Stochastic Gradient Descient

CPSC 383: Explorations in Artificial Intelligence and Machine Learning Fall 2025

Jonathan Hudson, Ph.D. Associate Professor (Teaching) Department of Computer Science University of Calgary

August 27, 2025

Copyright © 2025



Review

Training a neural net is is the process of adjusting _____to minimize across all of the _____ data.



Training



Definition

Optimization is a set of problems that involve selecting the "best" element from some set of values.

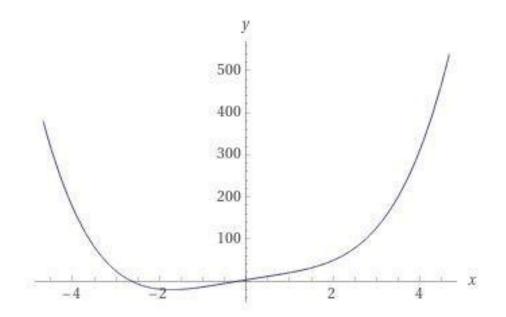
Usually values maximizing or minimizing some objective function

Optimization problems can be discrete or continuous.



Example

• Suppose you wanted to find the "best" value of x, i.e. the one that minimizes the cost function $f(x) = x^4 - x^2 + 17x + 3$ shown below.





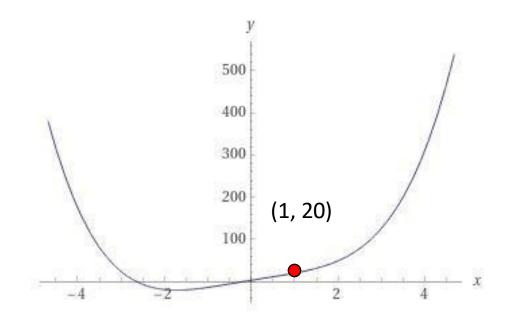
In general, optimization problems are **intractable**, i.e. an exact solution is either computationally expensive or impossible to obtain.

However, in practice we are usually happy with a "good enough" solution that is easier to find.



Idea: what if you just want to find something better than what you currently have?

$$f(x) = x^4 - x^2 + 17x + 3$$



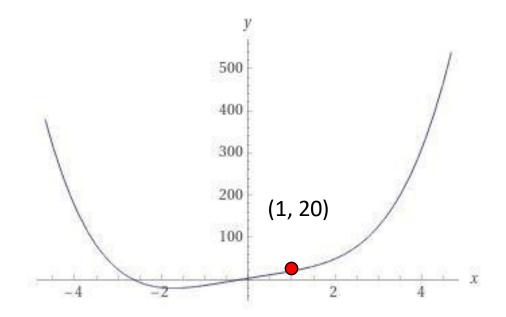


Idea: what if you just want to find something better than what you

currently have?

$$f(x) = x^4 - x^2 + 17x + 3$$

 $f'(x) = 4x^3 - 2x + 17$



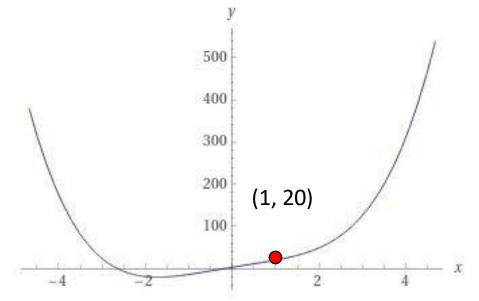


Idea: what if you just want to find something better than what you

currently have?

$$f(x) = x^4 - x^2 + 17x + 3$$

 $f'(x) = 4x^3 - 2x + 17$



f'(1) = 19 > 0 so move left!



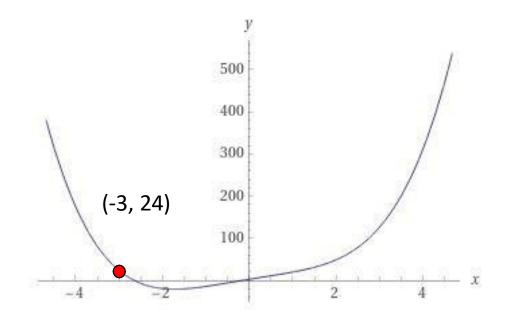
Idea: what if you just want to find something better than what you

currently have?

$$f(x) = x^4 - x^2 + 17x + 3$$

 $f'(x) = 4x^3 - 2x + 17$

$$f'(-3) = -85 < 0$$
 so move right!





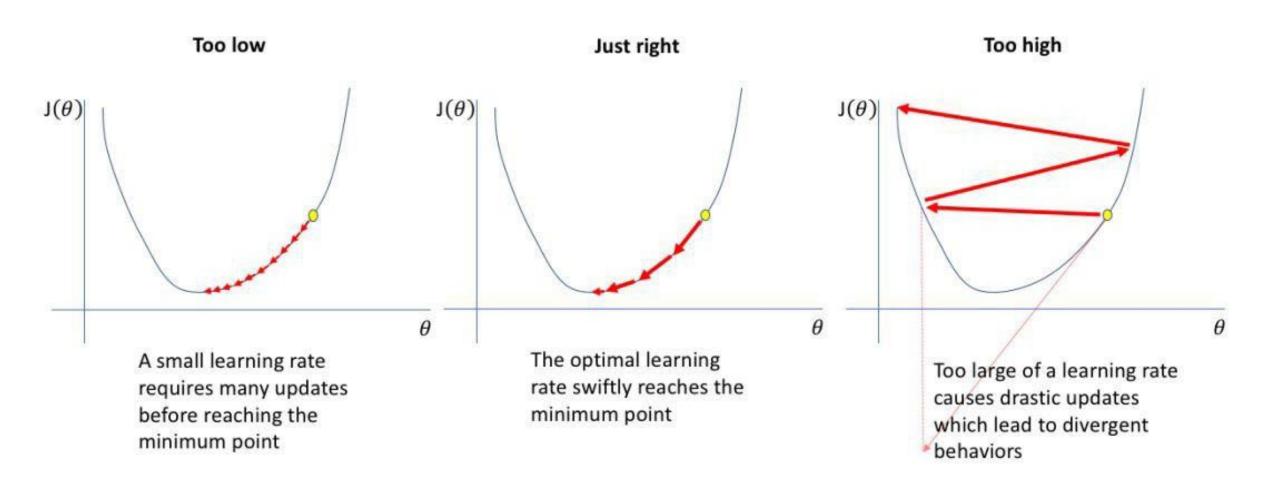
Learning rate

The step size in this approach is called the **learning rate**, η .

$$x = x - \eta * sign(f'(x))$$
 to minimize

$$x = x + \eta * sign(f'(x))$$
 to maximize





Gradients

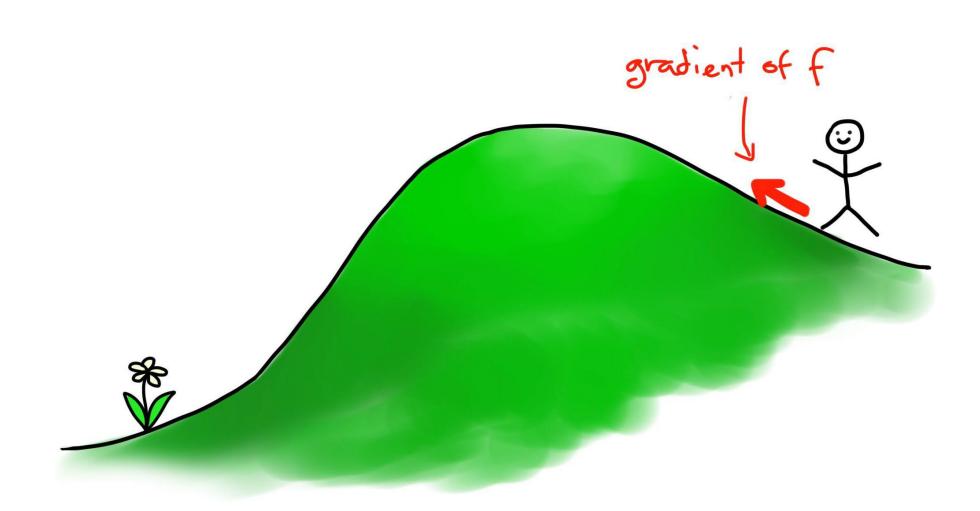


Gradient

The derivative at a point in k dimensions is called the **gradient**, which is a vector that points in the direction of steepest ascent.

$$f'(x_1, x_2, ... x_k) = \left(\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, ... \frac{\partial y}{\partial x_k}\right)$$





Definition

Gradient descent is an iterative algorithm for minimizing a differentiable multivariable function.

- 1. Pick a random starting point
- 2. Try to find the minimum value by taking a step of size η in the direction of the gradient
- 3. Repeat step 2 until you run out of steps

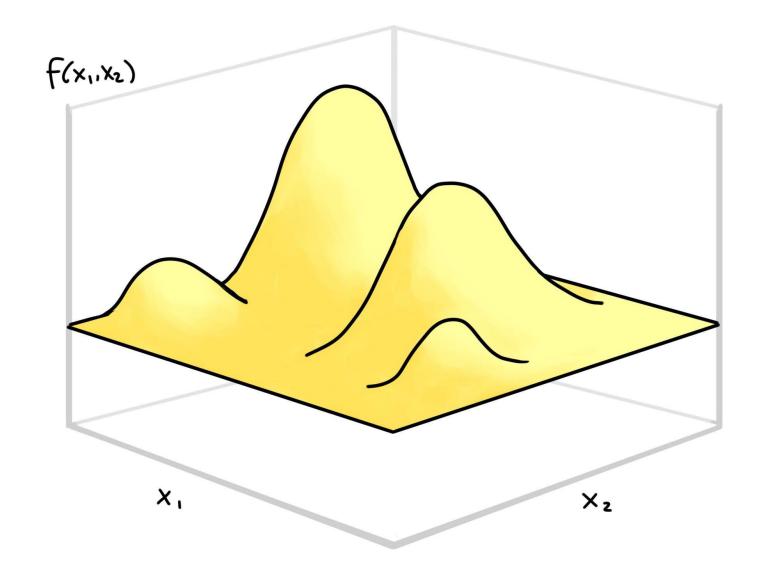


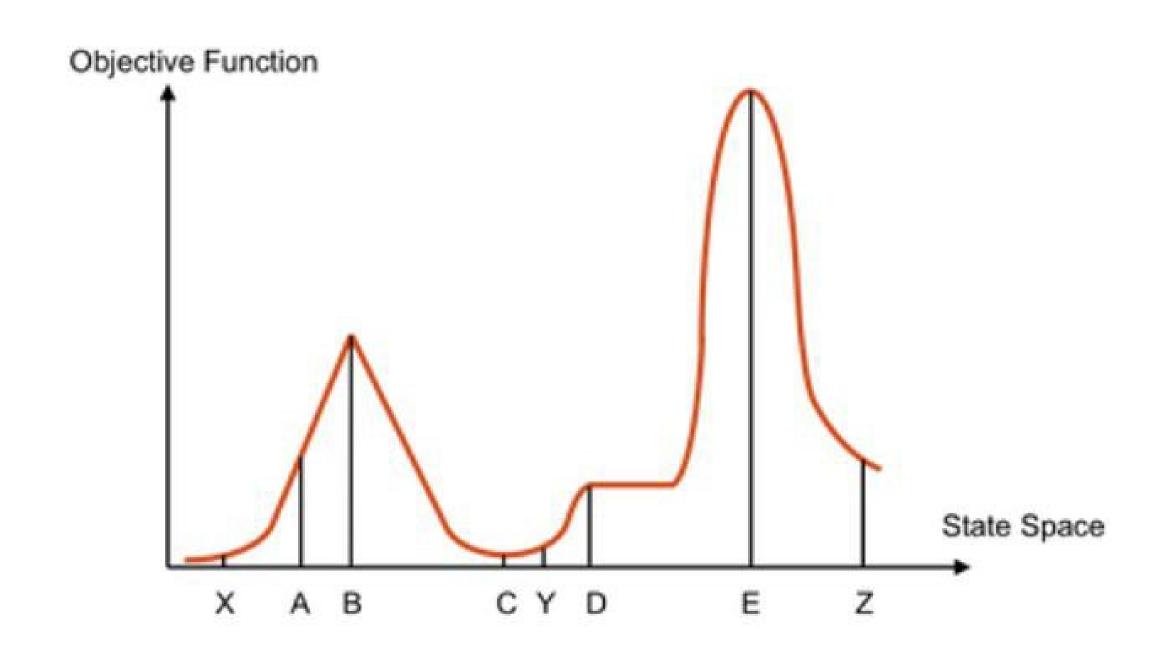
Gradient descent

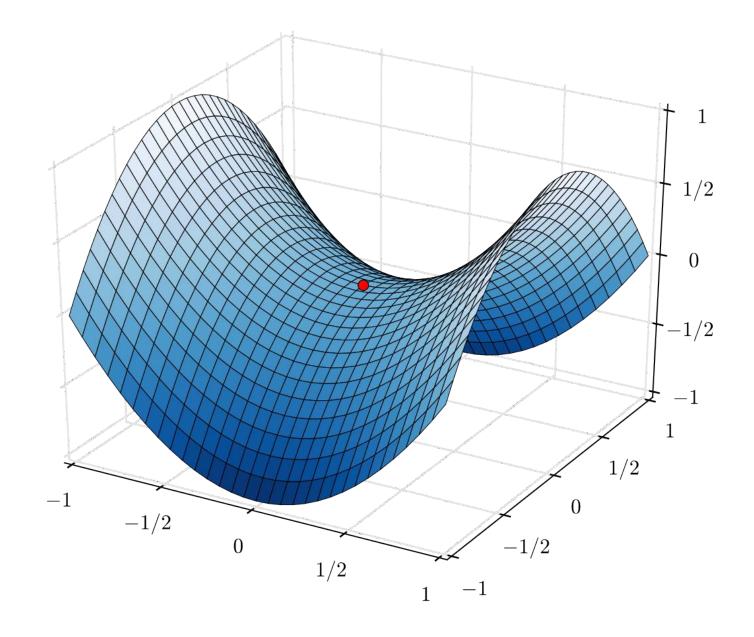
Gradient descent is very sensitive to both your random starting point and the shape (topography) of the space you are exploring.

...i.e. it doesn't always work well.









Note

Taking small steps against the gradient to try and find the global minimum is an example of following a **heuristic**.

And as we saw, there are always bad cases where any heuristic makes things worse, or cases where it is simply not helpful.



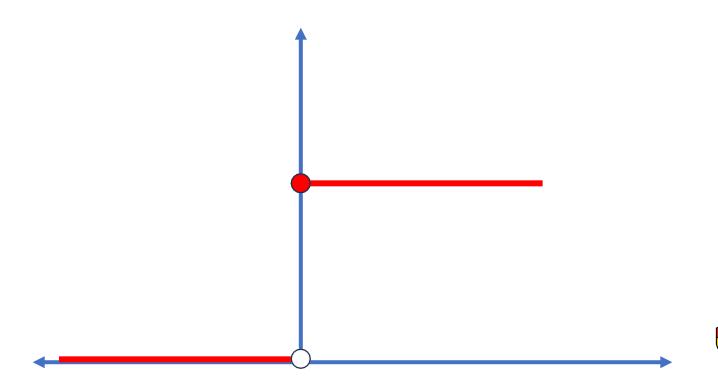
Activation Functions



Activation functions

One big problem with our current activation function is that it is not **continuous** (and the derivative is either flat or DNE).

This does not work well for our optimization!





Activation functions

To handle this, we have some alternative activation functions:

- Rectified linear unit (ReLU): $f(z) = \Phi_0^Z$ if z > 0 otherwise
- Sigmoid (logistic) function: $f(z) = \frac{1}{1+e^{-z}}$



Name \$	Plot	Function, $g(x)$ $\qquad \qquad \spadesuit$	Derivative of g , $g'(x)$ \Leftrightarrow	Range \$	Order of continuity
Identity		\boldsymbol{x}	1	$(-\infty,\infty)$	C^{∞}
Binary step		$\left\{egin{array}{ll} 0 & ext{if } x < 0 \ 1 & ext{if } x \geq 0 \end{array} ight.$	0	{0,1}	C^{-1}
Logistic, sigmoid, or soft step		$\sigma(x) \doteq rac{1}{1+e^{-x}}$	g(x)(1-g(x))	(0, 1)	C^{∞}
Hyperbolic tangent (tanh)		$ anh(x) \doteq rac{e^x - e^{-x}}{e^x + e^{-x}}$	$1-g(x)^2$	(-1,1)	C^{∞}
Soboleva modified hyperbolic tangent (smht)	5	$\mathrm{smht}(x) \doteq rac{e^{ax} - e^{-bx}}{e^{cx} + e^{-dx}}$		(-1, 1)	C^{∞}
Rectified linear unit (ReLU) ^[13]		$egin{aligned} (x)^+ &\doteq egin{cases} 0 & ext{if } x \leq 0 \ x & ext{if } x > 0 \ &= \max(0,x) = x 1_{x>0} \end{aligned}$	$\left\{egin{array}{ll} 0 & ext{if } x < 0 \ 1 & ext{if } x > 0 \end{array} ight.$	$[0,\infty)$	C^0

Activation functions

- A good activation function needs to be:
- Non-linear
- Easy to compute
- Continuous and differentiable (almost everywhere)
- Numerically stable, especially when combined with your loss function



Altogether



Neural nets

The goal of training is to find $w_1, w_2, ... w_n, b_1, b_2, ... b_m$ that minimize:

$$T(w_1, w_2, ... w_n, b_1, b_2, ... b_m) = \frac{1}{m} \sigma_{i=0}^m L(g_i, a_i),$$

where g_i is the guess output by our neural net for sample i.



Neural nets

The true objective function $T(w_1, w_2, ... w_n, b_1, b_2, ... b_m)$ depends simultaneously on <u>all</u> of the training samples.

• This is a lot of information/computation to factor in all at once



Neural nets

Instead, we usually prefer to improve $L(g_i, a_i)$ for a single <u>random</u> sample (x_i, y_i) at a time.

This approach is called **stochastic gradient descent**

Can prove this is almost as good as the original



Epochs

In practice, we usually shuffle the training data and then iterate through the samples, improving the loss on each one in turn.

Seems to be just as good as picking randomly each time

A single pass through the training data of this type is called an epoch.



Next...convolutional neural networks



