→ Optimization
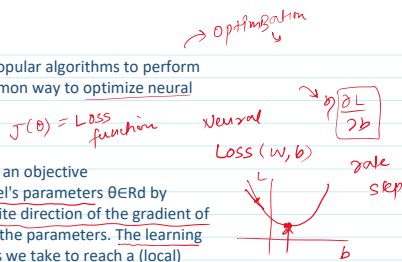
Gradient descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks.

$J(\theta) =$ Loss function    Neural
Loss $(w, b)$    $\boxed{\dfrac{\partial L}{\partial b}}$    rate step

Gradient descent is a way to minimize an objective function J(θ) parameterized by a model's parameters θ∈Rd by updating the parameters in the opposite direction of the gradient of the objective function ∇θJ(θ) w.r.t. to the parameters. The learning rate η determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley.

## Back propagation Algorithm

epochs = 5    ↙ STO    I point    → gradient descent ✓

for i in range(epochs):

    for j in range(X.shape[0]):

        → Select 1 row (random)

        → Predict (using Forward prop)

        → Calculate loss (using Loss function → mse)

        → Update weights and bias using GD

$$W_n = W_0 - \eta \frac{\partial L}{\partial w}$$

    → Calculate avg loss for the epoch
      L avg

{ Batch / Stochastic { mini batch }

There are three variants of gradient descent, which differ in how much data we use to compute the gradient of the objective function. Depending on the amount of data, we make a trade-off between the accuracy of the parameter update and the time it takes to perform an update.

$\dfrac{\partial L}{\partial w}$ } derivative
↓
3 types
accuracy → time → tradeoff.

→ Batch GD (Vanilla GD)

```
for i in range(nb_epochs):
    params_grad = evaluate_gradient(loss_function, data, params)
    params = params - learning_rate * params_grad
```

entire dataset ↓
update

epochs = # of updates.

Stochastic GD ✓
↓

```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for example in data:
        params_grad = evaluate_gradient(loss_function, example, params)
        params = params - learning_rate * params_grad
```

epoch → 10    (50 rows)    50 rows
for i in range(10):    → shuffle
    for i in range(X.shape[0])    $10 \times 50 = 500$ update

    ↳ 1 random point
    ↳ y_hat → forward
    ↳ loss
    ↳ w, b update → $W_n = W_0 - \eta \frac{\partial L}{\partial w}$

avg loss print → for the epoch

frequency of weight update higher

Batch GD → 1 loop    rows

epoch = ⑩    , epoch

for i in range(10):    , loop    → total → 10 times
    $y\_hat = np.dot(X, w) + b$    w, b update

50 values    dot → smart replacement → loops
Y = 50 values    ↳ vectorization ↳ faster loop

### Right column

50 points
1 epoch
    ↳ 50 times
      w, b update

Batch GD
epoch = 5    → current weight
    ↳ 50 points → predict ↓
      ↑
    dot product

$y\_hat = np.dot(X, w) + b$
    ↳ 50 predict

1 single    $Y = 50$ actual ocw.
(w, b)    y    y_hat    J loss

$$\int - \sum_{i=1}^{50}$$

$50 \times 10 = 500$    5 epochs
    ↳ 5 times
↓
w, b update

bias eliminate

50 values        $\underline{dot} \rightarrow$ smart replacement $\rightarrow$ loops
            $Y = 50$ values                                    $\downarrow$ vectorization $\curvearrowright$ faster loop
                $Y-hat, Y \rightarrow$ loss                                           $\uparrow$
                $\curvearrowleft$  W,b update    $W_n = W_0 - \eta \dfrac{\partial L}{\partial w}$    [Optimized]
        $\rightarrow$ loss

Which is faster ( given same no. of epochs)
    $\underset{10\eta}{\nearrow \uparrow}$                 $\overset{D}{\downarrow}$
            $320 \;\textcircled{$\eta$}\;\; 1$
                $\uparrow$
Which is the faster to converge (given same # epochs)
        $\nearrow$
    $\zeta$

Loss function $\rightarrow$ 2 Algo



$\rightarrow$ loss -3D
            Complex

$\;\swarrow$ 1D


Stochastic    Gradient

Stochastic $\nearrow$ 1
    $\hookrightarrow$ random $\rightarrow$ point $\downarrow$
                        updates

Batch $\downarrow$  $\rightarrow$ dataset
    $\hookrightarrow$ all point $\curvearrowright$ ↓update

Spikey SGD useful
    $\downarrow$         $\overline{\phantom{xx}}$
    Yes    No

SGD $\downarrow$
    help the algo
    to move out
    of local
        minima

    SGD
    BGD

local
minima

Exact
    solution
    $\hookrightarrow$
    Approximate
        diff
    $\rightarrow$ faster $\times$ loops
    $\searrow$ big dataset

Vectorization

    $\;\llcorner\;$ np. dot ( $\underline{X}$ , W) + b        $\overset{RAM}{\curvearrowright}$
            $\uparrow$
        $\;\zeta\;$ $x \rightarrow$ dataset of 10 crore
                            $\uparrow$

Mini Batch Gradient Descent $\nearrow$ ✓        $\boxed{BGD} \longleftrightarrow \boxed{SGD}$   Best of both
    $\uparrow$                                                worlds
                SGD $\rightarrow$ BGD                    $\downarrow$
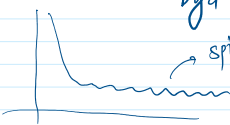
```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for batch in get_batches(data, batch_size=50):
        params_grad = evaluate_gradient(loss_function, batch, params)
        params = params - learning_rate * params_grad
```

                                                320 rows $\downarrow$
                                            In every epoch $\nwarrow$  batches
                                                10 batches         $\textcircled{32}$

$\dfrac{n}{x}$ = # of batches $\downarrow$                    $\downarrow$ 10 update
        $\zeta$  Keras  batch-size = $\underline{x}$
    $\downarrow$
    # updates/epoch                      for i in epochs $\rightarrow$
            $\downarrow$                            for j in num of batch
        bgd > mbgd > sgd                            1 batch
            con                          npdot    $\;\llcorner\;$ Y-pred (vector
        bgd < mbgd < sgd                          $\;\llcorner\;$ loss
MBGD                                              $\;\llcorner\;$ update
    $\downarrow$ spikey SGD $\searrow$
          Vectorization $\rightarrow$ $\zeta$ smaller $\zeta$
                                        batch

→ Why batch_size is provided in multiple of (2)?

2, 4, 8, 32, 64, 128, 256

10, 15

RAM effective → binary RAM → koras Optimization

faster


→ What if batch_size doesn't divide # rows properly

e.g    # of rows   $n = 400$
        batch_size = 150

       # of batch = $\dfrac{400}{150}$ = 2.66

   3    150, 150,  left 100
         ↑      ↑       ↑
      1 batch  2 batch  3rd batch