

- Machine Learning
  - ↳ Supervised Learning
    - ↳ Basic Models — Decision Trees
    - linear and non-linear classifiers

Inspired by Numerical Methods :-  
Regression Review.

Examples  $E \langle X_1, X_2, X_3 \dots X_n, Y \rangle$   $X_i \in \mathbb{R}$   
 $Y \in \mathbb{R}$

We are trying to learn  $\hat{Y} = (X_1, X_2, \dots, X_n) : Y$

Linear regression.

Assume  $\hat{Y} = w_0 + X_1 w_1 + X_2 w_2 + X_3 w_3 + \dots + X_n w_n$

where  $w_i$  are unknown.  
find such  $w_i$ !!

- Introduce  $X_0 = 1$  always. then:

$$\hat{Y} = \sum_{i=0}^n w_i X_i$$

Error Function :- to compare  $\hat{Y}$  and  $Y$

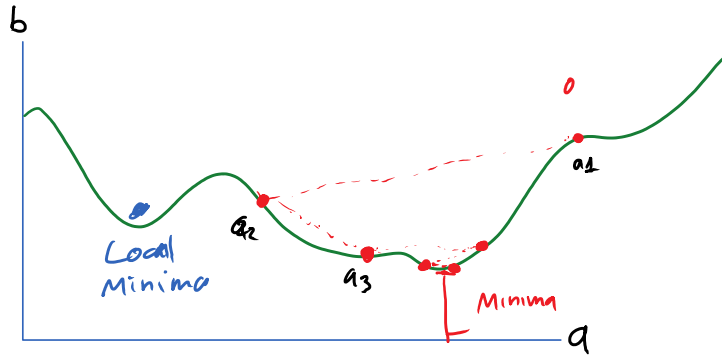
$$\begin{aligned} \text{Error}(E) &= \sum_{e \in E} (Y(e) - \hat{Y}(e))^2 \\ &= \sum_{e \in E} \left( Y(e) - \sum_{i=0}^n w_i X_i(e) \right)^2 \end{aligned}$$

Let's use sum of squares error

find  $w_i$ 's that minimize Error (E)?

Technique: Gradient Descent  
Iterative Method to find Minima of Functions

Cartoon Version



function must be differentiable.

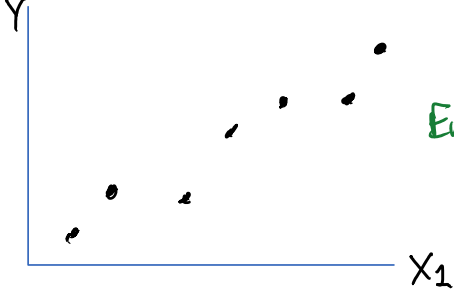
$$b = f(a)$$

$$a_{n+1} \leftarrow a_n + \eta \cdot f'(a)$$

$\uparrow$  learning rate       $\uparrow$  derivative of  $f'$

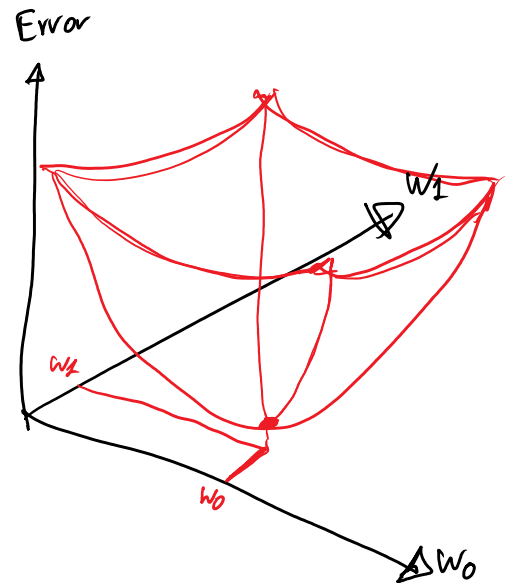
So: Apply Gradient Descent to our Error (E) function:

Examples:



$$\hat{Y} = w_0 + X_1 w_1$$

$$\text{Error} = \left( Y - \left( w_0 + X_1 w_1 \right) \right)^2$$



• update rule:

$$w_i \leftarrow w_i - \eta \cdot \frac{\partial}{\partial w_i} \text{Error}(E)$$

$\uparrow$  Learning Rate

if using Sum of square Errors

$$\frac{\partial}{\partial w_i} \text{Error}(E) = \sum_{e \in E} -2 \cdot \underbrace{(Y(e) - \hat{Y}(e))}_{\text{Error}} \cdot X_i(e)$$

$$\frac{\partial}{\partial w_i} \text{Error}(E) = \sum_{e \in E} -2 \cdot \underbrace{(Y(e) - \hat{Y}(e))}_{\delta(e)} \cdot X_i(e)$$

Compute for each  $w_i$   $\frac{\partial \text{Error}(E)}{\partial w_i} = \sum_{e \in E} -2 \cdot \delta(e) \cdot X_i(e)$   
 update each  $w_i$  by  $-\eta \cdot \frac{\partial \text{Error}(E)}{\partial w_i}$   
 Repeat until stop criteria reached.

-  $\delta(e)$ 's become small  
 - changes to  $w_i$  become small.

• Variant: incremental Gradient Descent.

- update  $w_i$  after each example.

$$w_i := w_i + \eta \cdot \delta(e) \cdot X_i(e)$$

-2 has been absorbed by  $\eta$



+ approaches solution faster

- does not converge.

- Stochastic = choose examples at random.

- Batched Gradient descent. - update of some number  $n_e$  of examples

① starts with  $n_e=1$  afterwards  $n_e$  increases until  $n_e=|E|$

① vary the size of  $\eta$  start with large value, decrease latter.

PROCEDURE LinearLearner ( E, eta )

- E : set of examples, each of the form  $\langle X_1, X_2, X_3, \dots, Y \rangle$
- eta : learning rate.

initialize  $w_0, \dots, w_n$  randomly

REPEAT

FOR EACH example e in E DO

Ycap :=  $\sum_i w_i \cdot X_i(e)$

delta :=  $Y(e) - Ycap$

update :=  $eta \cdot delta$

FOR EACH  $w_i$  DO

$w_i := w_i + update \cdot X_i(e)$

UNTIL some stop criteria is true

RETURN  $w_0, \dots, w_n$

from regression to Classification.

$X_0, X_1, \dots, X_n \in \mathbb{R}$   $Y \in \mathbb{R}$ .  $\leftarrow$  regression

$Y \in \{0, 1\}$   $\leftarrow$  classification

### Examples

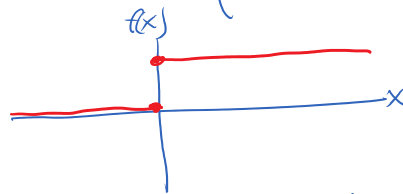
	Cats?	Celebrity?	Comedy?	Food?	Watch?
e1	True	False	False	True	Yes
e2	True	True	True	False	No
e3	False	False	True	False	Yes
e4	True	False	False	False	Yes
e5	True	False	True	False	No
e6	True	False	False	True	Yes
e7	False	False	True	True	Yes
e8	True	True	True	True	No

e9 T F F T ?

regression:  $\hat{Y} = \sum_{i=0}^n x_i w_i$   $\hat{Y} \in \mathbb{R}$

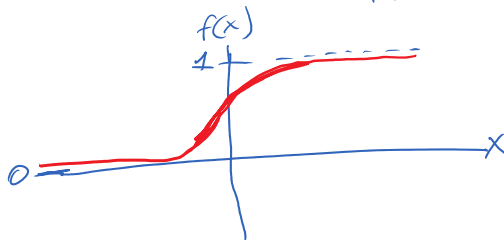
squash function  $\hat{Y} = f\left(\sum_{i=0}^n x_i w_i\right)$   $f$  squashes into  $[0, 1]$

one option for  $f$ : step function  $f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise.} \end{cases}$



second option for  $f$ : sigmoid function  $f(x) = \frac{1}{1 + e^{-x}} = \text{sig}(x)$

$$\text{sig}'(x) = \text{sig}(x) \cdot (1 - \text{sig}(x))$$



form of linear classification: "Logistic Regression"

$$\hat{Y}(e) = \text{sig}\left(\sum_{i=0}^n w_i X_i(e)\right)$$

$$\text{Error}(E) = \sum_{e \in E} (Y(e) - \hat{Y}(e))^2$$

$$\frac{\partial \text{Error}(E)}{\partial w_i} = \frac{\partial}{\partial w_i} \sum_{e \in E} \left( Y(e) - \text{sig}\left(\sum_{i=0}^n w_i X_i(e)\right) \right)^2$$

$$\begin{aligned} \frac{\partial \text{Error}(E)}{\partial w_i} &= \frac{\partial}{\partial w_i} \sum_{e \in E} \left( Y(e) - \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \right)^2 \\ &= \sum_{e \in E} 2 \cdot \left( \frac{\partial}{\partial w_i} \left( Y(e) - \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \right) \right) \\ &= \sum_{e \in E} 2 \cdot \left( Y(e) - \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \right) \cdot \frac{\partial}{\partial w_i} \left( \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \right) \\ &\quad \cdot \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \left( 1 - \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) \right) \cdot \frac{\partial}{\partial w_i} \sum_{r=0}^n w_r X_r(e) \\ &\quad \cdot X_i(e) \end{aligned}$$

$$p(e) = \text{sig} \left( \sum_{r=0}^n w_r X_r(e) \right) = \hat{Y}(e)$$

$$\delta(e) = Y(e) - \hat{Y}(e)$$

$$\frac{\partial \text{Error}(E)}{\partial w_i} = \sum_{e \in E} -2 \cdot \delta(e) \cdot p \cdot (1-p) \cdot X_i(e)$$

PROCEDURE LogisticRegression ( E, eta )

- E : set of examples, each of the form  $\langle X_1, X_2, X_3, \dots, Y \rangle$ 
  - Y is in  $\{0,1\}$
- eta : learning rate.

initialize  $w_0, \dots, w_n$  randomly

REPEAT

FOR EACH example e in E DO

$p := \text{sig} \left( \sum_i w_i * X_i(e) \right)$

$\text{delta} := Y(e) - p$

$\text{update} := \text{eta} * \text{delta} * p * (1 - p)$

  FOR EACH  $w_i$  DO

$w_i := w_i + \text{update} * X_i(e)$

UNTIL some stop criteria is true

RETURN  $w_0, \dots, w_n$

	Cats?	Celebrity?	Comedy?	Food?	Watch?
e1	True	False	False	True	Yes
e2	True	True	True	False	No
e3	False	False	True	False	Yes
e4	True	False	False	False	Yes
e5	True	False	True	False	No
e6	True	False	False	True	Yes
e7	False	False	True	True	Yes
e8	True	True	True	True	No

$$\hat{Y}(e) = \text{sig} \left( w_0 + w_1 \text{Cats}(e) + w_2 \text{Celebrity}(e) + w_3 \text{Comedy}(e) + w_4 \text{Food}(e) \right)$$

Example ②:

$$\langle X_1, X_2, X_3 \dots X_n \rangle \quad Y = \{\text{cat, dog, donut}\}$$

technique.

$$Y_{\text{cat}} = \{T, F\}$$

$$Y_{\text{dog}} = \{T, F\}$$

$$Y_{\text{donut}} = \{T, F\}$$

indicator variables



Step function:

$$f(x) = \begin{cases} 1 & \text{if } x > \theta \\ 0 & \text{otherwise.} \end{cases}$$



threshold.

Frank Rosenblatt : "Perceptron"  
circuit  $f(\sum w_i X_i)$

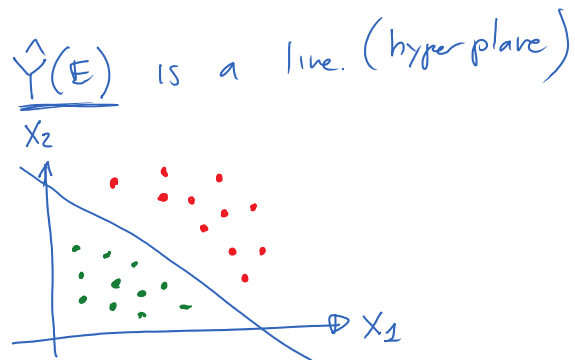
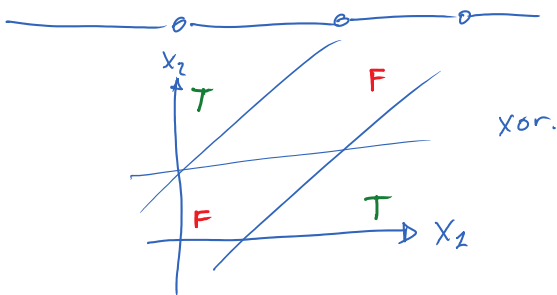
'58 "(Perceptrons) are the embryo of an electronic computer that will be able to talk, see, write, translate languages and reproduce itself and be conscious of its existence."

'69 = Marvin Minsky  
Seymour Papert

"Perceptrons"

xor function

= first AI winter.

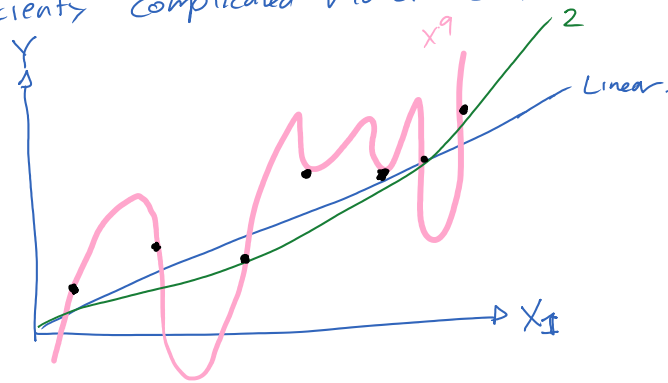


Data non-linearly separable

Why just  $\sum w_i X_i$  ???

a sufficiently complicated model can fit any data.

a sufficiently complicated model can fit any data.



Ockham's Razor: When faced with more than one explanation prefer the simplest one.

modify algorithm with a "regularizer"  
a component that rewards simplicity and punishes complexity

$$\hat{Y}(e) = \text{sig}\left(\sum_{i=0}^n w_i X_i\right)$$

$$\text{Error}(E) = \sum_e (Y(e) - \hat{Y}(e))^2 + \lambda \left( \sum_{i=0}^n |w_i| \right)$$

regularization parameter

regularizer L1

$$+ \lambda \left( \sum_{i=0}^n w_i^2 \right)$$

regularizer L2

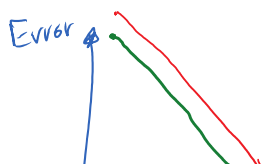
$$\frac{\partial}{\partial w} = \left( Y(e) - \hat{Y}(e) \right)^2$$

$$\frac{\partial}{\partial w} f(g(x)) = f'(g(x)) \cdot \frac{\partial}{\partial w} g(x)$$

## Errors & Pitfalls

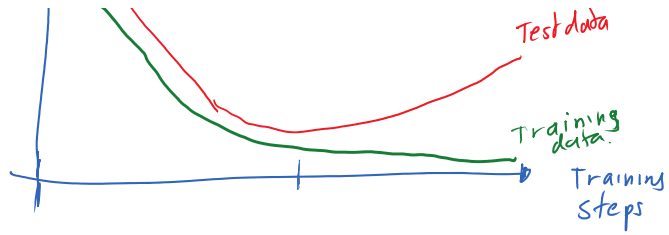
- Bias: Learner to find imperfect model
  - representation model not good enough
  - Search is not good enough
- Data is not good
  - Lack data
  - Data is noisy
- Overfitting:

Learner specializes to the training data.

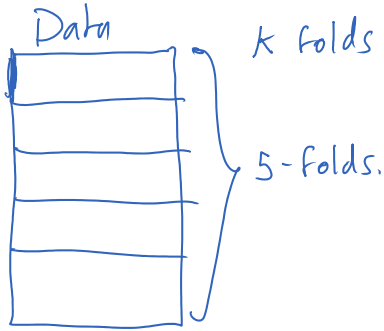


Test data





One technique to avoid over-fitting: Cross-Validation



pick 1 fold to test  
 train on remaining folds  
 Repeat K times.

Use all the data to train