

# INTRO TO MACHINE LEARNING

Slides: Andrew NG

(→ free course on Coursera)

In a problem of prediction, the outcome is what is to be predicted.

The outcome can be a continuous variable

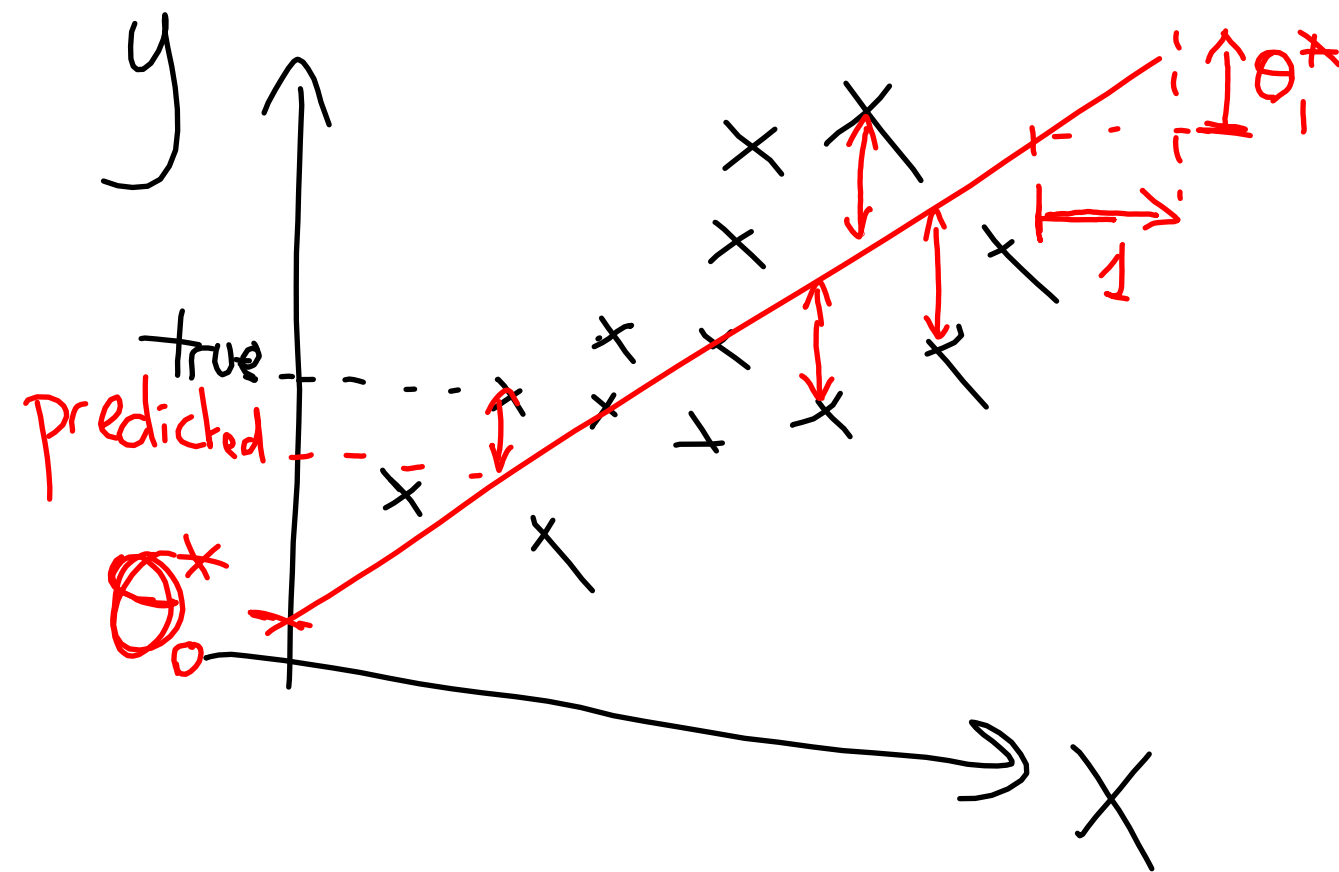
→ regression analysis

or a discrete / categorical variable

→ classification / clustering problem

# LINEAR REGRESSION (UNIVARIATE)

- \* predictors : one continuous variable ( $x$ )
- \* outcome : a continuous variable ( $y$ )
- \* training dataset :  $n$  points  $(x^{(i)}, y^{(i)})$   $i \in [1, n]$



we want to get to a  
model whereby  $y$  is a function  
of  $x$  :  $y = f(x)$

Linear regression means that our model must be of the form :  $y = \theta_0 + \theta_1 x$

$\Rightarrow$  the learning will consist in learning the best values for the 2 parameters  $\theta_0$  and  $\theta_1$

intercept  $+ \theta_2 x^2$  (quadratic term) slope

“the best”??  $\Rightarrow$  we will define a cost function

$J(\theta_0, \theta_1)$  to be minimized. Least-squares approach is

$$\text{to have } J(\theta_0, \theta_1) = \sum_i (\text{predicted}^{(i)} - \text{true outcome}^{(i)})^2$$

$$J(\theta_0, \theta_1) = \sum (\text{predicted}^{(i)} - \text{true outcome}^{(i)})^2$$

Means:  $J(\theta_0, \theta_1) = \sum_{i=1}^n (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$

There is an algebraic formula to calculate  $(\theta_0^*, \theta_1^*)$  that minimizes  $J(\theta_0, \theta_1)$ . Therefore we have our regression model:

$$y = \theta_0^* + \theta_1^* x$$

In R, linear regression is performed with the function `lm()` (think "linear model")

`mydata`

x	y

$\Rightarrow$  call: `lm(y ~ x, data = mydata)`

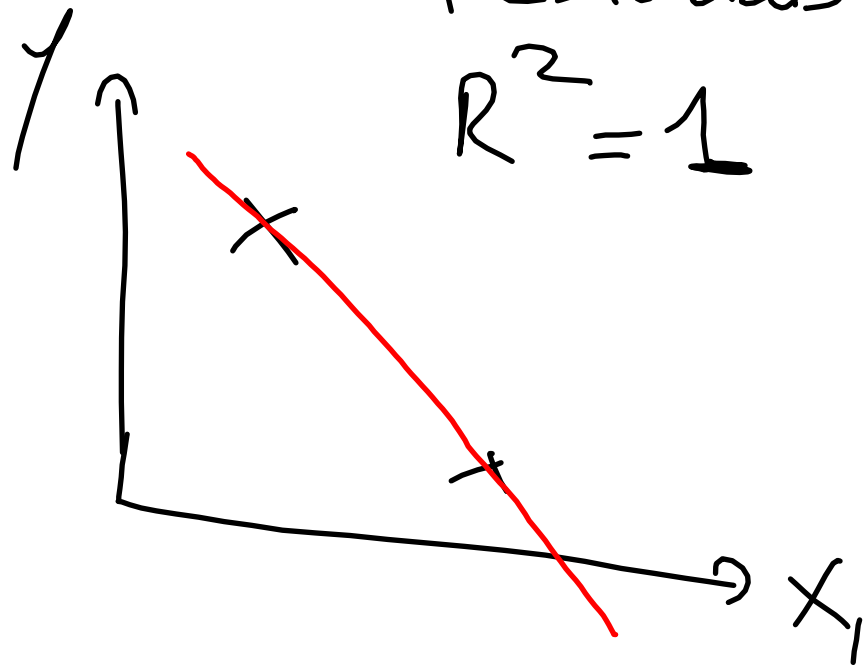
# MULTIVARIATE (MULTIPLE) LINEAR REGRESSION

- ⊗ still have a continuous outcome variable  $y$
- ⊗  $m > 1$  <sup>continuous or not</sup> predictor variables  $x_1, x_2, x_3, \dots, x_m$
- ⊗ one training datapoint  $\# i : (x_1^{(i)}, x_2^{(i)}, \dots, x_m^{(i)}, y^{(i)})$
- ⊗ model :  $y = \theta_0 + \sum_{j=1}^m \theta_j x_j$  ( $m+1$  parameters)
- ⊗ same Least Squares approach to minimize the sum of the squared prediction errors

# OVERFITTING

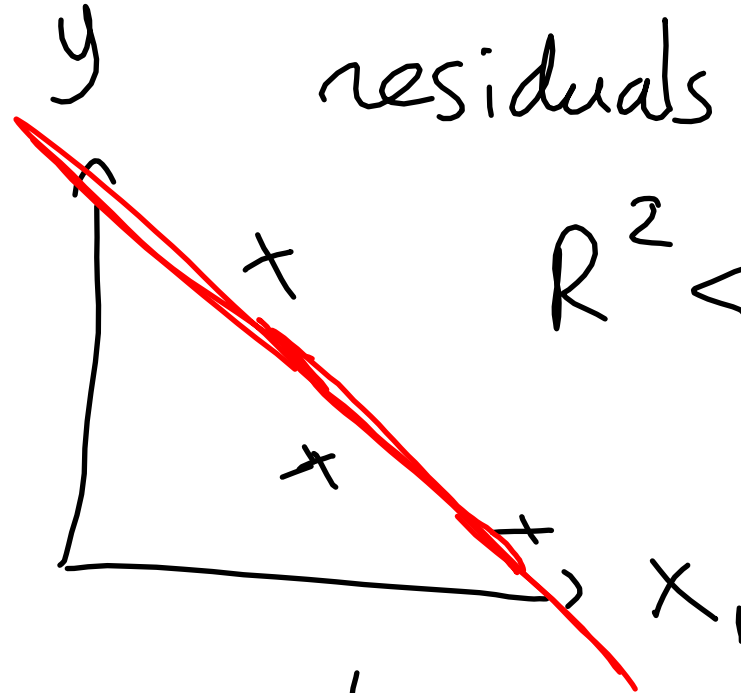
residuals = 0

$$R^2 = 1$$



residuals not zero

$$R^2 < 1$$



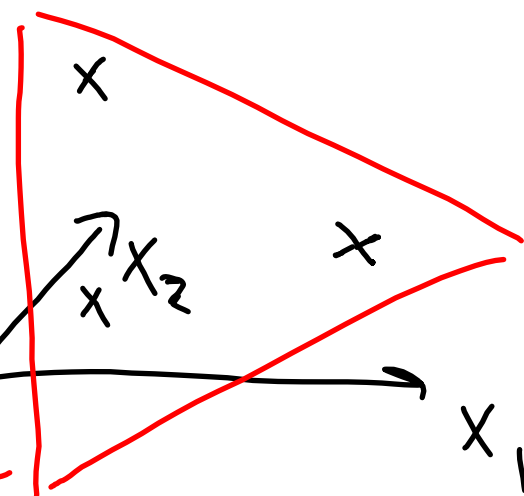
residuals = 0

$$R^2 = 1$$

$\text{plane}(\theta_0^*, \theta_1^*, \theta_2^*)$

contains all three training datapoints

outcome  
y



let us add one predictor

$$\text{model: } y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

defines a plane



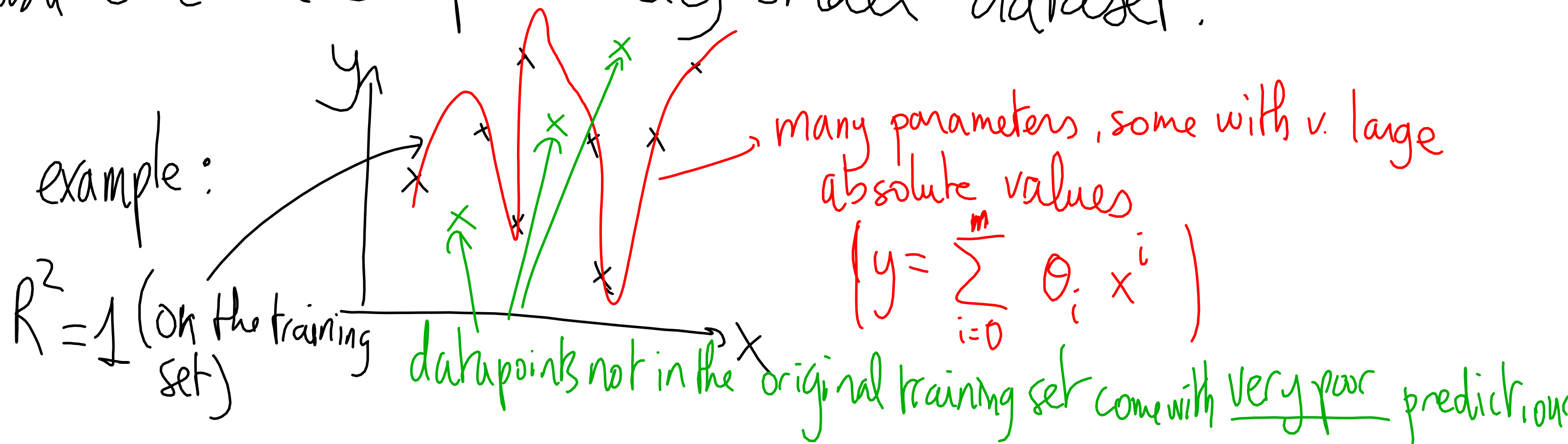
With a given Fixed training set of  $n$  datapoints, increasing the number of predictors automatically improves the  $R$  squared measure (better fit).

$\Rightarrow$  with  $(n-1)$  predictors, we reach  $R^2 = 1$  even if some individual predictors don't make sense.

The adjusted  $R^2$  includes a penalty  
in proportion with the number of  
predictors.

$\Rightarrow$  fairer to use the adjusted  $R^2$ .

Overfitting occurs when one designs a model of high complexity based on a comparatively small dataset.



to avoid overfitting, don't try  
and go for too complex a representation  
on "simple" data :

- check the adjusted  $R$  squared
- add a predictor only if it individually  
helps decrease significantly the value of the cost  
function

→ automatic procedures (based on the AIC/BIC values) to  
determine the best model, at the same time accurate and generic enough.