

Avoid local minima: basis of attraction.

Why do these local minima exist in deep neural network training?

If training is convex, there is no local minima (includes linear regression, logistic regression)

Neural network training is not convex.

- observe that DNNs have weight space symmetry
 ↳ permute all hidden neurons on a layer
 ↳ we obtain the same model

- if we average all permutations of the parameters, we would get a degenerate neural network where all hidden neurons are identical

- if the cost function for training the neural network was convex, the average network (defined above) should perform better than the original network

because when f is convex

$$f(\lambda_1 x_1 + \dots + \lambda_N x_N) \leq \lambda_1 f(x_1) + \dots + \lambda_N f(x_N) \quad \text{for } \lambda_i \geq 0, \sum \lambda_i = 1$$

loss function degenerate neural network

$\lambda_i = \frac{1}{N}$
 $x_i =$ model parameters for one of the permuted architectures

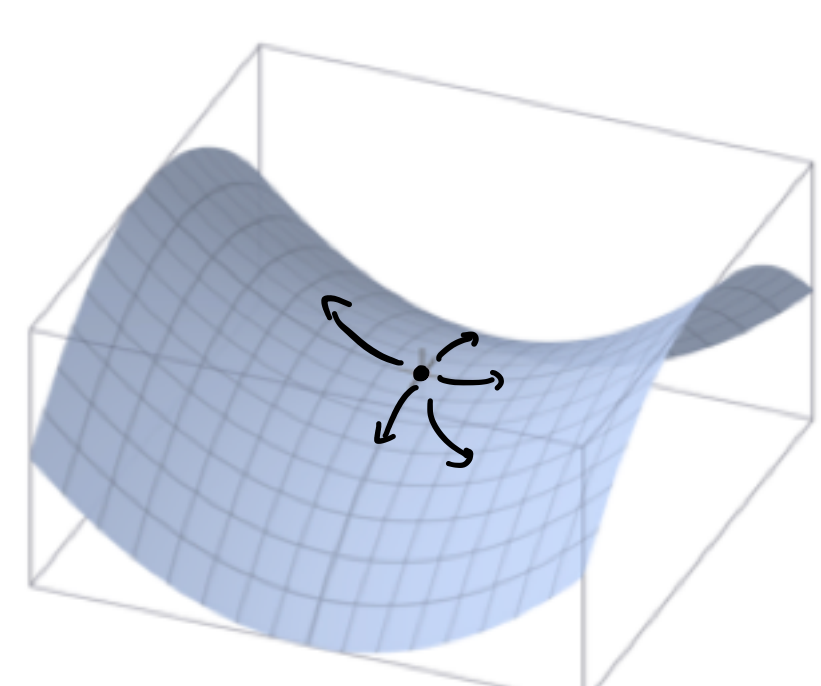
$f(x_i)$ is the loss of one of the permuted architectures

$$f(x_1) = f(x_2) = \dots = f(x_N)$$

$$\sum \lambda_i f(x_i) = f(x_i)$$

loss of our original model

Saddle points

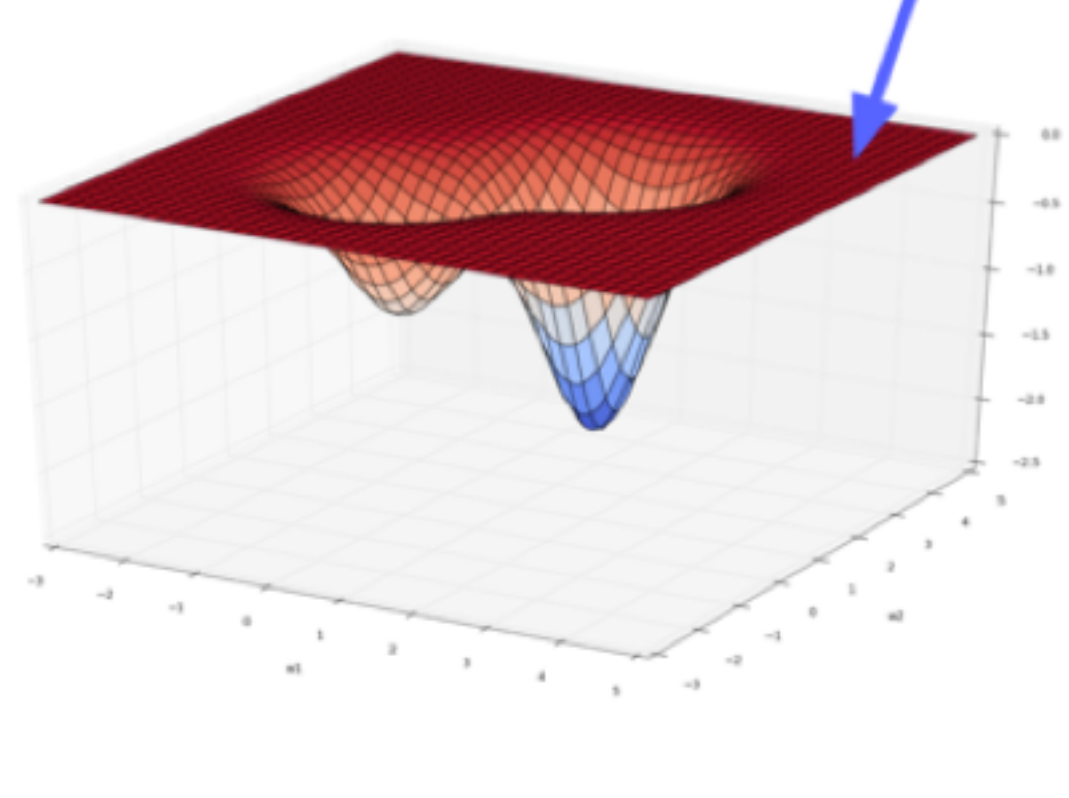


At a saddle point, we have $\frac{\partial L}{\partial w} = 0$ even though we are NOT at a minimum.

Exactly at the saddle point: gradient descent is stuck
 for a slight perturbation of the saddle point: gradient descent gets unstuck because we no longer have $\frac{\partial L}{\partial w} = 0$

When initializing the random weights of our model, we should avoid using weights that are initially set to 0 and instead we should prefer random small values.

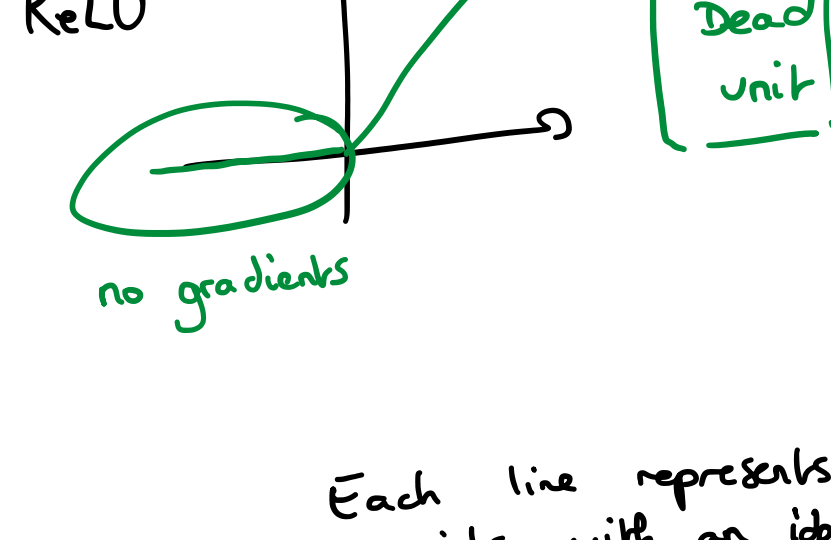
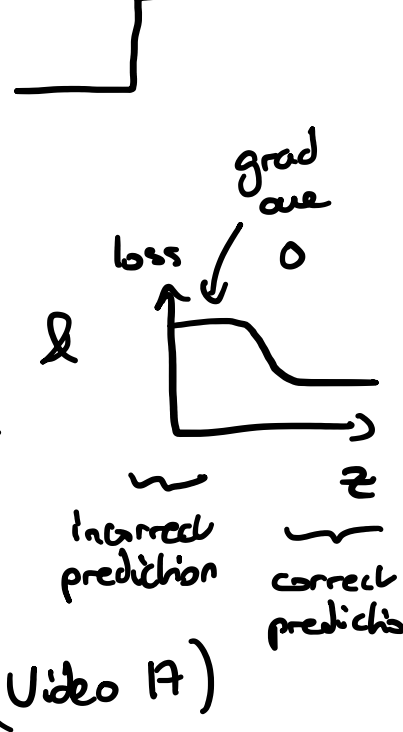
Plateaux



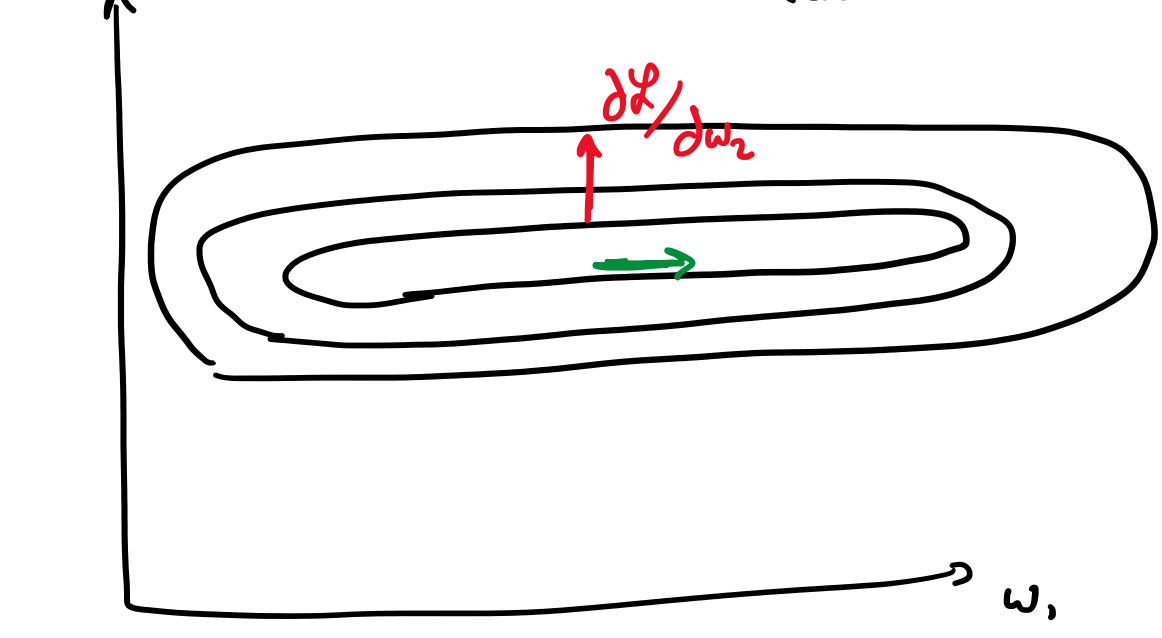
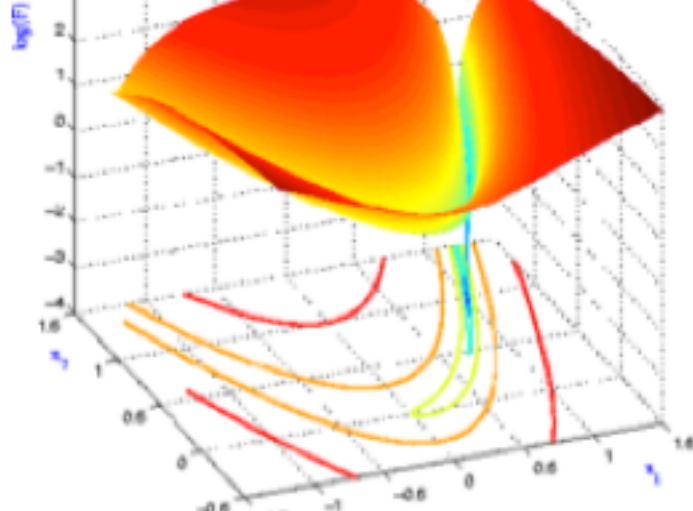
Flat region in the loss surface

Examples:

- 0-1 loss
- hard threshold activation
- logistic activations & least squares.
- saturated activation function (Video 11)



Ravines



Each line represents all points with an identical loss value

To avoid ravines: (A) make sure inputs are centered & have unit variance

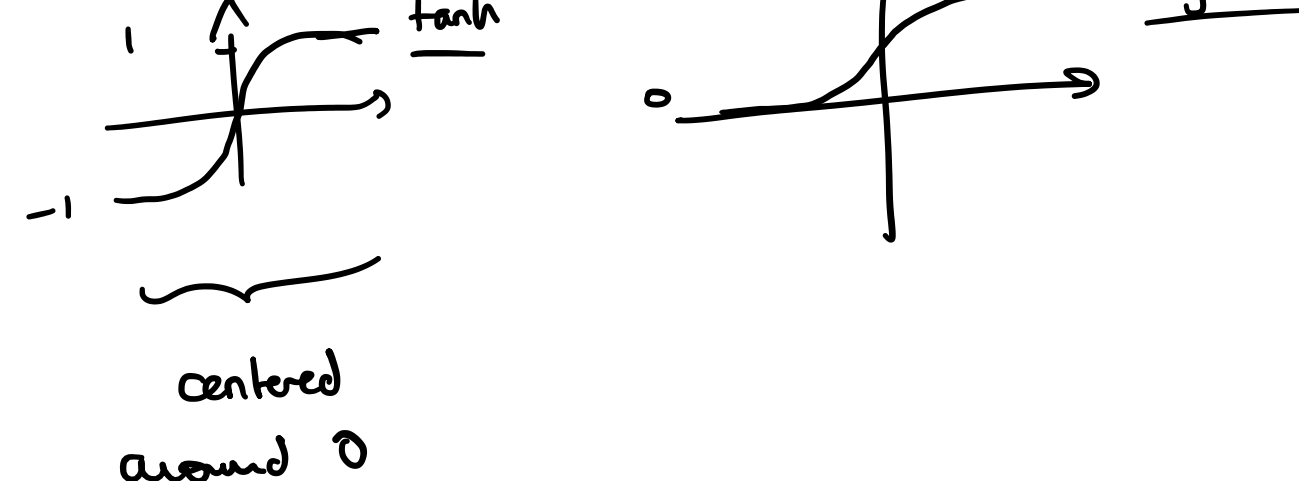
$$x_j \leftarrow \frac{x_j - \mu_j}{\sigma_j}$$

Especially important when features have \neq arbitrary units

(B) Hidden units/neurons should also have centered activations.

Harder to achieve:

→ use tanh activation instead of logistic activation



→ batch normalization: centers each hidden activation and can speed up training

(C) Momentum

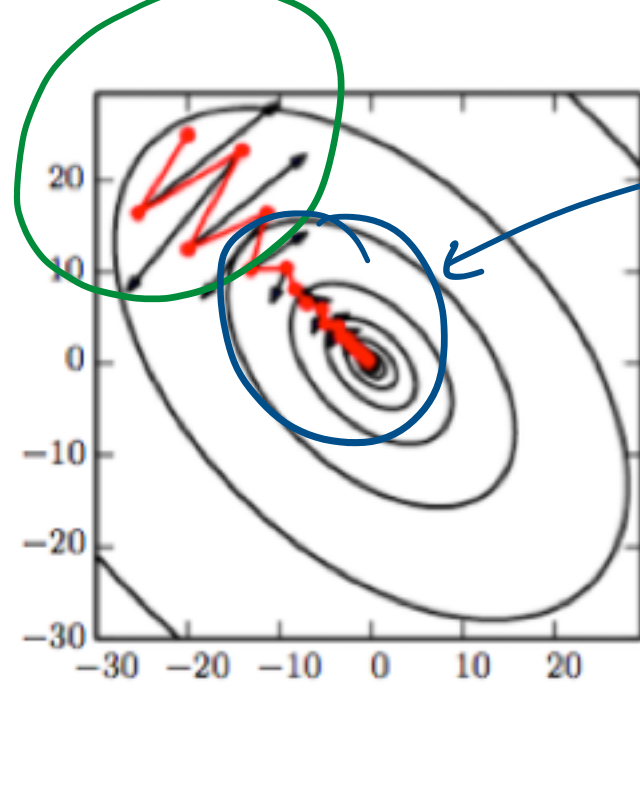
without momentum: $\vec{w} \leftarrow \vec{w} - \alpha \frac{\partial L}{\partial \vec{w}}$

with momentum: $\vec{p} \leftarrow \mu \vec{p} - \alpha \frac{\partial L}{\partial \vec{w}}$

$\vec{w} \leftarrow \vec{w} + \vec{p}$ (learning rate)

damping parameter (it should be slightly less than 1)

high curvature region gradients will cancel each other out
 ⇒ momentum will dampen the oscillations



low curvature region gradients all point in the same direction
 ⇒ momentum will accelerate the speed of learning

Second-order optimization

We need more information about the curvature of the loss function

We need to compute 2nd order derivatives

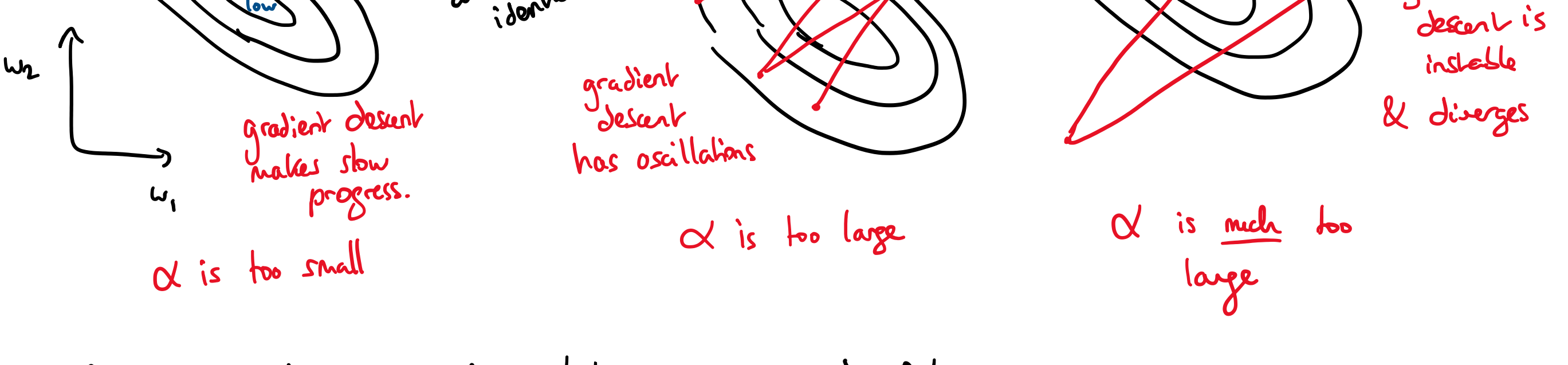
- ↳ complicated to compute
- ↳ difficult to scale them to DNNs (especially when training on large datasets)

One optimizer that uses a little bit of curvature information, yet is still practical for DNNs

Adam typically works better than gradient descent

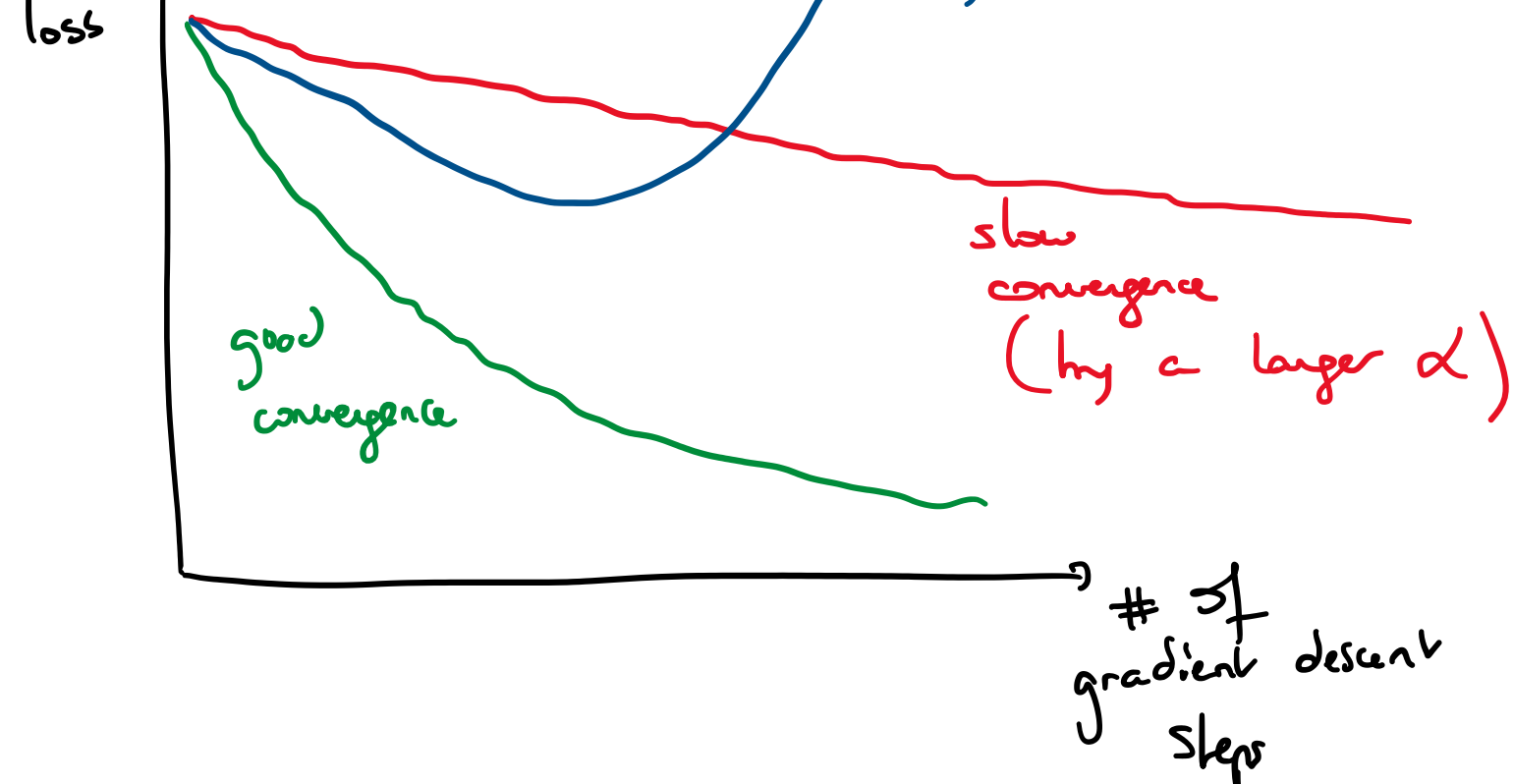
Learning rate

α is a hyperparameter that we need to tune using train/validation/test sets



Good values of α are often between 0.001 and 0.1

Training curves



⚠ training curves are useful to spot major problems with convergence, they do not guarantee convergence.