

# Gradient Descent and Siblings

Ananda Biswas

✉ ami.ananda@bhu.ac.in

## Acknowledgements

- CS7015 Lectures<sup>a</sup> of Professor Mitesh Khapra
- Videos<sup>b</sup> of Ryan Harris on Backpropagation

---

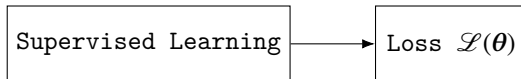
<sup>a</sup>YouTube

<sup>b</sup>YouTube

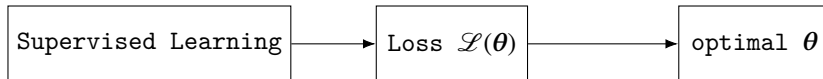
# Target : Learning Parameters

Supervised Learning

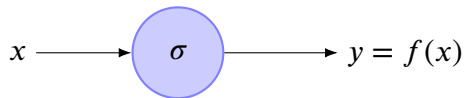
# Target : Learning Parameters



# Target : Learning Parameters



# Set-up



$$f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

## Input for Training

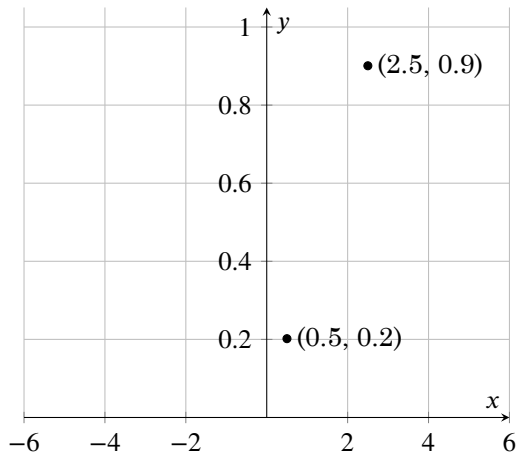
$(x_i, y_i)_{i=1}^n \rightarrow n$  pairs of  $(x, y)$

## Training Objective

Find  $w$  and  $b$  that minimizes

$$\mathcal{L}(w, b) = \sum_{i=1}^n (y_i - f(x_i))^2$$

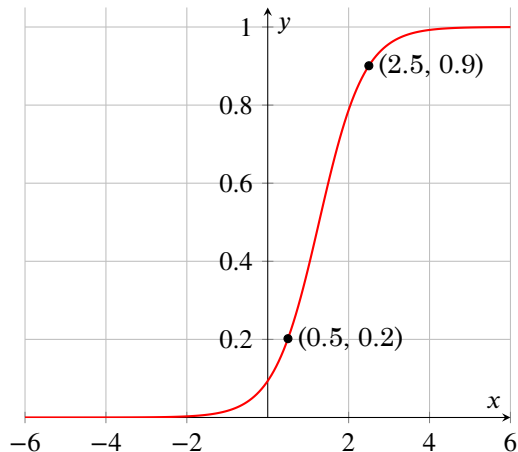
# Set-up



## What is training ?

At the end of the process we wish to get  $w^*$  and  $b^*$  so that  $f(0.5) \rightarrow 0.2$  and  $f(2.5) \rightarrow 0.9$ .

# Set-up



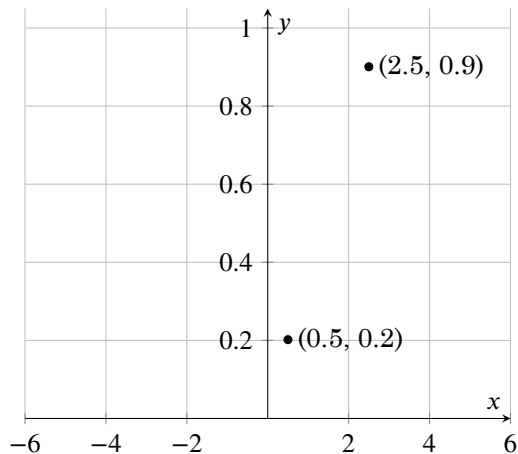
## What is training ?

At the end of the process we wish to get  $w^*$  and  $b^*$  so that  $f(0.5) \rightarrow 0.2$  and  $f(2.5) \rightarrow 0.9$ .

In other words, we hope to find a sigmoid function that satisfies (0.5, 0.2) and (2.5, 0.9).

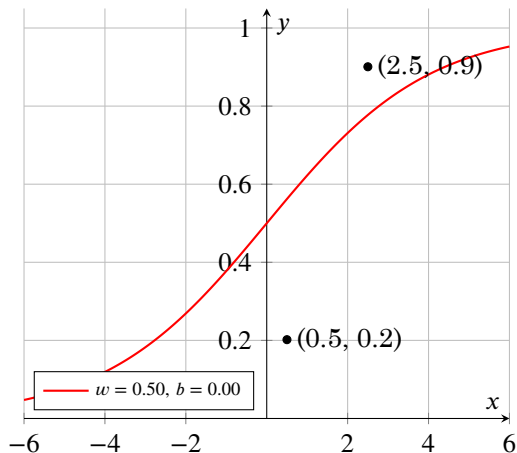


# Approach 1 : Guess Work



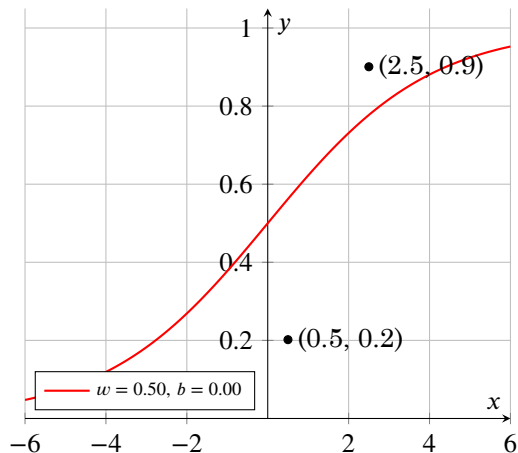
Let us try a random guess..

# Approach 1 : Guess Work



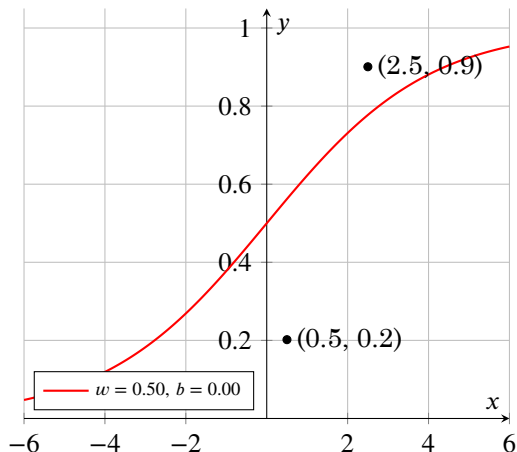
Let us try a random guess..  
say,  $w = 0.5, b = 0$

# Approach 1 : Guess Work



Let us try a random guess..  
say,  $w = 0.5, b = 0$   
Clearly not good, but how bad is it ?

# Approach 1 : Guess Work

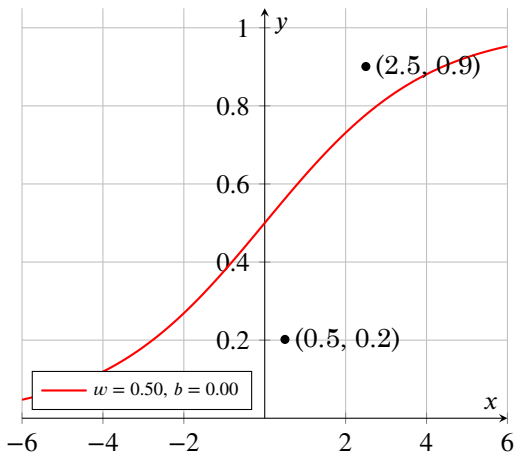


Let us try a random guess..

say,  $w = 0.5, b = 0$

Clearly not good, but how bad is it ?

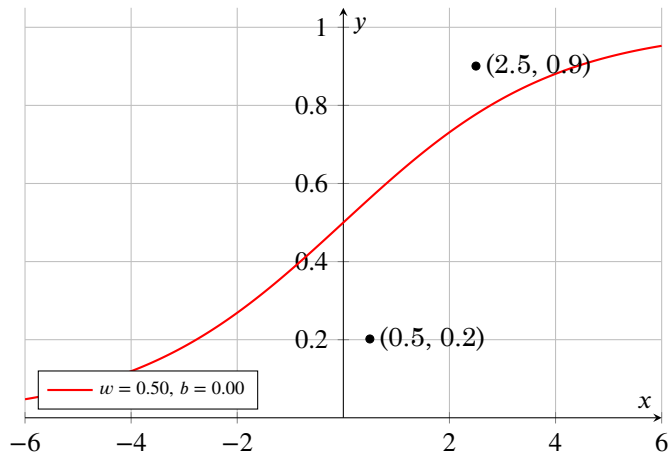
$$\mathcal{L}(w, b) = \sum_{i=1}^n (y_i - f(x_i))^2 \text{ will tell us.}$$



$$\begin{aligned}
 \mathcal{L}(w, b) &= \sum_{i=1}^n (y_i - f(x_i))^2 \\
 &= \frac{1}{2} \left[ (y_1 - f(x_1))^2 + (y_2 - f(x_2))^2 \right] \\
 &= \frac{1}{2} \left[ (0.2 - f(0.5))^2 + (0.9 - f(0.2))^2 \right] \\
 &= 0.073
 \end{aligned}$$

We want  $\mathcal{L}(w, b) = \sum_{i=1}^n (y_i - f(x_i))^2$  to as close to 0 as possible.

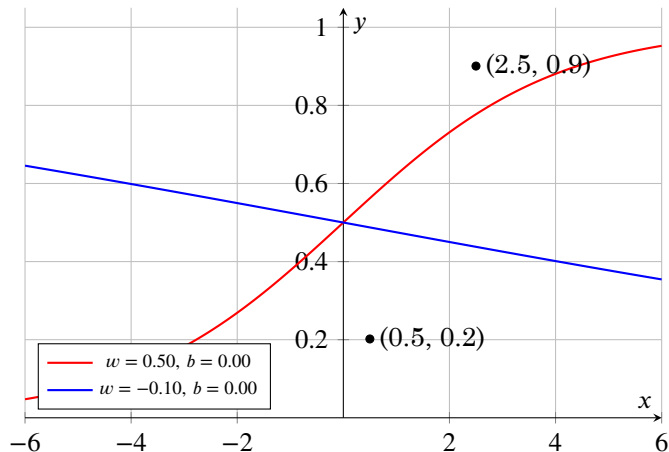
# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730

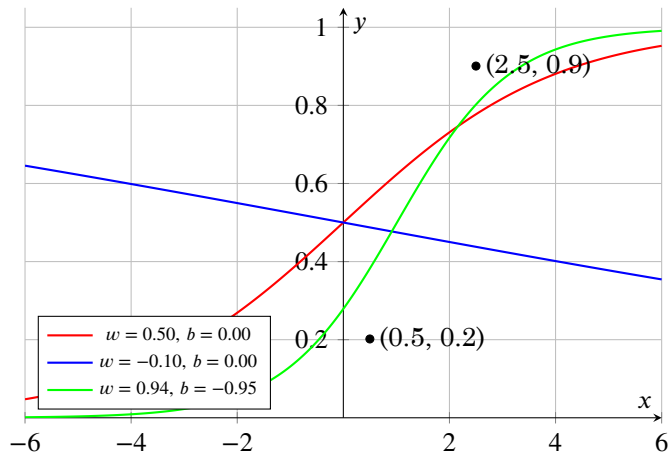
# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481

# Approach 1 : Guess Work

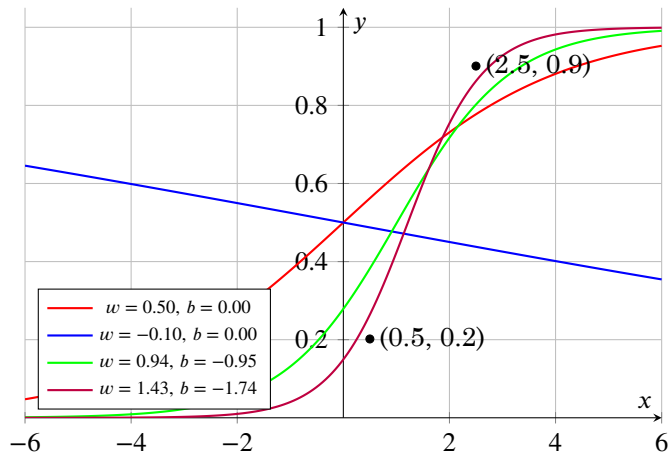


We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214



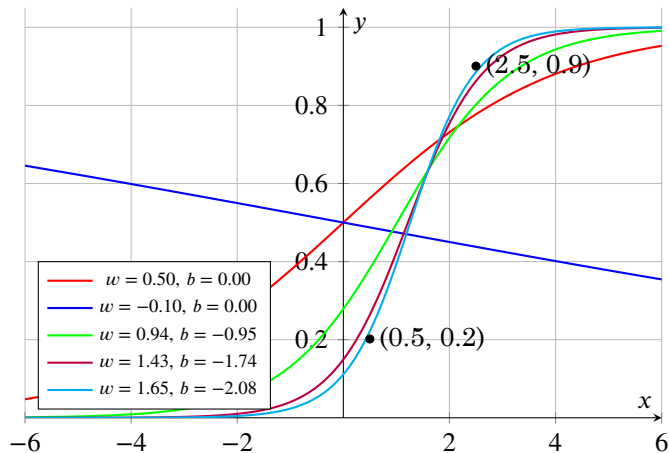
# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028

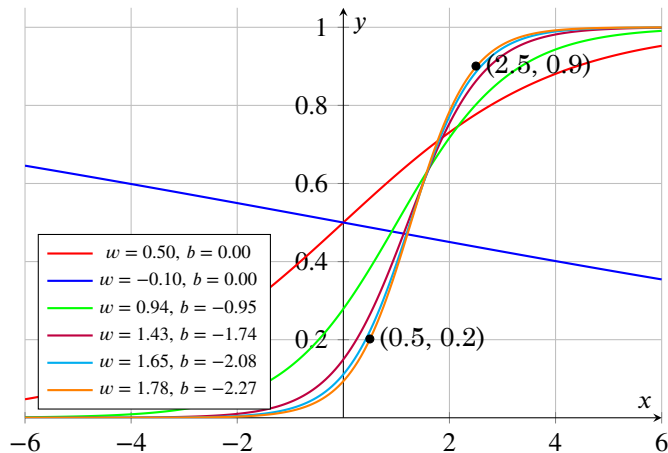
# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003

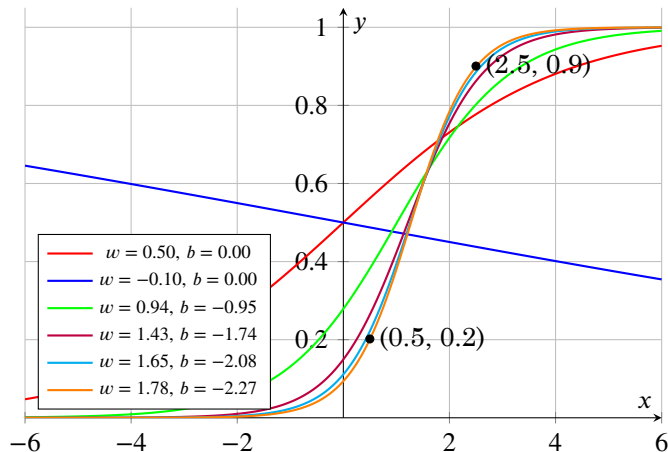
# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

# Approach 1 : Guess Work

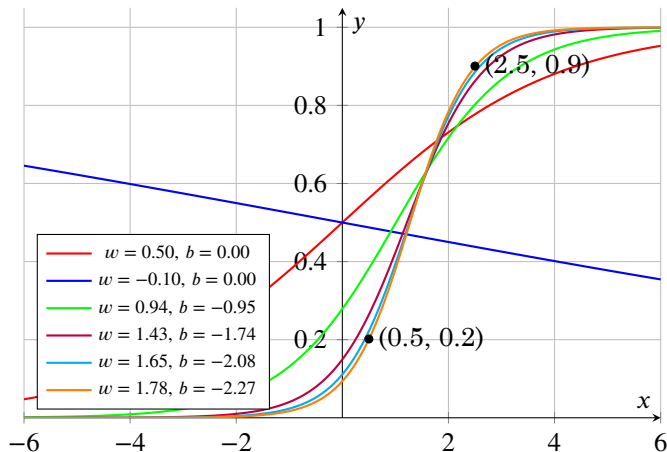


We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

Job done ! But

# Approach 1 : Guess Work



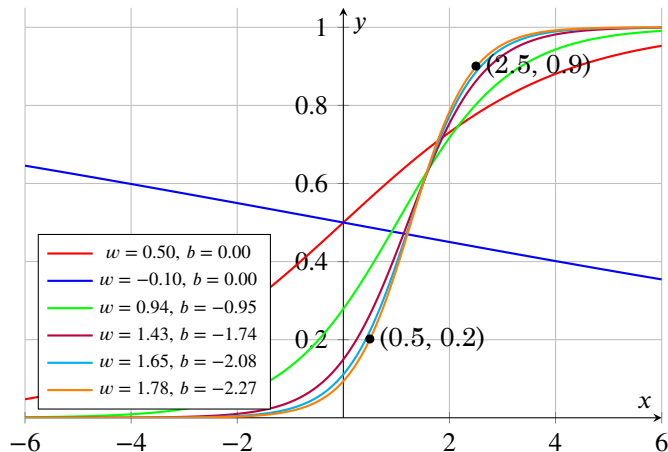
We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

Job done ! But

- Infeasible

# Approach 1 : Guess Work



We try with some other values of  $w$  and  $b$ .

$w$	$b$	$\mathcal{L}(w, b)$
0.50	0.00	0.0730
-0.10	0.00	0.1481
0.94	-0.94	0.0214
1.42	-1.73	0.0028
1.65	-2.08	0.0003
1.78	-2.27	0.0000

Job done ! But

- Infeasible
- Does not guarantee correct solution

# Approach 2 : Brute-force Search

Compute  $\mathcal{L}(w, b)$  for all possible  $w$  and  $b$

# Approach 2 : Brute-force Search

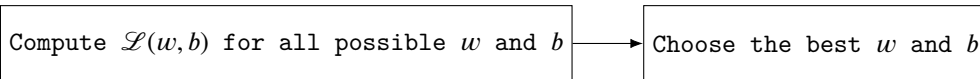
Compute  $\mathcal{L}(w, b)$  for all possible  $w$  and  $b$



Choose the best  $w$  and  $b$

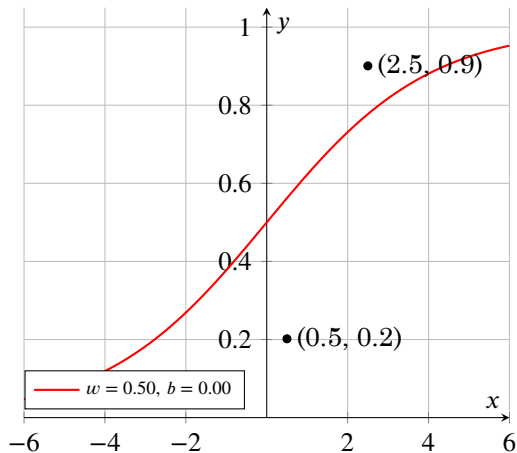


# Approach 2 : Brute-force Search

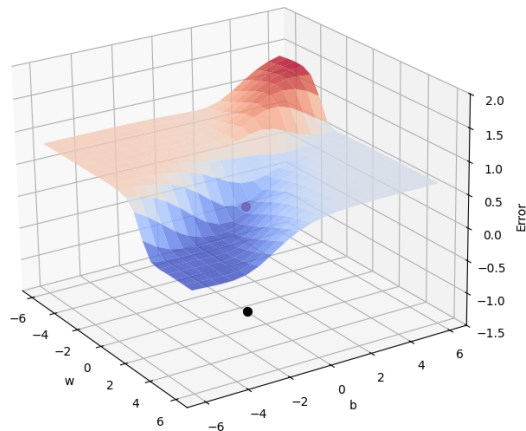


- Computationally infeasible

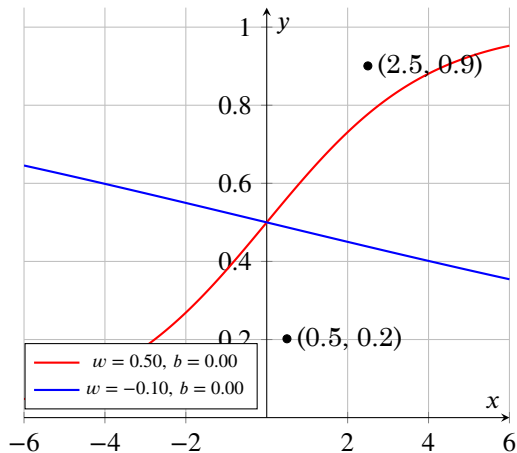
# Revisit to Approach 1



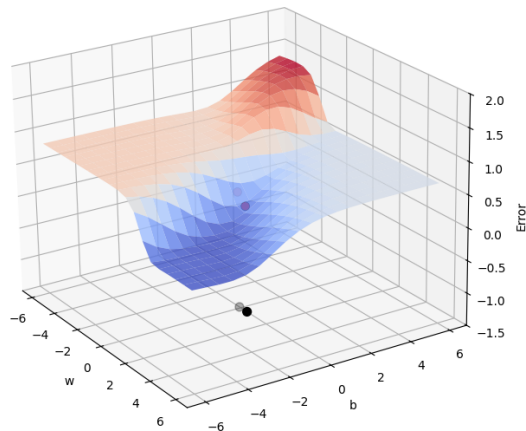
What we were actually doing !



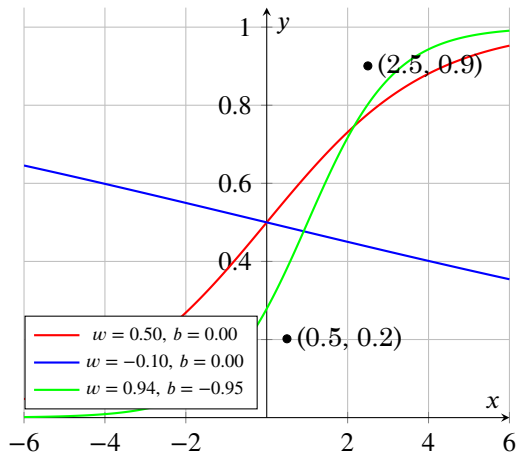
# Revisit to Approach 1



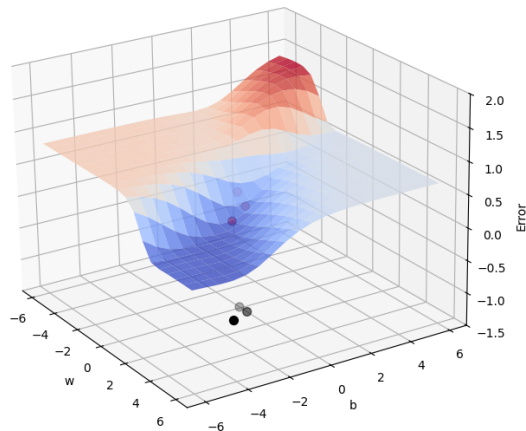
What we were actually doing !



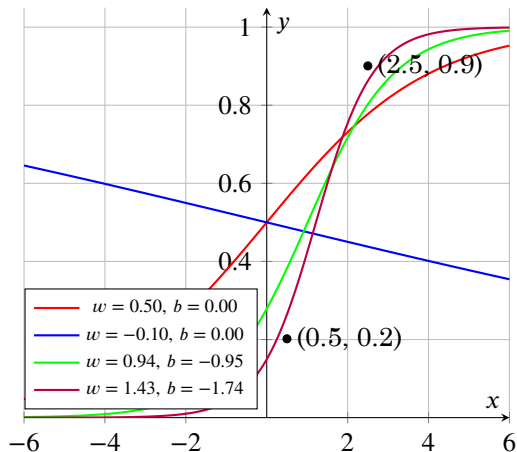
# Revisit to Approach 1



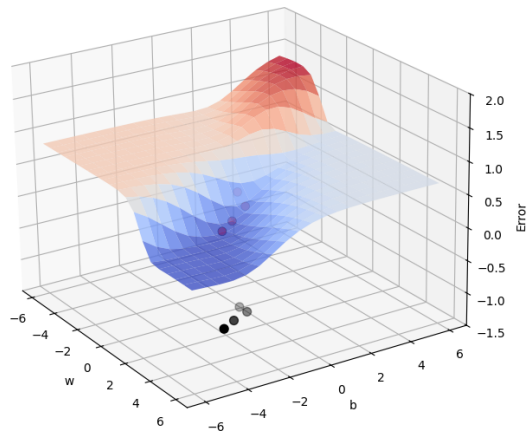
What we were actually doing !



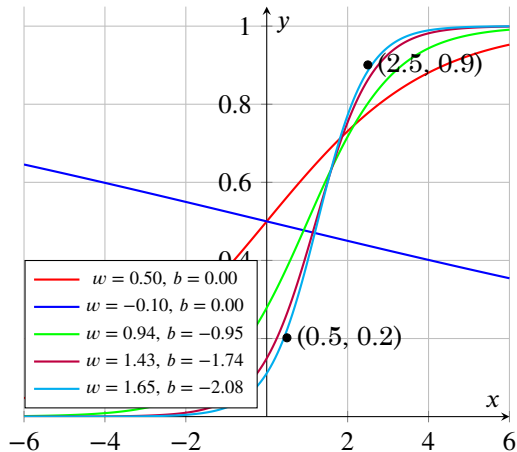
# Revisit to Approach 1



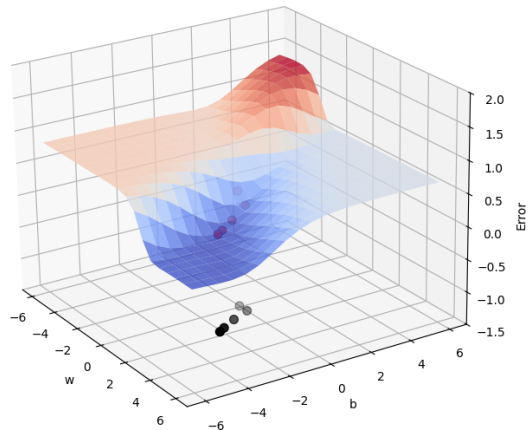
What we were actually doing !



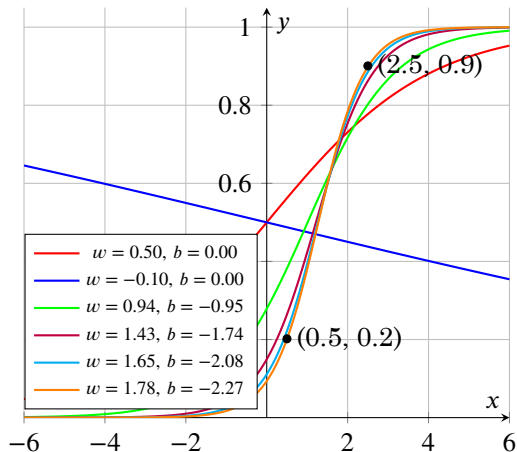
# Revisit to Approach 1



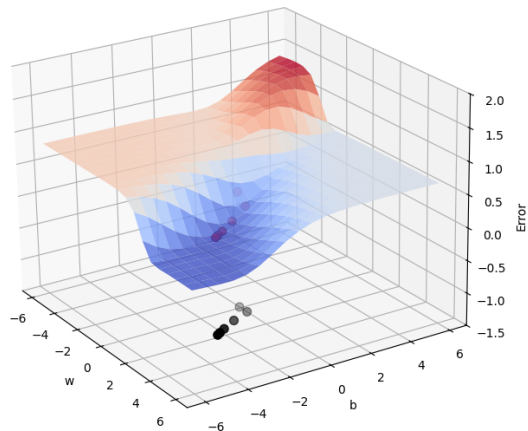
What we were actually doing !



# Revisit to Approach 1



What we were actually doing !



So far ...

We were traversing the error surface by fiddling with  $w$  and  $b$  with mere intuition and guess work.



So far ...

We were traversing the error surface by fiddling with  $w$  and  $b$  with mere intuition and guess work.

Up next ...

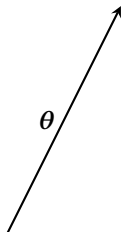
Looking for a more efficient and principled way of doing this.



## Goal ahead

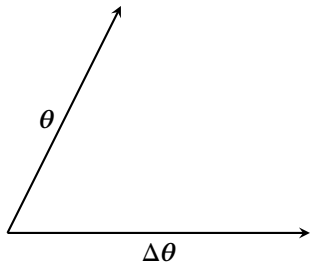
To find a better way of traversing the error surface so that we can reach the minimum value quickly without resorting to guess work or brute force search which is any how infeasible.

# Approach 3 : A Principled Way



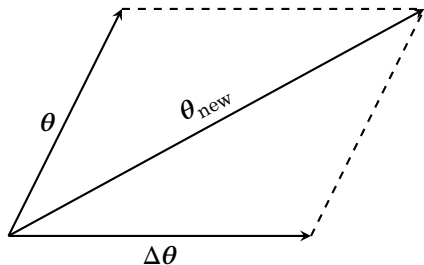
- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .

# Approach 3 : A Principled Way



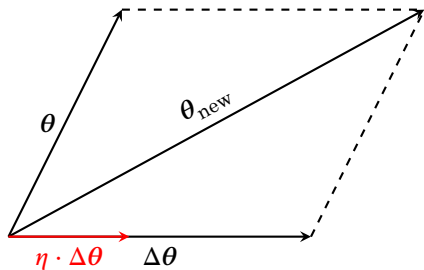
- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.

# Approach 3 : A Principled Way



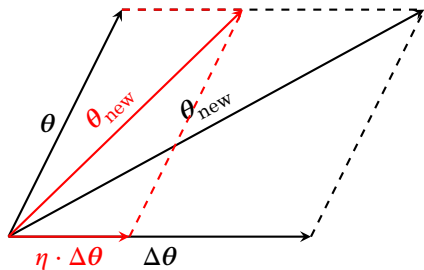
- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.
- So we move in the direction of  $\Delta\theta$ .

# Approach 3 : A Principled Way



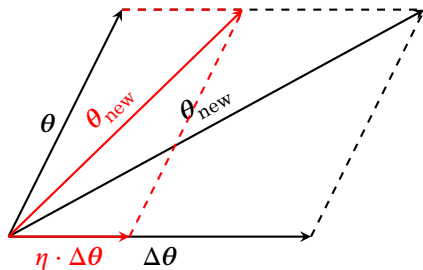
- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.
- So we move in the direction of  $\Delta\theta$ .
- Let us be a bit conservative: we move only by a small amount  $\eta > 0$ .

# Approach 3 : A Principled Way



- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.
- So we move in the direction of  $\Delta\theta$ .
- Let us be a bit conservative: we move only by a small amount  $\eta > 0$ .
- $\theta_{new} = \theta + \eta \cdot \Delta\theta$

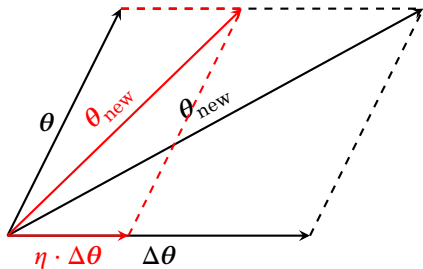
# Approach 3 : A Principled Way



- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.
- So we move in the direction of  $\Delta\theta$ .
- Let us be a bit conservative: we move only by a small amount  $\eta > 0$ .
- $\theta_{new} = \theta + \eta \cdot \Delta\theta$
- Now the question is what is the right  $\Delta\theta$  to use.



# Approach 3 : A Principled Way



- Suppose we have a randomly initialized vector of parameters  $\theta = [w, b]$ .
- Let  $\Delta\theta = [\Delta w, \Delta b]$  denote change in the values of  $w$  and  $b$  respectively.
- So we move in the direction of  $\Delta\theta$ .
- Let us be a bit conservative: we move only by a small amount  $\eta > 0$ .
- $\theta_{new} = \theta + \eta \cdot \Delta\theta$
- Now the question is what is the right  $\Delta\theta$  to use.
- The answer comes from Taylor Series.

## Taylor Series

For a function  $\mathcal{F}(x)$  expanded around  $a \in \mathbb{R}$ , we have,

$$\mathcal{F}(x) = \mathcal{F}(a) + (x - a)\mathcal{F}'(a) + \frac{1}{2!}(x - a)^2\mathcal{F}''(a) + \dots$$

## Taylor Series

For a function  $\mathcal{F}(x)$  expanded around  $a \in \mathbb{R}$ , we have,

$$\mathcal{F}(x) = \mathcal{F}(a) + (x - a)\mathcal{F}'(a) + \frac{1}{2!}(x - a)^2\mathcal{F}''(a) + \dots$$

## Multivariate Taylor Series

For a function  $\mathcal{F}(\mathbf{x})$  expanded around  $\mathbf{a} \in \mathbb{R}^n$ , we have,

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!}(\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \dots$$

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}$$

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}(\theta + \eta \mathbf{u}) =$$

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}(\theta + \eta \mathbf{u}) = \mathcal{L}(\theta) +$$



$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}(\theta + \eta \mathbf{u}) = \mathcal{L}(\theta) + (\eta \mathbf{u})'$$

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}(\theta + \eta \mathbf{u}) = \mathcal{L}(\theta) + (\eta \mathbf{u})' \nabla \mathcal{L}(\theta) +$$

$$\mathcal{F}(\mathbf{x}) = \mathcal{F}(\mathbf{a}) + (\mathbf{x} - \mathbf{a})' \nabla \mathcal{F}(\mathbf{a}) + \frac{1}{2!} (\mathbf{x} - \mathbf{a})' \nabla^2 \mathcal{F}(\mathbf{a}) (\mathbf{x} - \mathbf{a}) + \dots$$

For ease of notation, take  $\Delta\theta = \mathbf{u}$ .

With  $\mathcal{F} = \mathcal{L}$  (our Loss Function),

$$\mathbf{x} = \theta + \eta \mathbf{u},$$

$\mathbf{a} = \theta$ ; from the Taylor Series we have,

$$\mathcal{L}(\theta + \eta \mathbf{u}) = \mathcal{L}(\theta) + (\eta \mathbf{u})' \nabla \mathcal{L}(\theta) + \frac{1}{2!} (\eta \mathbf{u})' \nabla^2 \mathcal{L}(\theta) (\eta \mathbf{u}) + \dots$$

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + (\eta \mathbf{u})' \nabla \mathcal{L}(\boldsymbol{\theta}) + \frac{1}{2!} (\eta \mathbf{u})' \nabla^2 \mathcal{L}(\boldsymbol{\theta}) (\eta \mathbf{u}) + \dots$$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) &= \mathcal{L}(\boldsymbol{\theta}) + (\eta \mathbf{u})' \nabla \mathcal{L}(\boldsymbol{\theta}) + \frac{1}{2!} (\eta \mathbf{u})' \nabla^2 \mathcal{L}(\boldsymbol{\theta}) (\eta \mathbf{u}) + \dots \\ &= \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta}) + \frac{\eta^2}{2!} \mathbf{u}' \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{u} + \dots\end{aligned}$$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) &= \mathcal{L}(\boldsymbol{\theta}) + (\eta \mathbf{u})' \nabla \mathcal{L}(\boldsymbol{\theta}) + \frac{1}{2!} (\eta \mathbf{u})' \nabla^2 \mathcal{L}(\boldsymbol{\theta}) (\eta \mathbf{u}) + \dots \\ &= \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta}) + \frac{\eta^2}{2!} \mathbf{u}' \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{u} + \dots \\ &= \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta}) \left[ \eta \text{ typically being small } \eta^2, \eta^3, \dots \rightarrow 0 \right]\end{aligned}$$

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

Recall that our move was  $(\eta \mathbf{u})$ .



$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

Recall that our move was  $(\eta \mathbf{u})$ .

The move  $(\eta \mathbf{u})$  would be favourable only if,

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

Recall that our move was  $(\eta \mathbf{u})$ .

The move  $(\eta \mathbf{u})$  would be favourable only if,

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) - \mathcal{L}(\boldsymbol{\theta}) < 0$$

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

Recall that our move was  $(\eta \mathbf{u})$ .

The move  $(\eta \mathbf{u})$  would be favourable only if,

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) - \mathcal{L}(\boldsymbol{\theta}) < 0$$

*i.e.* if the new loss is less than the previous loss.

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) = \mathcal{L}(\boldsymbol{\theta}) + \eta \mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$$

Recall that our move was  $(\eta \mathbf{u})$ .

The move  $(\eta \mathbf{u})$  would be favourable only if,

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{u}) - \mathcal{L}(\boldsymbol{\theta}) < 0$$

*i.e.* if the new loss is less than the previous loss.

This implies

$$\mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta}) < 0$$

- More the negative  $\mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$  is,

- More the negative  $\mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$  is,

the more favourable is the move  $\eta \mathbf{u}$ .

- More the negative  $\mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$  is,  
the more favourable is the move  $\eta \mathbf{u}$ .
- Now let us find out the range of  $\mathbf{u}' \nabla \mathcal{L}(\boldsymbol{\theta})$ .

- Let  $\beta$  be the angle between  $\mathbf{u}$  and  $\nabla_{\theta}\mathcal{L}(\theta)$ , then we know that,

$$-1 \leq \cos(\beta) = \frac{\mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta)}{\|\mathbf{u}\| \cdot \|\nabla_{\theta} \mathcal{L}(\theta)\|} \leq 1$$



- Let  $\beta$  be the angle between  $\mathbf{u}$  and  $\nabla_{\theta}\mathcal{L}(\theta)$ , then we know that,

$$-1 \leq \cos(\beta) = \frac{\mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta)}{\|\mathbf{u}\| \cdot \|\nabla_{\theta} \mathcal{L}(\theta)\|} \leq 1$$

- Multiplying throughout by  $k = \|\mathbf{u}\| \cdot \|\nabla_{\theta} \mathcal{L}(\theta)\|$ :

$$-k \leq k \cdot \cos(\beta) = \mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta) \leq k$$

- Let  $\beta$  be the angle between  $\mathbf{u}$  and  $\nabla_{\theta}\mathcal{L}(\theta)$ , then we know that,

$$-1 \leq \cos(\beta) = \frac{\mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta)}{\|\mathbf{u}\| \cdot \|\nabla_{\theta} \mathcal{L}(\theta)\|} \leq 1$$

- Multiplying throughout by  $k = \|\mathbf{u}\| \cdot \|\nabla_{\theta} \mathcal{L}(\theta)\|$ :

$$-k \leq k \cdot \cos(\beta) = \mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta) \leq k$$

- Thus,

$$\mathcal{L}(\theta + \eta \mathbf{u}) - \mathcal{L}(\theta) = \mathbf{u}^T \nabla_{\theta} \mathcal{L}(\theta) = k \cdot \cos(\beta)$$

will be most negative when  $\cos(\beta) = -1$ , i.e., when  $\beta$  is  $180^\circ$ .

## Best Update Rule

So our best move is at  $180^\circ$  w.r.t. the gradient  $\nabla \mathcal{L}(\boldsymbol{\theta})$ .

## Best Update Rule

So our best move is at  $180^\circ$  w.r.t. the gradient  $\nabla \mathcal{L}(\theta)$ .

In other words, we should move in a direction opposite to the gradient *i.e.*

$$\eta \Delta \theta = -\eta \nabla \mathcal{L}(\theta)$$

## Best Update Rule

So our best move is at  $180^\circ$  w.r.t. the gradient  $\nabla \mathcal{L}(\theta)$ .

In other words, we should move in a direction opposite to the gradient *i.e.*

$$\eta \Delta \theta = -\eta \nabla \mathcal{L}(\theta)$$

## Parameter Update Equations

$$w_{t+1} = w_t - \eta \nabla w_t$$

$$b_{t+1} = b_t - \eta \nabla b_t$$

$$\text{where, } \nabla w_t = \left. \frac{\partial \mathcal{L}(w, b)}{\partial w} \right|_{w=w_t, b=b_t}, \quad \nabla b_t = \left. \frac{\partial \mathcal{L}(w, b)}{\partial b} \right|_{w=w_t, b=b_t}$$

# A Formal Definition : Gradient Descent

## Definition

# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a

# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a **first-order**



# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a **first-order iterative** algorithm for

# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a **first-order iterative** algorithm for **minimizing**

# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a **first-order iterative** algorithm for **minimizing** a **differentiable** multivariate function.

# A Formal Definition : Gradient Descent

## Definition

Gradient Descent is a **first-order iterative** algorithm for **minimizing** a **differentiable** multivariate function.

The idea is to take repeated steps in the opposite direction of the gradient (or approximate gradient) of the function at the current point, because this is the direction of steepest descent.

# Batch / Vanilla Gradient Descent

Let's formulate the algorithm.

---

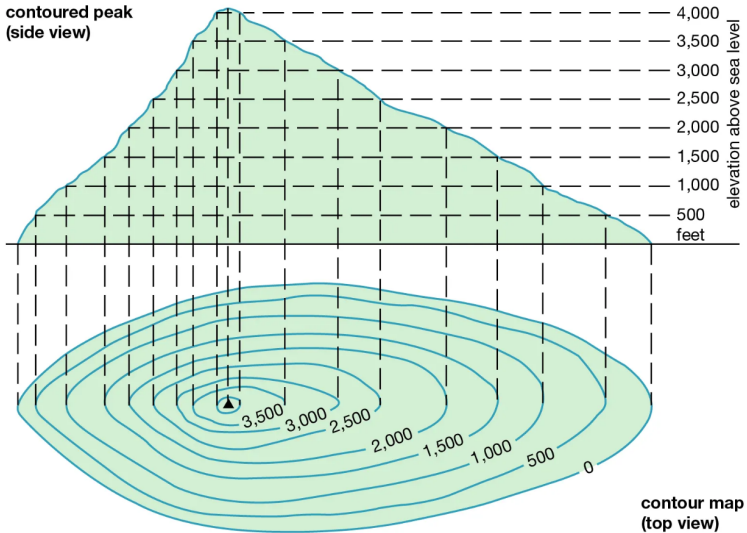
**Algorithm 1:** `gradient_descent()`

---

```
1  $t \leftarrow 0$ ;  
2  $max\_iterations \leftarrow 1000$ ;  
3 while  $t < max\_iterations$  do  
4    $w_{t+1} \leftarrow w_t - \eta \nabla w_t$ ;  
5    $b_{t+1} \leftarrow b_t - \eta \nabla b_t$ ;  
6    $t \leftarrow t + 1$ ;  
7 end
```

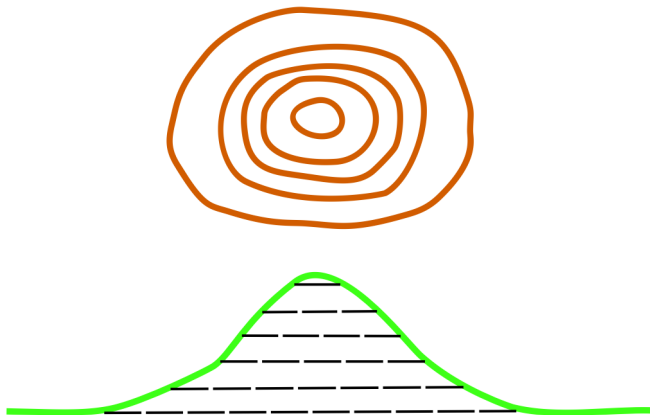
---

# A Digression to Contours

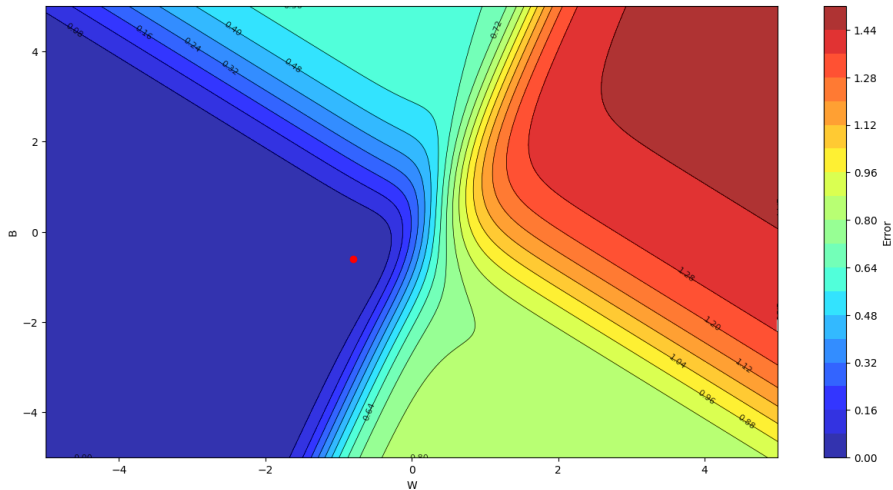


© 2011 Encyclopædia Britannica, Inc.

# A Digression to Contours

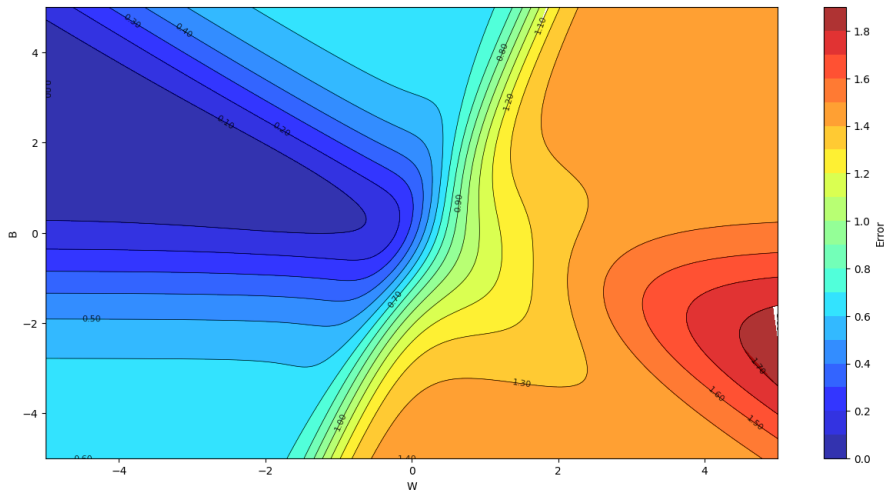


# A Digression to Contours

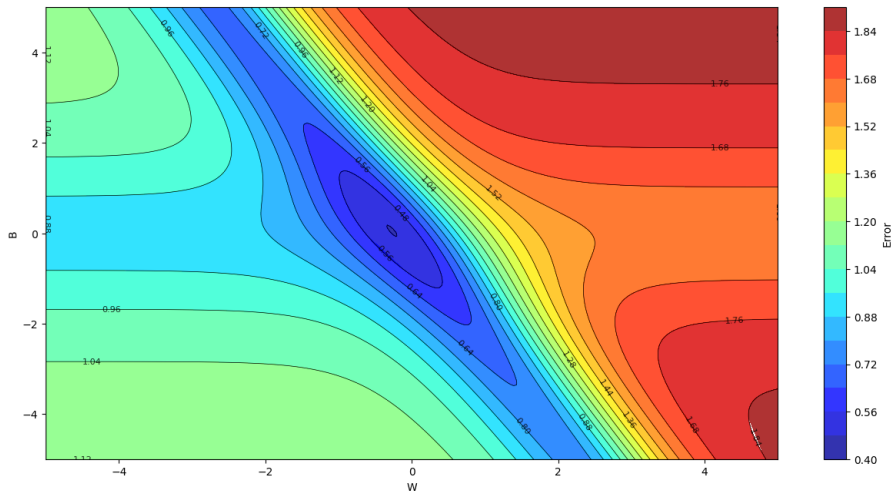




# A Digression to Contours



# A Digression to Contours



## Some Observations in Vanilla Gradient Descent

- It takes a lot of time to navigate regions having a gentle slope.

## Some Observations in Vanilla Gradient Descent

- It takes a lot of time to navigate regions having a gentle slope.
- This is because the gradient in these regions is very small.

## Some Observations in Vanilla Gradient Descent

- It takes a lot of time to navigate regions having a gentle slope.
- This is because the gradient in these regions is very small.
- We have to do something better.

## Intuition

- If I am repeatedly being asked to move in the same direction then I should probably gain some confidence and start taking bigger steps in that direction.

## Intuition

- If I am repeatedly being asked to move in the same direction then I should probably gain some confidence and start taking bigger steps in that direction.
- Just as a ball gains momentum while rolling down a slope

## A New Update Rule

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \cdot \nabla w_t$$

$$w_{t+1} = w_t - \text{update}_t$$



## A New Update Rule

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \cdot \nabla w_t$$

$$w_{t+1} = w_t - \text{update}_t$$

- We have a similar update rule for  $b$ .

## A New Update Rule

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \cdot \nabla w_t$$

$$w_{t+1} = w_t - \text{update}_t$$

- We have a similar update rule for  $b$ .
- In addition to the current update, also look at the history of updates.

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = \gamma \cdot \text{update}_0 + \eta \nabla w_1 = \eta \nabla w_1$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = \gamma \cdot \text{update}_0 + \eta \nabla w_1 = \eta \nabla w_1$$

$$\text{update}_2 = \gamma \cdot \text{update}_1 + \eta \nabla w_2 = \gamma \cdot \eta \nabla w_1 + \eta \nabla w_2$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = \gamma \cdot \text{update}_0 + \eta \nabla w_1 = \eta \nabla w_1$$

$$\text{update}_2 = \gamma \cdot \text{update}_1 + \eta \nabla w_2 = \gamma \cdot \eta \nabla w_1 + \eta \nabla w_2$$

$$\begin{aligned}\text{update}_3 &= \gamma \cdot \text{update}_2 + \eta \nabla w_3 = \gamma(\gamma \cdot \eta \nabla w_1 + \eta \nabla w_2) + \eta \nabla w_3 \\ &= \gamma^2 \cdot \eta \nabla w_1 + \gamma \cdot \eta \nabla w_2 + \eta \nabla w_3\end{aligned}$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = \gamma \cdot \text{update}_0 + \eta \nabla w_1 = \eta \nabla w_1$$

$$\text{update}_2 = \gamma \cdot \text{update}_1 + \eta \nabla w_2 = \gamma \cdot \eta \nabla w_1 + \eta \nabla w_2$$

$$\begin{aligned}\text{update}_3 &= \gamma \cdot \text{update}_2 + \eta \nabla w_3 = \gamma(\gamma \cdot \eta \nabla w_1 + \eta \nabla w_2) + \eta \nabla w_3 \\ &= \gamma^2 \cdot \eta \nabla w_1 + \gamma \cdot \eta \nabla w_2 + \eta \nabla w_3\end{aligned}$$

$$\text{update}_4 = \gamma \cdot \text{update}_3 + \eta \nabla w_4 = \gamma^3 \cdot \eta \nabla w_1 + \gamma^2 \cdot \eta \nabla w_2 + \gamma \cdot \eta \nabla w_3 + \eta \nabla w_4$$



$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

$$\text{update}_0 = 0$$

$$\text{update}_1 = \gamma \cdot \text{update}_0 + \eta \nabla w_1 = \eta \nabla w_1$$

$$\text{update}_2 = \gamma \cdot \text{update}_1 + \eta \nabla w_2 = \gamma \cdot \eta \nabla w_1 + \eta \nabla w_2$$

$$\begin{aligned}\text{update}_3 &= \gamma \cdot \text{update}_2 + \eta \nabla w_3 = \gamma(\gamma \cdot \eta \nabla w_1 + \eta \nabla w_2) + \eta \nabla w_3 \\ &= \gamma^2 \cdot \eta \nabla w_1 + \gamma \cdot \eta \nabla w_2 + \eta \nabla w_3\end{aligned}$$

$$\text{update}_4 = \gamma \cdot \text{update}_3 + \eta \nabla w_4 = \gamma^3 \cdot \eta \nabla w_1 + \gamma^2 \cdot \eta \nabla w_2 + \gamma \cdot \eta \nabla w_3 + \eta \nabla w_4$$

$$\vdots$$

$$\begin{aligned}\text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_t \\ &= \gamma^{t-1} \cdot \eta \nabla w_1 + \gamma^{t-2} \cdot \eta \nabla w_2 + \cdots + \eta \nabla w_t\end{aligned}$$

## Some observations and questions

- Even in the regions having gentle slopes, momentum based gradient descent is able to take large steps because the momentum carries it along.

## Some observations and questions

- Even in the regions having gentle slopes, momentum based gradient descent is able to take large steps because the momentum carries it along.
- Is moving fast always good? Would there be situations where momentum would cause us to run pass our goal?

- Momentum based gradient descent oscillates in and out of the minima valley as the momentum carries it out of the valley.

- Momentum based gradient descent oscillates in and out of the minima valley as the momentum carries it out of the valley.
- Takes a lot of u-turns before finally converging.

- Momentum based gradient descent oscillates in and out of the minima valley as the momentum carries it out of the valley.
- Takes a lot of u-turns before finally converging.
- Despite these u-turns it still converges faster than Vanilla Gradient Descent.

## Question

- Can we do something to reduce these oscillations ?

👉 Of Course !

## Recall in Momentum Gradient Descent

- We had  $w_{t+1} = w_t - \text{update}_t$  where



## Recall in Momentum Gradient Descent

- We had  $w_{t+1} = w_t - \text{update}_t$  where

$$\text{update}_t = \underbrace{\gamma \cdot \text{update}_{t-1}} + \underbrace{\eta \nabla w_t}$$

## Recall in Momentum Gradient Descent

- We had  $w_{t+1} = w_t - \text{update}_t$  where

$$\text{update}_t = \underbrace{\gamma \cdot \text{update}_{t-1}}_{\text{momentum derived from past updates}} + \underbrace{\eta \nabla w_t}$$

## Recall in Momentum Gradient Descent

- We had  $w_{t+1} = w_t - \text{update}_t$  where

$$\text{update}_t = \underbrace{\gamma \cdot \text{update}_{t-1}}_{\text{momentum derived from past updates}} + \underbrace{\eta \nabla w_t}_{\text{push from the current gradient}}$$

## Recall in Momentum Gradient Descent

- We had  $w_{t+1} = w_t - \text{update}_t$  where

$$\text{update}_t = \underbrace{\gamma \cdot \text{update}_{t-1}}_{\text{momentum derived from past updates}} + \underbrace{\eta \nabla w_t}_{\text{push from the current gradient}}$$

- So we know that we are going to move by at least by  $\text{update}_{t-1}$  and then a bit more by  $\eta \nabla w_t$ .

## Intuition

- Look before you leap.

## Intuition

- Look before you leap.
- If we already know the direction we're moving in, why not “look ahead” to where the momentum will take us before calculating the gradient?

## Intuition

- Look before you leap.
- If we already know the direction we're moving in, why not “look ahead” to where the momentum will take us before calculating the gradient?
- $w_{\text{look\_ahead}} = w_t - \gamma \cdot \text{update}_{t-1}$

## Intuition

- Look before you leap.
- If we already know the direction we're moving in, why not “look ahead” to where the momentum will take us before calculating the gradient?
- $w_{\text{look\_ahead}} = w_t - \gamma \cdot \text{update}_{t-1}$
- Why not calculate the gradient at this partially updated value of  $w$  *i.e.*  $w_{\text{look\_ahead}}$  instead of calculating it using the current value  $w_t$  ?



# Nesterov Accelerated Gradient Descent

## A New Update Rule

$$\begin{aligned}w_{\text{look\_ahead}} &= w_t - \gamma \cdot \text{update}_{t-1} \\ \text{update}_t &= \gamma \cdot \text{update}_{t-1} + \eta \nabla w_{\text{look\_ahead}} \\ w_{t+1} &= w_t - \text{update}_t\end{aligned}$$

We have similar update rule for  $b$ .