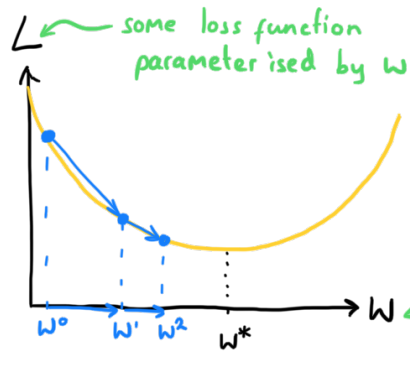


Gradient descent



At the heart of it:

Iteratively move down the slope

follow this heuristic,
the gradient

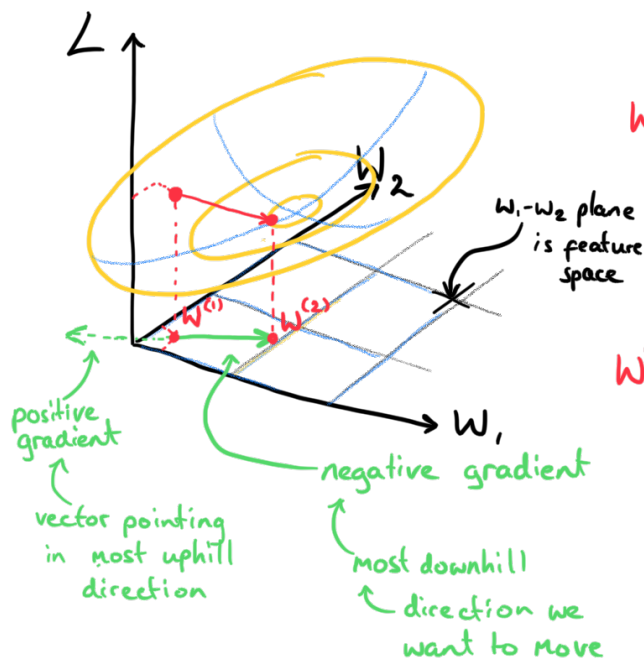
We update the model parameters as follows

$$W \leftarrow W - \alpha \frac{\partial L}{\partial W}$$

gradient descent
update rule

"set the new param equal to the old param value shifted in the direction of the negative gradient"

Most models of interest have multiple params, so let's get an idea of how this looks for a higher dimensional feature space



$$W^{(1)} = \begin{bmatrix} W_1^{(1)} \\ W_2^{(1)} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} W_1^{(2)} \\ W_2^{(2)} \end{bmatrix} = W^{(1)} - \alpha \frac{\partial L}{\partial W} \Big|_{W=W^{(1)}}$$

evaluated where $W = W^{(1)}$

$$= \begin{bmatrix} W_1^{(1)} \\ W_2^{(1)} \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial L}{\partial W_1} \\ \frac{\partial L}{\partial W_2} \end{bmatrix} = \begin{bmatrix} W_1^{(1)} - \alpha \frac{\partial L}{\partial W_1} \\ W_2^{(1)} - \alpha \frac{\partial L}{\partial W_2} \end{bmatrix}$$

Notice that the gradient is a vector $\rightarrow \frac{\partial L}{\partial W}$

↳ it's always the same shape as the model parameters

↳ WHY? \rightarrow because we need an update for each parameter

As long as we can compute $\frac{\partial L}{\partial W}$ then we can do gradient descent to descend the loss function in parameter space & optimise the params

So how does that work in the case of linear regression?

Our model: $\hat{y} = XW + b$

$$\begin{aligned}\text{Our loss: } L &= \frac{1}{n} \sum_{i=1}^n (\hat{y} - y)^2 \\ &= \frac{1}{n} \sum (XW + b - y)^2 \\ &= \mathbb{E} (XW + b - y)^2\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial W} &= \frac{\partial}{\partial W} \left(\frac{1}{n} \sum (XW + b - y)^2 \right) \\ &= \frac{1}{n} \sum \frac{\partial}{\partial W} \left((XW + b - y)^2 \right)\end{aligned}$$

gradient of mean
=
mean of gradient

$$\begin{aligned}&= 2 (XW + b - y) X \\ &= \frac{1}{n} \sum 2 \underbrace{(x^{(i)} W + b - y^{(i)})}_{\in \mathbb{R}} \underbrace{x^{(i)}}_{\in \mathbb{R}^n}\end{aligned}$$

$\in \mathbb{R}^n$
vector of gradients

average over each example

Can you figure out how to compute this in vector/matrix form?

Same procedure for bias:

$$\begin{aligned}\frac{\partial L}{\partial b} &= \frac{\partial}{\partial b} \left(\frac{1}{n} \sum (xw + b - y)^2 \right) \\ &= \frac{1}{n} \sum \underbrace{\frac{\partial}{\partial b} \left((xw + b - y)^2 \right)}_{= 2(xw + b - y)}\end{aligned}$$

gradient of mean
=
mean of gradient

$$= \frac{1}{n} \sum 2 \underbrace{(xw + b - y)}_{\in \mathbb{R}^n}$$

average over
the vector