



Northeastern University
Khoury College of
Computer Sciences

Recap

DS 4400 | Machine Learning and Data Mining I

Zohair Shafi

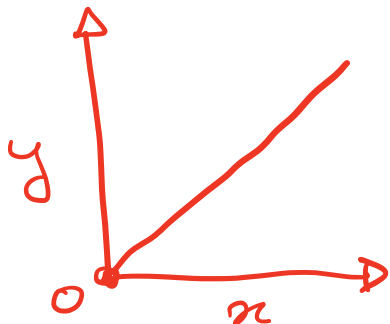
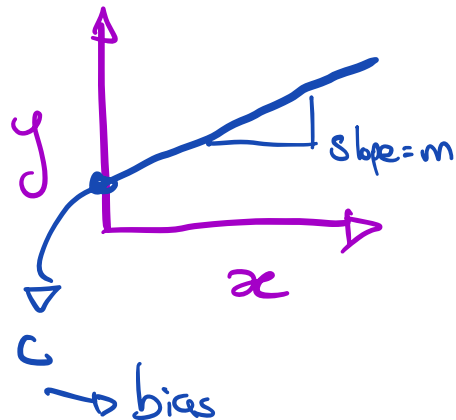
Spring 2026

Monday | January 26, 2026

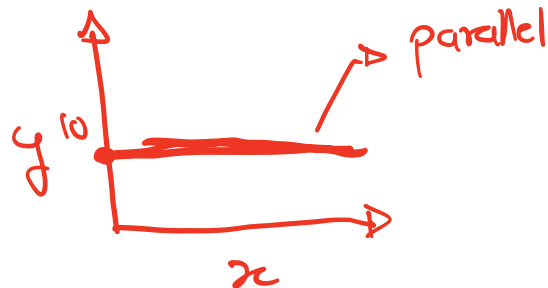
Linear Regression

$$y = mx + c$$

↓ ↓
Slope intercept



$$y = mx + c$$
$$(1)x + 0$$



$$y = mx + c$$
$$y = 0 \cdot x + c$$

Linear Regression

$$y = mx + c$$

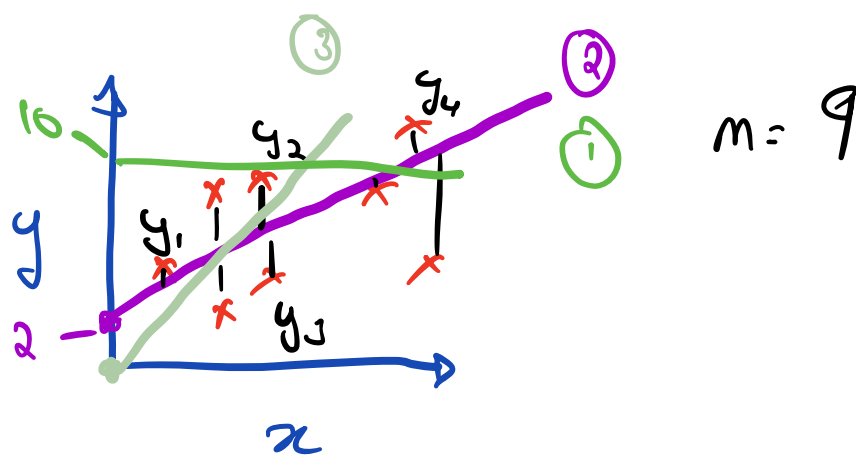
$$\hat{y} = \theta_1 x + \theta_0$$

input

predicted/
output

$$\text{Loss} = \frac{1}{M} \sum_{i=0}^{M-1} (y_i - \hat{y}_i)^2$$

Cost Error \mathcal{L}, J, C



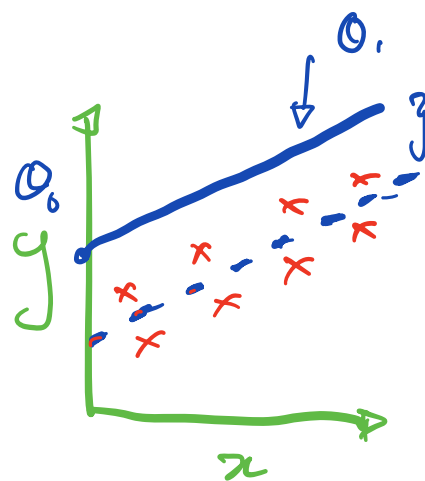
- ① $\hat{y} = \theta_1 x + \theta_0$ — $\theta_1 = 0$
 $\theta_0 = 10$
- ② $\hat{y} = \theta_1 x + \theta_0$ — $\theta_1 = 2$
 $\theta_0 = 2$
- ③ $\hat{y} = \theta_1 x + \theta_0$ — $\theta_1 = 4$
 $\theta_0 = 0$

Linear Regression

Model: Line \rightarrow

$$\hat{y} = \theta_0 + \theta_1 \cdot x$$

function $f_{\theta}(x)$
 θ parameters
 x input



Minimize

Loss:

$$L_{\theta}(x) = \frac{1}{n} \sum_{i=0}^n (y - \hat{y})^2$$

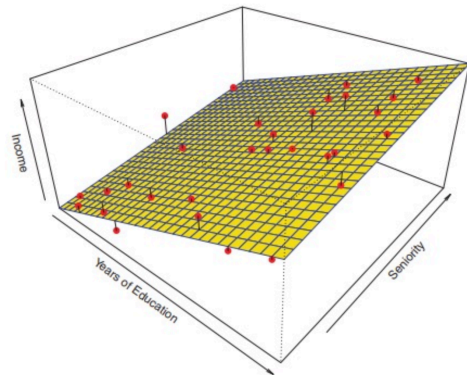
$$= \frac{1}{n} \sum_{i=0}^n (y - (\theta_0 + \theta_1 x))^2$$

- ① Compute derivative
- ② Set derivative to zero
- ③ Solve for θ

Linear Regression

- Linear Model

$$f_{\theta}(x) = \theta_0 + \theta_1 x_0 + \theta_2 x_1$$

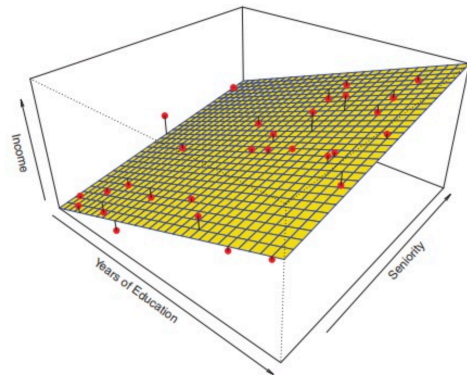


Linear Regression

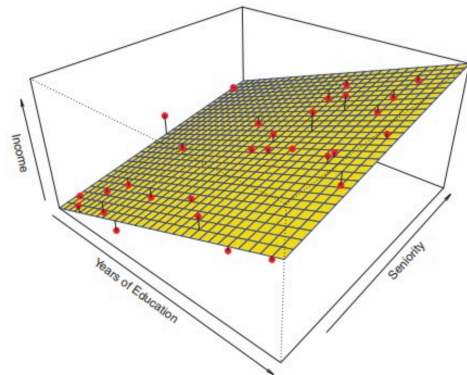
- Linear Model

$$f_{\theta}(x) = \theta_0 + \theta_1 x_0 + \theta_2 x_1$$

Learnable parameters



Linear Regression



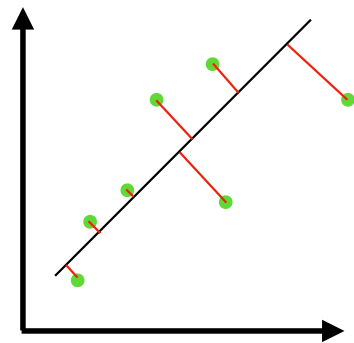
- Linear Model

$$f_{\theta}(x) = \theta_0 + \theta_1 x_0 + \theta_2 x_1$$

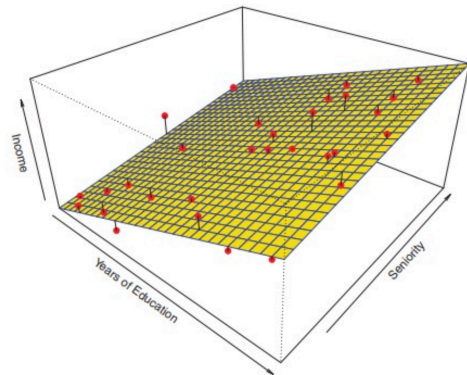
- Loss Functions (also called Cost Functions)

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [f_{\theta}(x_i) - y_i]^2 - \text{Mean Squared Error}$$

The red lines are called **residuals**



Linear Regression



- Linear Model

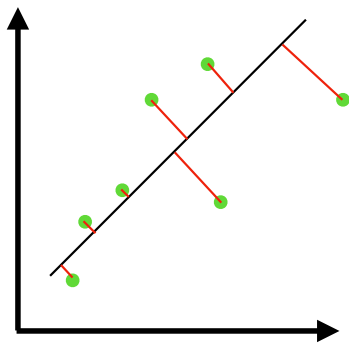
$$f_{\theta}(x) = \theta_0 + \theta_1 x_0 + \theta_2 x_1$$

- Loss Functions (also called Cost Functions)

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [f_{\theta}(x_i) - y_i]^2 - \text{Mean Squared Error}$$


$$L(\theta) = \sum_{i=1}^m [f_{\theta}(x_i) - y_i]^2 - \text{Residual Sum of Squares}$$

The red lines are called **residuals**



Linear Regression

- Linear Model

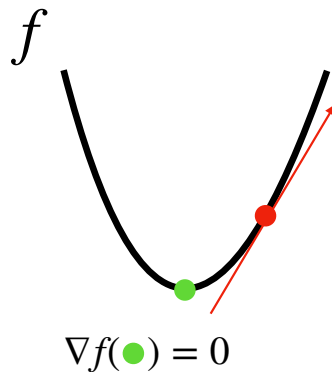
$$f_{\theta}(x) = \theta_0 + \theta_1 x$$


$\nabla f(\bullet)$ points in direction of steepest ascent

- How do we find the solution to this? How do we find the optimal θ ?
 - We optimize θ to minimize the loss function

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [f_{\theta}(x_i) - y_i]^2$$

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [\theta_0 + \theta_1 \cdot x - y_i]^2$$



Linear Regression

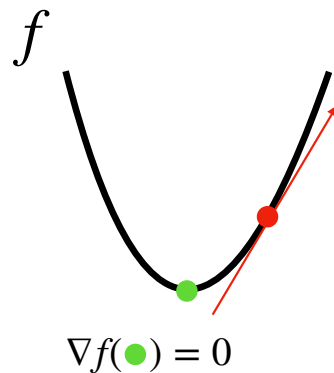
$$L_{\theta}(x) = \frac{1}{n} \sum (y - (\theta_0 + \theta_1 x))^2$$

$$\frac{\partial}{\partial \theta_1} L(\theta_1, x) = x$$

$$\frac{\partial L}{\partial \theta_0} = 2 \cdot \frac{1}{n} \sum (y - (\theta_0 + \theta_1 x))(-1)$$

$$\frac{\partial L}{\partial \theta_1} = 2 \cdot \frac{1}{n} \sum (y - (\theta_0 + \theta_1 x))(\underline{x})$$

$\nabla f(\bullet)$ points in direction of steepest ascent



Linear Regression

$$y = \boxed{\theta_0} + \boxed{\theta_1 x}$$

- How do we find the solution to this? How do we find the optimal θ ?
- We optimize θ to minimize the loss function

$$\theta_0 \rightarrow \frac{\partial L}{\partial \theta_0} = 0$$

$$\theta_1 \rightarrow \frac{\partial L}{\partial \theta_1} = 0$$

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [f_{\theta}(x_i) - y_i]^2$$

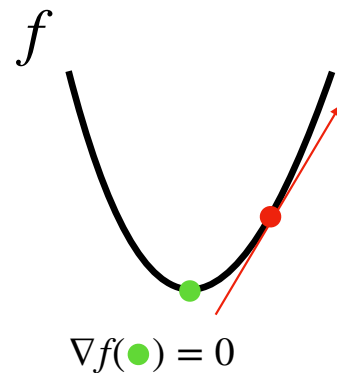
$$L(\theta) = \frac{1}{m} \sum_{i=1}^m [\theta_0 + \theta_1 \cdot x_i - y_i]^2$$

Find the point where $\nabla L(\theta) = 0$

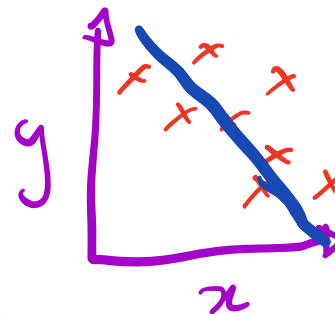
$$\frac{\partial L(\theta)}{\partial \theta_0} = \frac{2}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i) = 0$$

$$\frac{\partial L(\theta)}{\partial \theta_1} = \frac{2}{m} \sum_{i=1}^m x_i \cdot (\theta_0 + \theta_1 x_i - y_i) = 0$$

$\nabla f(\bullet)$ points in direction of steepest ascent



Linear Regression

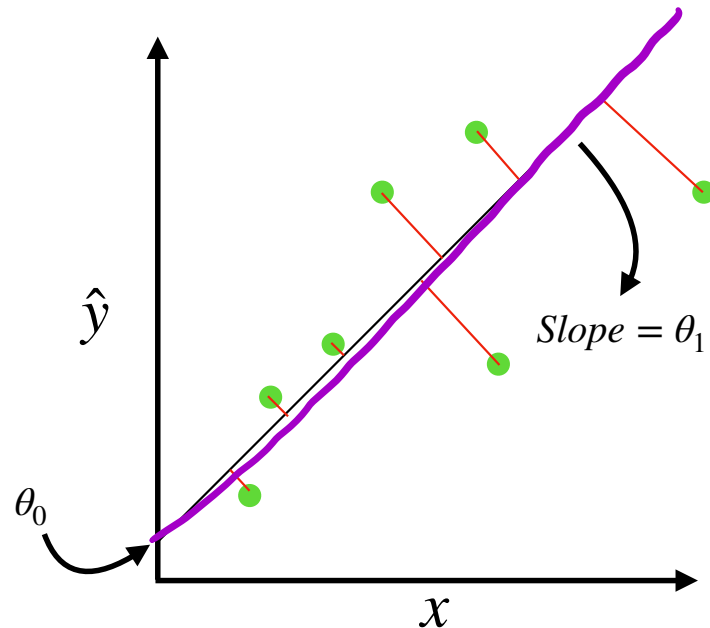


$$\begin{aligned}\theta_0 &= \bar{y} - \theta_1 \bar{x} \\ \theta_1 &= \frac{\text{Cov}(x, y)}{\text{Var}(x)}\end{aligned}$$

The slope $\theta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$ makes sense:

- If x and y covary strongly (move together), the slope is steeper
- If x has high variance (spread out), the slope is gentler
- The **sign** of covariance determines if the line goes up or down

$$f_{\theta}(x) = \theta_0 + \theta_1 x$$



$$Loss \rightarrow \frac{1}{n} \sum (y - \hat{y})^2$$

$$\frac{1}{n} \sum (y - x\theta)^2$$

Linear Regression

Solutions in Matrix Form

Predict y	Input x
10	3
15	5

Loss

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\hat{y} = \theta_0 + \theta_1 (3) = 10$$

$$\hat{y} = \theta_0 + \theta_1 (5) = 15$$

$$\begin{cases} \hat{y} = \theta_0 (1) + \theta_1 (3) = 10 \\ \hat{y} = \theta_0 (1) + \theta_1 (5) = 15 \end{cases}$$

$$\theta_0 (1) + \theta_1 (3) = 10$$

$$\theta_0 (1) + \theta_1 (5) = 15$$

Train data.

θ ~ # param $n \times 1$

$$\begin{bmatrix} 1 & 3 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} 10 \\ 15 \end{bmatrix}$$

$X \in n \times n$
train # param

$y = \# \text{ train } \times 1$
 $n \times 1$

Linear Regression

Solutions in Matrix Form

- Let's look at the matrix formulation of the same problem

$$L(\theta) = \frac{1}{m} \sum_i (y_i - \hat{y}_i)^2$$

But in matrix form, $f_{\theta}(x) = \hat{Y} = X\theta$, where $X \in \mathbb{R}^{m \times d}$ has m rows of data and d columns of features and $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \in \mathbb{R}^{d \times 1}$

$$L(\theta) = \frac{1}{m} \sum (Y - X\theta)^2$$

(think back to system of equations for why this is true)

Quick Recap

Systems of Linear Equations - Linear Regression Example

- Consider the equation $y = w_0 x_0 + w_1 x_1$

y x_0 x_1

Price	# Rooms	Sq. Ft.
2000	1	450
2100	1	510
2400	2	980
3000	3	1500

$$\begin{aligned} (1) \cdot w_0 + (450) \cdot w_1 &= 2000 \\ (1) \cdot w_0 + (510) \cdot w_1 &= 2100 \\ (2) \cdot w_0 + (980) \cdot w_1 &= 2400 \\ (3) \cdot w_0 + (1500) \cdot w_1 &= 3000 \end{aligned}$$



$$\begin{aligned} X &\in \mathbb{R}^{4 \times 2} & W &\in \mathbb{R}^{2 \times 1} & y &\in \mathbb{R}^{4 \times 1} \\ \begin{bmatrix} 1 & 450 \\ 1 & 510 \\ 2 & 980 \\ 3 & 1500 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} &= \begin{bmatrix} 2000 \\ 2100 \\ 2400 \\ 3000 \end{bmatrix} \end{aligned}$$

$$1(w_0) + (450)(w_1) = 2000$$

Linear Regression

Solution

$$L_{\theta} = \frac{1}{n} \sum (y - x\theta)^2$$

$$\nabla L_{\theta}$$

We want to find the minimum so set gradient to zero

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$\nabla L(\theta) = - \boxed{2X^T Y + 2X^T X \theta = 0}$$

$$2X^T X \theta = 2X^T Y$$

$$X^T X \theta = X^T Y$$

$$\theta_1 = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

If $X^T X$ is invertible, then

$$\boxed{\theta = (X^T X)^{-1} X^T Y}$$

Practical Issues in Linear Regression

Multicollinearity

- When two features are highly correlated or are linearly dependent on each other

$$5 \cdot \frac{1}{5} = 1$$


$$A \cdot A^{-1} = I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathcal{O} = \underbrace{(x^T x)}_{\substack{\text{full rank matrix} \\ \text{No linearly dependent} \\ \text{rows or columns.}}} \boxed{-1} x^T y$$

Practical Issues in Linear Regression

Multicollinearity

x_0 x_1

- When two features are highly correlated or are linearly dependent on each other
- Why it's a problem: $\theta = (X^T X)^{-1} X^T Y$
 - $X^T X$ becomes nearly singular (ill-conditioned) 
 - Small changes in data cause huge changes in coefficients
 - Coefficients become unreliable and hard to interpret
 - Standard errors blow up

cannot compute
inverse.

Practical Issues in Linear Regression

Multicollinearity

$x_0 \rightarrow x_{10}$

- When two features are highly correlated or are linearly dependent on each other



- Why it's a problem:

$$\theta = (X^T X)^{-1} X^T Y$$

Simple Detection:
If correlation between features ≥ 0.8

- $X^T X$ becomes nearly singular (ill-conditioned)
- Small changes in data cause huge changes in coefficients
- Coefficients become unreliable and hard to interpret
- Standard errors blow up

Practical Issues in Linear Regression

Quick Aside

$$\theta = (X^T X)^{-1} X^T Y$$

When else is this not going to be invertible?

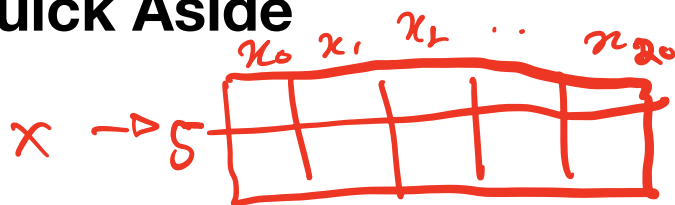
$$X \in \mathbb{R}^{m \times n}$$

m : Number of training examples

n : Number of parameters in the model

Practical Issues in Linear Regression

Quick Aside

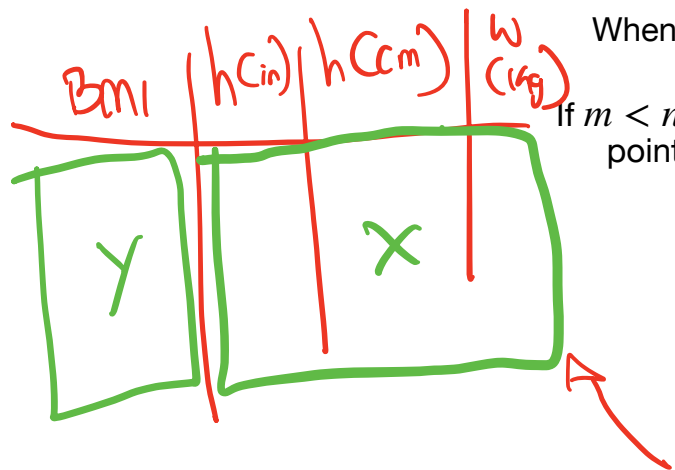


$$X = 5 \times 20$$

$$\theta = (X^T X)^{-1} X^T Y$$

When else is this not going to be invertible?

If $m < n$, then $\text{rank}(X) < m$, so need **more** data points than number of parameters to get a unique set of parameters



$$X \in \mathbb{R}^{m \times n}$$

m : Number of training examples

n : Number of parameters in the model

$$\text{rank}(X) = \min(m, n)$$

Lin = 2.54 cm

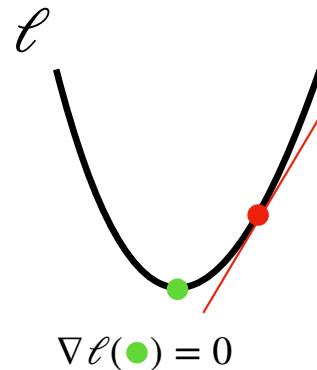
$y = 2.54 \text{ cm}$

in

Gradient Descent: Optimizing Loss Functions

- For any loss function $\ell(\theta)$
 - To find minimum, set $\nabla \ell = 0$ and solve for θ

$\nabla \ell(\bullet)$ points in direction of steepest ascent

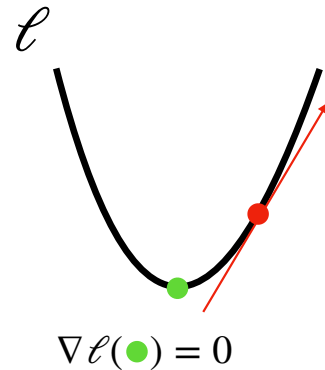


Optimizing Loss Functions

$$\underbrace{(X^T X)^{-1}} \rightarrow O(n^3)$$

- For any loss function $\ell(\theta)$
 - To find minimum, set $\nabla \ell = 0$ and solve for θ
 - This is called the **closed form solution**
 - But it's not always possible to find closed form solutions, especially when there are a large number of parameters
 - Inverting a matrix is a costly operation - most common methods have complexity $O(n^3)$

$\nabla \ell(\bullet)$ points in direction of steepest ascent



Optimizing Loss Functions

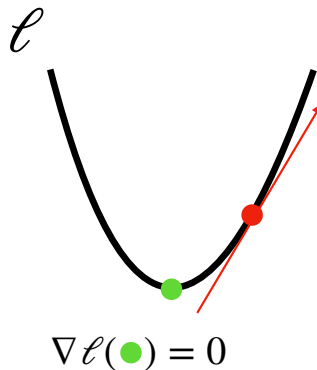
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$O(mTn) \begin{matrix} \text{\# Params} \\ \text{\# data} \quad \text{\# iterations} \end{matrix}$$

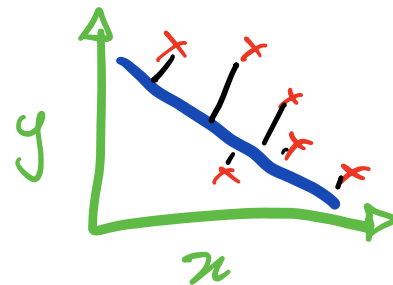
n^3

- This is where Gradient Descent comes in
- Practical and efficient - has $O(mTn)$ where m is number of training points, T is number of epochs and n is number of features
- Generally applicable to different loss functions
- Convergence guarantees for certain types of loss functions (e.g., convex functions)

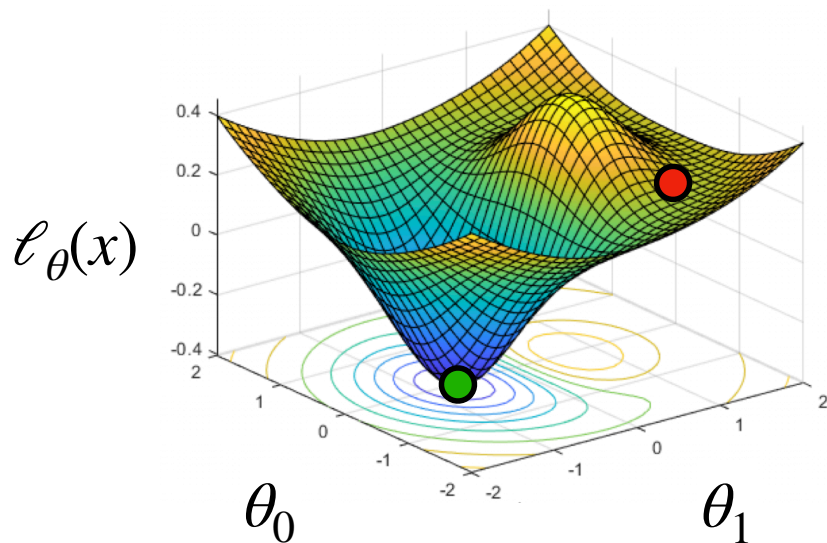
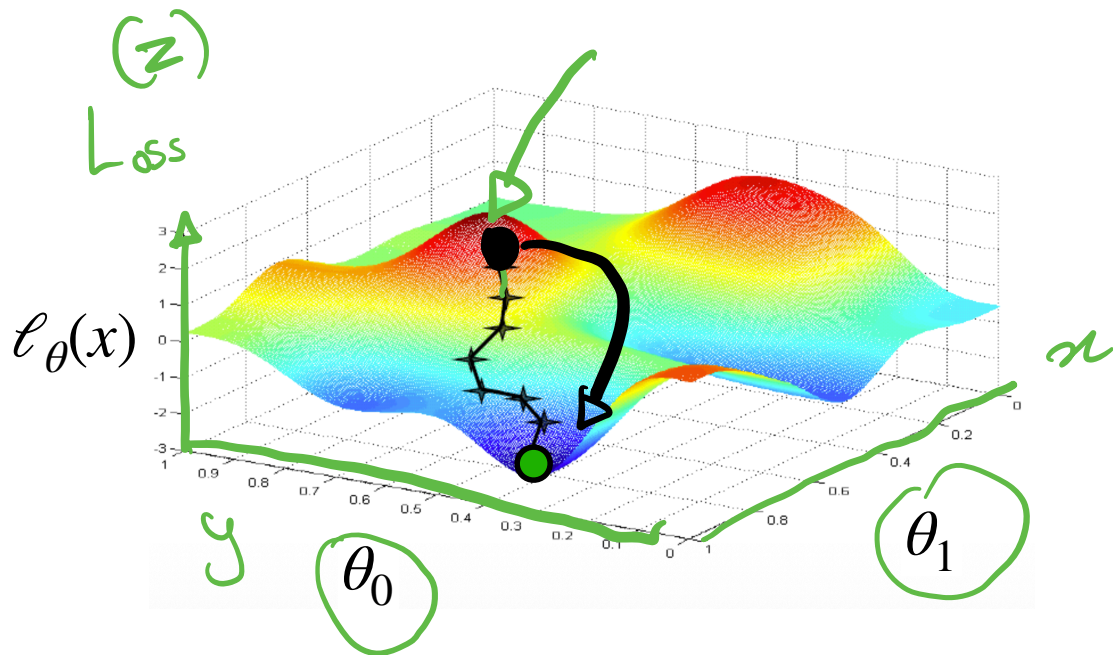
$\nabla \ell(\bullet)$ points in direction of steepest ascent



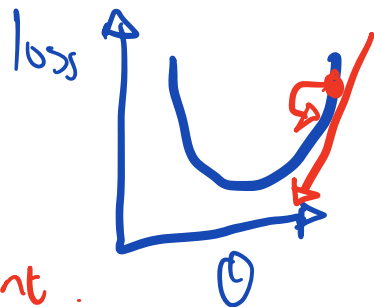
Optimizing Loss Functions



- What does the loss landscape look like with multiple learnable parameters?



Optimizing Loss Functions



- ① Define loss function
- ② Compute derivative.
- ③ Take a single step in the direction of the gradient.

$$① L_{\theta}(x) = \frac{1}{n} \sum (y - (\theta_0 + \theta_1 x))^2$$

$$② \frac{\partial L_{\theta}(x)}{\partial \theta} = \nabla_{\theta} L(x) \begin{cases} \nabla_{\theta_0} L(x) = -\frac{2}{n} \sum (y - (\theta_0 + \theta_1 x)) \\ \nabla_{\theta_1} L(x) = -\frac{2}{n} \sum (y - \theta_0 + \theta_1 x) \cdot x \end{cases}$$

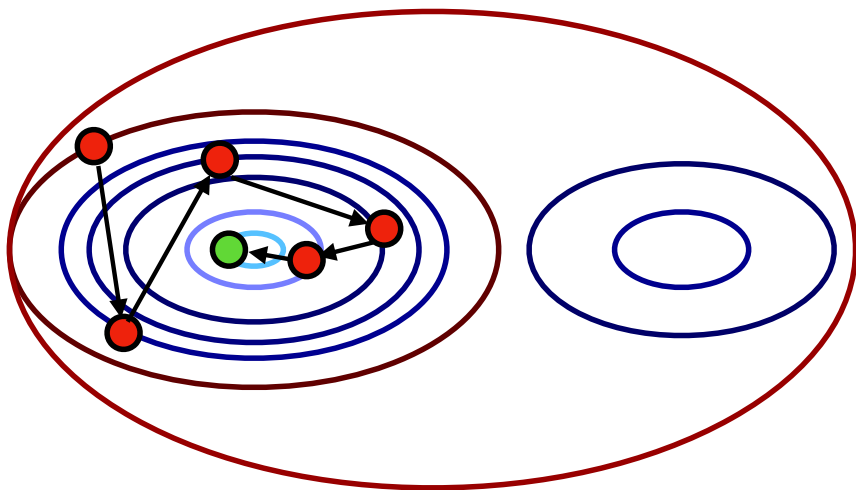
$$③ \begin{aligned} \theta_0 &\leftarrow \theta_0 - \alpha \cdot \nabla_{\theta_0} L(x) \\ \theta_1 &\leftarrow \theta_1 - \alpha \cdot \nabla_{\theta_1} L(x) \end{aligned}$$

for t in range $(0, 100)$
theta: theta - alpha - deriv.

$$\begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \leftarrow \theta_t \leftarrow \theta_{t-1} - \alpha \cdot \nabla_{\theta_{t-1}} L(x)$$

Optimizing Loss Functions

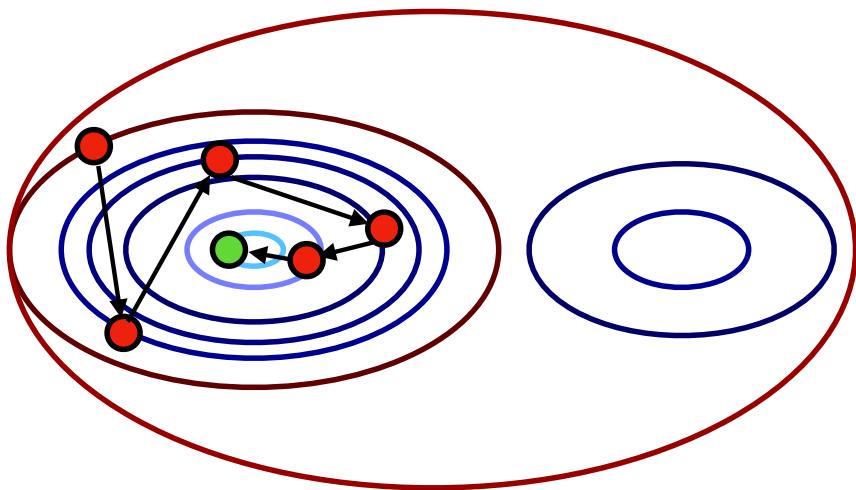
Gradient Descent - Formulation



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Optimizing Loss Functions

Gradient Descent - Formulation

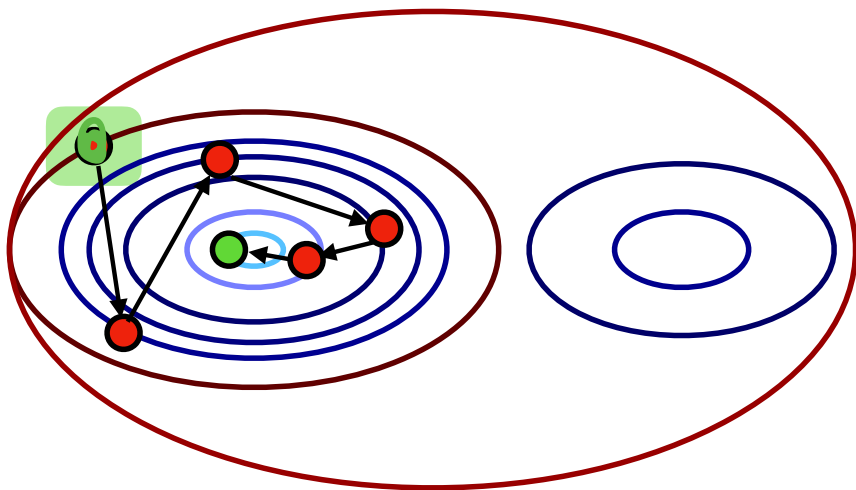


$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \overset{0.5}{\theta_0} - \overset{10}{\theta_1} x_i)^2$$

Step 1: Initialize θ_0, θ_1

Optimizing Loss Functions

Gradient Descent - Formulation



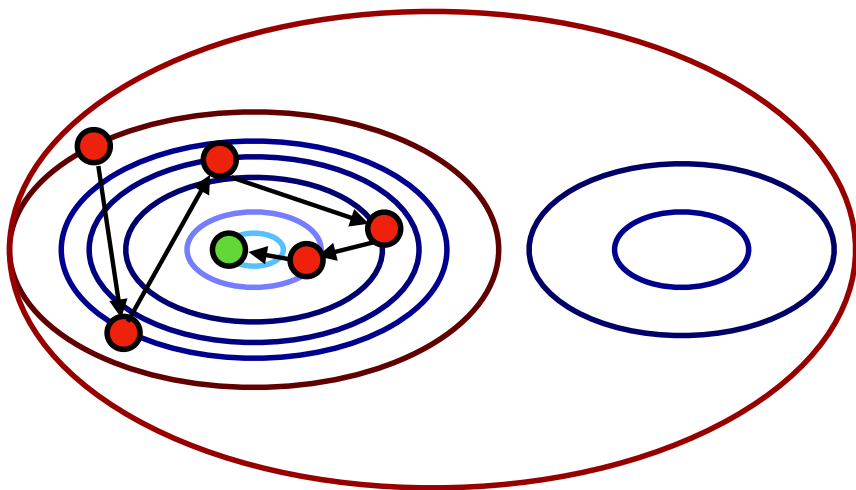
$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

This is going to be your “starting point” on the loss landscape

Optimizing Loss Functions

Gradient Descent - Formulation



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

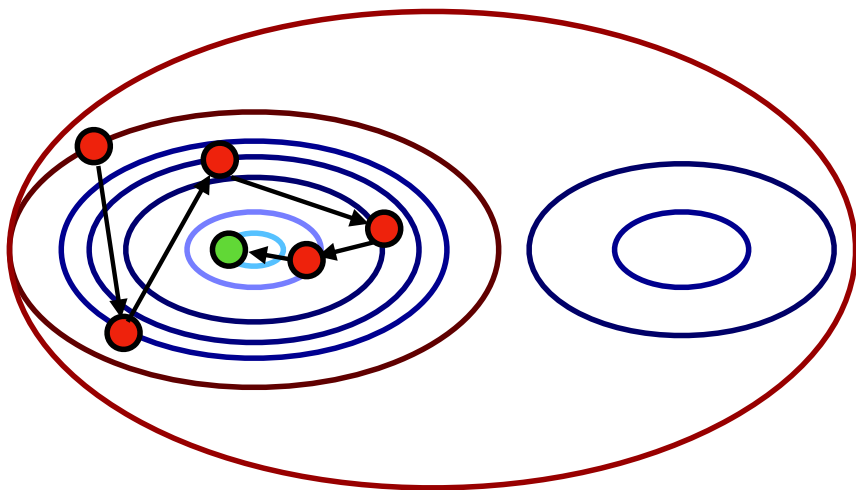
Step 1: Initialize θ_0, θ_1

Step 2: Repeat Until Convergence

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

Optimizing Loss Functions

Gradient Descent - Formulation



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

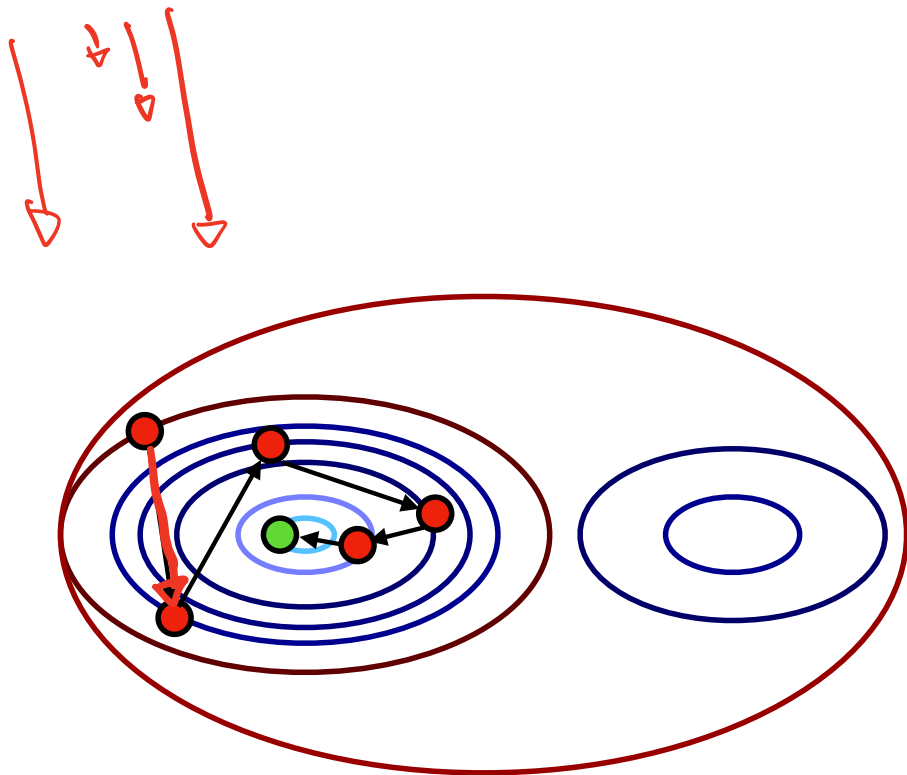
Step 2: Repeat Until Convergence

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

Negative of partial derivative points
in the direction of steepest descent

Optimizing Loss Functions

Gradient Descent - Formulation



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

Step 2: Repeat Until Convergence

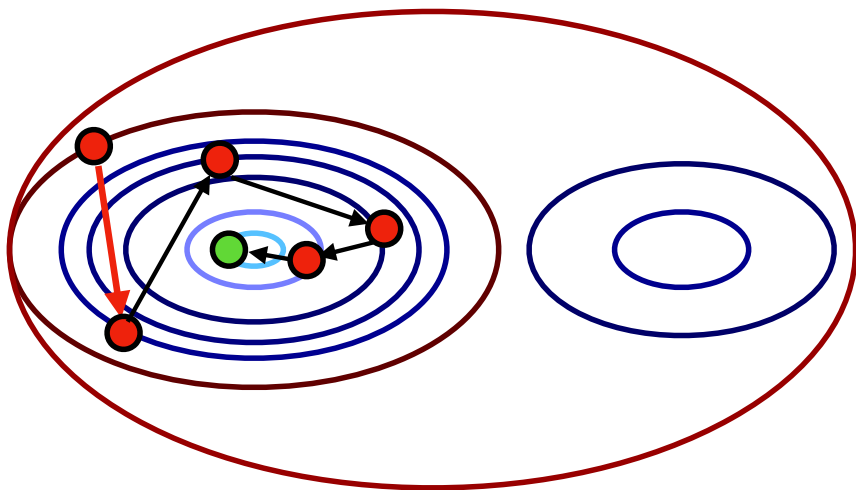
$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

α : Learning Rate

Optimizing Loss Functions

Gradient Descent - Formulation

α controls how big a step to take



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

Step 2: Repeat Until Convergence

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

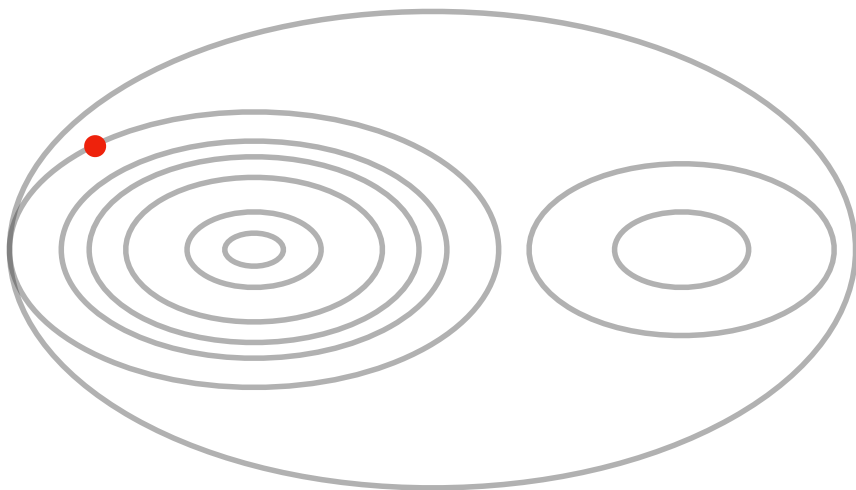
α : Learning Rate

Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

What happens when α is too small?

Say $\alpha = 10^{-5}$



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

Step 2: Repeat Until Convergence

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

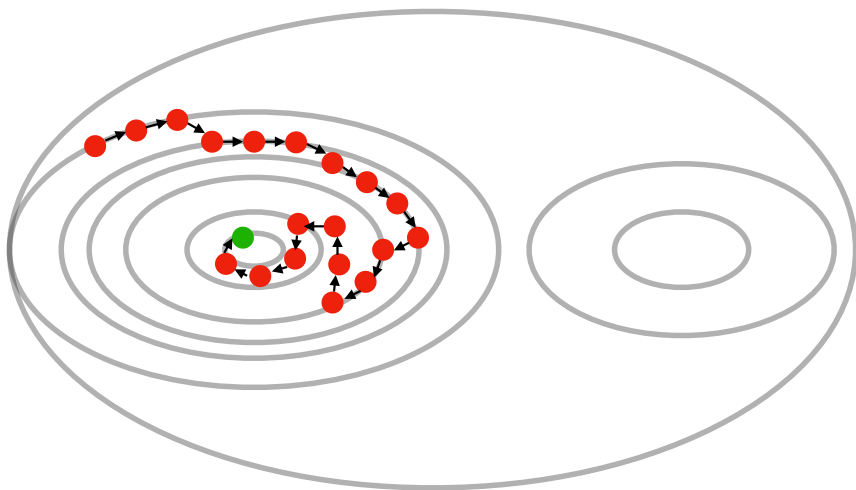
α : Learning Rate

Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

What happens when α is too small?

Say $\alpha = 10^{-5}$



$$\ell_{\theta}(x) = \frac{1}{m} \sum_i (y_i - \theta_0 - \theta_1 x_i)^2$$

Step 1: Initialize θ_0, θ_1

Step 2: Repeat Until Convergence

$$\theta_j \leftarrow \theta_j - \alpha \cdot \frac{\partial \ell_{\theta}(x)}{\partial \theta_j}$$

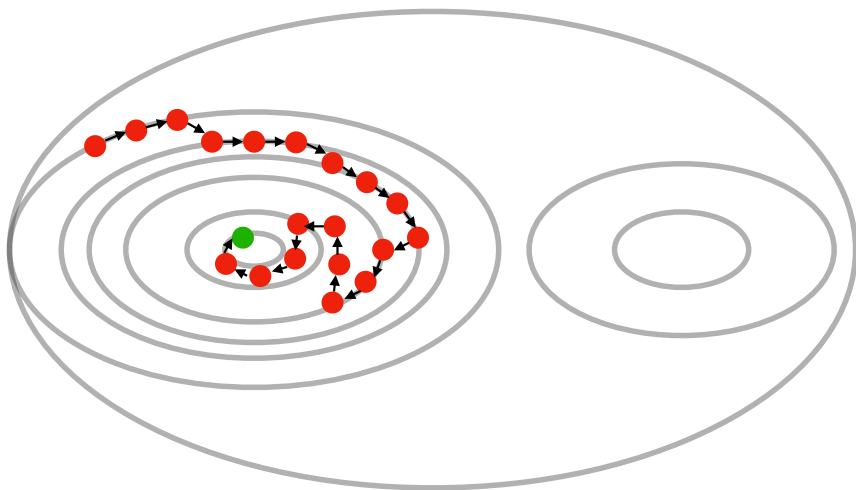
α : Learning Rate

Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

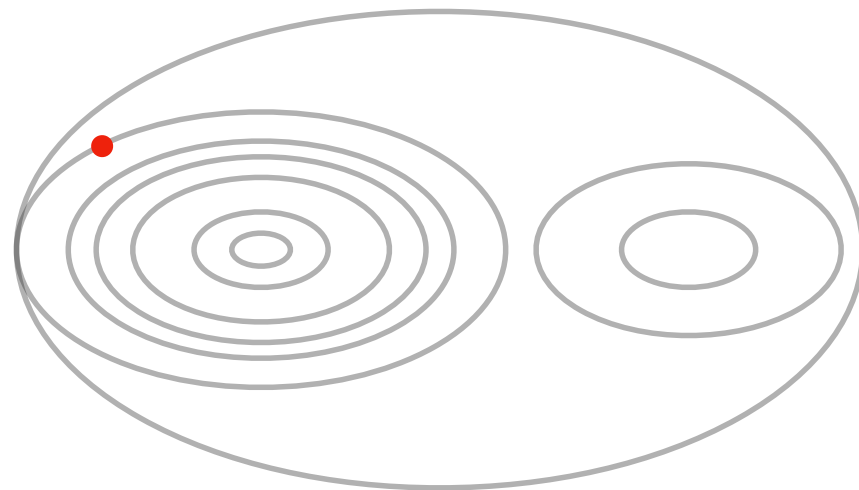
What happens when α is too small?

Say $\alpha = 10^{-5}$



What happens when α is too large?

Say $\alpha = 10$

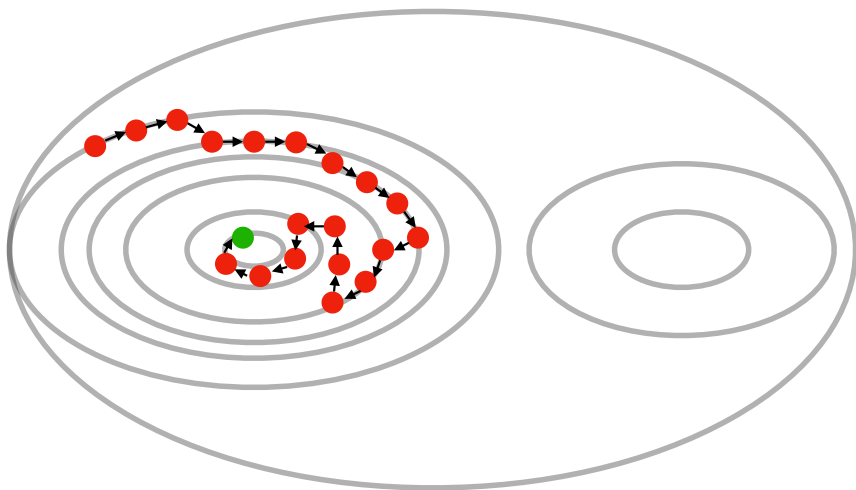


Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

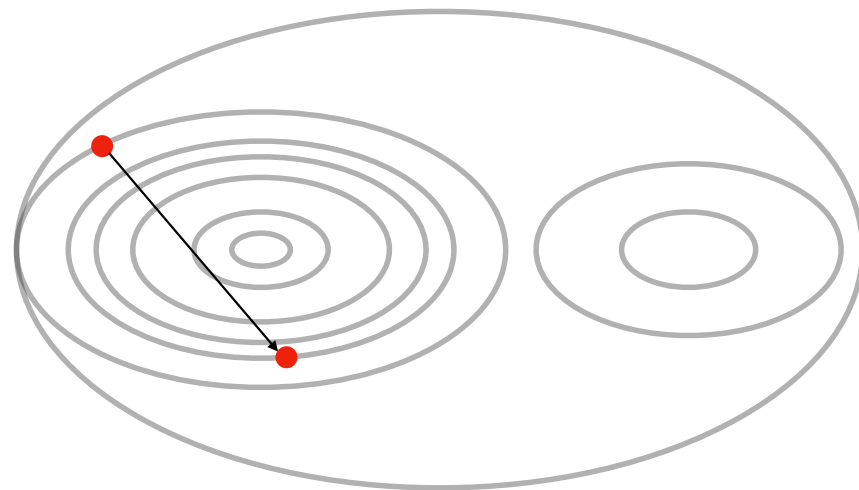
What happens when α is too small?

Say $\alpha = 10^{-5}$



What happens when α is too large?

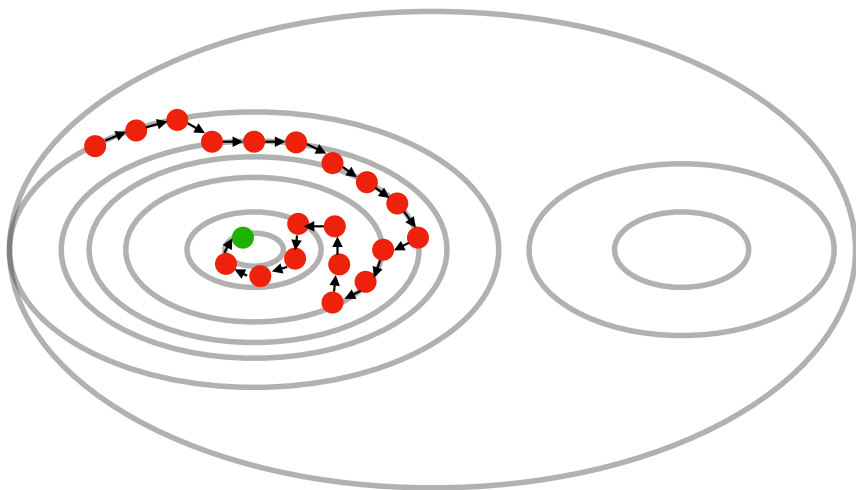
Say $\alpha = 10$



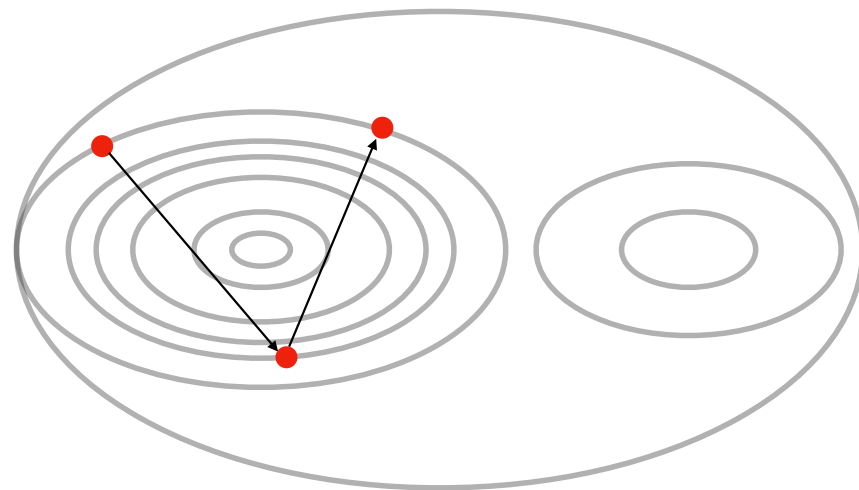
Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

What happens when α is too small?
Say $\alpha = 10^{-5}$



What happens when α is too large?
Say $\alpha = 10$

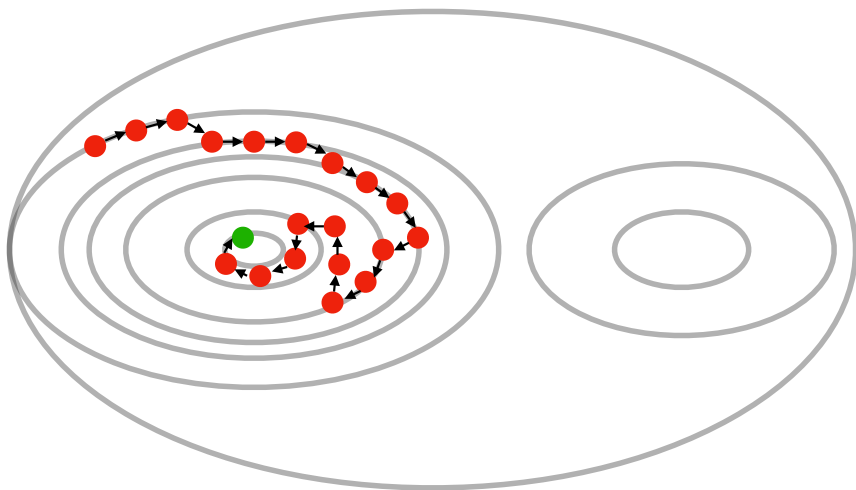


Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

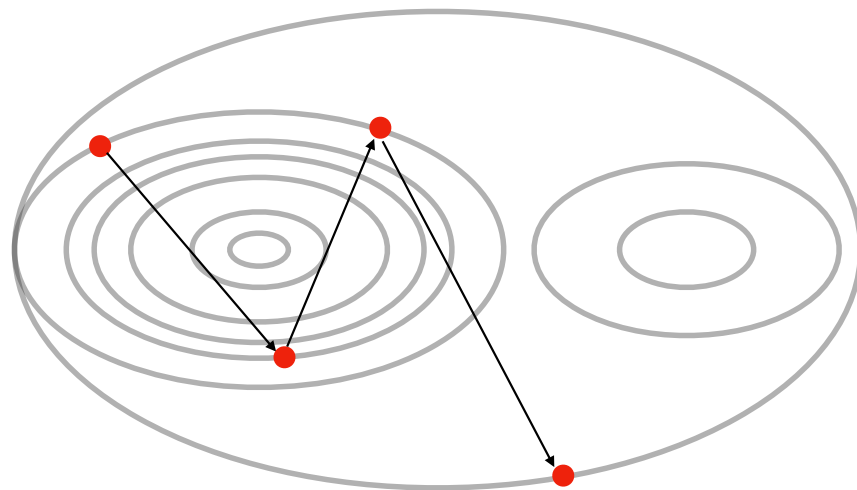
What happens when α is too small?

Say $\alpha = 10^{-5}$



What happens when α is too large?

Say $\alpha = 10$



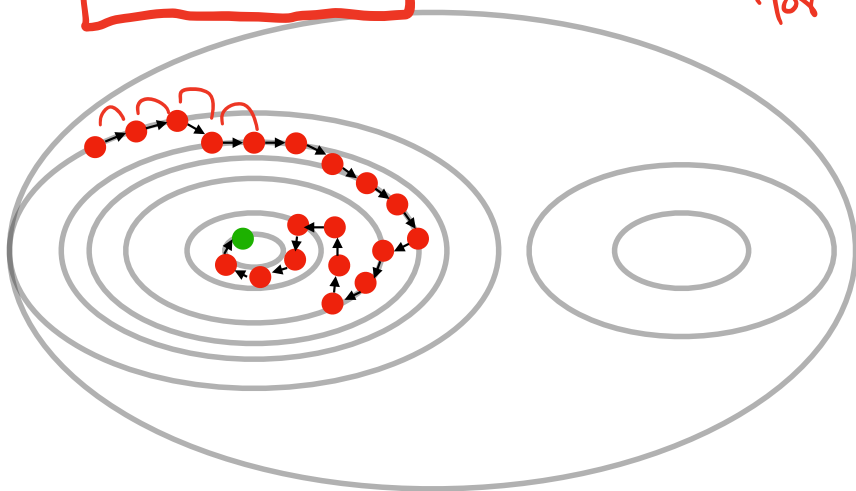
Optimizing Loss Functions

Gradient Descent - Effect of Learning Rate

What happens when α is too small?

Say $\alpha = 10^{-5}$

With a small learning rate α , if the loss function is convex, the optimization will eventually **converge**

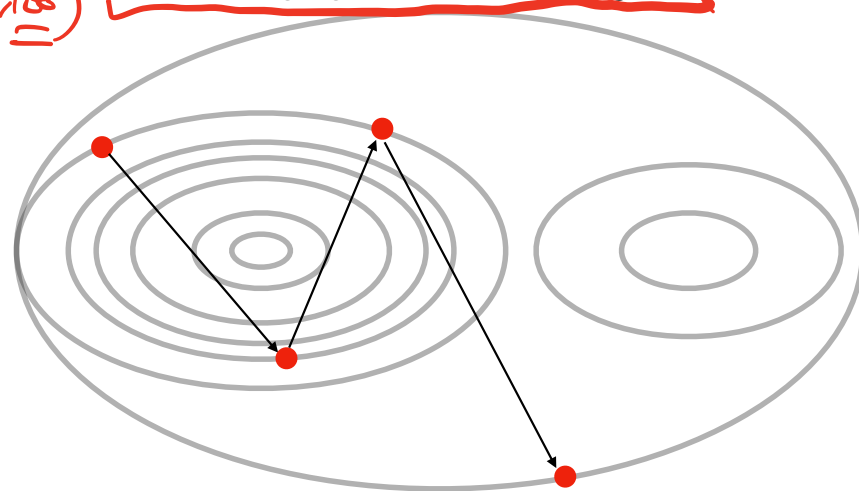


What happens when α is too large?

Say $\alpha = 10$

With a large learning rate α , if the loss function is convex, the optimization could possibly start **diverging** and never converge

for i in (0, 100)



Optimizing Loss Functions

Gradient Descent - Stopping Criterion

$$\theta_t \leftarrow \theta_{t-1} - \underbrace{\alpha}_{10^{-3}} \cdot \underbrace{\nabla L(\theta_{t-1})}_{\text{gradient}}$$

Maximum Iteration

Gradient Norm Threshold

Function Value Change

Parameter Value Change

for i in range(1000)

epochs

1 whole pass through
dataset.

$$\|\nabla L(\theta_{t-1})\| < \epsilon$$

\downarrow
 10^{-2}

Loss.

$$|L(\theta_t) - L(\theta_{t-1})| \leq \epsilon$$

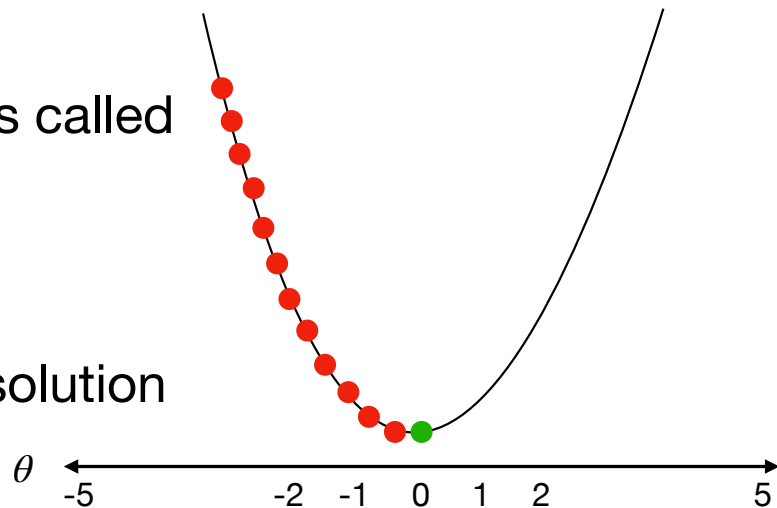
\downarrow
 10^{-2}

$$|\theta_t - \theta_{t-1}| \leq \epsilon$$

Optimizing Loss Functions

Gradient Descent - Stopping Criterion

- When do you stop your iterations?
 - Maximum Iteration]
 - Each iteration through the training dataset is called an “epoch”
 - Terminate after a fixed number of epochs
 - Simple, but provides no guarantees about solution quality



Optimizing Loss Functions

Gradient Descent - Stopping Criterion

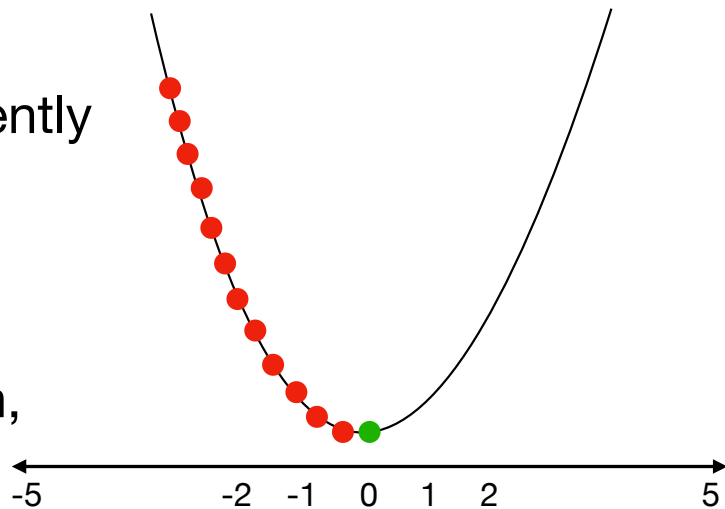
- When do you stop your iterations?

- Gradient Norm Threshold

- Terminate when the gradient becomes sufficiently small

$$\|\nabla \ell_{\theta}(x)\|_2 \leq \epsilon$$

- At this point, if the gradients are small enough, the parameters won't move much anyway

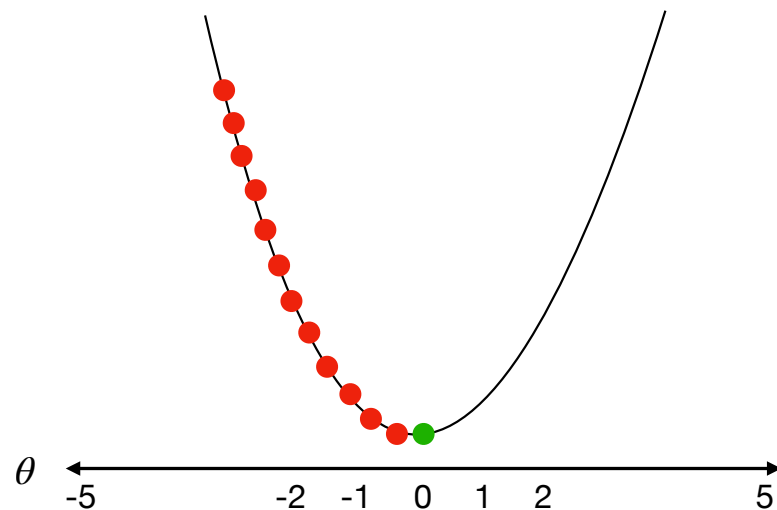
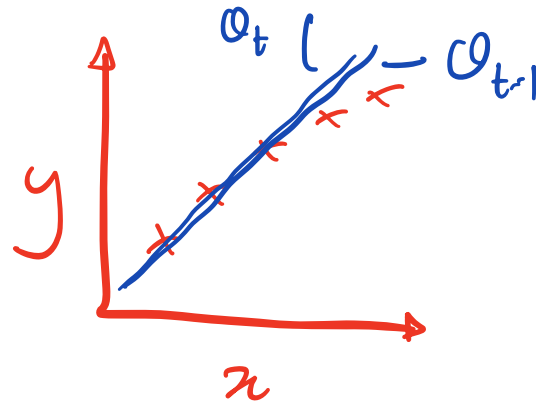


Optimizing Loss Functions

Gradient Descent - Stopping Criterion

- When do you stop your iterations?
 - Function Value Change
 - Terminate when the loss stops changing meaningfully

$$|\ell_{\theta_t}(x) - \ell_{\theta_{t-1}}(x)| \leq \epsilon$$

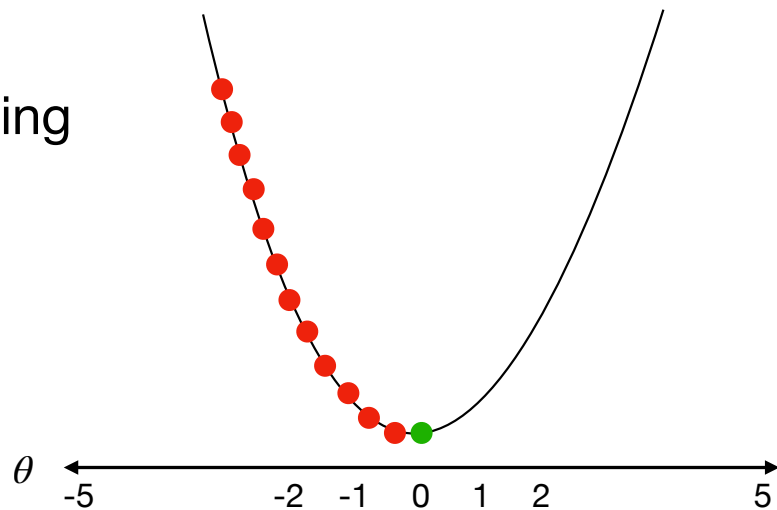


Optimizing Loss Functions

Gradient Descent - Stopping Criterion

- When do you stop your iterations?
 - Parameter Value Change
 - Terminate when the parameters stop changing meaningfully

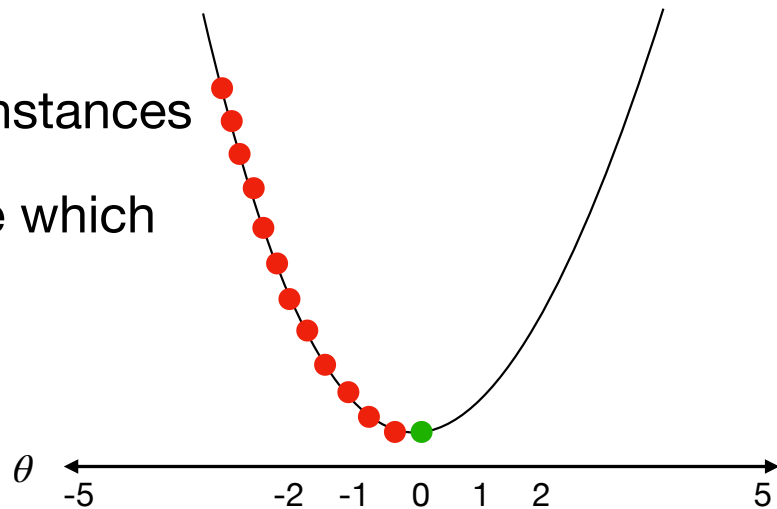
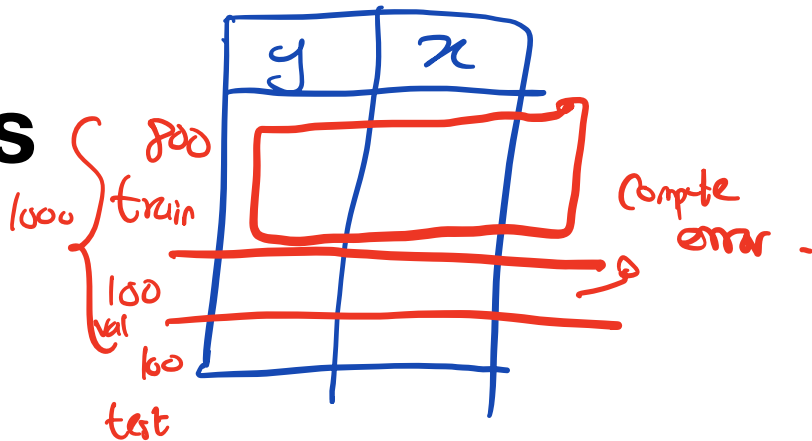
$$|\theta_t - \theta_{t-1}| \leq \epsilon$$



Optimizing Loss Functions

Gradient Descent - Stopping Criterion

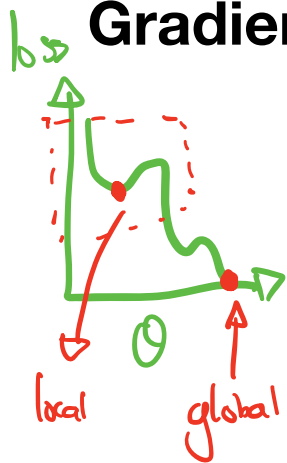
- When do you stop your iterations?
- Validation Based Stopping (Early Stopping)
 - Monitor performance on a validation set of instances
 - Stop when validation loss begins to increase which signals overfitting
 - Serves as both stopping criterion and regularization



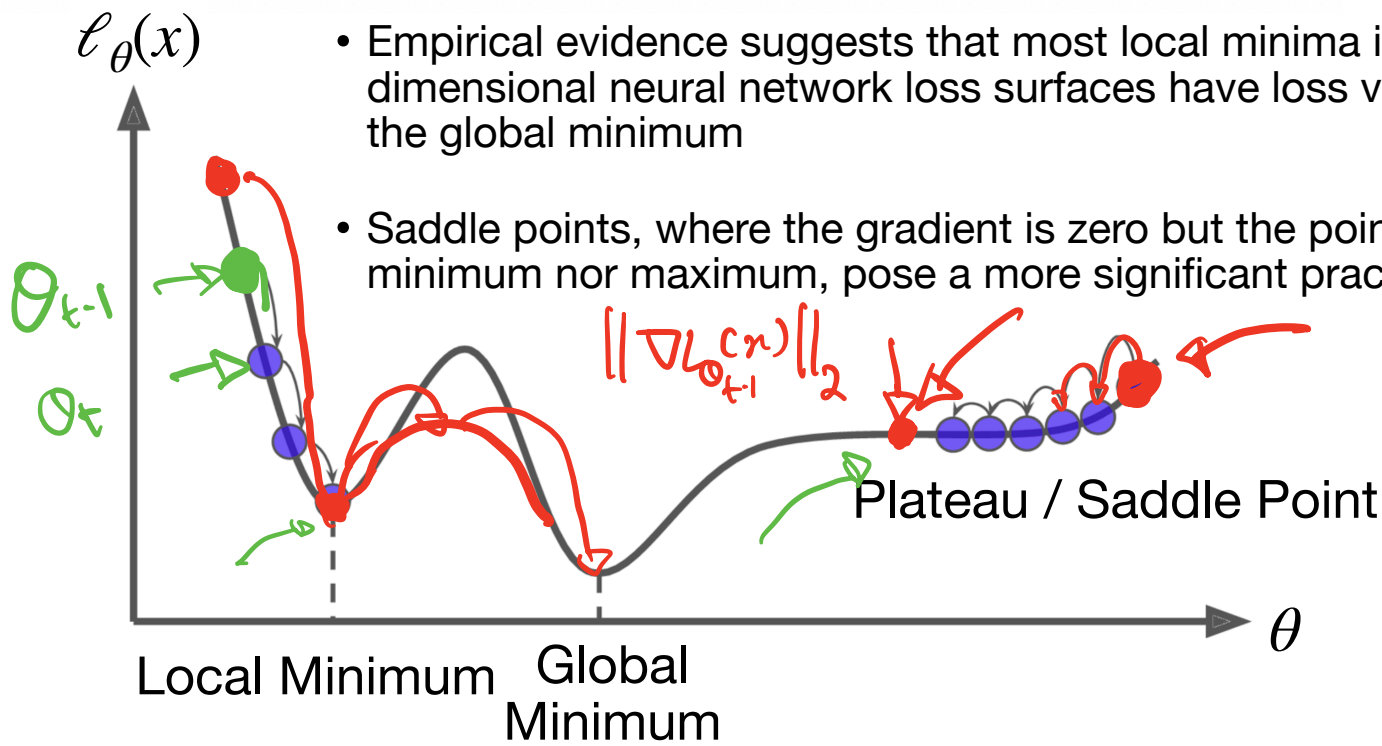
$$\theta_t \leftarrow \boxed{\theta_{t-1}} - \alpha \cdot \nabla_{\theta_{t-1}} L(n)$$

Optimizing Loss Functions

Gradient Descent - More Complicated Functions



- Most deep learning models however have **highly non-convex** loss landscapes
- Empirical evidence suggests that most local minima in high-dimensional neural network loss surfaces have loss values close to the global minimum
- Saddle points, where the gradient is zero but the point is neither a minimum nor maximum, pose a more significant practical challenge.



Optimizing Loss Functions

Gradient Descent - Momentum

① G.D $\rightarrow \theta_t \leftarrow \theta_{t-1} - \alpha \cdot \nabla_{\theta} L(x)_{t-1}$ — gets too small

② Problem $\rightarrow \nabla_{\theta} L(x)_{t-1}$ gets too small.

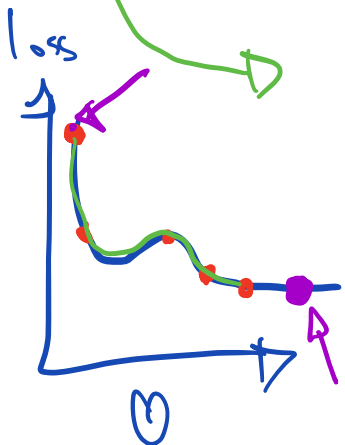
③ Fix \rightarrow Remember previous values and velocities.

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot v_t$ — does not get too small.

$v_t \leftarrow \beta v_{t-1} + \nabla_{\theta} L(x)_{t-1}$

Sum of all gradients seen so far

$\beta = 0.9 \rightarrow$ hyper parameter



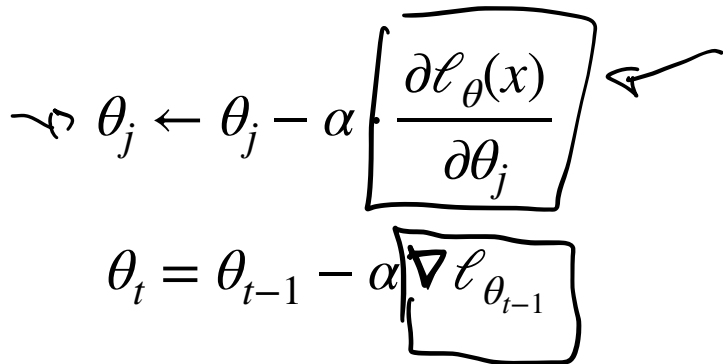
Optimizing Loss Functions

Gradient Descent - Momentum

Optimizing Loss Functions

Gradient Descent - Momentum

- Standard gradient descent can oscillate in ravines
 - Areas where the surface curves **more steeply in one dimension** than another
 - Or they can get stuck in plateau / saddle points
- Momentum helps accelerate gradient descent by accumulating velocity in **directions of consistent gradient**

$$\leadsto \theta_j \leftarrow \theta_j - \alpha \left[\frac{\partial \ell_{\theta}(x)}{\partial \theta_j} \right]$$
$$\theta_t = \theta_{t-1} - \alpha \left[\nabla \ell_{\theta_{t-1}} \right]$$


Optimizing Loss Functions

Gradient Descent - Momentum

- Momentum helps accelerate gradient descent by accumulating velocity in **directions of consistent gradient** *out of bounds.*

$$\rightarrow \theta_t = \theta_{t-1} - \boxed{\alpha \nabla \ell_{\theta_{t-1}}}$$

With Momentum

$$\rightarrow v_t = \beta \cdot v_{t-1} + \nabla \ell_{\theta_{t-1}}$$

$$\rightarrow \theta_t = \theta_{t-1} - \alpha \cdot \boxed{v_t}$$

Optimizing Loss Functions

Gradient Descent - Momentum

- Momentum helps accelerate gradient descent by accumulating velocity in **directions of consistent gradient**

$$\theta_t = \theta_{t-1} - \alpha \nabla \ell_{\theta_{t-1}}$$

With Momentum

Velocity Vector $v_t = \beta \cdot v_{t-1} + \nabla \ell_{\theta_{t-1}}$

$$\theta_t = \theta_{t-1} - \alpha \cdot v_t$$

Optimizing Loss Functions

Gradient Descent - Momentum

- Momentum helps accelerate gradient descent by accumulating velocity in **directions of consistent gradient**

$$\theta_t = \theta_{t-1} - \alpha \nabla \ell_{\theta_{t-1}}$$

With Momentum

β is the momentum coefficient, typically set to 0.9

$$v_t = \beta \cdot v_{t-1} + \nabla \ell_{\theta_{t-1}}$$

$$\theta_t = \theta_{t-1} - \alpha \cdot v_t$$

Optimizing Loss Functions

Gradient Descent - Momentum

- Momentum helps accelerate gradient descent by accumulating velocity in **directions of consistent gradient**

$$\theta_t = \theta_{t-1} - \alpha \nabla \ell_{\theta_{t-1}}$$

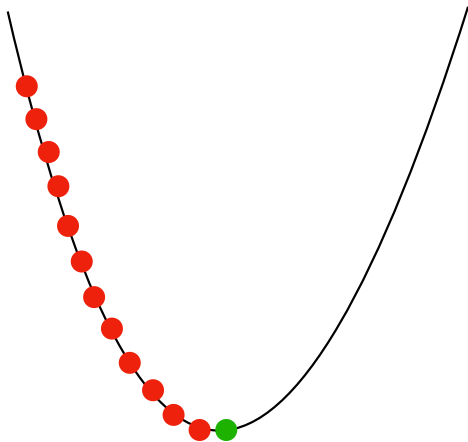
With Momentum

If $\beta = 0$, you get back standard gradient descent $v_t = \beta \cdot v_{t-1} + \nabla \ell_{\theta_{t-1}}$

$$\theta_t = \theta_{t-1} - \alpha \cdot v_t$$

Optimizing Loss Functions

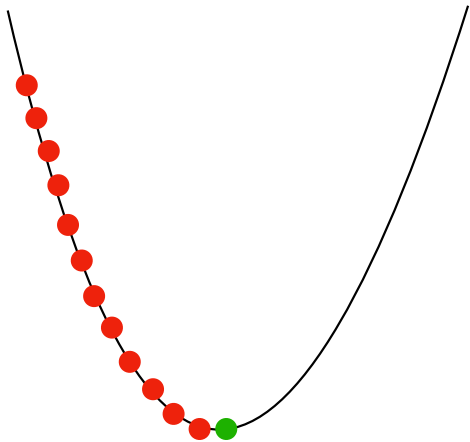
Gradient Descent - Adaptive Step Sizes



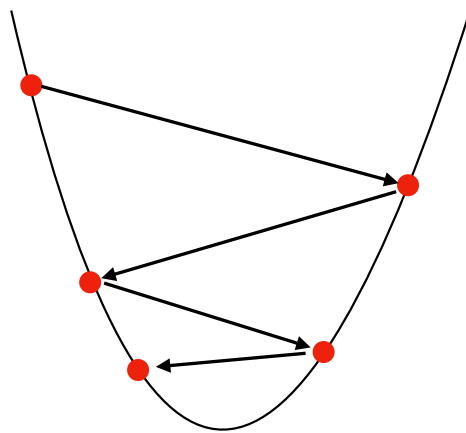
α is too small
Finds the optimal but too slow

Optimizing Loss Functions

Gradient Descent - Adaptive Step Sizes



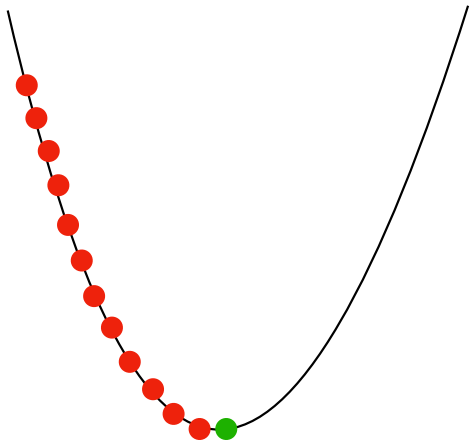
α is too small
Finds the optimal but too slow



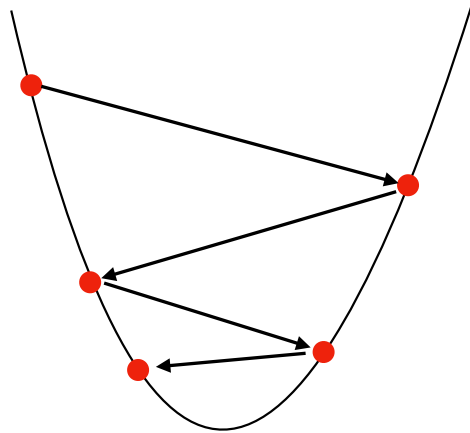
α is too large
Might not find optimal
Could even begin to diverge

Optimizing Loss Functions

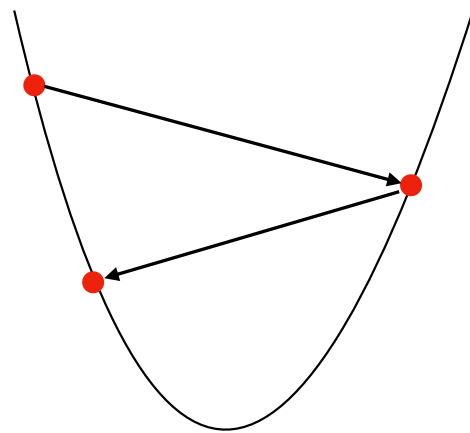
Gradient Descent - Adaptive Step Sizes



α is too small
Finds the optimal but too slow



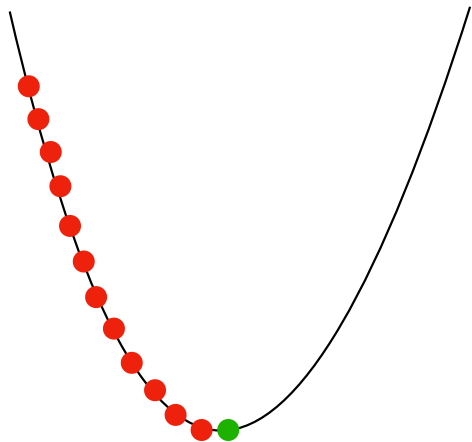
α is too large
Might not find optimal
Could even begin to diverge



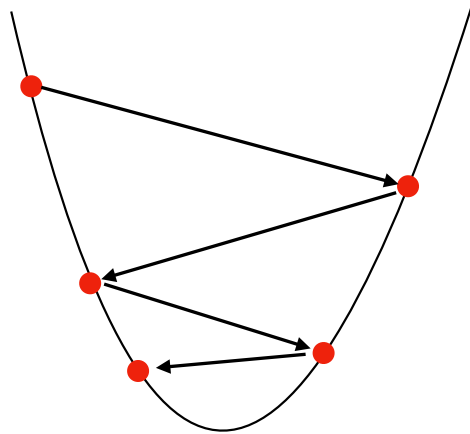
What if you set α to be large
initially?

Optimizing Loss Functions

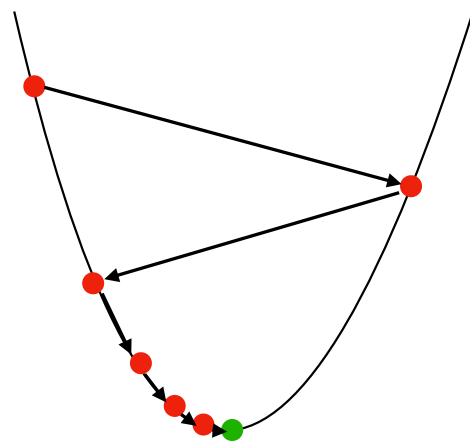
Gradient Descent - Adaptive Step Sizes



α is too small
Finds the optimal but too slow



α is too large
Might not find optimal
Could even begin to diverge




And keep reducing α as
number of epochs increases?

Optimizing Loss Functions

Gradient Descent - Per Parameter Adaptive Learning Rates

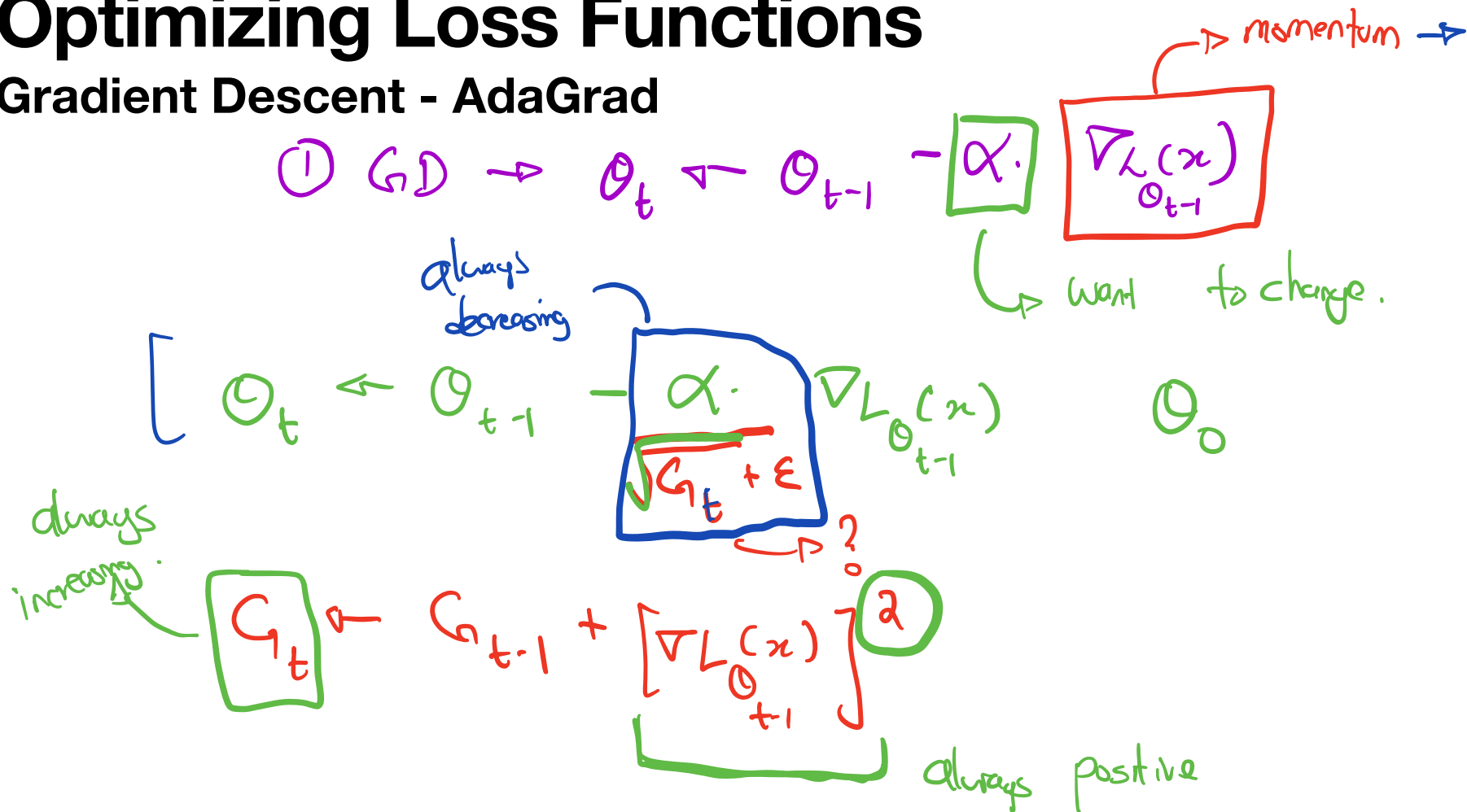
- A single global learning rate may be suboptimal
 - Some parameters might benefit from larger updates while others need smaller ones.
 - Adaptive methods adjust the learning rate for each parameter individually based on historical gradient information.

$$\hat{y} = \theta_0 + \theta_1 x$$


A handwritten diagram in red ink. It shows the equation $\hat{y} = \theta_0 + \theta_1 x$. Below the θ_0 term, there is an upward-pointing arrow from the Greek letter α . A curved arrow also originates from the α and points towards the θ_1 term, indicating that the learning rate α is applied to both parameters.

Optimizing Loss Functions

Gradient Descent - AdaGrad



Optimizing Loss Functions

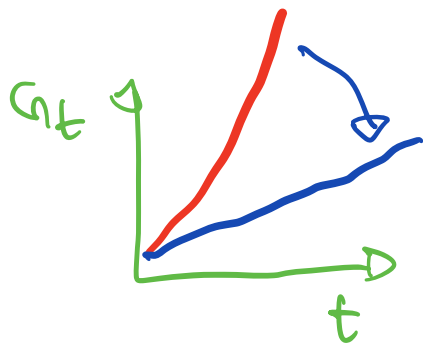
Gradient Descent - AdaGrad + RMSProp

$$\theta_t \leftarrow \theta_{t-1} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \nabla_{\theta_{t-1}} L(x)$$

$$G_t \leftarrow \beta G_{t-1} + (1-\beta) [\nabla_{\theta_{t-1}} L(x)]^2$$

↳ decay rate = 0.9


$$G_t \leftarrow 0.9 \cdot G_{t-1} + (0.1) \boxed{[\nabla_{\theta_{t-1}} L(x)]^2}$$




Optimizing Loss Functions

Gradient Descent - AdaGrad

- AdaGrad adapts the learning rate for **each parameter** based on the sum of squared historical gradients

$$G_t = G_{t-1} + (\nabla \ell_{\theta_{t-1}})^2$$


$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \cdot \nabla \ell_{\theta_{t-1}}$$


- Parameters with large historical gradients receive smaller updates
- Parameters with small historical gradients receive larger updates
- The limitation is that the accumulated sum G_t grows monotonically, eventually making the learning rate vanishingly small.

Optimizing Loss Functions

Gradient Descent - AdaGrad

- AdaGrad adapts the learning rate for **each parameter** based on the sum of squared historical gradients

$$G_t = G_{t-1} + (\nabla \ell_{\theta_{t-1}})^2$$
$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \cdot \nabla \ell_{\theta_{t-1}}$$

- Parameters with large historical gradients receive smaller updates
- Parameters with small historical gradients receive larger updates
- The limitation is that the accumulated sum G_t grows monotonically, eventually making the learning rate vanishingly small.

Optimizing Loss Functions

Gradient Descent - RMSProp

- RMSprop addresses AdaGrad's diminishing learning rate by using an exponentially decaying average of squared gradients

$$G_t = \rho \cdot G_{t-1} + (1 - \rho) (\nabla \ell_{\theta_{t-1}})^2$$
$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \cdot \nabla \ell_{\theta_{t-1}}$$

- The decay rate ρ is typically set to 0.9.
- This prevents the learning rate from decaying to zero while still adapting to the gradient scale.

Optimizing Loss Functions

Gradient Descent - ADAM

- Adam (**A**daptive **M**oment Estimation) combines the benefits of **momentum** (first moment) with the adaptive learning rates of **RMSProp** (second moment)

$$GD \rightarrow \theta_t \leftarrow \theta_{t-1} - \boxed{\alpha} \boxed{\nabla_{\theta_t} L(x)}$$

want decay

$$\theta_t \leftarrow \theta_{t-1} - \frac{\alpha}{\sqrt{G_t} + \epsilon} \cdot v_t$$

Want momentum

$$v_t \leftarrow \beta_1 v_{t-1} + (1 - \beta_1) \nabla_{\theta_t} L(x)$$
$$G_t \leftarrow \beta_2 G_{t-1} + (1 - \beta_2) \left[\nabla_{\theta_t} L(x) \right]^2$$