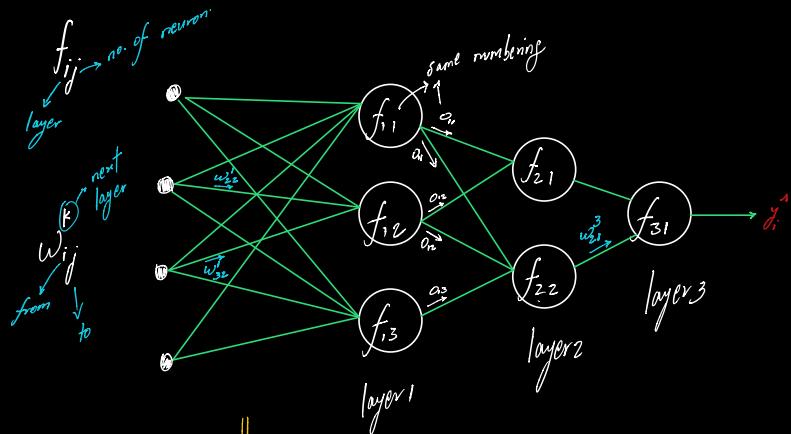


Fully connected NN

$$D = \{x_i, y_i\} ; \quad x_i \in \mathbb{R}^4 \\ y_i \in \mathbb{R}$$



\Downarrow

$$\left(\begin{array}{ccc} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{array} \right)$$

$\xrightarrow{\text{next layer}}$

$\xrightarrow{\text{previous layer}}$

The math becomes easy
once the notation is clear.

Quick Note

Linear Regression

$$\min_{w_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{regularization}$$

NN

$$D = \{x_i, y_i\}_{i=1}^n$$

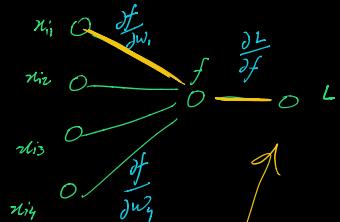
$$x_i \in \mathbb{R}^4$$

$$y_i \in \mathbb{R}$$

① \Rightarrow loss function.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}$$

for all data



② \Rightarrow optimization.

$$\min_{w_i} \sum_{i=1}^n (y_i - f(w^T x_i))^2 + \text{regularization}$$

$$w^* = \arg \min_w \sum_{i=1}^n (y_i - f(w^T x_i))^2 + \text{regularization}$$

$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial w_4} \right]^T$$

gradient of
L w.r.t. w.

if we know

follow the connecting
path and keep differentiating
until input is reached

$$\begin{aligned} \text{chain rule} \\ \left(\frac{\partial L}{\partial f} \right) \left(\frac{\partial f}{\partial w_i} \right) &= \frac{\partial f(w^T x_i)}{\partial w_i} \\ &= x_i \end{aligned}$$

this, we can then
change w accordingly
to minimize loss.

$$\underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\partial f} + \text{(reg)} \xrightarrow{\text{ignore for simplicity}}$$

$$= -2 \sum_{i=1}^n (y_i - \hat{y}_i) \frac{\partial \hat{y}_i}{\partial f}$$

$$= -2 \sum_{i=1}^n (y_i - f(w^T x_i)) \cdot (-1) \underbrace{\frac{\partial f}{\partial f}}_{\rightarrow 1}$$

$$\therefore \frac{\partial L_i}{\partial w_j} = -2 x_j (y_i - \hat{y}_i)$$

$$\frac{\partial L}{\partial w_j} = -2 \sum_{i=1}^n x_j (y_i - \hat{y}_i)$$

Similarly we can get stuff for $\frac{\partial L}{\partial w_k}, \dots$

$\therefore \nabla_w L \rightarrow$ ^{# Gradient} how much should we change to move the value of L

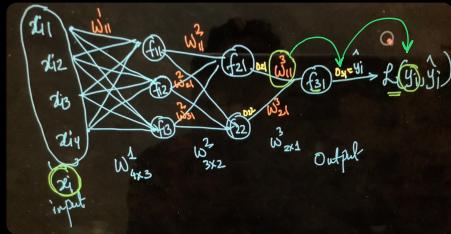
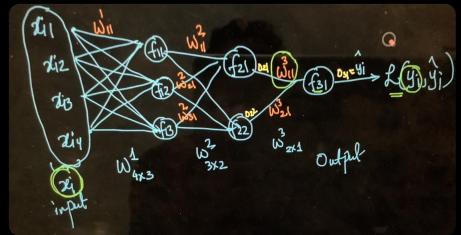
$$\therefore w_{\text{new}} = w_{\text{old}} - \eta \left[\nabla_w L \right]_{w_{\text{old}}}$$

Training regularization = $\sum_{i,j} (w_{ij}^k)^2$ norm^{L2}

Gradient Descent

① Initialize w_{ij}^k [lots of techniques out there]

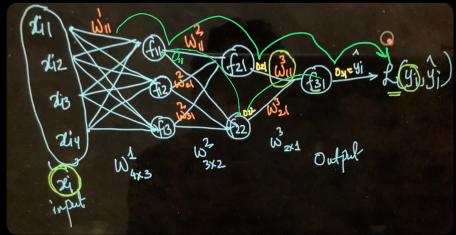
② $w_{ij}^k_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}^k}$



$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}$$

similarly

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{21}^2}$$



→ multiple paths
so adding them up

$$\frac{\partial L}{\partial w'_{ii}} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w'_{ii}}$$

$$= \frac{\partial L}{\partial O_{31}} \cdot \left[\frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w'_{ii}} + \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial w'_{ii}} \right]$$

$$\frac{\partial L}{\partial w'_{ii}} = \frac{\partial L}{\partial O_{31}} \cdot \left[\frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{ii}} + \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{ii}} \right]$$

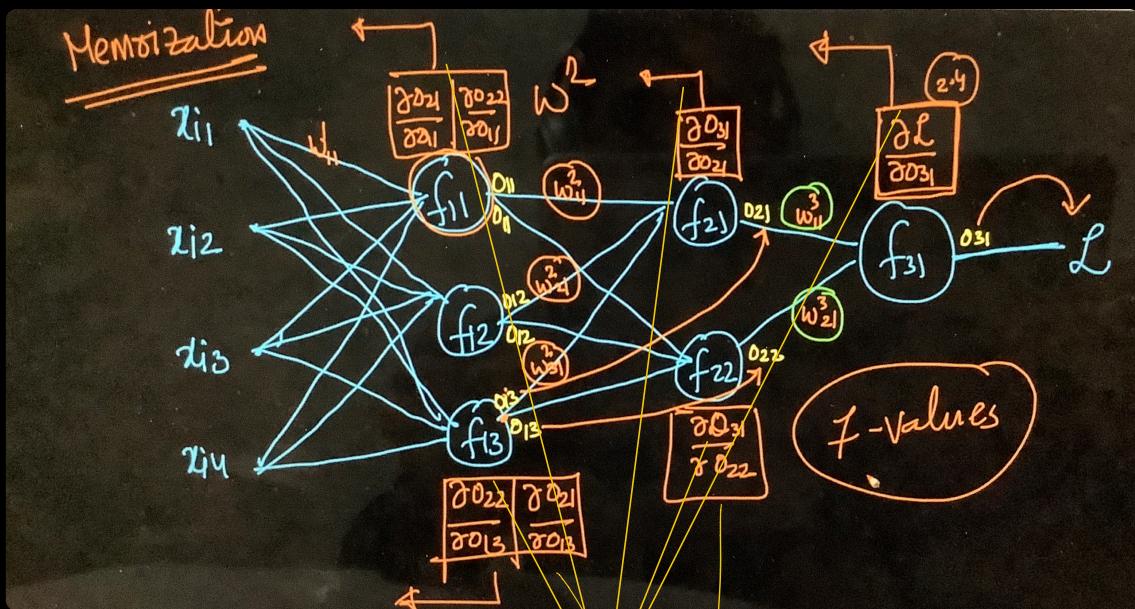
→ This involves computing a lot of gradients.

There's a surf that one gradient will be computed many times. This is why

we use memoization → store the gradients in a hash table and use those instead of calculating [→ only works for the current iteration.]

Dynamic Programming

then again again
for entire ∇L



Store them like this.
It helps a ton with
speed. Uses only
slightly more memory

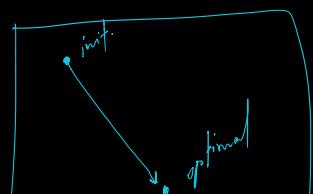
Chain rule + memorization

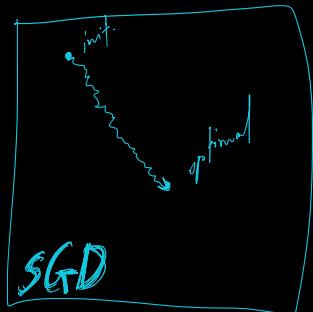
↓
boom!

↓
Backpropagation!

- ① get gradient, update weight matrix
 - ② update matrix closest to output first, and then the one at the back
 - ③ do this for every $D = \{x_1, x_2, x_3, \dots\}$
 - ④ stop when $(w_{ij}^k)_{\text{new}} \approx (w_{ij}^k)_{\text{old}}$
↳ convergence.
- * Backprop only works if the activation functions are differentiable. Otherwise the whole idea fails.
- Computing the derivatives is the most time consuming step. So the activation functions need to be easily differentiable.
- Instead of doing backprop for every point, Σ loss over a batch of points and then do backprop.
- kinda the same thing but a ton faster.

Here, the term "stochastic" comes from the fact that the gradient based on a single training sample is a "stochastic approximation" of the "true" cost gradient. Due to its stochastic nature, the path towards the global cost minimum is not "direct" as in Gradient Descent, but may go "zig-zag" if we are visualizing the cost surface in a 2D space. However, it has been shown that Stochastic Gradient Descent almost surely converges to the global cost minimum if the cost function is convex (or pseudo-convex)[1].



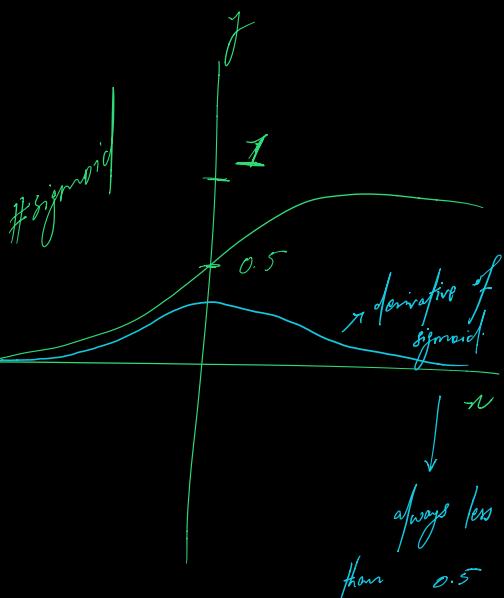


GD

Activation functions

- Differentiable
- Easy to differentiate

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

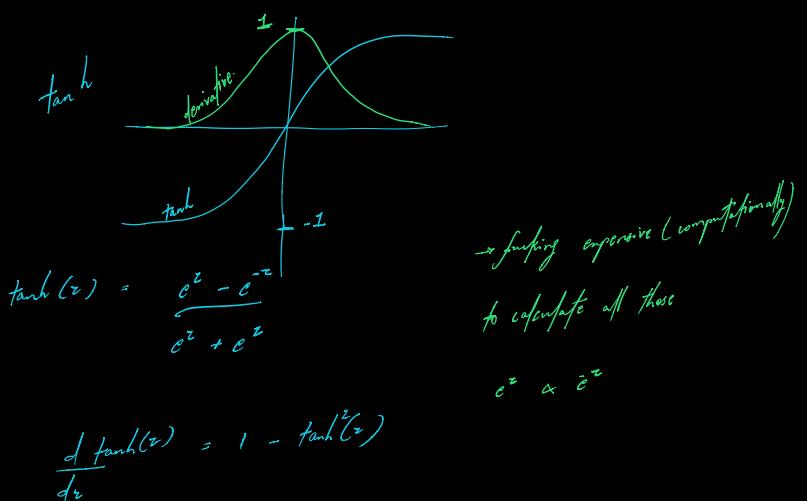


$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

↓
vanishing
gradient!!

Derivative is expressed in terms

of the f^k itself. So it can
be used for forward as well as backprop.



- ① Range is -1 to 1 compared to 0 to 1 for sigmoid
- ② Derivative still less than $1 \Rightarrow$ vanishing gradient

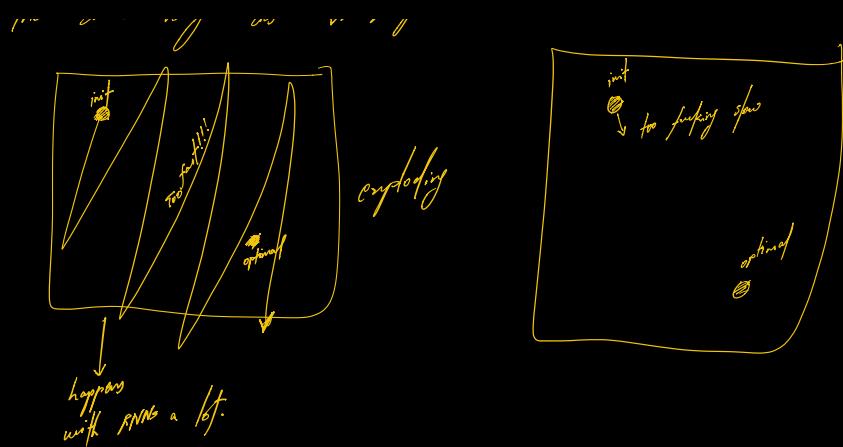
* Vanishing gradient:

Derivatives $\leftarrow 1$

\therefore chain rule makes them $\ll 1$

$\therefore w_{new} \approx w_{old}$
 kinda remains the same \rightarrow no receptive gradient with more hidden layers. \rightarrow because multiple weights \downarrow
 \therefore These activations are bad! \downarrow gradient cannot propagate down the network.

There's also something called exploding gradient, which also fucks up our weights the same way as vanishing



→ more parameters
allow to fit more complex stuff

we need to have the right amount
of parameters to prevent over/underfit

→ with RNNs it's always overfitting.

↓
we can use dropout,
Batch Normalization,
 L_1/L_2 regularization, etc. or all of them
together.

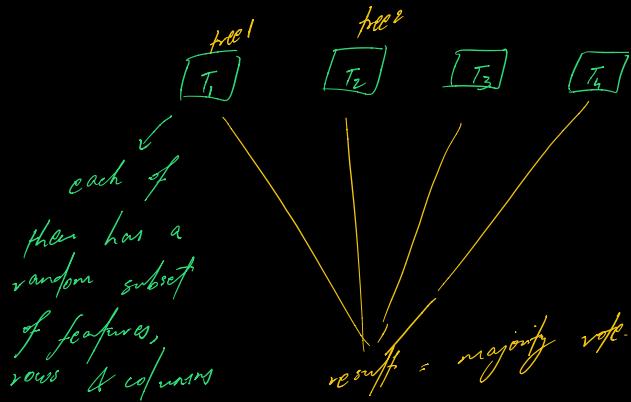
L_1 regularization creates sparsity in weights

Most deep learning is empirical.
 ↳ experiments first → the engineer's way
 ↳ theory will be built later on.

↓ ,

too much
fine, too slow

\Rightarrow Random Forest.

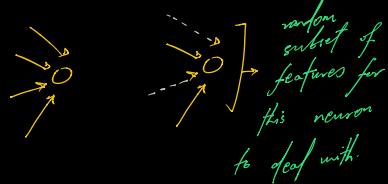


\rightarrow Each tree may overfit but a collection is used to regularize the DT

Dropout is like applying the same philosophy to NNs.

$p = 0.2$ *chosen randomly for each iteration.*
That means 20% of the neurons are inactive for each layer.

\rightarrow This is like trying out a random subset of features



\rightarrow There is a mathematical proof of why this thing works. \rightarrow not worth the time investment. But we care about the results for now.

\rightarrow Dropout is only for training time, not during question.