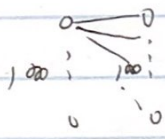


nodes in hidden layers }
hidden layers } examples of
activation fns } "hyperparameters"
↓
not trainable parameters θ .

• NN basically repeated matrix multiplications
+ vectorized activations
→ perfect for GPU acceleration!!

• Counting parameters

let's say $d=100$, $L=10$, $n_h=1000$



$$\dim W^{(n)} \sim (10^3)^2 \sim 10^6$$

• $\times 10$ hidden layers

$\sim 10^7$ parameters

→ # of NN parameters can easily get very large!

↳ can still be fit to data using powerful algorithms & hardware

→ 1st demo nb: constructing a simple MLP
in Keras & in pytorch

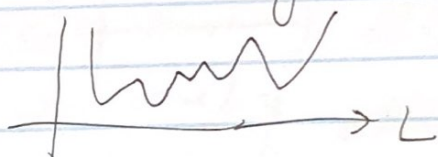
Training NNs

How to fit $f(x; \theta)$ w/ many (eg 1M) parameters?

Want to minimize $L(\theta)$ wrt θ

$$\frac{\partial L}{\partial \theta} \Big|_{\theta_*} = 0 \quad + \text{(ideally) global minimum}$$

L might have many local minima, or saddles



↓
don't want to get stuck

Many approaches...

— one idea: Newton method

$\theta = \theta_0$ initial guess

$$L(\theta) = L(\theta_0) + \delta\theta \cdot \frac{\partial L}{\partial \theta} \Big|_{\theta_0} + \frac{1}{2} \delta\theta \cdot \delta\theta \cdot \frac{\partial^2 L}{\partial \theta \partial \theta} \Big|_{\theta_0} + \dots$$

— Minimize 2nd order approx $\rightarrow \delta\theta$

— Repeat for $\theta + \delta\theta$

— Hopefully converge to ~~the~~ ^{correct} θ_* .

— Why is this infeasible!?

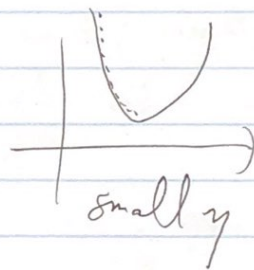
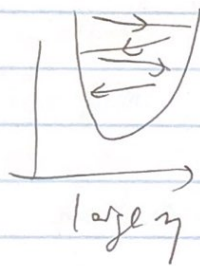
$\rightarrow \frac{\partial^2 L}{\partial \theta_i \partial \theta_j}$ too high dim & have to invert it!

• Instead, NNs are optimized w/ 1st order methods:

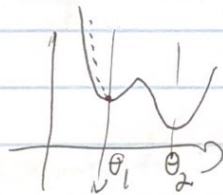
$$\delta\theta = \eta \frac{\partial L}{\partial \theta} \Big|_{\theta_0} \quad \text{"gradient descent"}$$

idea: move in direction where L decreases

η = "learning rate" — important hyperparameter!
often need to tune for best perf.



• Problem: GD can get stuck in local minima



$$\frac{\partial L}{\partial \theta} \Big|_{\theta_1} = 0 \Rightarrow \delta\theta = 0 \quad \text{no update}$$

get stuck at $\theta = \theta_1$

→ "Stochastic gradient descent"

noisy gradients are better than true gradients!

how to get noisy gradients?