

Computer Vision: Fall 2022 — Lecture 9

Dr. Karthik Mohan

Univ. of Washington, Seattle

October 28, 2022

Check-In

- 1 Mini Project 1 Assigned

Check-In

- ① Mini Project 1 Assigned
- ② Identify your team mate through the spreadsheet

Check-In

- ① Mini Project 1 Assigned
- ② Identify your team mate through the spreadsheet
- ③ Other thoughts/questions?

References

- ① Good Book for Machine Learning Concepts
- ② Deep Learning Reference

Today

- ① Neural Networks/Deep Learning
 - ② Gradient Descent
 - ③ Back propogation
 - ④ Overfitting in Deep Learning
- } mechanics of DL

Introduction to Deep Learning

Deep Learning

- 1 Lot of buzz around Deep Learning in the past decade!

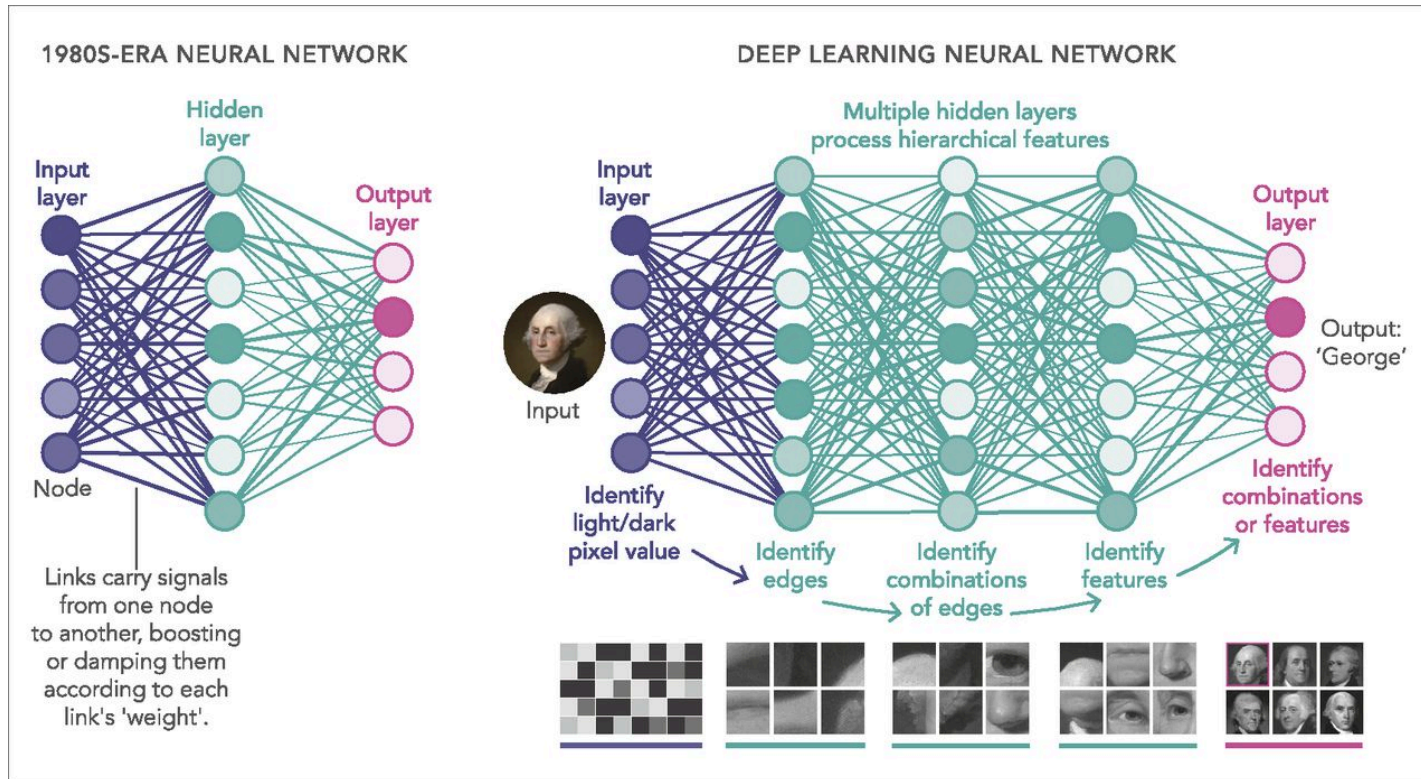
Introduction to Deep Learning

Deep Learning

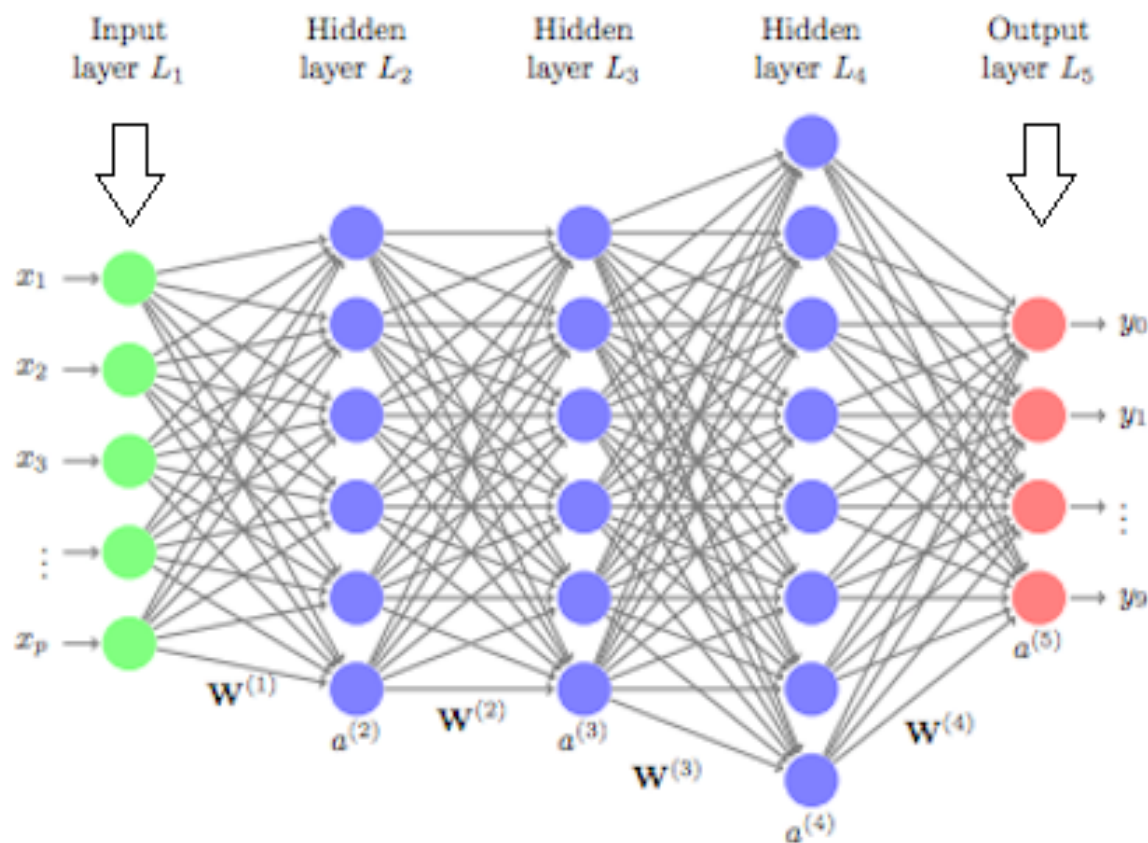
- ① Lot of buzz around Deep Learning in the past decade!
- ② Deep Learning refers to Neural Networks that is a loose approximation of how the brain works

Feed-forward Deep Learning Architecture Example

Feed Forward
→



Feed-forward Deep Learning Architecture Example



Other Neural Network Architectures

Neural Networks Zoo

- Input Cell
- Backfed Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Gated Memory Cell
- Kernel
- Convolution or Pool

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Lejnen asimovinstitute.org

Perceptron (P)



Feed Forward (FF)



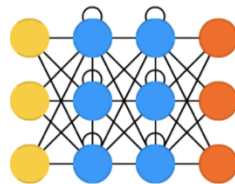
Radial Basis Network (RBF)



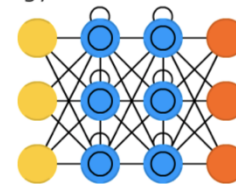
Deep Feed Forward (DFF)



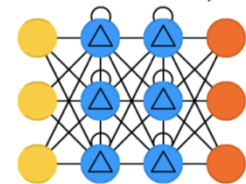
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denosing AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Hyper-parameters in Deep Learning

Hyper-parameters

- ① Learning rate
- ② Number of Hidden Layers
- ③ Number of neurons per hidden layer

Hyper-parameters in Deep Learning

Hyper-parameters

- ① Learning rate
- ② Number of Hidden Layers
- ③ Number of neurons per hidden layer
- ④ Type of non-linear activation function used ✓

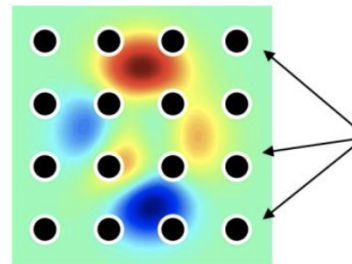
Hyper-parameters in Deep Learning

Hyper-parameters

- ① Learning rate
- ② Number of Hidden Layers
- ③ Number of neurons per hidden layer
- ④ Type of non-linear activation function used
- ⑤ Anything else?

Hyper-parameter tuning methods

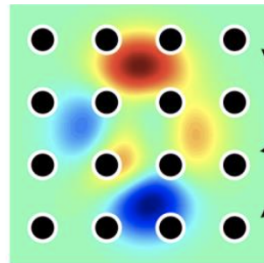
Grid search:



Hyperparameters
on 2d uniform grid

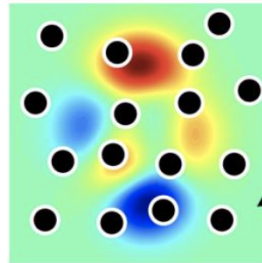
Hyper-parameter tuning methods

Grid search:



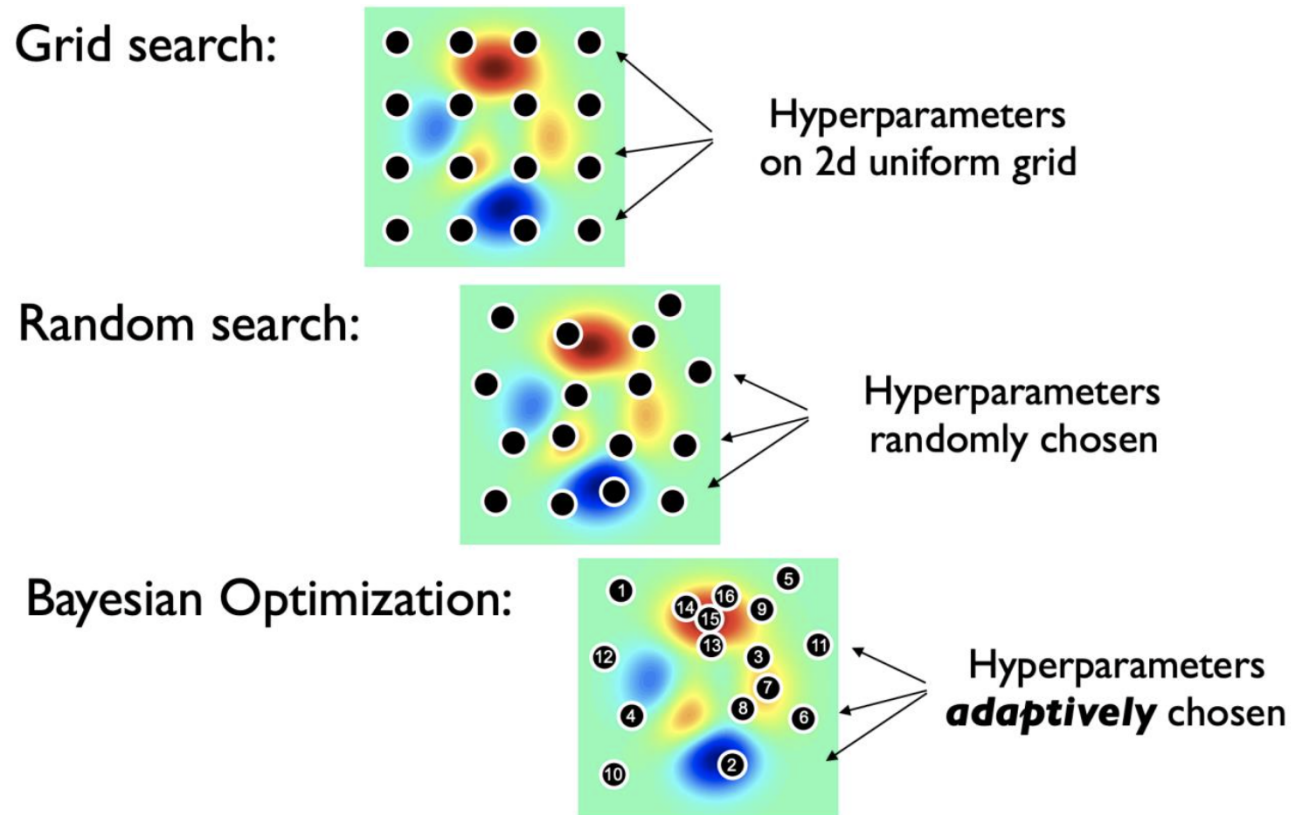
Hyperparameters
on 2d uniform grid

Random search:



Hyperparameters
randomly chosen

Hyper-parameter tuning methods



ICE #0

Compute the number of parameters in DNN model

Consider a DNN model with 3 hidden layers where each hidden layer has 1000 neurons. Let the input layer be raw pixels from a 100x100 image and the output layer has 10 dimensions, let's say for a 10 class image classification example. How many total parameters exist in the DNN model?

- ① 10 million parameters
- ② 11 million parameters
- ③ 12 million parameters ✓
- ④ 13 million parameters

Over-fitting in DNNs

How to handle over-fitting in DNNs

- 1 A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.

Over-fitting in DNNs

How to handle over-fitting in DNNs

- 1 A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.
- 2 Weight regularization can help - l_1, l_2

$$\min \ell(w) + \lambda \underbrace{\|w\|_1 / 2}_{\text{Regularization}}$$

Over-fitting in DNNs

How to handle over-fitting in DNNs

- ① A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.
- ② Weight regularization can help - ℓ_1, ℓ_2
- ③ More common over-fitting strategy for DL?

Over-fitting in DNNs

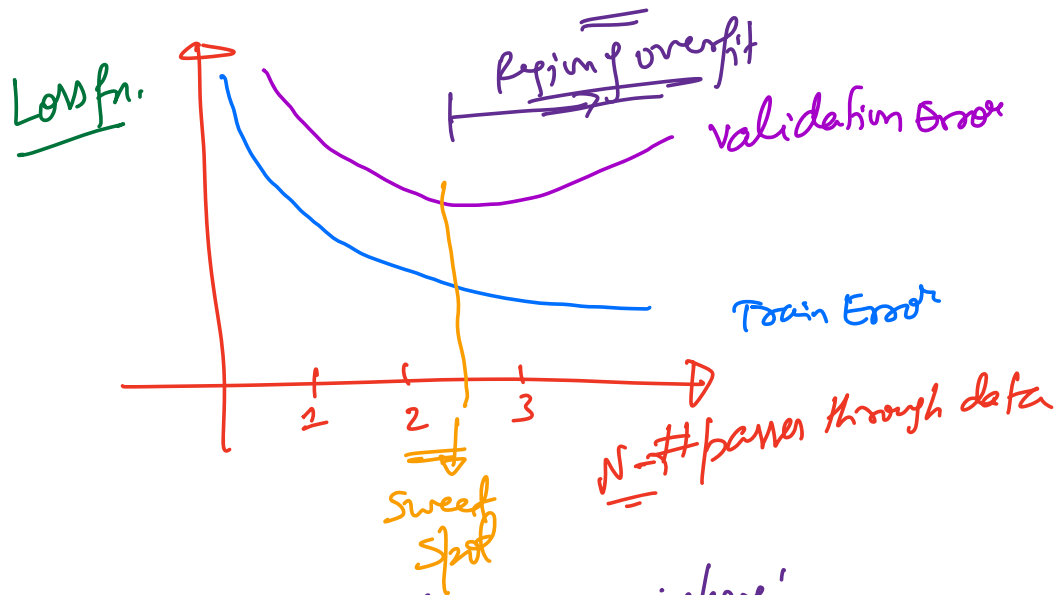
How to handle over-fitting in DNNs

- ① A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.
- ② Weight regularization can help - ℓ_1, ℓ_2
- ③ More common over-fitting strategy for DL?
- ④ Dropouts!

Over-fitting in DNNs

How to handle over-fitting in DNNs

- 1 A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.
- 2 Weight regularization can help - ℓ_1, ℓ_2
- 3 More common over-fitting strategy for DL?
- 4 Dropouts!
- 5 Early stopping is also a great strategy! Stop training the DL model when the validation error starts increasing. How's this different from regular validation we were doing earlier??



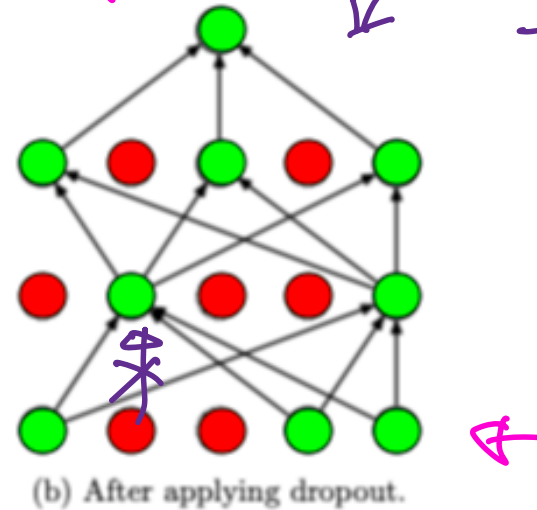
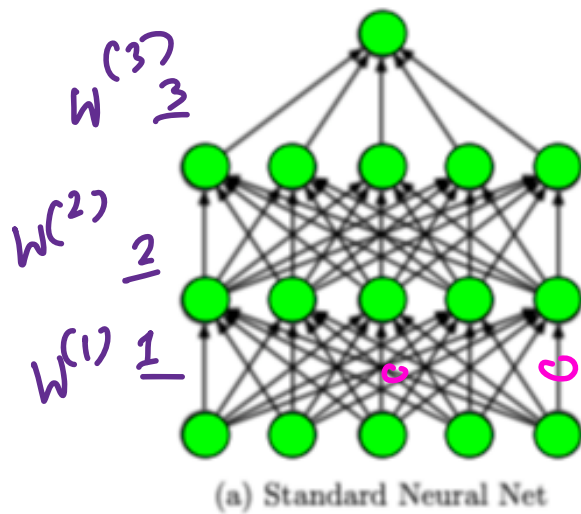
↳ stop training here!
 = "Early Stopping"

Over-fitting in DNNs

How to handle over-fitting in DNNs

- 1 A DNN model with 100 million parameters and only 100k data points or even a million data points will overfit unless we take care of over-fitting.
- 2 Weight regularization can help - ℓ_1, ℓ_2
- 3 More common over-fitting strategy for DL?
- 4 Dropouts!
- 5 Early stopping is also a great strategy! Stop training the DL model when the validation error starts increasing. How's this different from regular validation we were doing earlier??
- 6 Book by Yoshua Bengio has tons of details and great reference for Deep Learning!

Taking care of Over-fitting: Dropouts



Dropout helps overfit
- Makes training faster
- Ensemble of models
Specific kind of network

==

==

Hyper parameters \leftarrow dropout = 0.6

ICE #1

Size of Dataset \sim Size of parameters / weights?

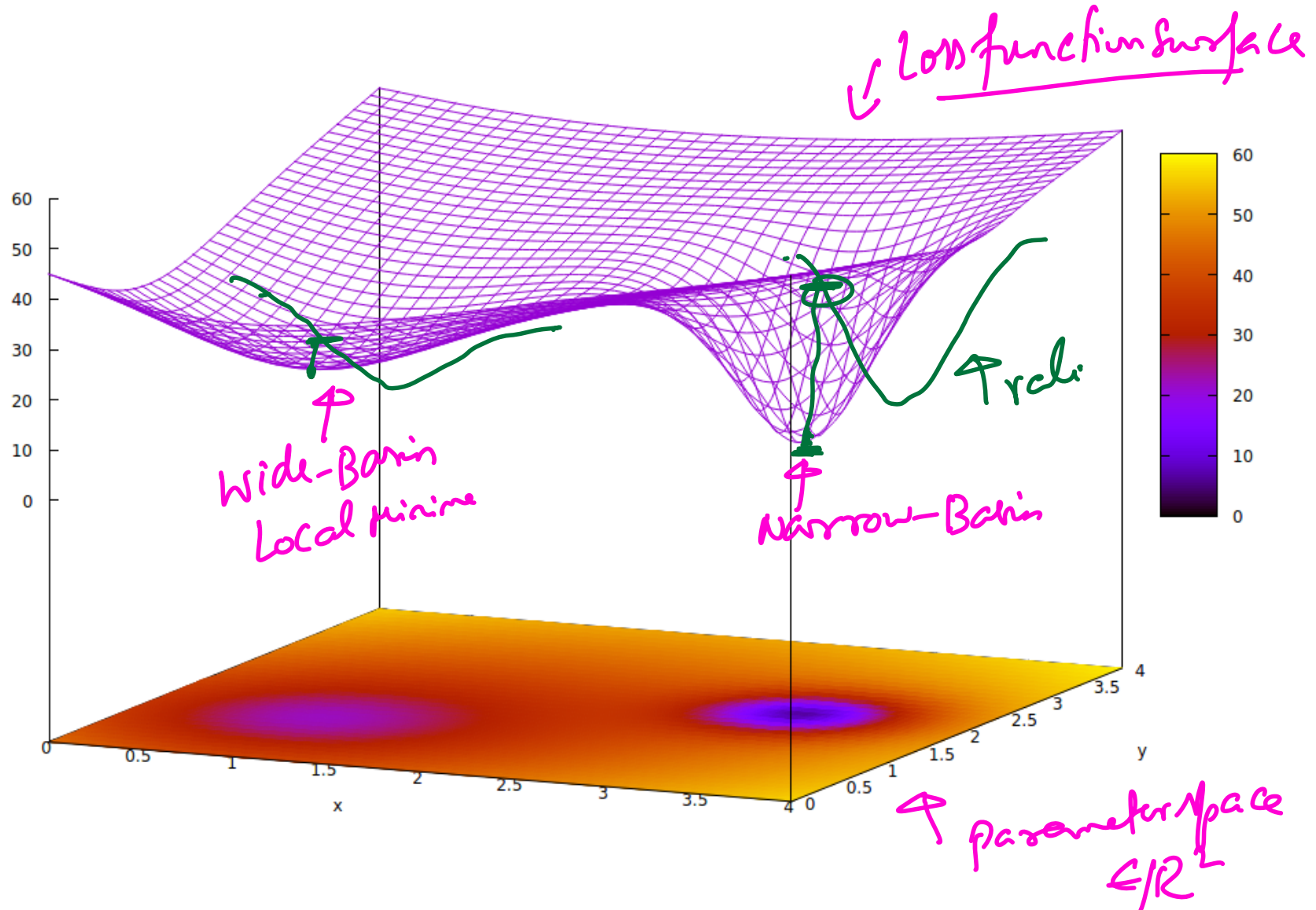
Overfitting

Say you trained your deep learning model on an image data set to classify images as a cat, dog or a human. You realize you are probably overfitting as you are doing much better on training set as compared to the validation and test set. Which of these strategies would probably help reduce overfitting?

- Data Augmentation
- 1 Create variations of the images in the training set using image processing techniques and adding it to the training set ✓
 - 2 Using ℓ_1 regularization on the weights ✓
 - 3 Down-sampling the training data set and retraining ✗
 - 4 Error analysis to understand which parts of the validation data does the model perform poorly on? ✓

Good vs Bad Local minima

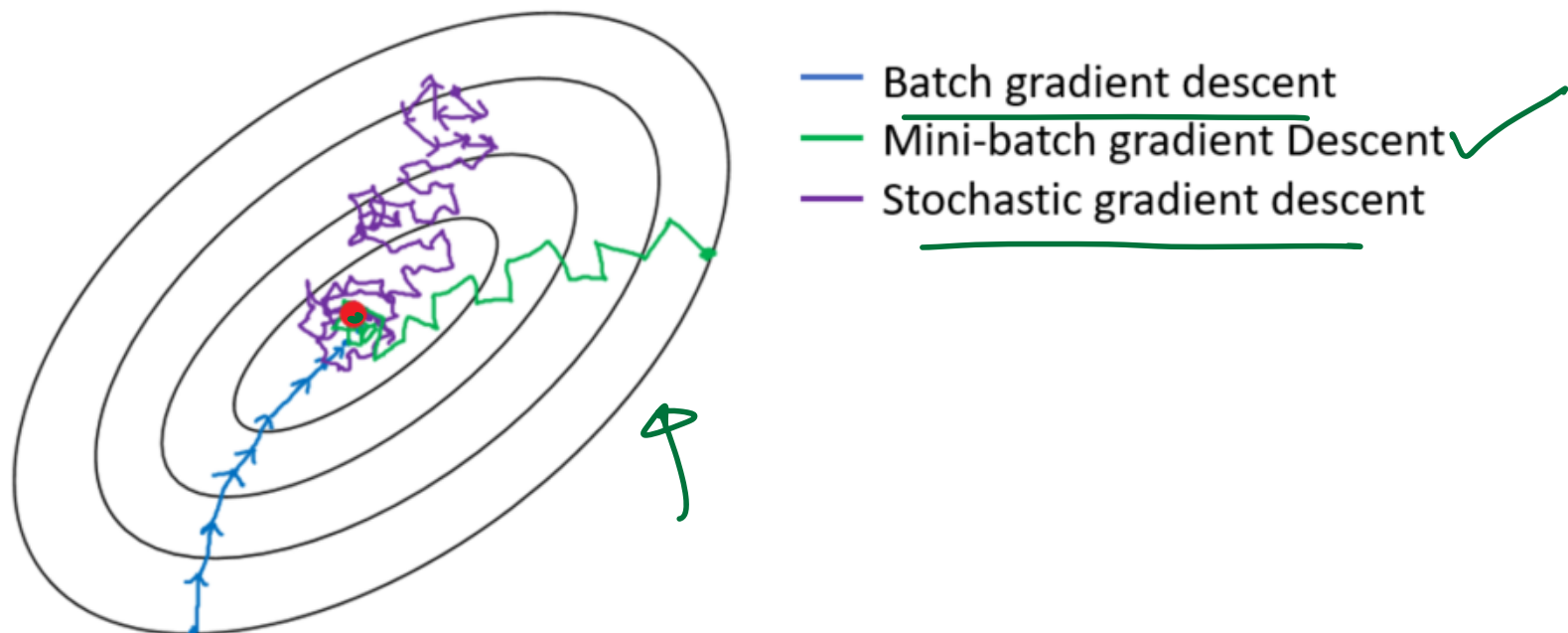
(Mechanics of DL)



Algorithmic foundations to Machine Learning

Underlying Engine behind ML Training

(Mini-batch) Stochastic Gradient Descent Almost every model and problem-space in ML uses SGD of some kind - Clustering, Regression, Deep Learning, Computer Vision and NLP to name a few. Almost every algorithm in every library - Scikit-learn, Keras, Pytorch, etc uses **mini-batch SGD under the hood**.



So what is Gradient Descent?

Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .

So what is Gradient Descent?

Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .

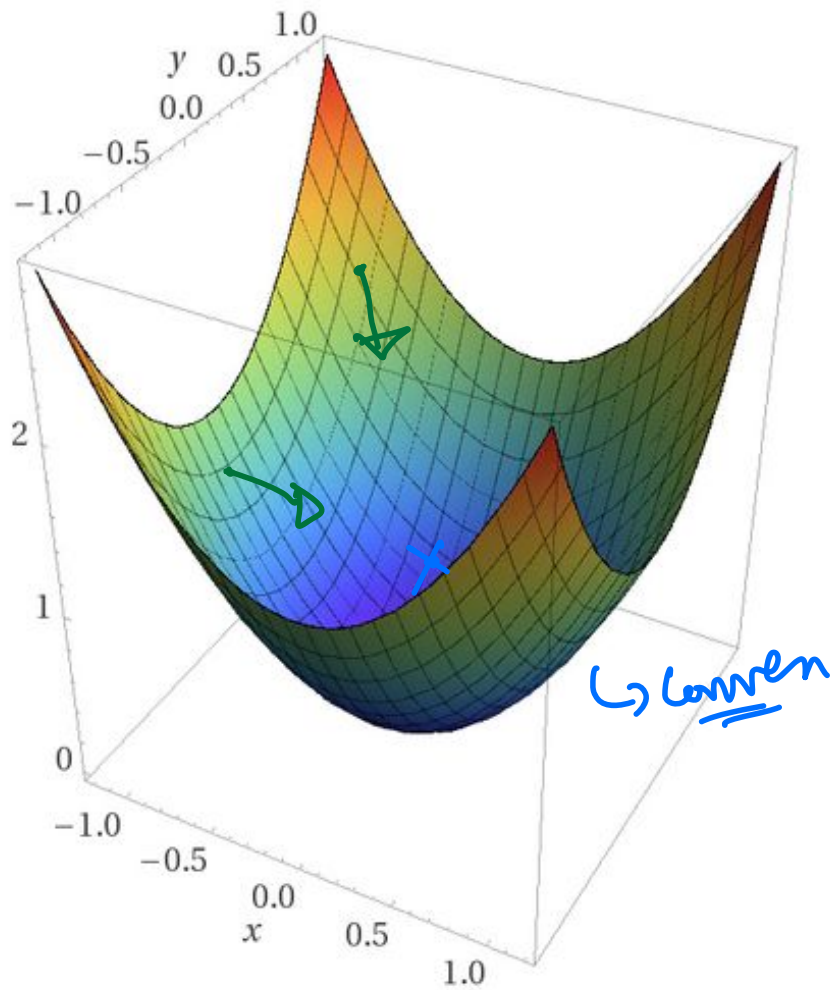
Why is this important?

Consider the Logistic Regression problem. \hat{w} is a local optimum (and global optimum) to the function

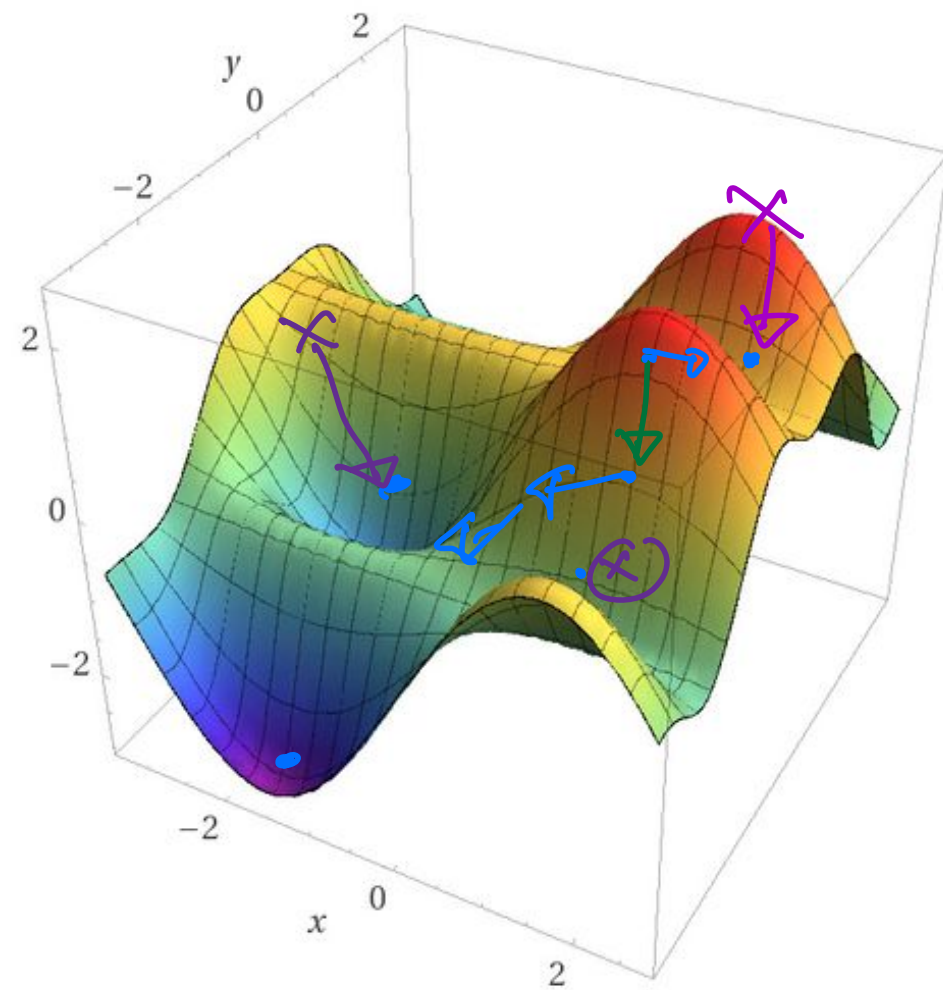
$$l(w) = \sum_i y_i \log(1 + e^{-w^T x^i}) + \sum_i (1 - y_i) \log(1 + e^{-w^T x^i})$$

Binary Cross-Entropy

Negative Gradient helps you view the direction of descent



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

Negative Gradient = Direction of Steepest Descent

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

① **Initialize** $w = w_0$ (maybe randomize)

② **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$
↳ Learning Rate / step size

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)
- 2 **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$
- 3 **Iterate** Repeat step 2 until w converges, i.e.

$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

ICE #2

Derivative (1 min)

Find the derivative of w^2

- a) $2w$ ✓
- b) w
- c) $0.5w$
- d) 0

ICE #3

Pytorch: Automatic Differentiation

$$\lambda(w_1^2 + w_2^2 + \dots + w_d^2)$$

↑

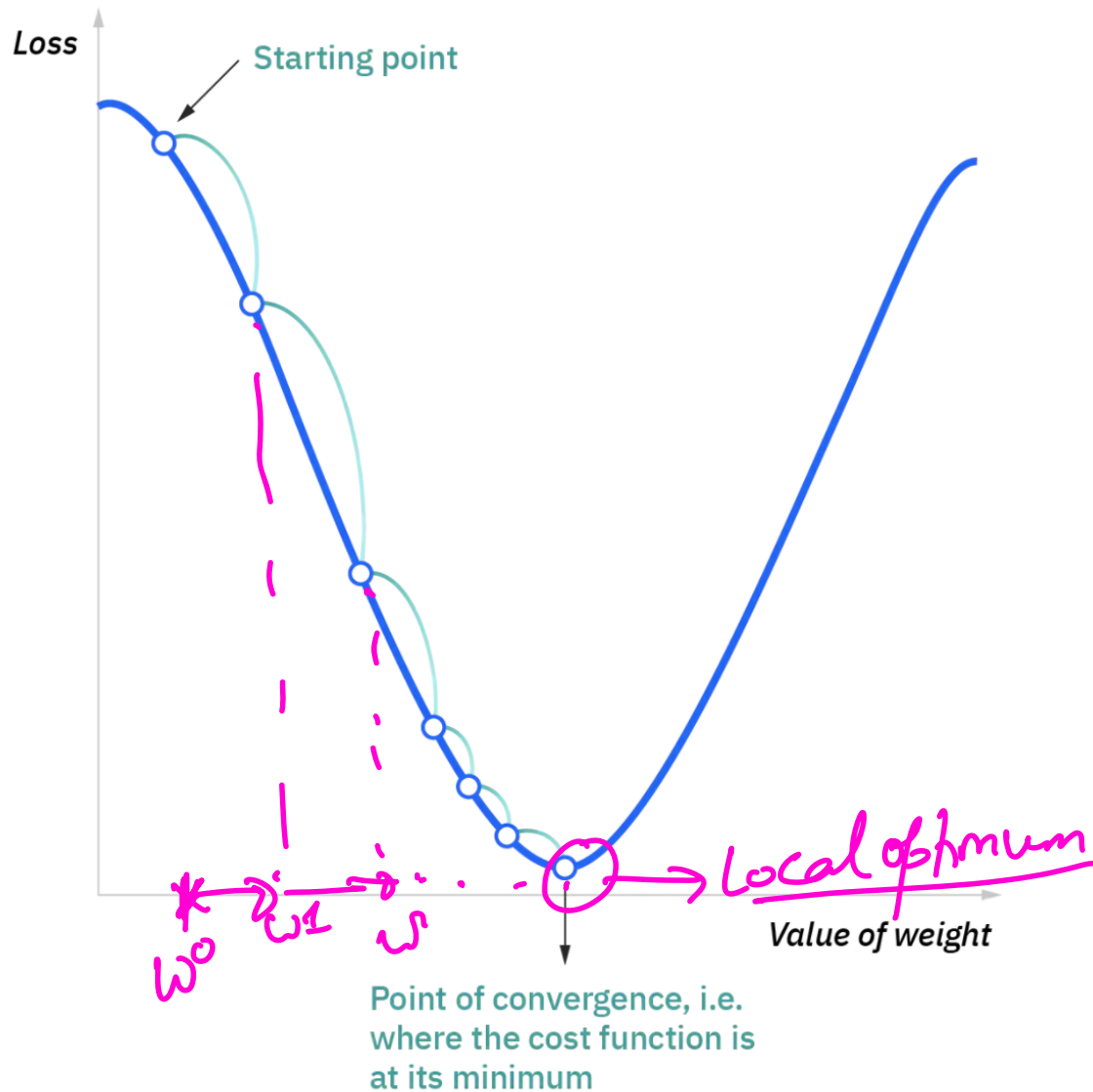
Gradient of Ridge Regularizer (2 mins)

Find the gradient of the regularization function, $R(w) = \lambda \|w\|_2^2$. I.e. obtain the expression for, $\nabla_w R(w)$?

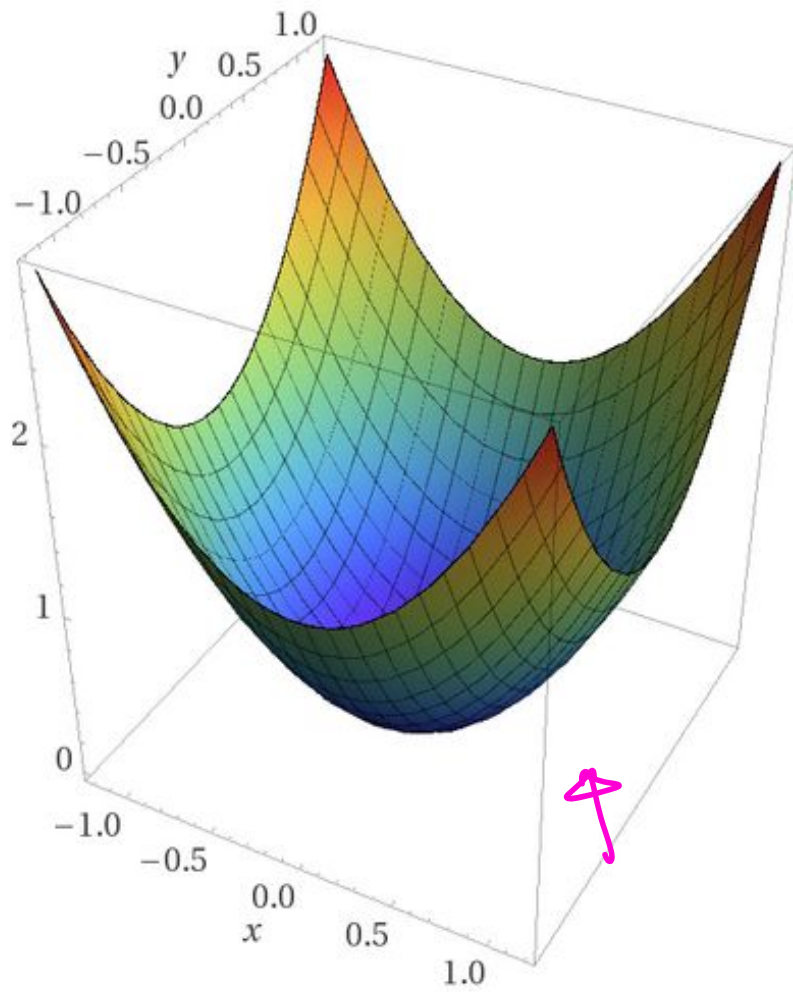
- a) $2\lambda \|w\|_2$
- b) $\lambda \|w\|_2 w$
- c) $2\lambda w$
- d) $2\lambda \|w\|_2 w$

$$\begin{aligned} \frac{\partial}{\partial w_1} R(w) &= \frac{\partial}{\partial w_1} (\lambda w_1^2) = 2\lambda w_1 \\ \frac{\partial}{\partial w_j} R(w) &= 2\lambda w_j \end{aligned}$$
$$\begin{bmatrix} 2\lambda w_1 \\ 2\lambda w_2 \\ \vdots \\ 2\lambda w_d \end{bmatrix} = 2\lambda w$$

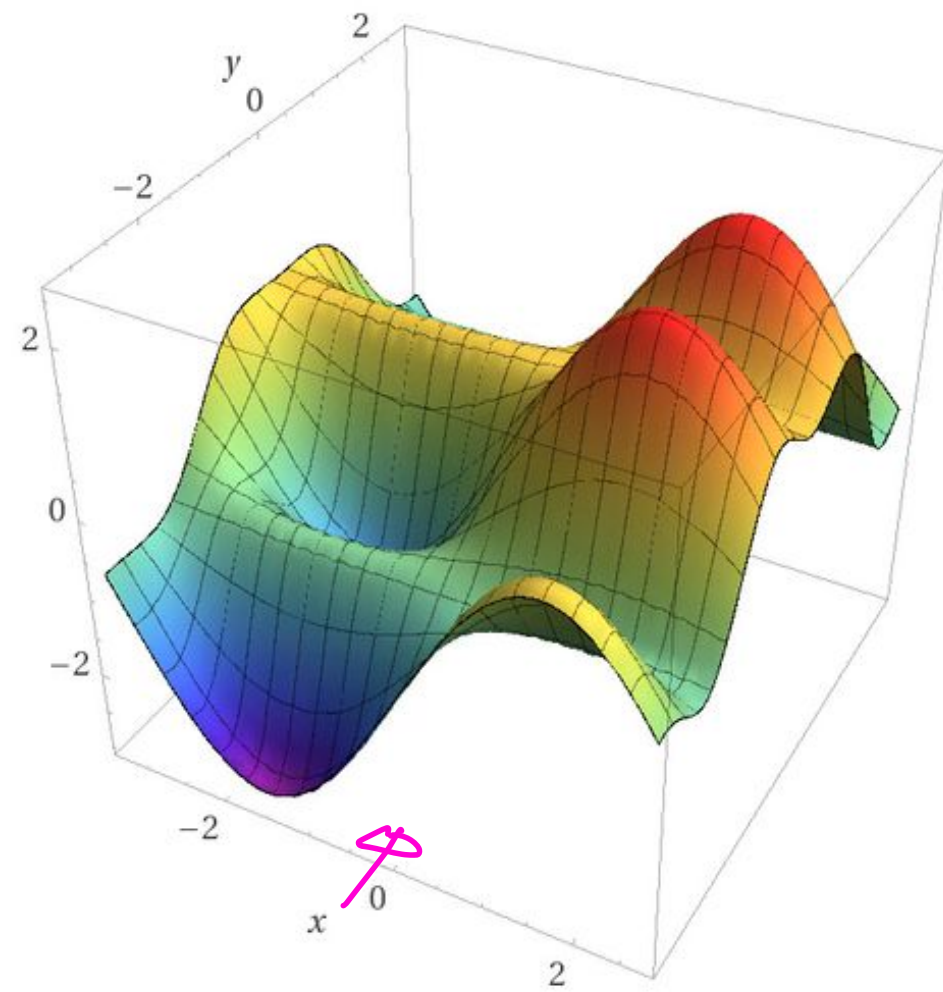
GD in one dimension



Loss function in 2 dimensions



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

Gradient Descent Properties

- 1 Gradient Descent converges to a local minimum

Gradient Descent Properties

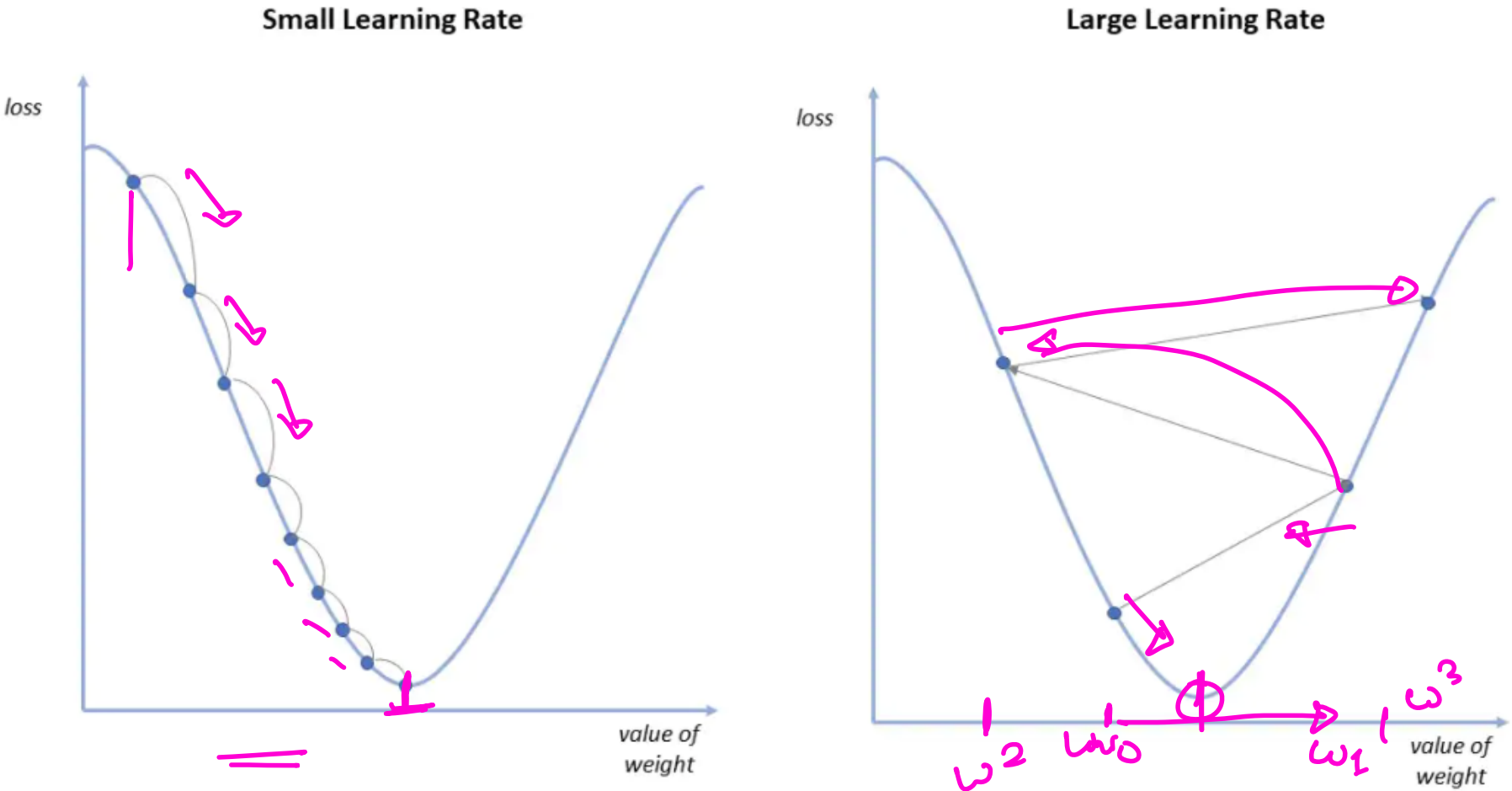
- ① Gradient Descent converges to a local minimum
- ② If L is a convex function, all local minima become a global minima!

Gradient Descent Properties

- 1 Gradient Descent converges to a local minimum
- 2 If L is a convex function, all local minima become a global minima!
- 3 Wherever we start, gradient descent usually finds a local minima closest to the start.

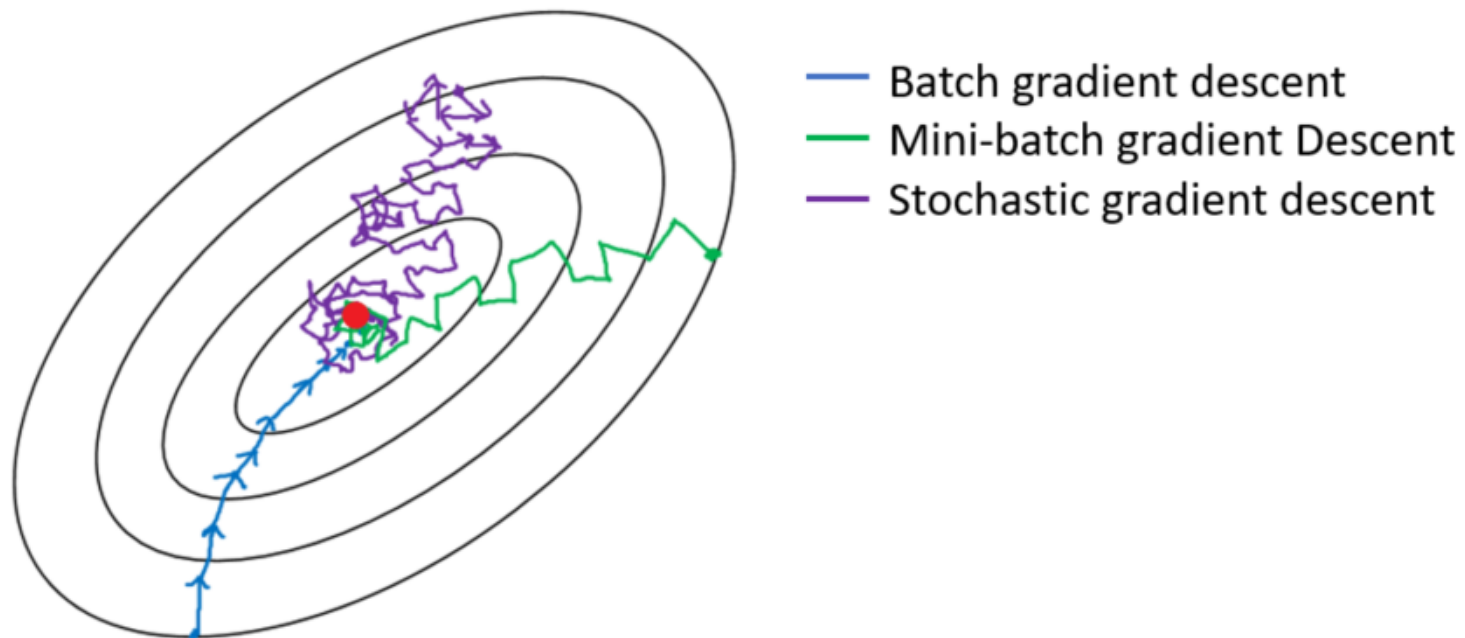
↳ Pick starting point at Random

Effect of Learning Rate



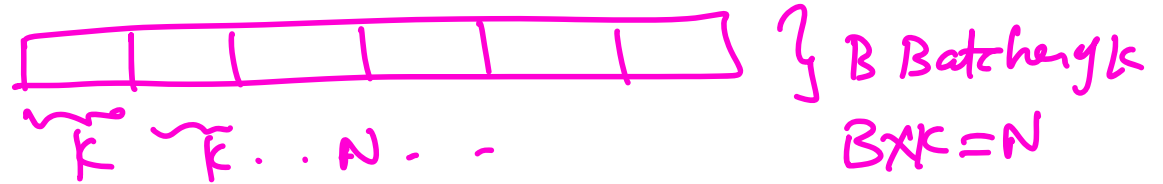
(Adaptive Learning Rate)

SGD behavior in search space



SGD in practice - mini-batch SGD!

mini-batch SGD



Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

$k = 128$ or 256

- 1 **Initialize** $w = w_0$ (randomize)

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$

SGD in practice - mini-batch SGD!

mini-batch SGD

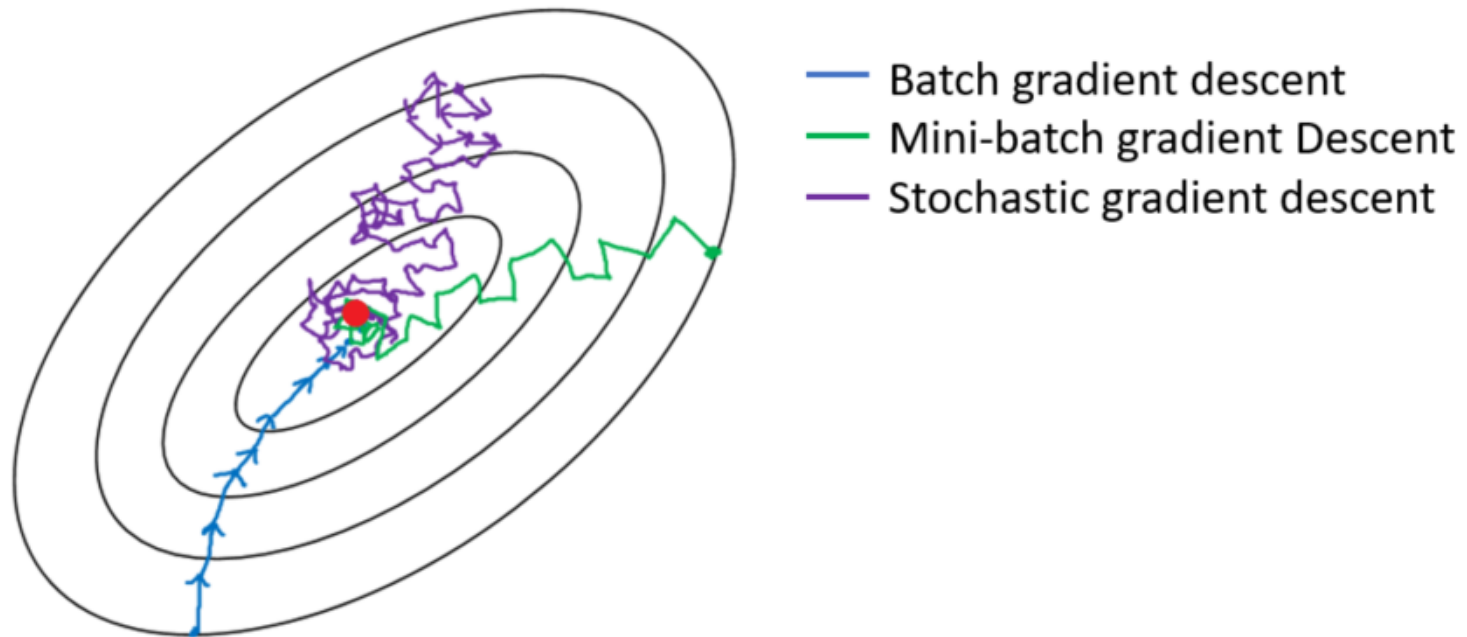
} *How many passes through the data?*

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

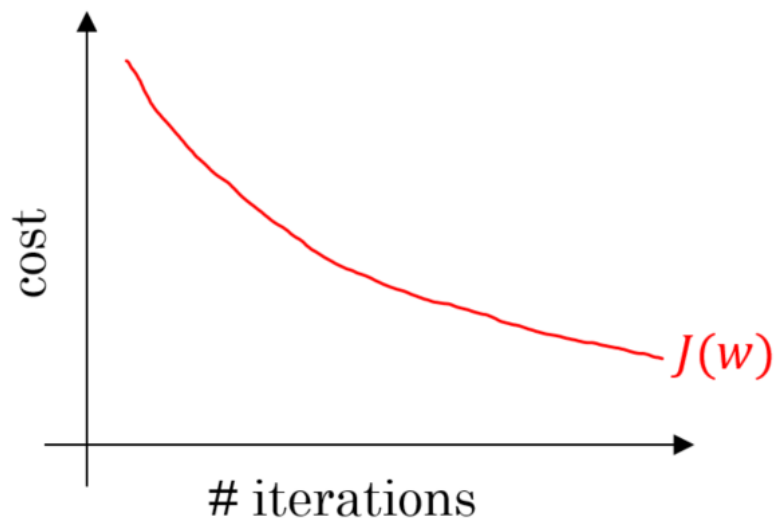
$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

GD behavior in the search space

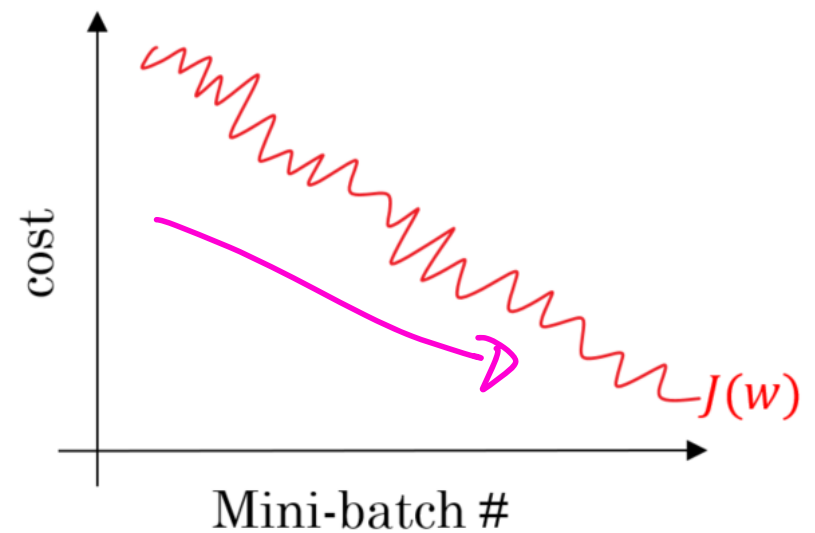


GD vs Mini-batch convergence behavior

Batch gradient descent



Mini-batch gradient descent

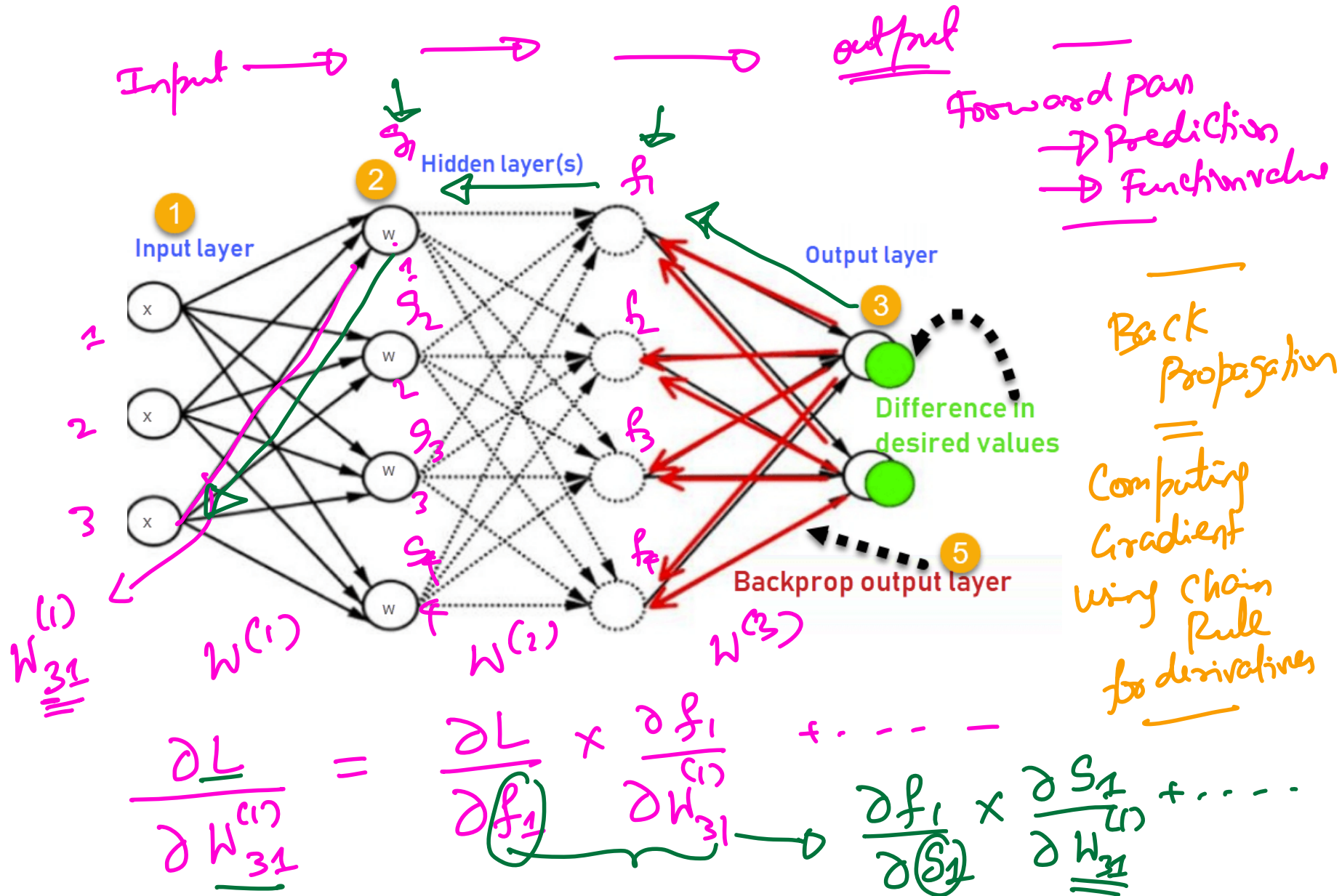


GD vs mini-batch SGD

→ In practice

Factor	<u>GD</u>	<u>Mini-batch SGD</u>
<u>Data</u>	All per iteration	Mini-batch (usually <u>128</u> or 256)
Randomness	<u>Deterministic</u>	Stochastic
Error reduction	Monotonic	<u>Stochastic</u>
Computation	<u>High</u>	Low
<u>Memory big data</u>	Intractable	<u>Tractable</u>
Convergence	Low relative error	<u>Few "passes" on data</u>
Local Minima traps	Yes	<u>No</u>

Forward Propagation vs Back-propagation in NN



ICE #4

Back Prop

How is back-prop related to Gradient Descent?

- 1 Back-propagation is an alternative to Gradient Descent for Neural Networks
- 2 Back-propagation computes the gradient that can then be used in gradient descent
- 3 Back-prop is the same as gradient descent for neural networks
- 4 Back-prop is a different concept from Gradient Descent

Mini-Project Pointers

- ① Form a team of 2 (today if you can!) Share your team on the spreadsheet in discord


Mini-Project Pointers

- ① Form a team of 2 (today if you can!) Share your team on the spreadsheet in discord
- ② Play with the architecture details in your modeling process. Start simple and add more layers, more neurons per layer if it's help your validation metrics. Hyper-parameters are tuned on validation set

Mini-Project Pointers

- ① Form a team of 2 (today if you can!) Share your team on the spreadsheet in discord
- ② Play with the architecture details in your modeling process. Start simple and add more layers, more neurons per layer if it's help your validation metrics. Hyper-parameters are tuned on validation set
- ③ **Two Deadlines for Mini-Project:** First one is November 6th as a check-point with deliverables including baseline and your first NN model. Second includes full report, best Kaggle submission, all metrics and CNN model that is due November 13th.

Mini-Project Pointers

- ① Form a team of 2 (today if you can!) Share your team on the spreadsheet in discord
 - ② Play with the architecture details in your modeling process. Start simple and add more layers, more neurons per layer if it's help your validation metrics. Hyper-parameters are tuned on validation set
 - ③ **Two Deadlines for Mini-Project:** First one is November 6th as a check-point with deliverables including baseline and your first NN model. Second includes full report, best Kaggle submission, all metrics and CNN model that is due November 13th.
 - ④ Will work with PyTorch for this Mini-Project - Get yourself familiar with tutorials on this (Will be covered in quiz section as well)
- 

Summary

- 1 Introduction to Neural Networks
- 2 Neural Network Architecture and its components
- 3 Gradient Descent
- 4 Backpropagation in Neural Networks
- 5 Overfitting in Neural Networks