

Backpropagation

An algorithm that trains a neural network.

weights / bias : optimum value is found for it.

seen

Epoch: An epoch simply means the model has every sample in your dataset once

epochs = 5

) One complete pass through all training data

Batch: If there are 4 rows

1 epoch, It sees all 4 rows at once (parallel)

5 epochs: 20 times, Batch = 4

Here batch size =

for i in range (epoch):

for j in range (4):

select 1 row (random)

Predict (using forward prop)

Calc loss (mse)

Update weights = $w_0 - n \frac{\partial L}{\partial w}$

Calculate avg loss for every epoch

Learning = minimizing mistakes

If changing a weight reduces error, change it in that direction.

→ Loss function is a function of all trainable parameters.

→ Concept of gradient:

↳ fancy word for derivative

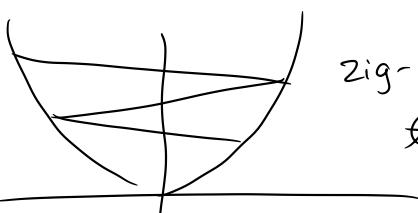
$y = f(x)$ Dependent

$$\frac{dy}{dx} = 2 \quad \text{So, upon doing 1 change in } x \text{ would bring 2 change in } y$$

We don't do $\frac{\partial L}{\partial w} = 0$ because neural networks

$$w_n = w_n - n \frac{\partial L}{\partial w}$$

contain millions of parameters



Zig-zag motion

that's why keep learning rate

They contain nonlinear activations

(ReLU, sigmoid, tanh)

No analytic solution exists for them

$$w \leftarrow w - n \nabla L(w)$$

What is convergence?

We keep on updating our weights till

$$w_n' = w_n - \frac{\eta \delta L}{\partial w_n} \approx 0$$

$\frac{\delta L}{\partial w_n} \approx 0$ means $w_n' \approx w_n$
we have hit a flat surface.

For Regression problems, output: linear activation

Why do we need weights & bias ??

Neural networks learn a function that maps input & output
We don't know the fn explicitly, so the network learns it
from data

Weighted sum:

$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 - \dots$$

Weights represent importance of inputs.

large weight \rightarrow input matters a lot

Add bias:

Bias shift the decision:

\rightarrow Controls when the neuron activates.

\rightarrow like an adjustable threshold.

Activation:

Without it, the neural networks become just a linear model.

$$y = w_2 (w_1 x) = (w_2 w_1) x$$

add non-linearity

2) Ability to learn complex functions.

Multiple layers
Each layer learns features of increasing complexity

Is this info important enough to read

pars nothing?

pars little?

pars truly or very

pars strongly ?

Memoization: optimization technique from Dynamic programming (DP) that speeds up computations by caching techniques from Dynamic Programming (DP) that speeds

(let's say Fibonacci Recursion)

$$f(n) = f(n-1) + f(n-2)$$

Stores outputs of f^n s with specific inputs (e.g., $f(n)$) in a table

There are 3 variants of gradient descent, which differ in how much data we use to compute the gradient of the objective f^n .

Trade off b/w time & accuracy

Batch GD (Vanilla GD)

entire dataset update

50 pts \rightarrow predict

5 epochs

50 predictions (Compare it with y)

epoch = 10

for in range(10):

$y_{\text{hat}} = \text{np.dot}(x, w) + b$.

50 values.

$y = 50$ values

10 times update

w, b update

no. of epochs =

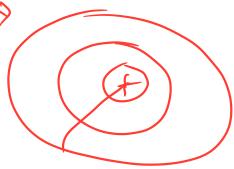
no. of updates

Which is faster (given same no. of epochs): Batch

Which is faster to converge (given same no. of epochs): SGD

Stochastic: We keep pt. in random and do

Batch: We take the whole dataset into account.



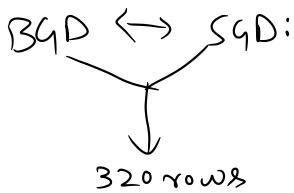
The spiky behaviour has good as well bad effects
Batch might fall into local minima
Sgd might jump out of it

Vectorization in Batch

$$\text{np. dot}(x, w) + b$$

(For small datasets its very fast)

② → dataset of 10 crore
problem when the dataset is very large.



Mini-batch Gradient Descent

No. of batches = 10 updates

$$\text{no. of batches} = \text{no. of updates} = \frac{\text{no. of rows}}{\text{Batch size}}$$

fast : bgd > mbgd > sgd

conv : bgd < mbgd < sgd

,

Why is batch_size is provided in multiple of 2?

Optimization technique 4, 8, 32, 128

RAM takes binary values
faster

What if batch_size doesn't divide # no. rows properly?

Rows = 800 Batch size = 150

B1: 150, B2: 150

So, 3 updates per epoch.

Nanishing Gradient Problem

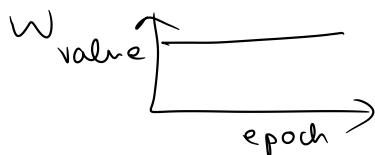
$$\frac{d(\text{sigmoid})}{dx} = \boxed{0-0.5}$$

The neural network might stop training.

Major reasons why NNs failed in so

How to recognize:

- 1) focus on loss: epoch → no changes → VGP
- 2) weights → graph values



How to handle Vanishing Gradient Problem?

- 1) Reduce model complexity. Not applicable because we want the model to identify complex patterns.
- 2) Using ReLU \rightarrow derivative $\rightarrow 0$ or 1
dying relu: $\text{relu} = 0$
 \hookrightarrow then makes it zero.
- 3) Proper Weight Initialization (Glorot, XAVIER)
- 4) Batch Normalization
- 5) Residual Network

Sigmoid, tanh leads to Vanishing gradient problem.

Solved by ReLU, leaky ReLU

Exploding Gradient Problem

RNNs $10, 10, 10, 10 = 10^4$ linear or exponential

$$\frac{\partial L}{\partial w} = d_1 \cdot d_2 \cdot d_3 \quad \text{tanh}$$

$$w_h = w_o - n \frac{\partial L}{\partial w}$$

How to improve a neural network??

- 1) Fine tuning NN hyper parameters:
 - No. of hidden layers.
 - no. of neurons per layer
 - learning rate
 - Optimizer
 - Batch size.
 - Activation function
 - no. of epochs.

2) By solving problems:

Vanishing gradients / Exploding
Not enough data
Slow training
Overfitting

Fine tuning Hyperparameters:

- 1) No. of hidden layers
Multiple hidden layers with fewer \rightarrow neurons \rightarrow 1 hidden layer with large no. of neurons.

Representation learning.
Initial layers capture primitive features. } Hierarchy learning

Transfer learning

stop adding layers till overfitting starts.

2) Neuron / layer



primitive features are captured early on & they are more in number.
Start with more & then reduce it based on overfitting

3) Batch Size

Batch Stochastic Min-Batch

Smaller
better results

large
training faster (warning the learning rate)
learning rate scheduler (start with small value & then raise it)

4) Epochs: $\xrightarrow{100}$ $\xrightarrow{500}$ $\xrightarrow{1000}$

Early Stopping
Keras \rightarrow Stabilize

Keras
 \downarrow
call back.

- Vanishing & exploding gradients
- ↳ weight init
- ↳ Activation fⁿ
- ↳ Batch Norm
- ↳ Gradient Clipping

- Not enough data
 - ↳ Transfer learning
 - ↳ Unsupervised pre-trained

Slow training
Optimizers
learning rate scheduler

Overfitting
 L_1 & L_2
reg
Dropouts

Problems with Neural Networks

Early Stopping

monitor = "val-loss"
min-delta = 0.0001

Tells where performance on validation set reduces

To what epochs we should continue

Normalizing Inputs

(Age) x_1 , x_2 (Salary)
 w_1 w_2 Magⁿ is high

$$w_2 = w_2 - n \frac{\partial L}{\partial w_2}$$