

TRAINING LOOP

1 Initialize parameters to random values

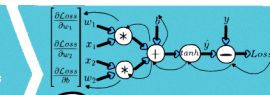


2 Forward pass



Predicted labels \hat{y} True labels y

3 Calculate Loss



4 Backward pass

CONTINUE...

LEARNING RATE

$$w_1 = w_1 - \eta \left(\frac{\partial \text{Cost}}{\partial w_1} \right)$$

0.1 0.001 ...

ΔLoss

5 Backpropagation

Small updates

UNTIL CONVERGENCE

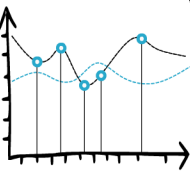


SPEED
Stability

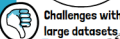
- More precise
- Slower
- Faster
- Right overshoot the minimum

FULL-BATCH

Use average gradient over entire dataset to update parameters in a single step



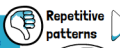
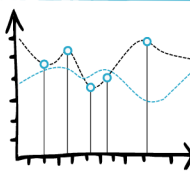
Less noisy



Challenges with large datasets

INCREMENTAL

Update parameters after each training example is processed



Repetitive patterns



Smaller updates



Random samples



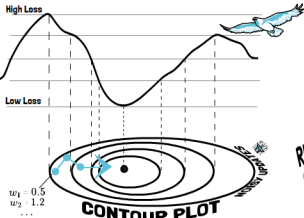
Noisy

STOCHASTIC GRADIENT DESCENT (SGD)

More robust

SGD
Full Batch
Mini-Batch
Batch Size (32, 64, 128, ...)
Efficient
Less noisy

GRADIENT DESCENT OPTIMIZATIONS



MOMENTUM
Keep moving in consistent direction

$\Delta \text{Loss} + \text{MEMORY}$
RUNNING AVERAGE = MEMORY

ROOT MEAN SQUARE PROPAGATION (RMSProp)

Adjust learning rate for each parameter

Boost
Smaller updates

ADAM
Adaptive Moment Estimation

Adaptive learning rate

YouTube [youtube.com/@donatocapitella](https://www.youtube.com/@donatocapitella)
<https://llm-chronicles.com/>

