

Optimization:

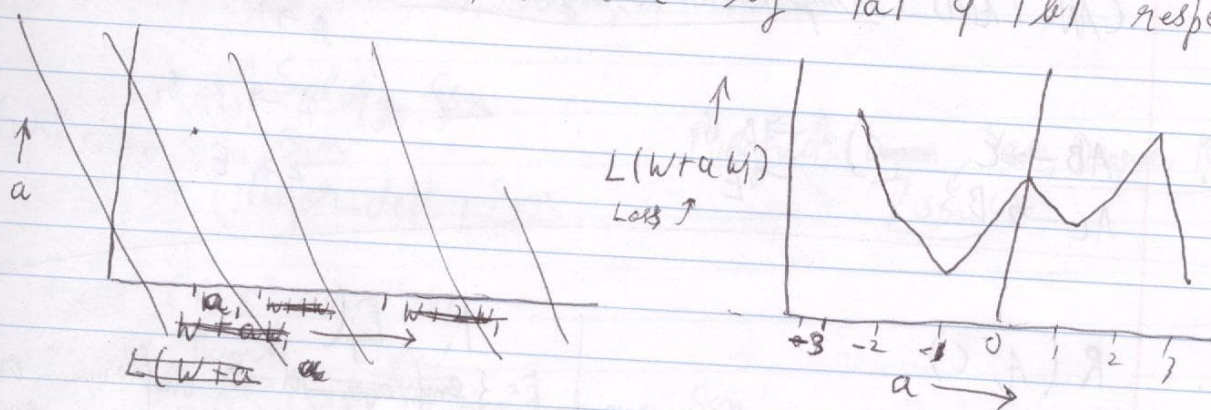
★ This is where we calculate W !

• Visualising loss function: Although this seems confusing, the plot of a loss function is just:

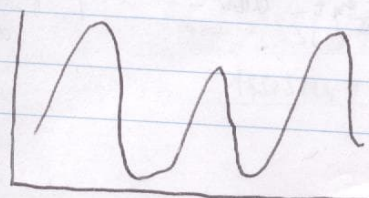
- i) $W + aW_1$ (1-D line)
- ii) $W + aW_1 + bW_2$ (2-D, colorful heatmap-like plot)

W is the initial weight. It can be given ~~represented~~ by a single point in space (just have it correspond to some point in space through a bijective function.)

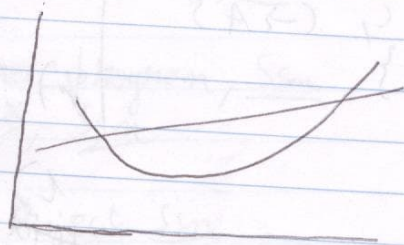
W_1 & W_2 are matrices of same dimensions as W with values that we add $(+a, +b)$ or subtract $(-a, -b)$ from W , scaled by $|a|$ & $|b|$ respectively.



Convex Function (Side Note)



Non-Convex



A line segment is always above or on the graph

→ Has global minimum only 1

★ Optimization Strategies:

1. Random Search : Randomly search for W

repeat $1..n$

$W = \text{randomWeights}()$

$\text{loss} = \text{Loss}(W)$

If $\text{loss} < \text{bestLoss}$:

$\text{bestLoss} = \text{loss}$

$\text{best } W = W$

2. Random Local Search : Go in Random direction from W ,
 $W = \text{randomWeights}()$ update if W_{new} is lower loss.

repeat $1..n$

$W_1 = \text{randomWeights()} * \text{stepSize}$ (stepSize = 0.0001)

$W_{\text{new}} = W + W_1$

$\text{loss} = \text{Loss}(W_{\text{new}})$

If $\text{loss} < \text{bestLoss}$

$\text{bestLoss} = \text{loss}$

$W = W_{\text{new}}$

3. Gradient Descent : Mathematically compute best direction along which W should be changed \rightarrow Find direction of steepest Descent of Loss function.

Time for some calculus, kids!

- ★ Gradient: Vector of partial derivatives of a vector of numbers.
 • Slopes in each of the dimensions.

$$\frac{d}{dx} f(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

★ Calculating gradient:

1. For all values in a vector, add a small value h , see if change is negative or positive. (Numerical / Analytical Method)

gradient(x, f):
 original = f(x) // f is our loss function.
 grad = []
 h = 0.00001 // h has to be very, very small. (h → 0)

```
for i in 1...len(x):
    old_value = x[i]
    f_xh = f
    x[i] = x[i] + h
    f_xh = f(x) // Updated vector with x[i] + h
    grad[i] = (f_xh - original) / h
    x[i] = old_value // Revert
return grad
```

(it is W, weights, we update.)

- ★ Note: x is not input data here, just some variable. Usually, ↑
 • Note how each time we update only one value in vector x by h , and then get $f(x)$ with this updated vector. This is equivalent to partial derivative of that value - ~~update~~ slope for that value, loss function along in that dimension.

★ How to Use:

$W = \text{random_weights}()$
 $df = \text{gradient}(W, \text{loss_func})$ // From prev. page

$W_{\text{new}} = W - \text{StepSize} * df$ // StepSize is variable
// df tells direction in which loss func increases, so we move in opposite direction

★ StepSize - also known as "Learning Rate"

Hyperparameter we tune.

Gradient tells us direction of steepest increase, but not how far we must step in that direction.

2. Calculus (Analytical Method)

$$\frac{d}{dw_j} L_i$$

i = data point ~~for~~ we are calculating loss for.
 y_i = true class

w_{y_i} = Weights for true class.

Example: SVM Hinge Loss Gradient:

$$\begin{aligned} \nabla_{w_{y_i}} L_i &= \frac{d}{dw_j} \left[\sum_{j \neq y_i} \max(0, w_j^T x_i - w_{y_i}^T x_i + \Delta) \right] \\ &= - \left(\sum_{j \neq y_i} \mathbb{I}(w_j^T x_i - w_{y_i}^T x_i + \Delta > 0) \right) x_i \end{aligned}$$

Mini-batch Gradient Descent: Compute gradient over batches / samples of training data, in large-scale applications.

→ May also be referred to as "Stochastic Gradient Descent"