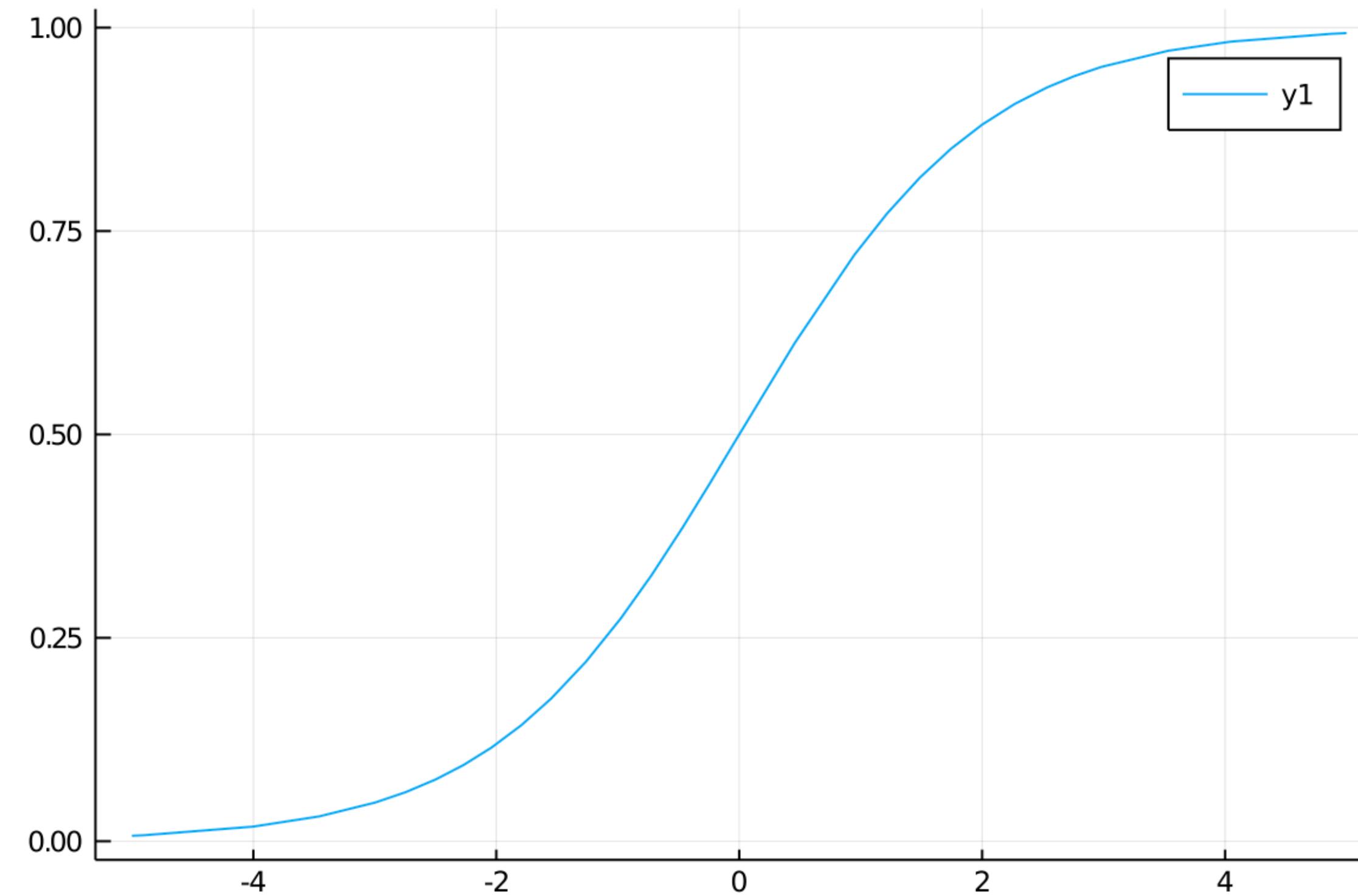
$$X = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \cdots & x_1^n \\ 1 & x_2^1 & x_2^2 & \cdots & x_2^n \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_m^1 & x_m^2 & \cdots & x_m^n \end{bmatrix}$$

- **sigmoid** function

$$h(x) = \frac{1}{1 + \exp^{-z(x)}}$$

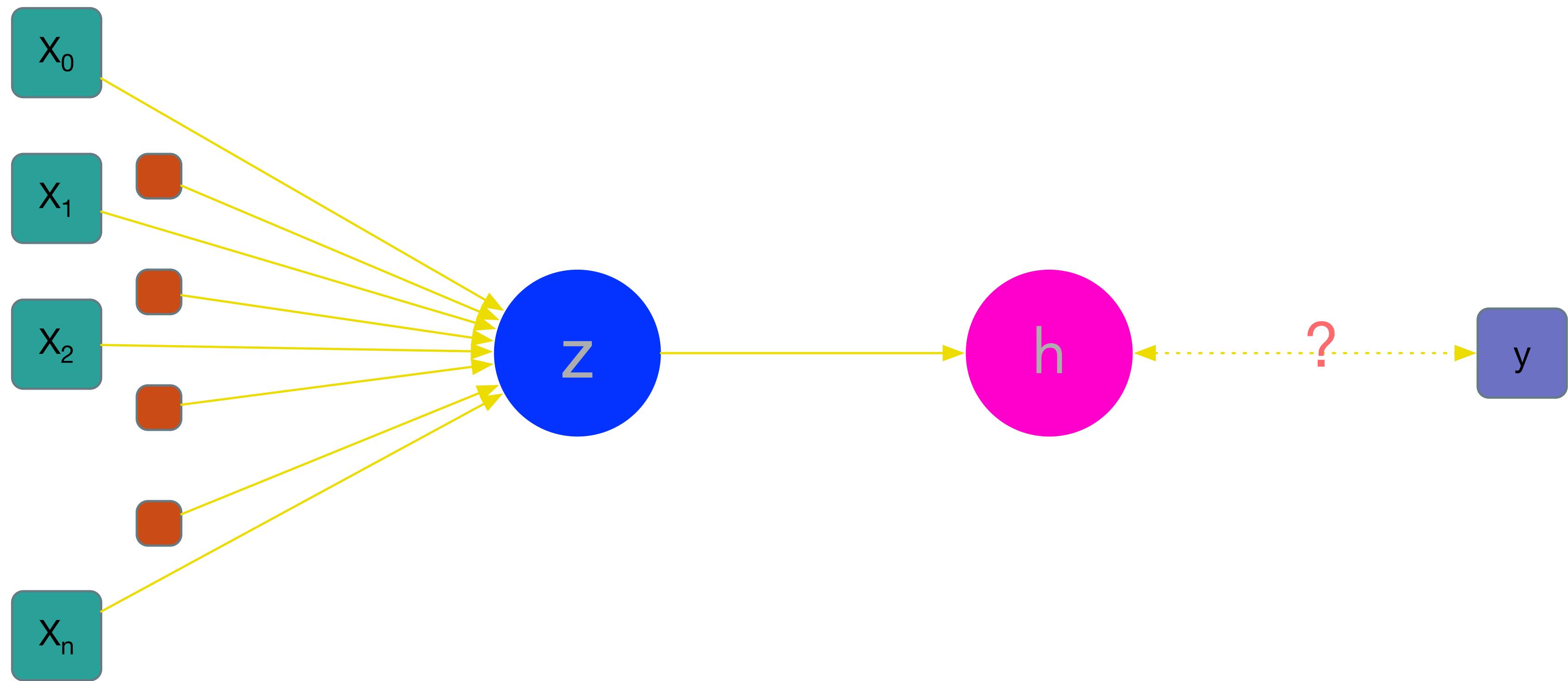
$$z = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n$$

$$h(x) = p(y=1|x; \theta)$$



$$h(x) = p(y = 1|x; \theta)$$

$$\begin{cases} h(x) \geq 0.5 \implies y = 1 \\ h(x) < 0.5 \implies y = 0 \end{cases}$$



$$\mathcal{L}(h(x), y) = \begin{cases} -\log(h(x)) & \text{if } y = 1 \\ -\log(1 - h(x)) & \text{if } y = 0 \end{cases}$$

$$\mathcal{L}(h(x), y) = -y \log(h(x)) - (1 - y) \log(1 - h(x))$$

$$J(\Theta) = \frac{1}{2m} \sum_{i=1}^m (-y^i \log(h(x^i)) - (1 - y^i) \log(1 - h(x^i))) + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2$$

$$argmin \quad \frac{1}{2m} \sum_{i=1}^m (-y^i \log(h(x^i)) - (1 - y^i) \log(1 - h(x^i))) + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2$$

$$\theta'_j = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(X^i) - Y^i) X_j^i + \frac{\lambda}{n} \theta_j \right]$$

for  $j > 0$

$$\theta'_0 = \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h(X^i) - Y^i) X_0^i$$

- metrics
- confusion matrix

	$h = 0$	$h = 1$
$y = 0$	true negative (TN)	false positive (FP)
$y = 1$	false negative (FN)	true positive (TP)

- accuracy
  - overall, how often is the classifier correct?
- precision
  - when it predicts “yes”, how often is it correct?

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TN} + \text{FP} + \text{TP} + \text{FN}}$$

$$\text{precision} = \frac{\text{TP}}{\text{FP} + \text{TP}}$$

- error rate
- overall, how often is it wrong?
- true positive rate (recall)
- when it is actually “yes”, how often does it predict “yes”?

**error rate** = 1 – accuracy

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- true negative rate (specificity)
  - when it is actually “no”, how often does it predict “no”
- and more...
  - e.g., F-score

$$\text{specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

# **multinomial classification**

- first strategy

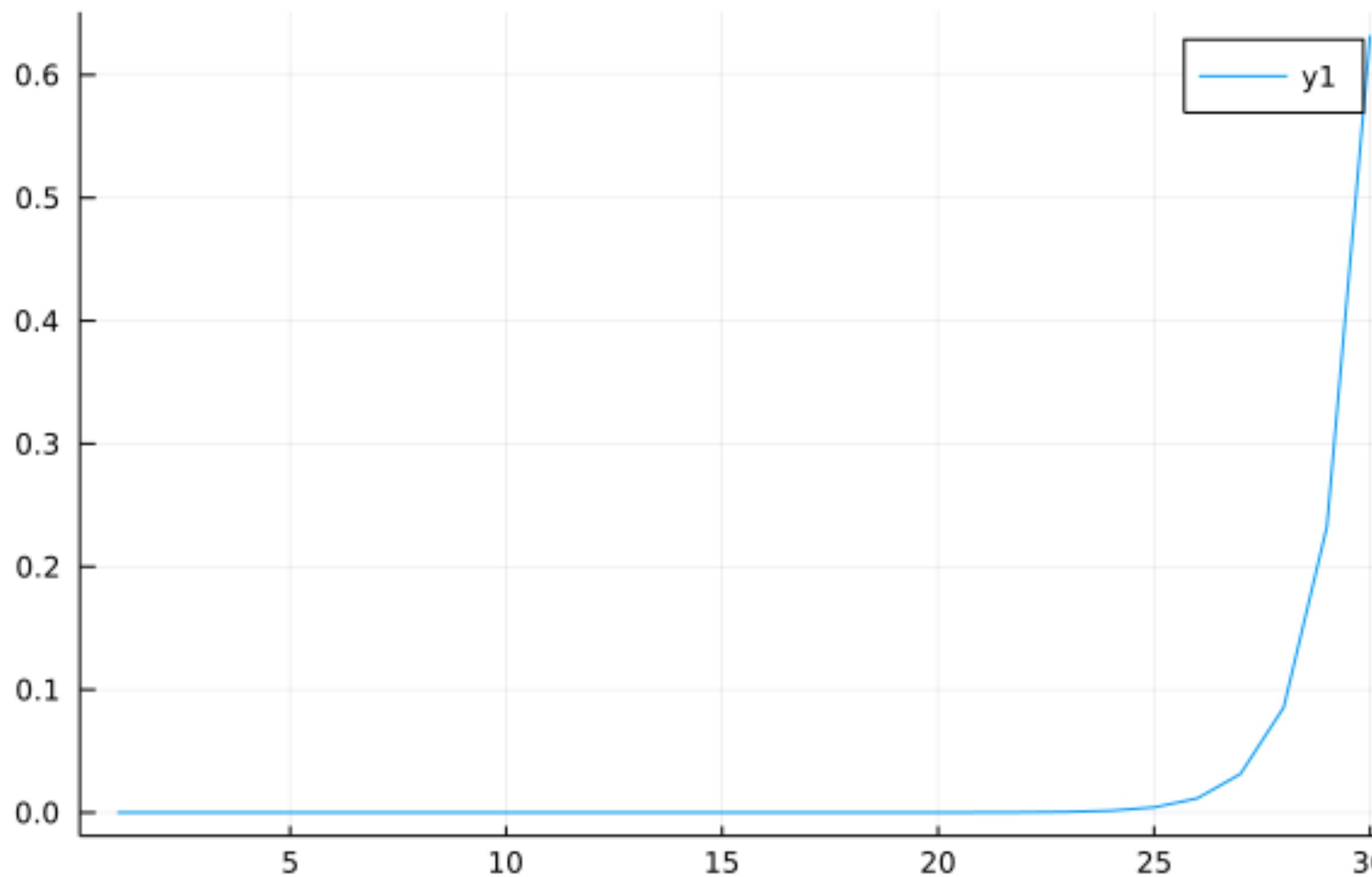
- use the training set for multiple separate binomial classifications
  - e.g., for a 3-class classifier  $c_1, c_2, c_3$ 
    - $c_1$  vs  $c_2 + c_3$ ;  $c_2$  vs  $c_1 + c_3$  and  $c_3$  vs  $c_1 + c_2$
  - train the logistic classifier for each class to predict the probability of the outcome in a given class
  - pick the class that maximises its probability

- second strategy

- use the **softmax** function instead of the sigmoid

$$h(x) = \frac{\exp(Z_i(X))}{\sum_{j=1}^k \exp Z_j(X)} \quad 1 \leq i \leq k$$

## Softmax function



**in short...**

- linear regression
- logistic regression
- regularisation