

## Simple Regression

Linear regression with one input

$$(x_1 = \text{sq.ft}, y_1 = \$)$$

$$(x_2 = \text{sq.ft}, y_2 = \$)$$

$$(x_3 = \text{sq.ft}, y_3 = \$)$$

$$(x_4 = \text{sq.ft}, y_4 = \$)$$

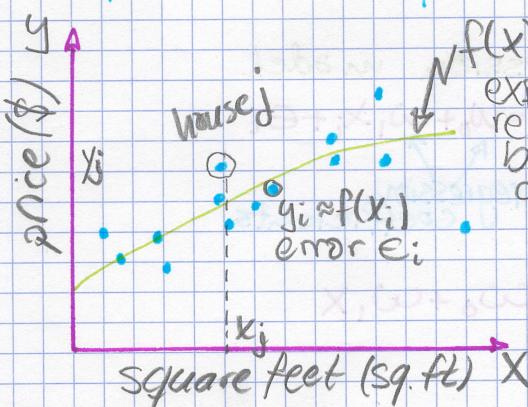
$$(x_5 = \text{sq.ft}, y_5 = \$)$$

Input

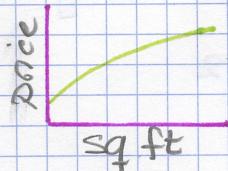
Output

### Input vs. Output

- $y$  is the quantity of interest
- assume  $y$  can be predicted from  $x$



Which Model  $f(x)$ ?



Regression model:

$$y_i = f(x_i) + \epsilon_i$$

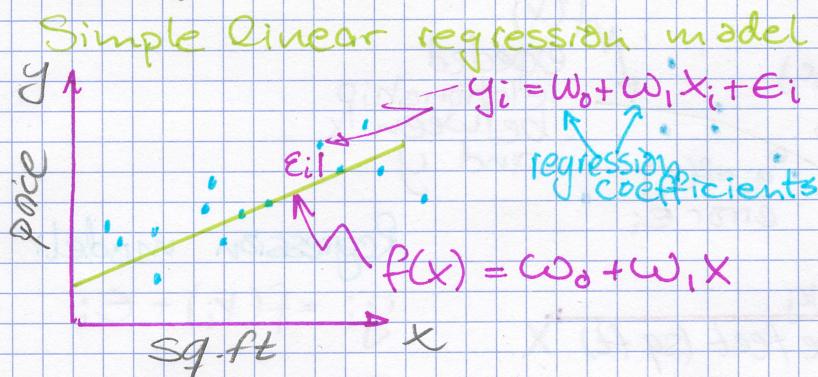
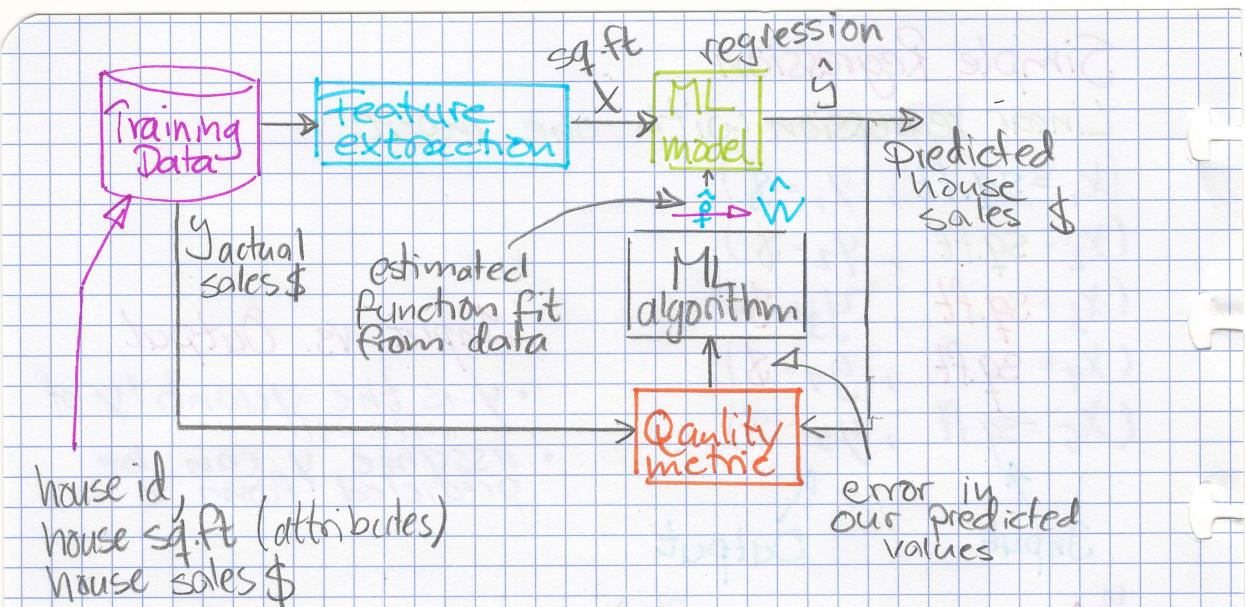
$$\mathbb{E}[\epsilon_i] = 0$$

expected value

equally likely  
that error is  
+ or -



$y_i$  is equally likely to be above or below  $f(x_i)$



Cost of using a given line  
Residual sum of squares (RSS)

$$\text{RSS}(w_0, w_1) = \left( \$_{\text{house}} - [w_0 + w_1 \text{sq.ft}_{\text{house}}] \right)^2 + \dots$$

+ [include all training houses]

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N \left( y_i - [w_0 + w_1 x_i] \right)^2$$

Regression model

$$y_i = \omega_0 + \omega_1 x_i + \epsilon_i$$

parameters (unknown variables)

Estimated parameters

$$\hat{\omega}_0, \hat{\omega}_1$$

take actual values

Minimizing the cost

$$\min_{(\omega_0, \omega_1)} \sum_{i=1}^N (y_i - (\omega_0 + \omega_1 x_i))^2$$

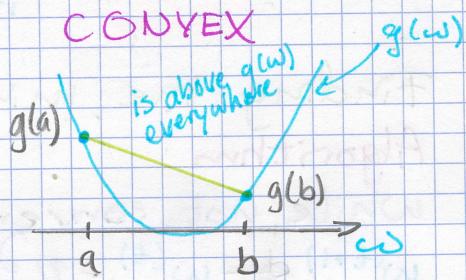
RSS( $\omega_0, \omega_1$ ) is a function  
of 2 variables =  $g(\omega_0, \omega_1)$

Finding maxima or minima analytically  
Convex/concave functions

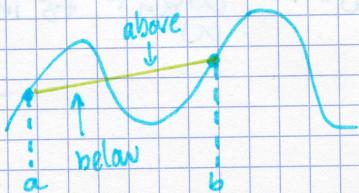
CONCAVE



CONVEX



NEITHER



## Example

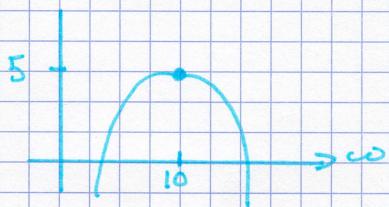
$$g(w) = 5 - (w - 10)^2$$

$$\frac{\partial g(w)}{\partial w} = 0 - 2(w - 10) \cdot 1 = -2w + 20$$

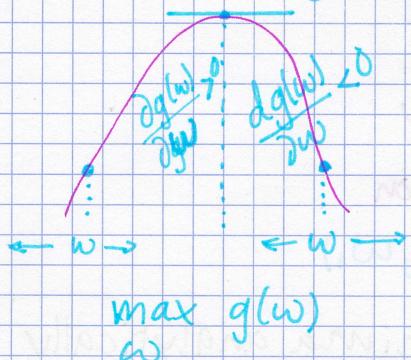
set derivative = 0:

$$-2w + 20 = 0$$

$$w = 10$$



Finding the max via hill climbing  
derivative = 0



How do we know whether to move  $w$  to right or left?  
(Inc or dec. the value of  $w$ ?)

while not converged

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \frac{\partial g(w)}{\partial w}$$

iteration  $t$

step size

Finding the min via hill descent

Algorithm

while not converged

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{\partial g(w)}{\partial w}$$

$w^{(t)}$

When derivative is positive, we want to decrease  $w$  and when derivative is negative, we want to increase  $w$ .

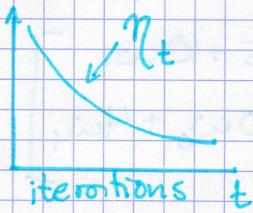
## Choosing the stepsize - fixed stepsize

$$\eta = 0.1$$

Decreasing stepsize or stepsize schedule  
Common choices:

$$\eta = \frac{\alpha}{t}$$

$$\eta = \frac{\alpha}{\sqrt{t}}$$



## Convergence criteria

For convex functions, optimum occurs when

$$\frac{\partial g(\omega)}{\partial \omega} = 0$$

In practice, stop when

$$\left| \frac{\partial g(\omega)}{\partial \omega} \right| < \epsilon_A$$

threshold to be set

## Gradient: derivatives in multiple dimensions

$$\nabla g(\omega) = \begin{bmatrix} \frac{\partial g}{\partial \omega_0} \\ \frac{\partial g}{\partial \omega_1} \\ \vdots \\ \frac{\partial g}{\partial \omega_p} \end{bmatrix} \quad \text{---} \quad \left. \begin{array}{l} \text{(p+1)-dimensional} \\ \text{vector} \end{array} \right\}$$

gradient

constants

partial derivatives is like a derivative with respect to  $\omega$ , treating all other variables as constants

## Gradient example

$$g(w) = 5w_0 + 10w_0w_1 + 2w_1^2$$

$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

$$\frac{\partial g}{\partial w_1} = 10w_0 + 4w_1$$

$$\nabla g(w) = \begin{bmatrix} 5 + 10w_1 \\ 10w_0 + 4w_1 \end{bmatrix}$$

## Gradient descent

### Algorithm

while not converged:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla g(w^{(t)})$$

Convergence:  $\|\nabla g(w)\| < \epsilon$

## Computing the gradient of RSS

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^n (y_i - [w_0 + w_1 x_i])^2$$

$$\frac{\partial}{\partial w} \sum_{i=1}^n g_i(w) = \sum_{i=1}^n \frac{\partial}{\partial w} g_i(w)$$

In our case

$$g_i(w) = (y_i - [w_0 + w_1 x_i])^2$$

$$\frac{\partial \text{RSS}(\omega)}{\partial \omega_0} = \sum_{i=1}^N \frac{\partial}{\partial \omega_0} (y_i - [\omega_0 + \omega_1 x_i])^2$$

same for  $\omega_1$

Taking the derivative w.r.t  $\omega_0$

$$\sum_{i=1}^N 2(y_i - [\omega_0 + \omega_1 x_i])' \cdot (-1) =$$

$$-2 \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i])$$

Taking the derivative w.r.t  $\omega_1$

$$\sum_{i=1}^N 2(y_i - [\omega_0 + \omega_1 x_i])' \cdot (-x_i) =$$

$$-2 \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i]) x_i$$

$$\checkmark \text{RSS}(\omega_0, \omega_1) = \left[ \begin{array}{l} -2 \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i]) \\ -2 \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i]) x_i \end{array} \right]$$

Approach 1

top line average house sales price

$$\hat{\omega}_0 = \frac{\sum_{i=1}^N y_i}{N} - \hat{\omega}_1 \frac{\sum_{i=1}^N x_i}{N}$$

estimate the slope

Bottom term

$$\sum y_i x_i - \hat{\omega}_0 \sum x_i - \hat{\omega}_1 \sum x_i^2 = 0$$

$$\hat{\omega}_1 = \frac{\sum y_i x_i - \frac{\sum y_i \sum x_i}{N}}{\sum x_i^2 - \frac{\sum x_i \sum x_i}{N}}$$

Note:

$$\sum y_i$$

$$\sum x_i$$

$$\sum y_i x_i$$

$$\sum_{i=1}^N x_i^2$$

Approach 2: Gradient descent

$$\nabla_{(w_0, w_1)} \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum y_i - \hat{y}_i(w_0, w_1) \\ -2 \sum (y_i - \hat{y}_i(w_0, w_1)) x_i \end{bmatrix}$$

while not converged:

$$\begin{aligned} w_0^{(t+1)} &= \underbrace{w_0^{(t)}}_{\text{initial}} - \eta \left[ -2 \sum (y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})) \right] \\ w_1^{(t+1)} &= \underbrace{w_1^{(t)}}_{\text{initial}} - \eta \left[ -2 \sum (y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})) x_i \right] \\ &= \begin{bmatrix} w_0^{(t+1)} \\ w_1^{(t+1)} \end{bmatrix} \end{aligned}$$

If overall underpredicting  $\hat{y}_i$ , then  $\sum [y_i - \hat{y}_i]$  is positive  $\rightarrow w_0$  is going to increase. Similar intuition for  $w_1$ , but multiply with  $x_i$

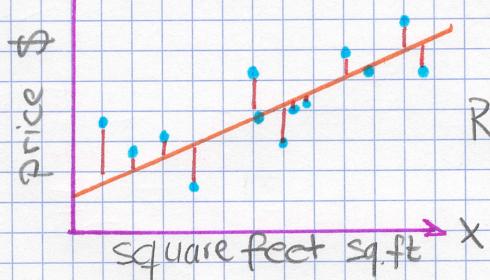
Comparing the approaches

- For most ML problems, cannot solve  $\text{gradient} = 0$
- Even if solving  $\text{gradient} = 0$  is feasible, gradient descent can be more efficient
- Gradient descent relies on choosing stepsize and convergence criteria

## Asymmetric cost functions

### Symmetric cost function

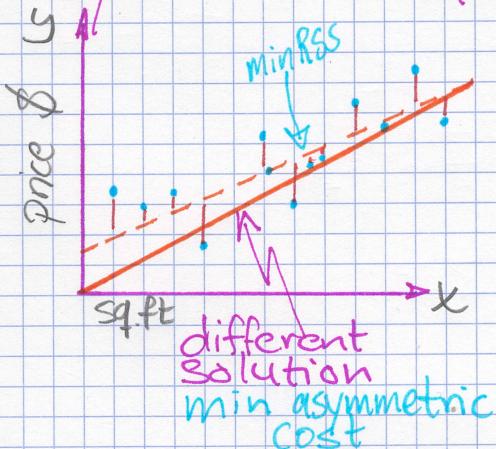
Residual sum of squares (RSS)



$$RSS(w_0, w_i) = \sum_{i=1}^N (y_i - [w_0 + w_i x_i])^2$$

Assumes cost of over-estimating sales price is same as under-estimating

## Asymmetric cost functions



What if cost of listing house too high has bigger cost?

Too high  $\rightarrow$  no offers \$=0  
Too low  $\rightarrow$  offers for lower