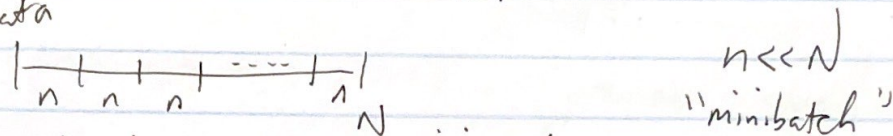


- Vanilla GD $L = \sum_{i=1}^N \mathcal{L}(f(x_i; \theta), x_i)$

- Stochastic GD $L = \sum_{i=1}^n \mathcal{L}(f(x_i; \theta), x_i)$

train data



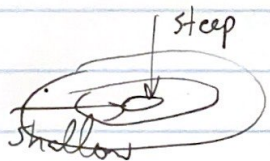
- divide up data into minibatches
- evaluate ~~gradient~~ loss & gradient on each
- ~~update~~ update parameters ~~at~~ minibatch by minibatch

minibatch size n is another key hyperparameter!
Some tradeoff btw n & η

- looping over full data once = "1 epoch"

Further improvements to SGD:

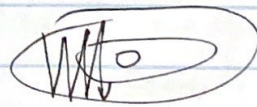
- "momentum" SGD w/ fixed η could mean



"ravine"

\rightarrow too large for large gradients
"small" "small"

• can lead to slow convergence
allow shallow direction



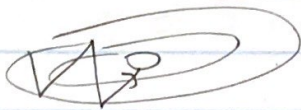
• idea: "momentum" like in physics \rightarrow acceleration

$$\delta\theta = \frac{\partial L}{\partial \theta} \left\{ \begin{array}{l} v_t = v_{t-1} - \eta \frac{\partial L}{\partial \theta} \\ \delta\theta = v_t \end{array} \right.$$

Push parameters in direction of previous gradient

- if gradients oscillate (bouncing around steep direction), momentum cancels out

- if gradients pointed along same dir $t-1, t$, then they add \rightarrow acceleration



- Adaptive learning rates

idea: large gradients \rightarrow slow down

small gradients \rightarrow speed up

+ diff learning rates \forall parameter

$$\delta\theta_i = \frac{\eta}{\sqrt{G_i}} \frac{\partial L}{\partial \theta_i} \quad G \sim \sum_{t \text{ in steps}} \left(\frac{\partial L}{\partial \theta_i} \right)^2$$

"Ada delta"

\downarrow
restrict G to finite
of time steps
or decaying avg

& "Adam" two popular examples

\downarrow
most common choice (~"Ada delta w/ momentum")

Backpropagation

A brief aside on how gradients of NNs are actually calculated
(motivated by explaining vanishing gradient problem)

NNs have a recursive structure
nested

(ignore L , as for simplicity)

$$f_{NN}(x) = A(w_L \cdot A(w_{L-1} \cdot A(\dots A(w_1 x) \dots)))$$

$$L = \sum \mathcal{L}(f_{NN}(x_i; w), x_i)$$

$$\frac{\partial L}{\partial w} \rightarrow \frac{\partial f_{NN}}{\partial w}$$

$$\left. \begin{array}{l} x_0 = x \\ x_1 = A(w_1 x_0) \\ x_2 = A(w_2 x_1) \\ \vdots \\ x_{L-1} = A(w_{L-1} x_{L-2}) \\ x_L = A(w_L x_{L-1}) \end{array} \right\} \begin{array}{l} \text{forward} \\ \text{pass} \end{array}$$

just chain rule!

$$1 \text{ mm} \quad \frac{\partial f_{NN}}{\partial w_L} = A'(w_L x_{L-1}) x_{L-1}$$

$$3 \text{ mm} \quad \frac{\partial f_{NN}}{\partial w_{L-1}} = A'(w_L x_{L-1}) w_L A'(w_{L-1} x_{L-2}) x_{L-2}$$

$$5 \text{ mm} \quad \frac{\partial f_{NN}}{\partial w_{L-2}} = A'(w_L x_{L-1}) w_L A'(w_{L-1} x_{L-2}) w_{L-1} A'(w_{L-2} x_{L-3}) x_{L-3}$$

To compute all gradients, need $\mathcal{O}(L^2)$ matrix multiplications
— expensive!

Backprop: reuse matrix multiplications $L \rightarrow L-1 \rightarrow L-2 \dots$

$$\text{Let } G_h = A'(w_L x_{L-1}) w_L A'(w_{L-1} x_{L-2}) \dots w_{h+1} A'(w_h x_{h-1})$$

$$\text{So } \frac{\partial f_{NN}}{\partial w_L} = G_L x_{L-1}$$

$$\frac{\partial f_{NN}}{\partial w_{L-1}} = G_{L-1} x_{L-2}$$

⋮

$$G_h = G_{h+1} w_{h+1} A'(w_h x_{h-1})$$

have it from forward pass

each step only needs 2 matrix multiplications $\rightarrow O(L)$ needed

\rightarrow Backprop was key step (seems obvious) in making ^{deep} NNs trainable!

\rightarrow also see gradients are products of many $A'(\dots)$ w.r.t pos. of derivatives!

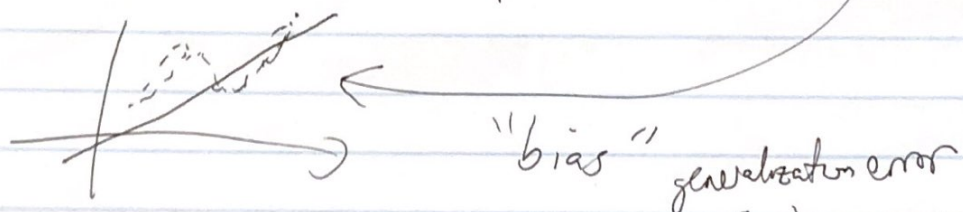
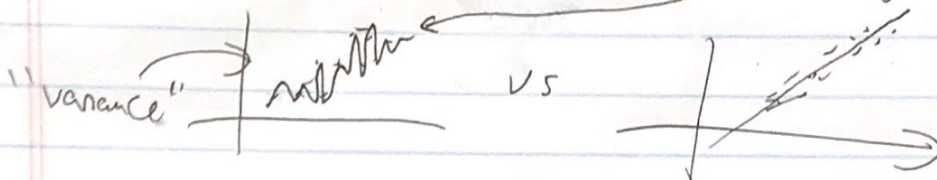
\rightarrow gradients of earlier layer weights can ~~be~~ be very suppressed or enhanced "vanishing/exploding gradient problem"

\rightarrow especially bad for sigmoid ameliorated by ReLU

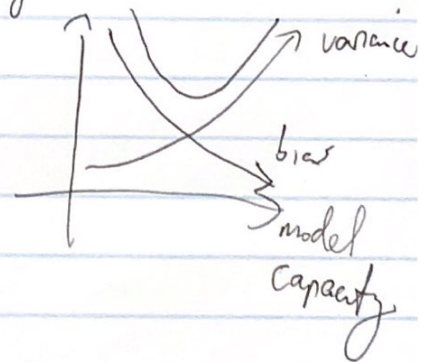
concept to introduce
One last thing to discuss before we can do
a demo of NN training:

overfitting & train/val/test split

How do we detect & prevent overfitting? Or underfitting?

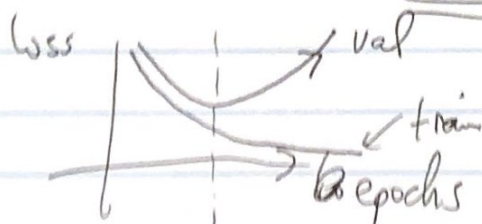


"bias-variance tradeoff"



→ want model to generalize well
to unseen data following same dist'n
as training data

→ set aside validation set



stop training! "early stopping" or "model selection"

Then possibly biased on val set

→ final metrics on another held out dataset
"test set"