

Artificial Intelligence/Machine Learning/Deep Learning: 'Bridging the Skills Gap'

Lesson 2: Gradient Descent

Hello, welcome back! This lesson is dedicated to the most widely used optimization technique: Gradient Descent!

1. Gradient Descent, Multivariate Gradient Descent, Stochastic Gradient Descent, Mini Batch Stochastic Gradient Descent,
2. Optimizers: Momentum, RMSprop, Adam, Adagrad

Prerequisite for this session:

Differentiation, Gradient (∇), Taylor Series, Vector. These items are explained in the math refresher on Calculus and Linear Algebra.

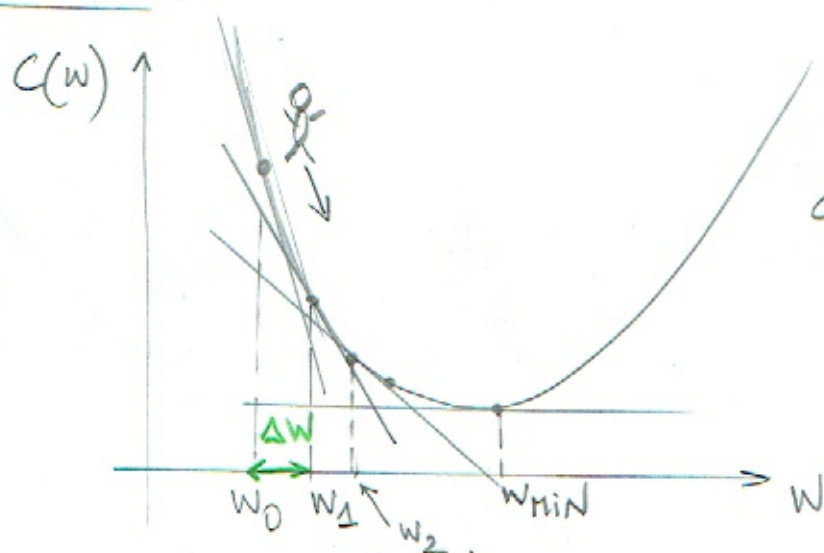
Last session we looked at a regression problem, trying to predict the house price given the size of the house. We introduced the concept of a Cost Function $C(w,b)$ and every time that $(y \neq y_p)$ the model incurs 'cost'.

Our goal is to 'learn' the model parameters w and b as to minimize that cost for a test point (x_{new}) \rightarrow minimization problem \rightarrow differentiation problem \rightarrow Gradient Descent

The cost function used for our single variate regression problem was MSE

$C(w,b) = \frac{1}{6} \sum_{i=1}^6 (y_p^i - y^i)^2$ with y_p the predicted housing price and y the actual housing price

$C(w,b)$ is quadratic function \rightarrow always positive and penalizes outliers



MSE = quadratic function

$$C(W, b) = \frac{1}{m} \sum_{i=1}^m (y_p^{(i)} - y^{(i)})^2$$

START ① pick w_0 randomly.

② calculate $\frac{dc}{dw}(w_0) \rightarrow \text{sign} \rightarrow -$

move downwards with step $\alpha \rightarrow w_1$

calculate $\frac{dc}{dw}(w_1) \rightarrow \text{sign} \rightarrow -$

move downwards with step $\alpha \rightarrow w_2$

Repeat \rightarrow till $\frac{dc}{dw} \approx 0 \rightarrow \text{MINIMUM}$

\hookrightarrow OR change in cost $< 10^{-3}$

$$w_{i+1} = w_i - \alpha \frac{dc}{dw}(w_i)$$

LEARNING Rate = hyperparameter!

$$i=0 \rightarrow w_1 = w_0 - \alpha \frac{dc}{dw}(w_0)$$

$$\uparrow \text{ if } \frac{dc}{dw}(w_0) < 0$$

$$-\alpha \frac{dc}{dw}(w_0) > 0$$

$w_1 = w_0 + \text{something} \rightarrow$ move to the right on W-axis

Chose AS A DATA Scientist

\rightarrow ANDREW NG
0.001, 0.003, 0.01, ...
 $\uparrow \times 3 \quad \uparrow \times 3$

α too small \rightarrow slow convergence

α too big \rightarrow Risk of oscillations \rightarrow NO convergence

So far \rightarrow 1 feature (size of the house)

Multivariate Regression

	Size	#rooms	AGE	PRICE
<i>m samples</i>				

→ d features ($d=3$)

$$y = wx + b \quad \xrightarrow{\text{set } x_0 = 1} \quad y = w_0 x_0 + w_1 x_1 + \dots + w_d x_d \quad (1)$$

1 feature d features

$$X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$
 feature vector

$$w^T X = (w_0 \ w_1 \ \dots \ w_d) \cdot \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = (1)$$

$$\text{MSE} = C(w) = \frac{1}{m} \sum_{i=1}^m (w^T x^{(i)} - y_i)^2$$

Now: $C(w_0 + \Delta w) \approx C(w_0) + \underbrace{\Delta w \cdot \nabla C(w_0)}_{\text{dot product} = \text{scalar}} + \dots$
 \downarrow
 Taylor

If Δw small enough \rightarrow linear!
 and minimum if Δw and $\nabla C(w_0)$ are opposite directions
 $\rightarrow \Delta w \sim -\nabla C(w_0) \rightarrow \frac{\Delta w}{\nabla C(w_0)} = -\alpha$
 $\xrightarrow{+w_0} w_0 + \Delta w = w_0 - \alpha \nabla C(w_0)$

(3)

$$W_{i+1} = W_i - \alpha \nabla C(W_i)$$

Gradient Descent
for multi variate regression!

example: 2 features \rightarrow 2 weights
 $x_1 \ x_2$ $w_1 \ w_2$

$$C(w) = w_0 + w_1 x_1 + w_2 x_2$$

\downarrow

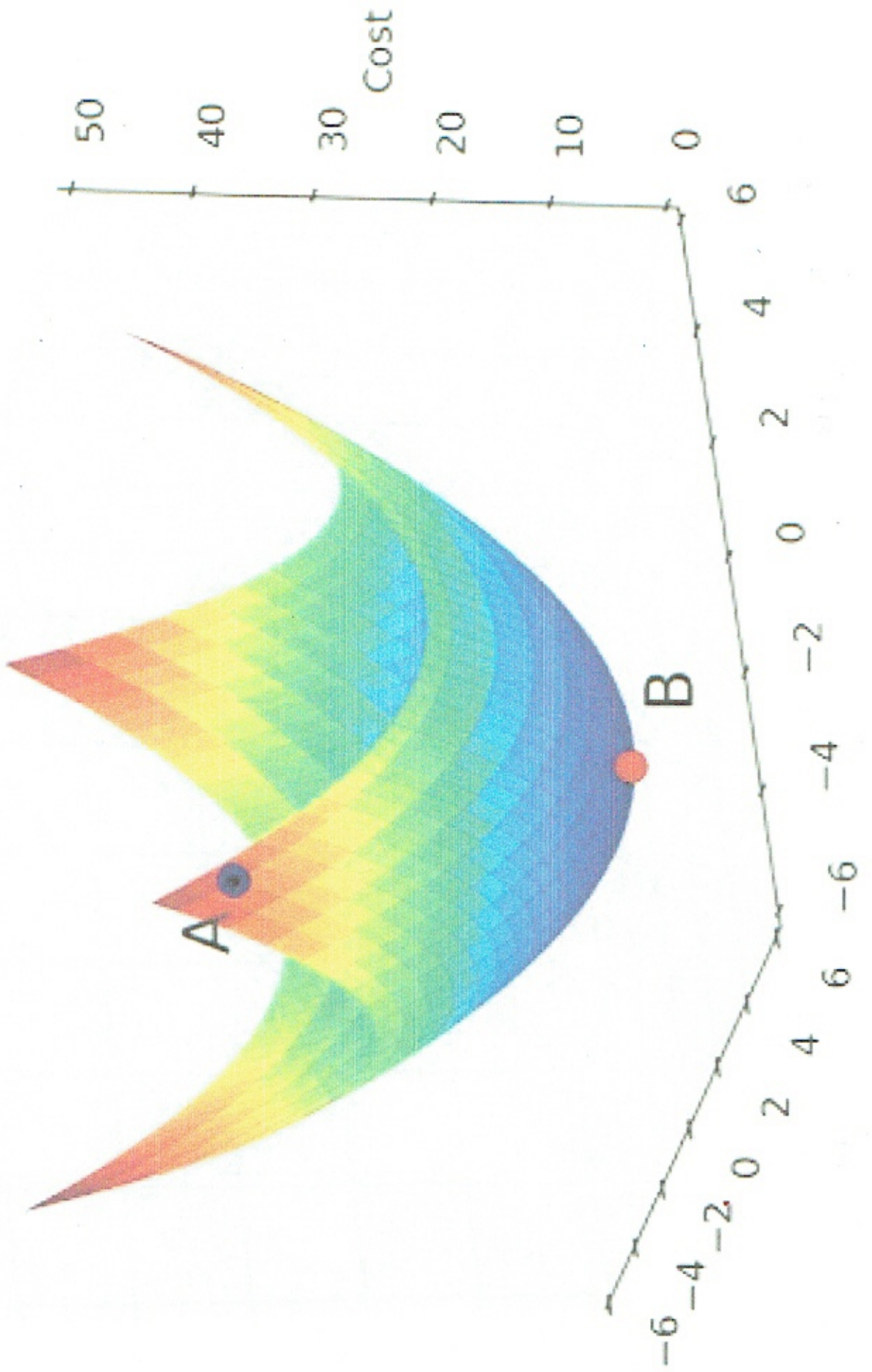
was BIAS b

\hookrightarrow

$$w_0 x_0 \dots$$

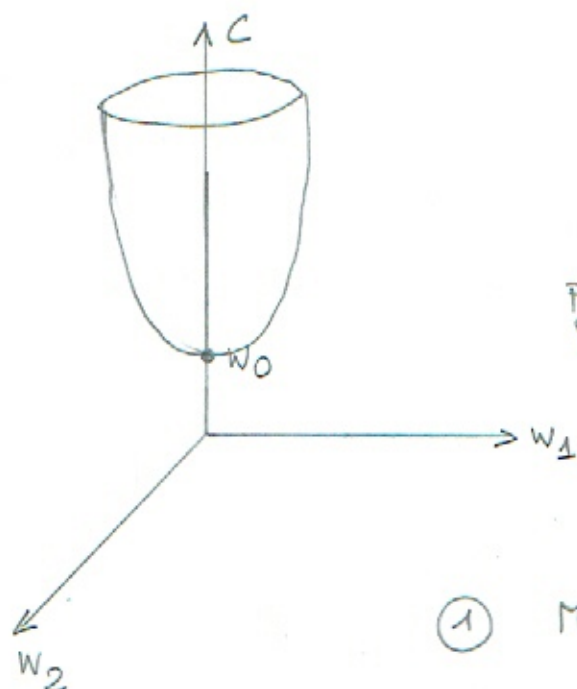
\downarrow

$$x_0 \stackrel{\text{set}}{=} 1.$$



PROJECTION ON w_1, w_2 plane \rightarrow CONTOUR PLOTS

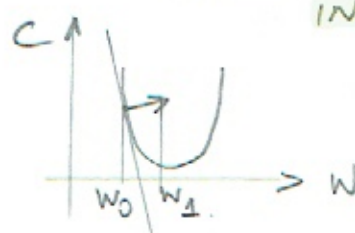
(5)



projection
on w_1, w_2 plane

CONTOUR PLOTS
showing points with
equal cost
see p 6

- ① MINIMIZATION \rightarrow gradient points inward.



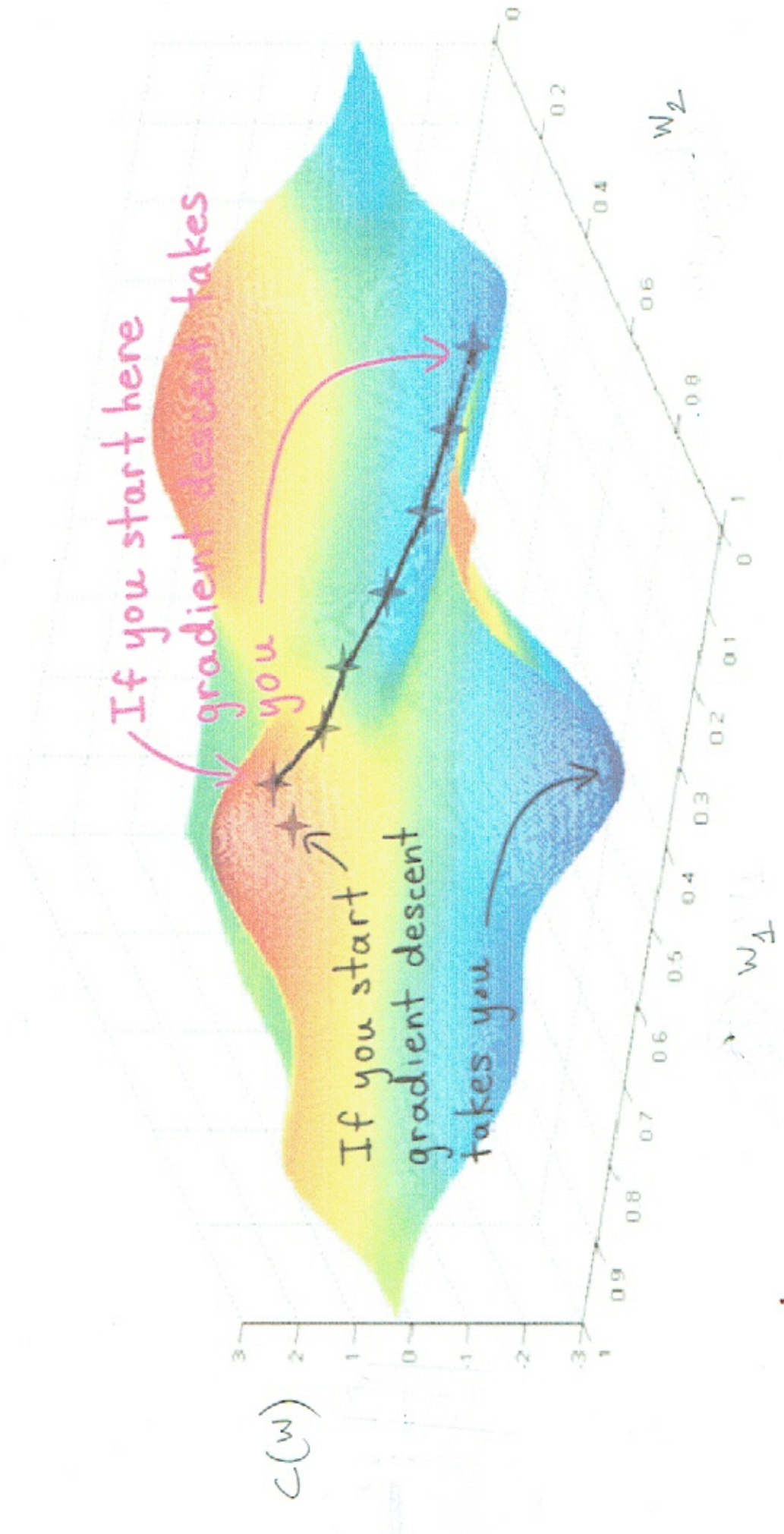
- ② Gradient is \perp ON CONTOUR lines and
is IN the direction of steepest descent.

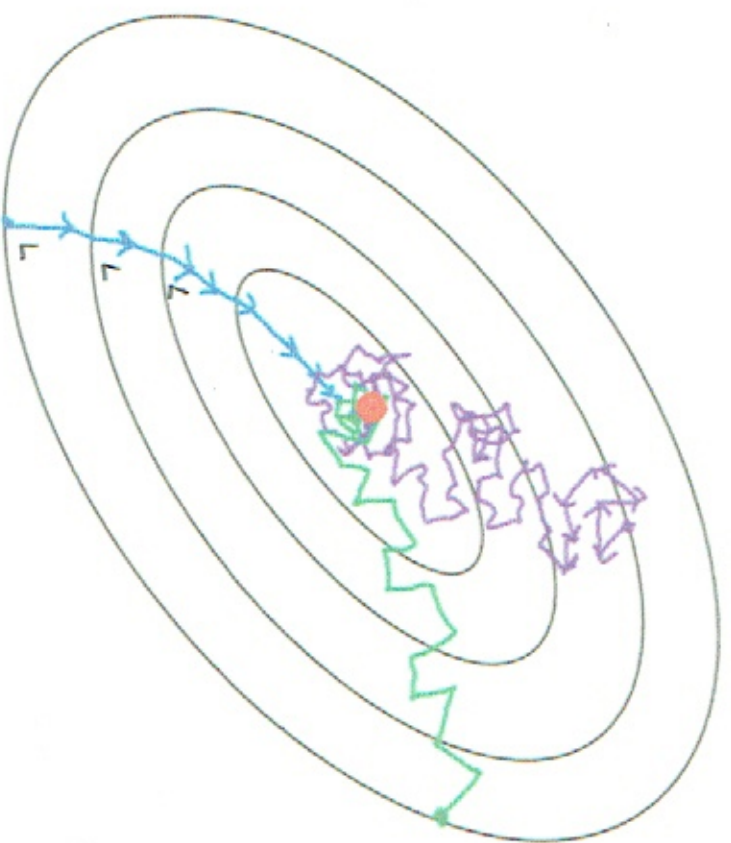


Intuitively
why gradient points INTO
Direction of Steepest Descent

convexity : NOT ALL cost functions are convex
ex: Neural Networks

\rightarrow you can get stuck IN local minimum.





- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

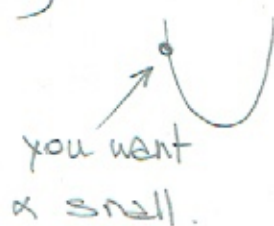
How to improve GD?

8

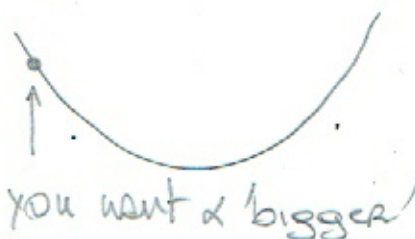
↳ 1) converge FASTER

2) AVOID getting stuck in local minima

3)



you want α small.



you want α 'bigger'

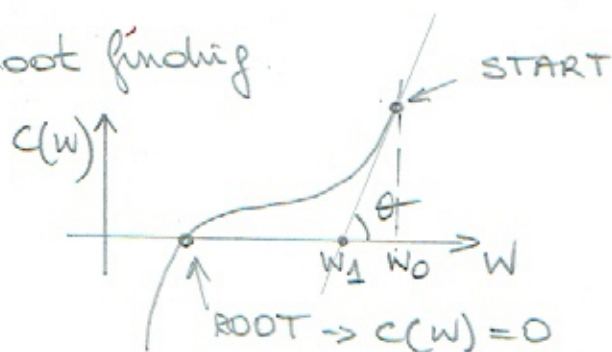
↓
ideally you have a different learning rate α for every feature!

→ you need info over curvature

→ Hessian

Newton's method: 2nd order optimization

↳ ① Root finding



$$\tan \theta = \frac{C(w_0)}{w_0 - w_1} \quad (1)$$

$$\downarrow$$
$$\frac{dC}{dw_0} = C'(w_0) \quad (2)$$

$$(1) \text{ and } (2) \rightarrow C'(w_0) = \frac{C(w_0)}{w_0 - w_1}$$

$$\rightarrow (w_0 - w_1) C'(w_0) = C(w_0)$$

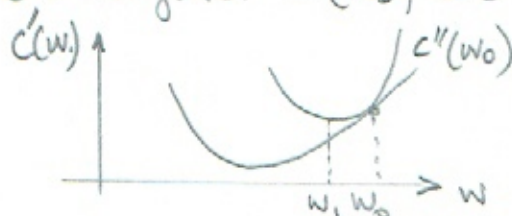
$$\rightarrow w_0 C'(w_0) - w_1 C'(w_0) = C(w_0) - w_0 C'(w_0)$$

$$\rightarrow w_1 C'(w_0) = w_0 C'(w_0) - C(w_0)$$

$$\rightarrow w_1 = \frac{w_0 C'(w_0) - C(w_0)}{C'(w_0)} \rightarrow \boxed{w_1 = w_0 - \frac{C(w_0)}{C'(w_0)}}$$

Analogy

→ we want to find $C'(w_0) = 0$



$$\rightarrow \boxed{w_1 = w_0 - \frac{C'(w_0)}{C''(w_0)}}$$

NEWTON converges faster!

computationally more expensive $\rightarrow \mathcal{O}^{-1}$
no hyperparameter \propto

↓
better when # features is small ($d < 1000$)

STOCHASTIC GRADIENT DESCENT

Full batch GD \rightarrow for every little step \rightarrow massive simult. calculations

$$MSE = c(w) = \frac{1}{m} \sum_{i=1}^m (w^T x^{(i)} - y^{(i)})^2$$

$\hookrightarrow m = 300 \cdot 10^6$ TRAINING samples

IDEA:

- ① shuffle dataset
- ② execute GD on 1 sample and update!
- ③ Repeat for next sample



↓ Noisy!

SGD doesn't really converge but ok if we are near minimum

Mini-batch Gradient Descent

\hookrightarrow GOTO method for large datasets

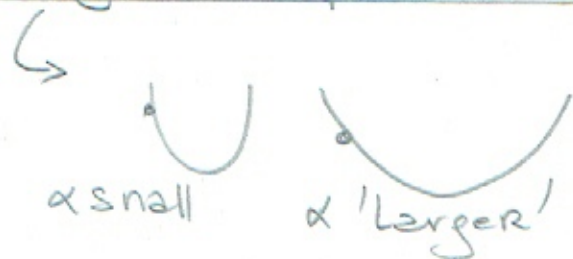
- ① split data in mini-batches and process batch per batch!

↓ size batch = hyperparameter!

typically 64, 128, 256, 512, ...

Optimizers

Adagrad: Adaptive Gradient (2011)



↳ dynamically adjust α for each update and for every weight individually

$$w_{i+1} = w_i - \frac{\alpha}{\sqrt{v_i + \epsilon}} \nabla C(w_i) \quad (1)$$

↳ Avoid division by zero

$$v_i = v_{i-1} + (\nabla C(w_i))^2$$

START: $v_0 = 0$

$$v_1 = v_0 + (\nabla C(w_1))^2 = (\nabla C(w_1))^2$$

$$v_2 = v_1 + (\nabla C(w_2))^2 = (\nabla C(w_1))^2 + (\nabla C(w_2))^2$$

$$v_i = (\nabla C(w_1))^2 + (\nabla C(w_2))^2 + \dots + (\nabla C(w_i))^2 \quad (2)$$

(1) \times (2) $\Rightarrow - \frac{\alpha}{\sqrt{v_i + \epsilon}}$

↳ keeps squares of past gradients

↳ aggressively decays learning rate α

RMSprop:

↳ Avoid rapid growth of denominator

$$w_{i+1} = w_i - \frac{\alpha}{\sqrt{v_i + \epsilon}} \nabla C(w_i)$$

EMA of squares of past gradients

$$v_i = \beta v_{i-1} + (1-\beta) (\nabla C(w_i))^2$$

$\beta = 0.9$

impact reduced!
→ less aggressive than ADAGRAD

MATH INTERPREZZO: Exponential weighted (moving) Average EMA ^②

$$EMA_{today} = \beta EMA_{yesterday} + (1-\beta) P$$

current value
ex: stock price

Example:

4 2 5 3 1
← past days
6
↓ today

β defines how much weight
is given to historical
DATA.

$$EMA_{today} = \frac{2+5+3+1+6}{5} = 17/5 = 3.4$$

$$EMA_{yesterday} = \frac{4+2+5+3+1}{5} = \frac{15}{5} = 3$$

$$\beta = 0.9$$

$$EMA_{today} = 2.7 + 0.1 \times 6 = \underline{\underline{3.3}}$$

Momentum:

uses EMA of past
gradients to
update weights



we want to clarpen sideways movements

we want faster learning along this direction

start, release ball.

2) ball gains speed and once it has velocity it doesn't
do steepest descent anymore because momentum
will keep ball going in previous direction!

3) momentum will allow ball to deal with 'smaller'
local minima.

ISSUE: overshoots minimum if momentum high! → could oscillate

Momentum

$$W_{i+1} = W_i - \underbrace{\beta \text{update}_{i-1}}_{\text{momentum coeff. move according to history}} - \underbrace{\alpha \nabla C(W_i)}_{\text{move according to current GRAD}}$$

0.9

learning rate

2 effects

momentum coeff.
move according
to history.

move according to
current GRAD

START: $update_0 = 0$

$$update_1 = \beta update_0 + \alpha \nabla C(w_1) = \alpha \nabla C(w_1)$$

$$update_2 = \alpha \beta \nabla C(w_1) + \alpha \nabla C(w_2)$$

$$update_i = \beta^{i-1} \alpha \nabla C(w_1) + \beta^{i-2} \alpha \nabla C(w_2) + \dots + \alpha \nabla C(w_i)$$

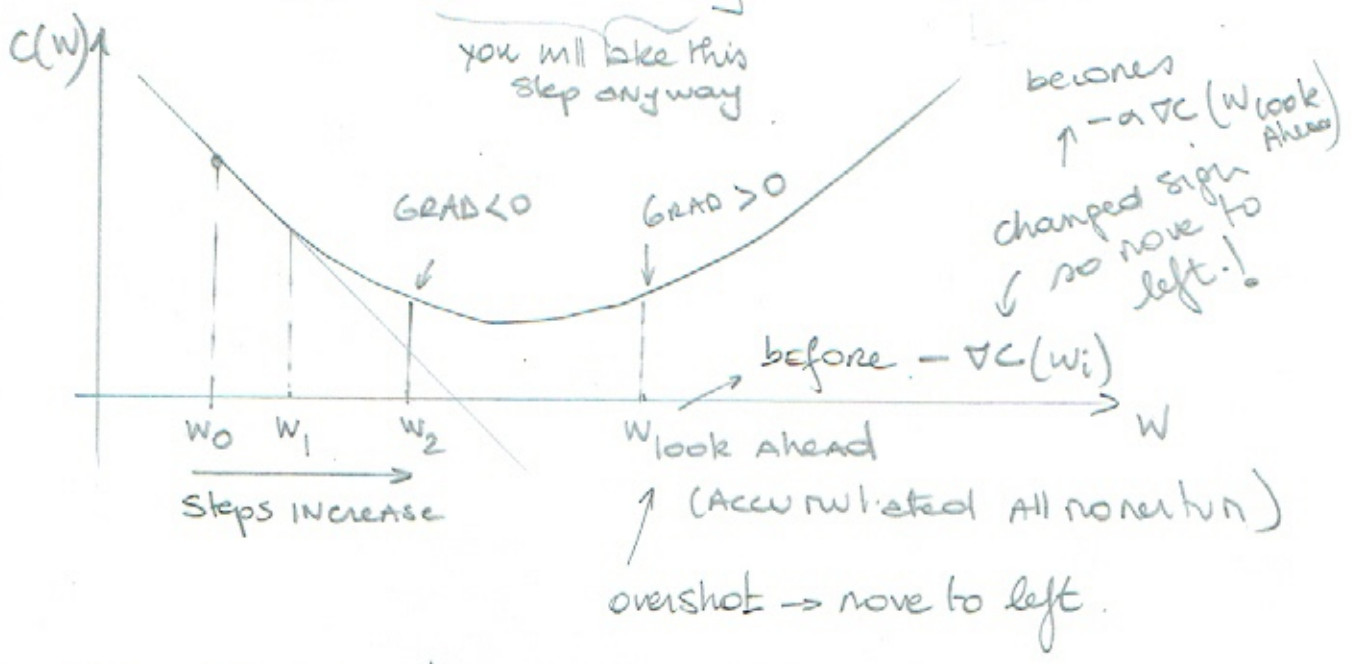
influence ↓ if we go along!

as you go faster

→ take larger steps!

NESTEROV

$$w_{i+1} = \underbrace{[w_i - \text{history}]}_{\text{momentum}} - \underbrace{\alpha \nabla C(w_i)}_{\text{current GRAD}}$$



ADAM: VERY POPULAR!

↳ momentum + RMS prop

EMAS → $m_i = \beta_1 m_{i-1} + (1-\beta_1) \nabla C(w_i)$ ← Accumulates history

$v_i = \beta_2 v_{i-1} + (1-\beta_2) \nabla C(w_i)$ ← takes care of α

$\beta_1 = 0.9$
 $\beta_2 = 0.999$
 $\epsilon = 10^{-8}$

$$w_{i+1} = w_i - \frac{\alpha}{\sqrt{\hat{v}_i} + \epsilon} \cdot \hat{m}_i$$

$$\hat{m}_i = \frac{m_i}{1-\beta_1^i}$$

$$\hat{v}_i = \frac{v_i}{1-\beta_2^i}$$

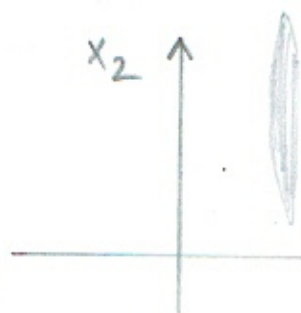
↳ hyperparameters

Feature Scaling / Normalizing

Python: sklearn \rightarrow StandardScaler

ISSUE: $x_1: 0 \rightarrow 5$

$x_2: 0 \rightarrow 20,000$



VARIANCE $x_1 \ll$ VARIANCE x_2

\rightarrow stretched contour lines

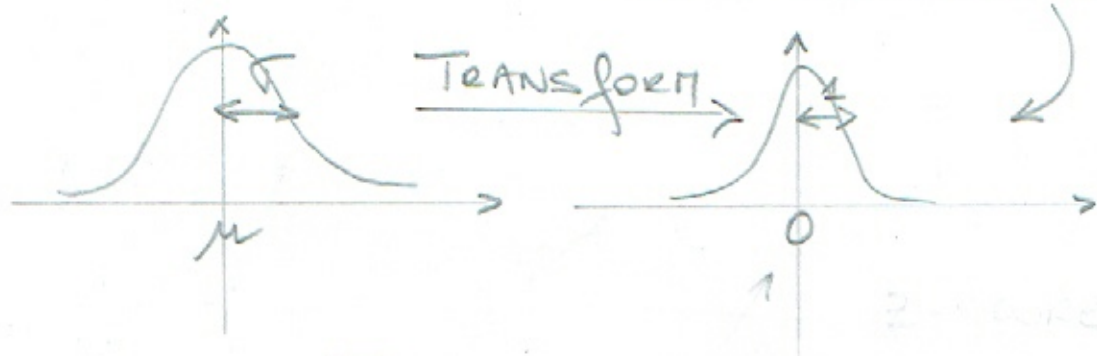
\rightarrow slows down convergence

\downarrow solution

good practice
in machine learning

\leftarrow NORMALIZATION

MEAN IS ZERO $= \mu$
VARIANCE IS 1 $= \sigma^2$



$$z = \frac{x - \mu}{\sigma}$$

\hookrightarrow use same μ, σ for
NORMALIZATION of test data