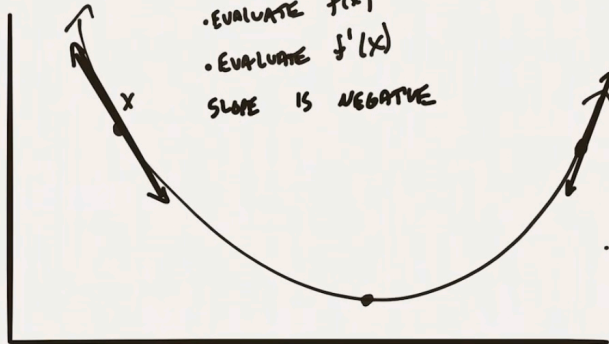


# GRADIENT DESCENT

GOAL: FIND MINIMUM  
OF THIS FUNCTION

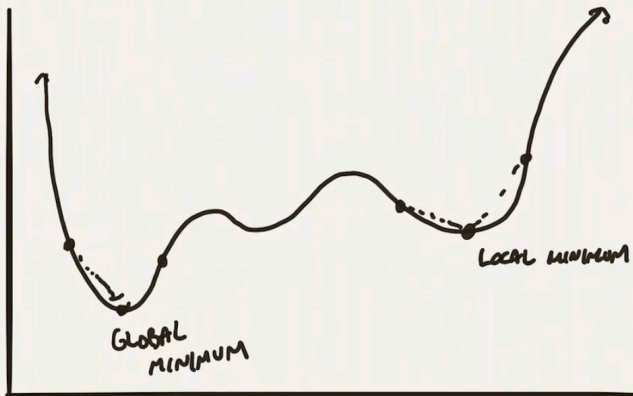


- EVALUATE  $f(x)$
- EVALUATE  $f'(x)$
- SLOPE IS NEGATIVE

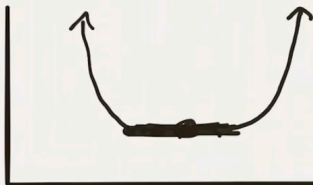
• IF  $f'(x) < 0$   
MOVE A BIT TO  
RIGHT

• IF  $f'(x) > 0$   
MOVE A BIT TO  
THE LEFT

• IF  $f'(x) \approx 0$   
STOP OUTPUT  $x$ .

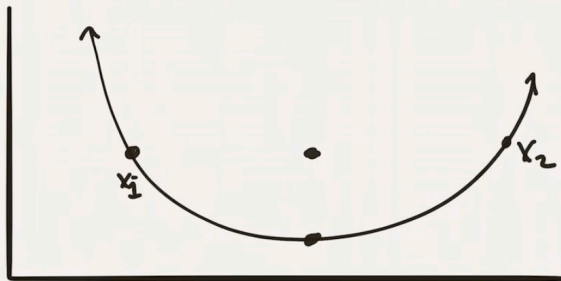


$f'(x) > 0$  MOVE LEFT  
 $f'(x) < 0$  MOVE RIGHT  
 O.W. STOP.



### CONVEXITY

A FUNCTION IS CONVEX IF  
 THE CHORD CONNECTING ANY 2  
 POINTS OF THE GRAPH LIES  
 ABOVE THE FUNCTION.



$$f\left(\frac{1}{2}x_1 + \frac{1}{2}x_2\right) \leq \frac{1}{2}f(x_1) + \frac{1}{2}f(x_2)$$

EQUIVALENT DEFINITION OF CONVEXITY: A FUNCTION IS CONVEX

$$\text{IF } f(ax_1 + (1-a)x_2) \leq af(x_1) + (1-a)f(x_2)$$

let's say  $x^*$  is GLOBAL MIN

WE'RE CURRENTLY AT  $x$

$$f(ax + (1-a)x^*) \leq af(x) + (1-a)f(x^*)$$

THUS FAR: OUR IDEA HAS BEEN TO  
LOOK AT <sup>TANGENT</sup> LINES AND THIS IDEA WORKS  
FOR SAY LINEAR FUNCTIONS AND SIMPLE CONVEX  
FUNCTIONS.

EVEN IF WE WANT TO MINIMIZE MORE COMPLICATED  
FUNCTIONS, ASSUME THEY ARE "LOCALLY" LINEAR.

$f$  <sup>FIXED</sup> AT POINT  $x$

$$f(x+\epsilon) = f(x) + \underbrace{\epsilon \cdot f'(x)}_{\text{LINEAR FUNCTION OF } \epsilon} + \underbrace{\frac{\epsilon^2}{2!} f''(x) + \frac{\epsilon^3}{3!} f'''(x) + \dots}_{\text{WHEN } \epsilon \text{ IS SMALL, THESE TERMS ARE NEGLIGIBLE.}}$$

EXPRESSION IN TERMS OF  $\epsilon$ .

TAYLOR'S THM ALSO HOLDS IN  $d$  DIMENSIONS

INSTEAD OF TAKING DERIVATIVES (UNIVARIATE CASE)

FOR HIGHER DIMENSIONS WE MUST LOOK AT GRADIENTS.

THE GRADIENT OF  $f$  AT POINT  $x$

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_d}(x) \right)$$

↙  $d$ -dimensional  
VECTOR.

$$f = w^T x + b \quad \frac{\partial f}{\partial x_i} = w_i \quad \nabla f = w$$

$$f(x) = x^T A x - b^T x$$

(n - dimensions)

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j - \sum_{i=1}^n b_i x_i$$

(\*) ..... .....

$$\frac{\partial f}{\partial x_k} = \sum_{j=1}^n \underline{a_{kj} x_j} + \sum_{i=1}^n \underline{a_{ik} x_i} - b_k$$

CASE 1:  $i=k$  kth row of A inner prod with x kth col of A inner prod with x

$$(when\ i=k\ and\ j=k\ in\ (k))\ a_{kk} x_k^2 \quad \frac{\partial (a_{kk} x_k^2)}{\partial x_k} = \underline{2 \cdot a_{kk} x_k}$$

CASE 2:  $i \neq k$  WE WON'T HAVE  $a_{kk} x_k$  TERM BECAUSE  $i \neq k$

ANSWER

$$\boxed{Ax + A^T x - b}$$

← GRADIENT AT POINT x  
IF A IS SYMMETRIC

$$2Ax - b$$

# DEFINE GRADIENT DESCENT

INITIALLY WE'LL CHOOSE  $w$  RANDOMLY  
(WANT TO MINIMIZE  $f(w)$ )

IF  $\|\nabla f(w)\|_2 < \epsilon$  STOP OUTPUT  $w$

OTHERWISE  $w_{\text{new}} = w_{\text{old}} - \eta \nabla f(w)$

↳ STEP-SIZE PARAMETER

$$w_j^{\text{new}} = w_j^{\text{old}} - \eta \frac{\partial f}{\partial w_j}(w)$$

IS USUALLY  
SET TO BE  
RELATIVELY SMALL.

# APPLY GRADIENT DESCENT TO LINEAR REGRESSION

$$h(x) = w^T x + b \quad (\text{SEARCHING FOR THIS FUNCTION})$$

(WE HAVE A TRAINING SET OF SIZE  $m$ )

$$\text{M.S.E.}(w) = \frac{1}{m} \sum_{j=1}^m \underbrace{(w^T x^j + b - y^j)}_{g_j}^2$$

$$\frac{\partial g_j}{\partial w_i}$$

$$\frac{\partial (h_1 + h_2)}{\partial z} = \frac{\partial h_1}{\partial z} + \frac{\partial h_2}{\partial z}$$



# APPLY GRADIENT DESCENT TO LINEAR REGRESSION

$$h(x) = w^T x + b \quad (\text{SEARCHING FOR THIS FUNCTION})$$

(WE HAVE A TRAINING SET OF SIZE  $m$ )

$$\text{M.S.E.}(w) = \frac{1}{m} \sum_{j=1}^m \underbrace{(w^T x^j + b - y^j)}_{g_j}^2$$

M.S.E.(w) IS  
A CONVEX FUNCTION.

$$\frac{\partial g_j}{\partial w_i} = 2 \cdot (w^T x^j + b - y^j) x_i^j$$

RUNNING TIME

$$O(m \cdot n)$$

$$\nabla g_j(w) = 2 (w^T x^j + b - y^j) x^j$$

$$\nabla \text{M.S.E.}(w) = \frac{2}{m} \cdot \sum_{j=1}^m (w^T x^j + b - y^j) \cdot x^j$$

## STOCHASTIC GRADIENT DESCENT

- PREVIOUSLY IN LINEAR REGRESSION EXAMPLE  
WE SUMMED OVER ALL POINTS IN TRAINING SET.

- CHOOSE AN INDEX  $j$  AT RANDOM; COMPUTE  
GRADIENT W.R.T. THIS POINT only

$$W_{\text{NEW}} = W_{\text{OLD}} - 2 \cdot \eta (W^T x^j + b - y^j) x^j$$

$$E[W_{\text{NEW}}] = W_{\text{OLD}} - 2\eta \cdot \frac{1}{m} \sum_{j=1}^m (W^T x^j + b - y^j) x^j$$

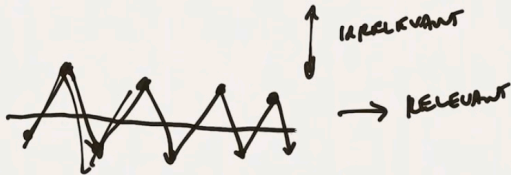
- USE "BATCHES" TO INTERPOLATE BETWEEN GRADIENT DESCENT  
AND PURE S.G.D. (SGD)

• HOW TO CHOOSE  $\eta$  THE STEP-SIZE

• MORE ART THAN SCIENCE; USE CROSS  
VALIDATION TO PICK  $\eta$

• MANY TECHNIQUES FOR ADAPTIVELY CHOOSING  $\eta$

• MOMENTUM



MOMENTUM HAS A "VELOCITY" VARIABLE  $V$

$$V_0 = 0$$

$$V_1 = -\eta g_1$$

$$V_i = \alpha \cdot V_{i-1} - \eta g_i$$

THIS TAKES A WEIGHTED MOVING AVERAGE OF  $-\eta g_i$ 's  
EXPONENTIAL " "

$$W_{\text{new}} = W_{\text{old}} + V_i$$

ACCELERATED GRADIENT DESCENT