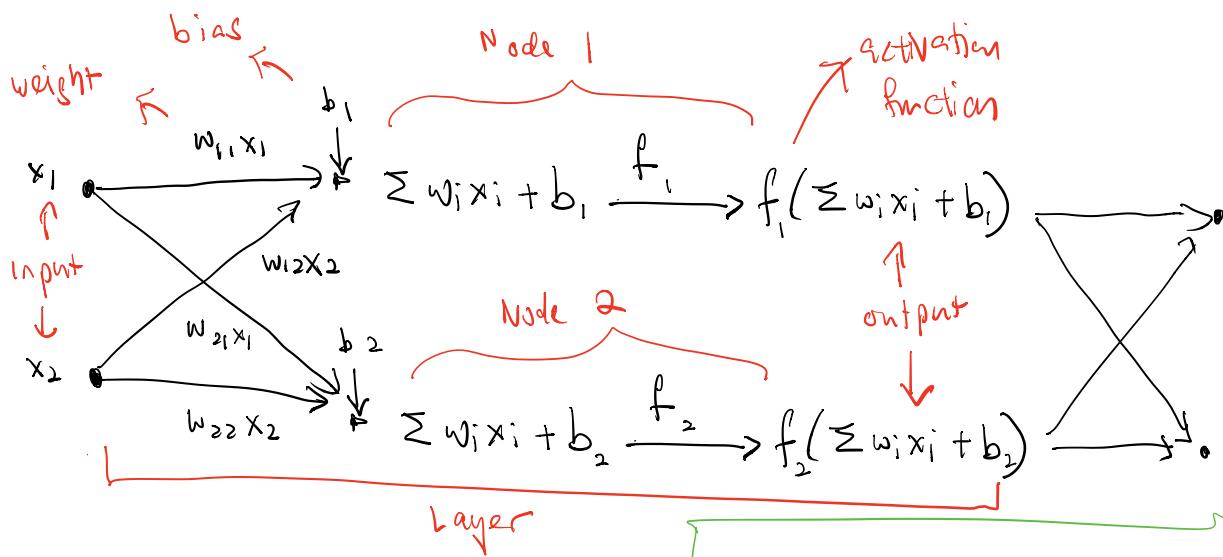
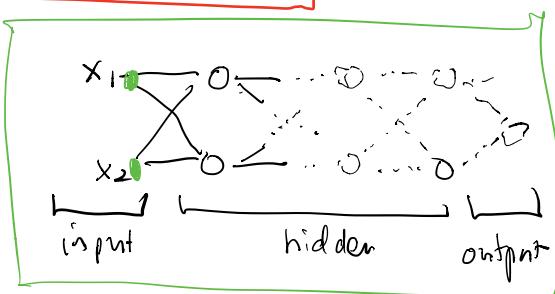


Anatomy



Each layer implements this function

$$[y] = f([w][x] + [b])$$



affine transformation (linear + translation)

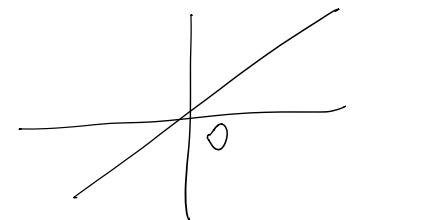
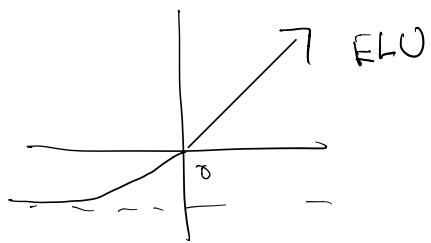
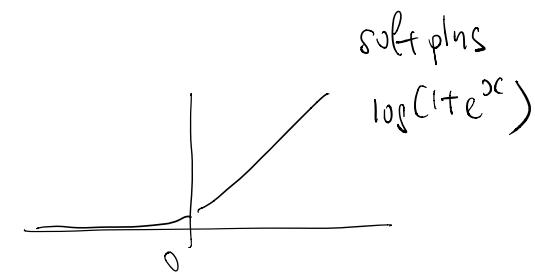
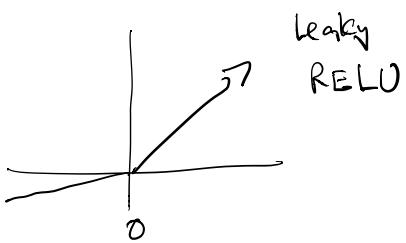
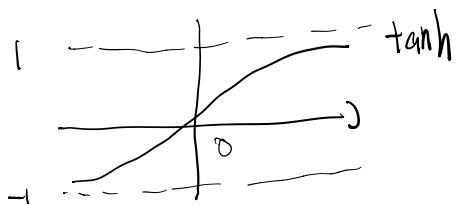
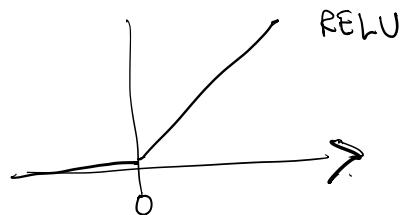
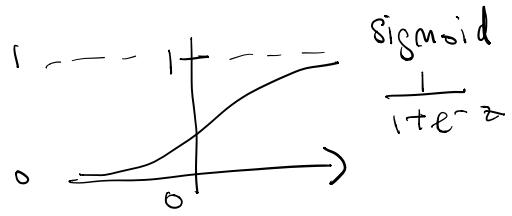
Affine can be expressed as an augmented matrix

$$\begin{bmatrix} ax_1 + bx_2 + c \\ dx_1 + ex_2 + f \\ 1 \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

If activation function is linear \rightarrow Matrix multiplication

$$\begin{aligned} y &= A_5 A_4 A_3 A_2 A_1 x \\ &= Bx \end{aligned}$$

Activation functions



Training

$\theta \rightarrow$ Collection of weights w_{ij}^k
3 lines b_i^k

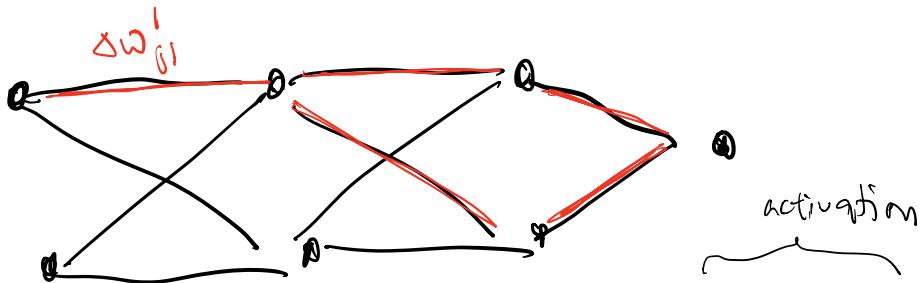
Loss with batch data 3 params

$$\theta_{k+1} = \theta_k - \alpha \frac{\partial L(x, \theta_k)}{\partial \theta}$$

This is just gradient descent

How is gradient calculated?

Sketch of main idea of backpropagation



$$\Delta L \approx \frac{\partial L}{\partial w_{jk}^l} \Delta w_{jk}^l$$

↓
 $x_m^l = f(wx + b)$
 layer = L
 node = m

Fw one path $\Delta L \approx \frac{\partial L}{\partial d_m^l} \frac{\partial d_m^l}{\partial x_n^{l-1}} \frac{\partial x_n^{l-1}}{\partial x_p^{l-2}} \dots \frac{\partial x_j^l}{\partial x_j^l} \frac{\partial x_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$

chain rule

Fw all paths

$$\Delta L \approx \sum \frac{\partial L}{\partial d_m^l} \frac{\partial d_m^l}{\partial x_n^{l-1}} \frac{\partial x_n^{l-1}}{\partial x_p^{l-2}} \dots \frac{\partial x_j^l}{\partial x_j^l} \frac{\partial x_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$

Expression for gradient

$$\frac{\partial L}{\partial w_{jk}^l} \approx \sum \frac{\partial L}{\partial d_m^l} \frac{\partial d_m^l}{\partial x_n^{l-1}} \frac{\partial x_n^{l-1}}{\partial x_p^{l-2}} \dots \frac{\partial x_j^l}{\partial x_j^l} \frac{\partial x_j^l}{\partial w_{jk}^l}$$

Automatic differentiation

Tensorflow uses reverse mode auto diff

Auto diff \rightarrow rules for differentiating simple expressions

by substitution e.g. $\frac{d}{dx} x^n \rightarrow nx^{n-1}$

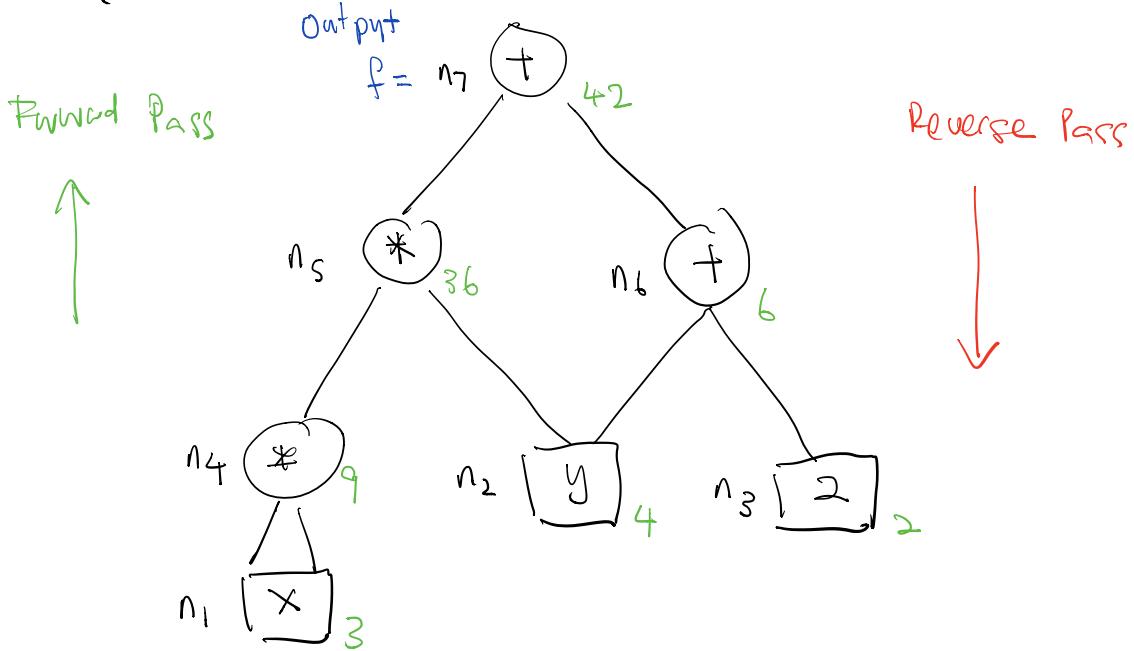
Example

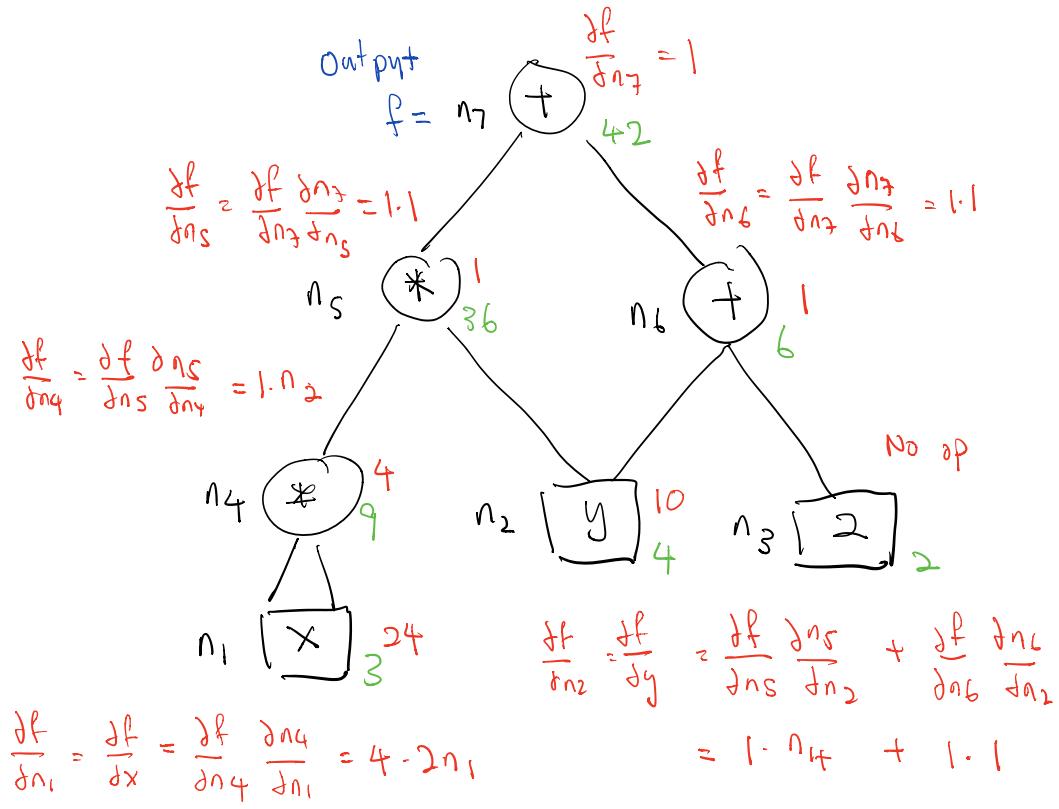
$$f(x, y) = x^2y + y + 2 \quad \text{Inputs } \begin{cases} x=3 \\ y=4 \end{cases} \quad \text{variables}$$

$$\text{Find } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \quad \begin{cases} 2=2 \end{cases} \quad \text{constant}$$

$$\frac{\partial f}{\partial x} = 2xy = 24 \quad \frac{\partial f}{\partial y} = x^2 + 1 = 10$$

Computation Graph

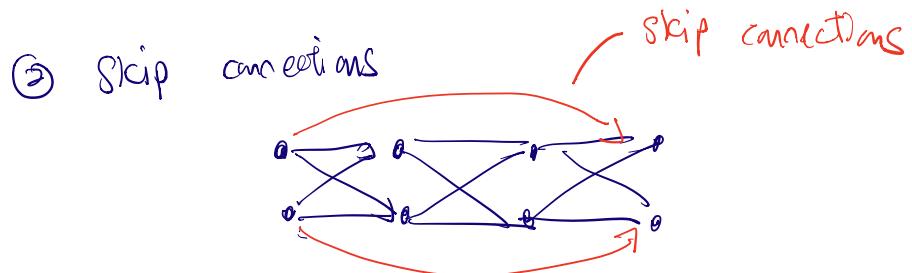
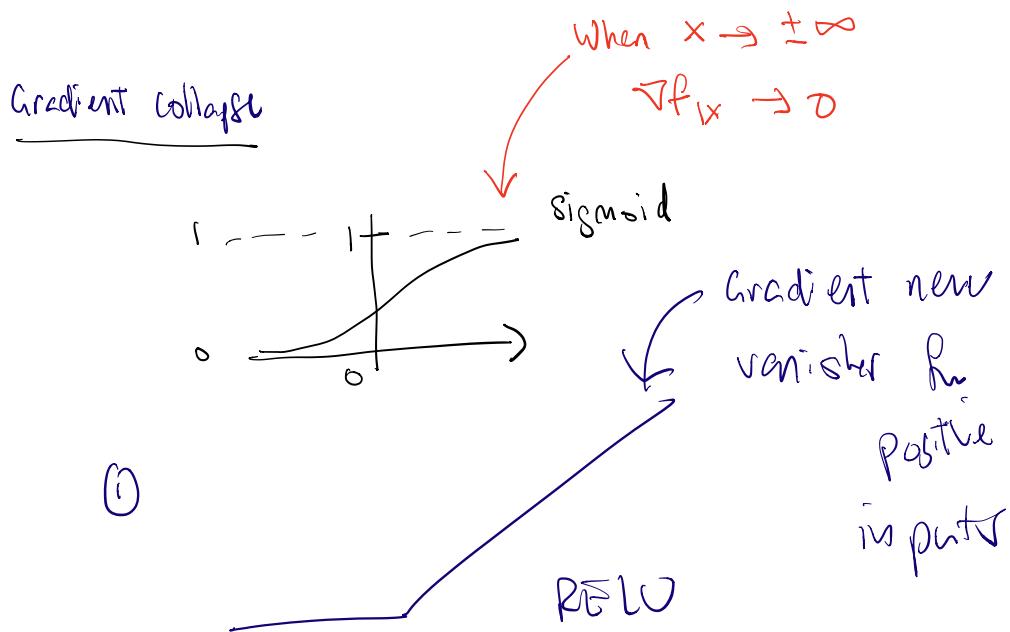




Regardless of # variables, only 2 passes needed
to calculate all partial derivatives.

Gradient issues

Deep networks have many layers.



Gradient explosion

- ① Gradient clipping
- ② Batch normalization

Batch normalization

$$\begin{aligned} \mu_B &= \text{batch mean} \\ \sigma_B &= \text{batch standard deviation} \end{aligned} \quad \left. \begin{array}{l} \text{estimated } \mu \\ \text{inputs} \end{array} \right\}$$

$$\hat{x} = \frac{x - \mu_B}{\sigma_B + \epsilon}$$

$$z = \gamma x + \beta$$

Learned during backprop

Learn optimal scaling for each layer.

Batch, mini-batch ? Stochastic gradient descent

Deep learning often uses large training sets - e.g. 10^6 images

Training size	10^6
Batch size	10^2
Iterations per epoch	10^4

} with batch size of 10^2 ,
we have 10^4 gradient updates
instead of 1 for same work

Gradient $\approx 10^2$ likely to be good approximation
of gradient $\approx 10^6$

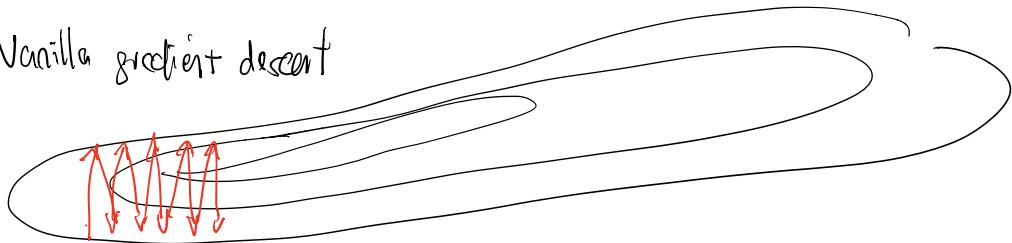
When batch size = 1 \longrightarrow stochastic gradient descent
SGD

Optimal batch size is tunable hyperparameter

depends on variability of training set,
hardware specs (e.g. amount of GPU RAM)

Optimizing gradient descent.

Vanilla gradient descent



① Momentum

$$\text{Vanilla } \theta_{k+1} = \theta_k - \alpha \nabla f(\theta_k)$$

* In physical systems \rightarrow gradient affects velocity or momentum, not position directly.

$$\begin{aligned} & \text{"momentum"} \quad \text{"friction"} \sim \beta = 0.9 \\ & \quad | \quad | \\ & m_{k+1} = \beta m_k - \alpha \nabla f(\theta_k) \\ & \theta_{k+1} = \theta_k + m_{k+1} \end{aligned}$$

(1) $\beta = 0.9$ accumulates history of gradients
 β opposite cancels out
 β some reinforces

Suppose ∇f is constant, say $\nabla f = 1$

$$m_0 = 0$$

$$m_1 = \beta \cdot 0 + \eta = \eta$$

$$m_2 = \beta \eta + \eta = \eta(1 + \beta)$$

$$m_3 = \beta^2 \eta + \beta \eta + \eta = \eta(1 + \beta + \beta^2)$$

$$m_\infty = \beta^2 \eta + \beta \eta + \eta = \eta(1 + \beta + \beta^2 + \beta^3)$$

$$m_\infty = \frac{1}{1-\beta} \eta$$

\downarrow oscillation

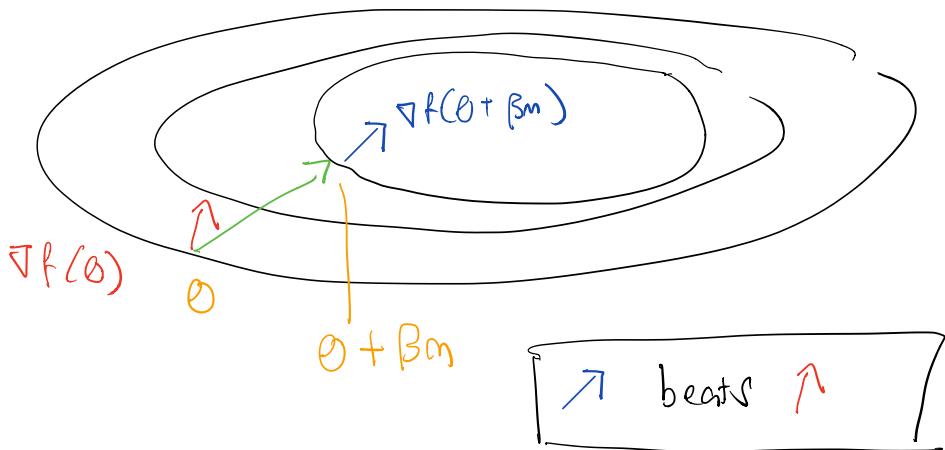
② If $\beta = 0.9$ terminal velocity is 10x more than vanilla.

Nesterov accelerated gradient

$$m_{k+1} = \beta m_k - \nabla f(\theta_k + \beta m_k)$$

$$\theta_{k+1} = \theta_k + m_{k+1}$$

Basic idea \rightarrow evaluate gradient not at current position, but somewhere along direction of momentum.



$$\begin{aligned}
 & \text{RMSProp} \quad \text{decay } \gamma = 0.9 \quad \text{element-wise} \\
 & \quad \quad \quad | \quad \quad \quad / \quad \quad \quad \text{multiplication} \\
 & s \leftarrow \gamma s + \nabla f(\theta) \otimes \nabla f(\theta) \\
 & \theta \leftarrow \theta - \alpha \nabla f(\theta) \odot s \\
 & \quad \quad \quad | \quad \quad \quad \text{element wise division}
 \end{aligned}$$

Basic idea

Reduce gradient in directions where
it is changing the fastest

Increase gradient in direction where

it is changing very slowly



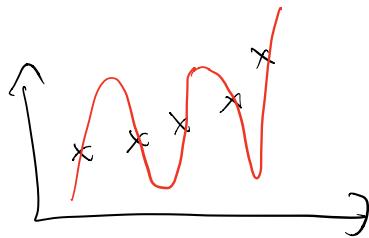
Helps escape long narrow valleys

ADAM (momentum + RMSprop)

NADAM (+ Nesterov)

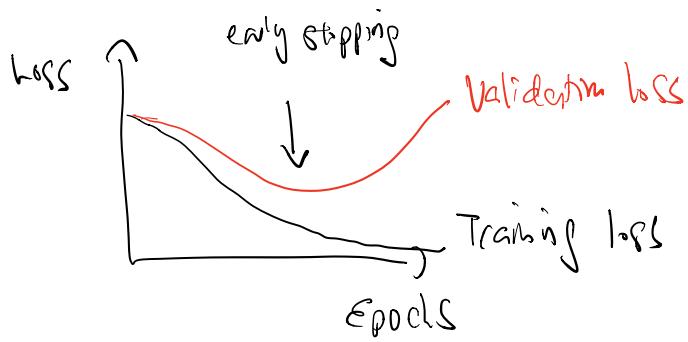
Adaptive Moment estimation

Regularization
Overfitting

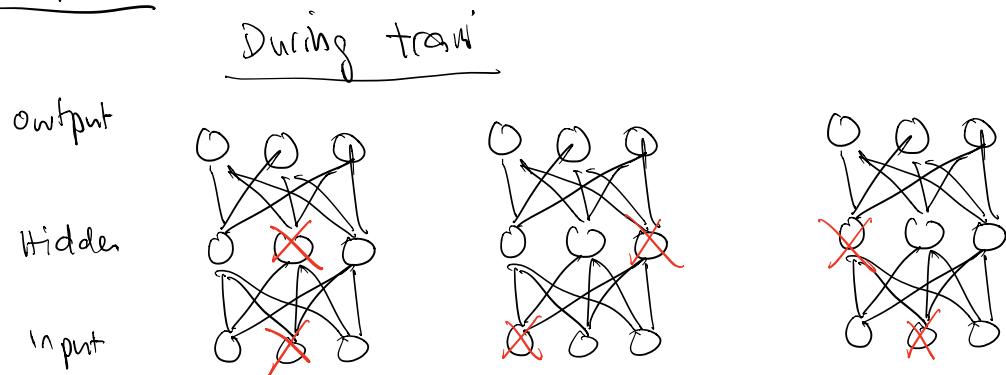


Remedies

- ① \uparrow data
- ② Early stopping
- ③ L₁ or L₂
- ④ Dropout



Dropout



At every step \rightarrow a fraction p of nodes is
inactivated (including input but
excluding output)

① Each step trains a DIFFERENT neural network

\rightarrow Ensemble perspective

② No single node can be very important

\rightarrow Regularization perspective

During test

All nodes are active

weights are multiplied by $1-p$

Review

- ① Anatomy of neural network
- ② What a simple node does $z = f(\sum \omega x + b)$
- ③ Update parameters by gradient descent
 - i) gradient calculated by backpropagation
 - ii) implemented by reverse mode auto differentiation
- ④ Collapsing ? exploding gradients
 - i) choice of activation function
 - ii) skip connections
 - iii) gradient clipping
 - iv) batch normalization
- ⑤ Optimizing gradient descent
 - i) Momentum
 - ii) NESTEROV
 - iii) RMSProp
 - iv) ADAM
- ⑥ Overfitting 3 regularization
 - i) Dropout