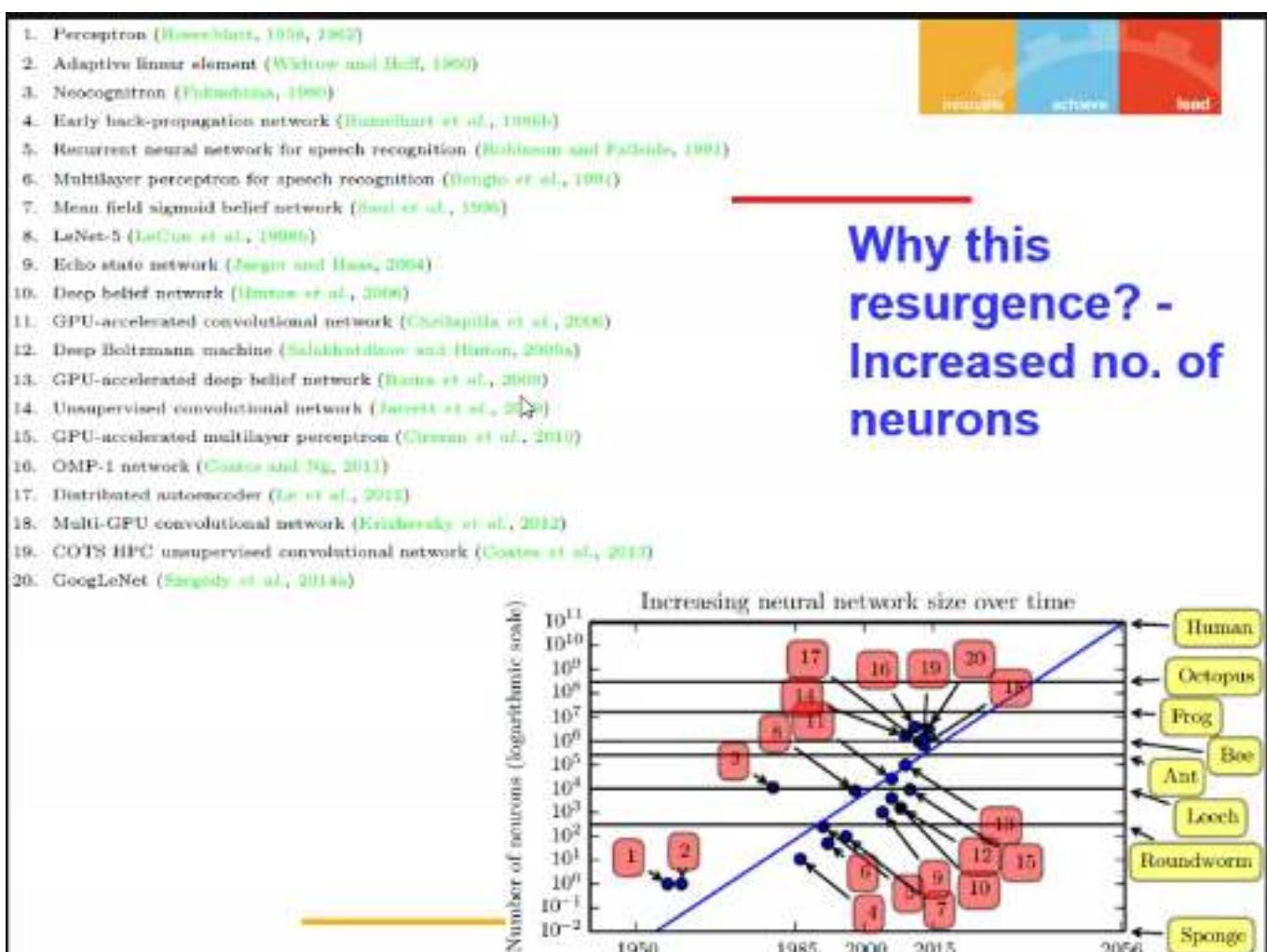


# L2 - I guess(?)

- Deep learning can help discern / disentangle the accent from the voice and get just the pronunciation, through having layers for representing the accent
- Three historical waves
  - Cybernetics (1970s)
  - Connectionism/ Neural networks (peaked in 1995)
  - Deep learning (2006+, layerwise training, big data)
- Hardware didn't increase immensely until 2015, so that's why they were able to increase suddenly only then
  - so we will need a immense development in hardware for a machine to reach the number of neurons in humans. And even then it may not think just as well.
  - DL is *not* simulating a brain. It is inspired from the brain, not imitation



- number of connections per neuron for a human is  $10^4$ , and we have reached that in a COTS HPC unsupervised CNN, but most of our models are still nowhere in comparison. Perhaps we don't even want to simulate the brain.

## **One learning algorithm Hypothesis**

---

- Most perception (input processing) in the brain may be due to one learning algorithm so even if the vision cortex is damaged, the auditory cortex is forced to learn and can make up for it

# L3

---

- Vision Cortex
- Auditory Cortex
- Somatosensory cortex
  - these guys can learn other senses when their primary functions are severed
    - So there may be just one algo behind this
- Perception based algorithms
  - Eg. recognising text and stuff
- focus onto the given problem is going down, but it's not
  - $\implies$  We can't say much about DL
  - We need to tune our DL algorithm, but we can't do a one-for-all algo

Pack learning , Bounce of learning etc.

No free lunch theorem

- averaging over all possible data generation distributions, Every classification algorithm has the same error rate, classifying previously unobserved points.
  - you should not look for a universal learning algorithm
- There is a proof

DL - did better than SVM for handwriting recognition

TL;DR - DL is kinda scam term only for nicer NN



$$y(x) = \sum_{i=1}^n a_i t_i k(x, x_i) + b$$

↓  
finding  
the decision  
boundary

↓  
you need  
to solve  
a function defined on x - testing  
& with training example

$$\max \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j k_i k_j k(x_i, x_j)$$

$$\text{st. } a_i \geq 0$$

$$\sum_{i=1}^m a_i t_i = 0$$

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

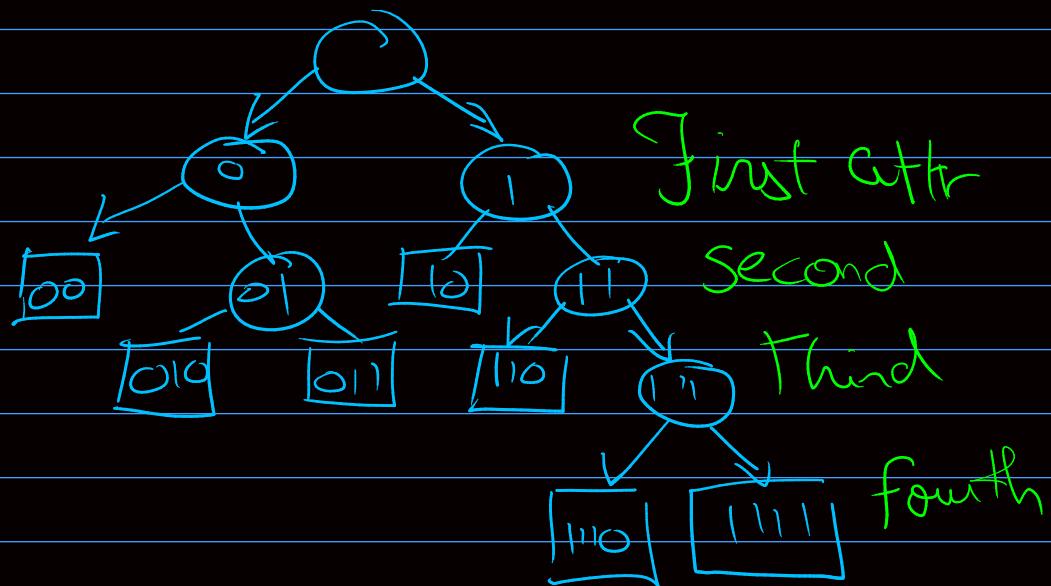
$$b = \frac{1}{N_S} \sum_{i \in S} \left( t_i - \sum_{j \in S} a_j k(x_i, x_j) \right)$$

$$e^{-\frac{\|x-x_j\|^2}{2\sigma^2}}$$

$k$  can be any famous kernel

disadvantage: Too many training examples - 1) or  
much to compute  $O(n^2)$   
no. of training examples

# Decision Tree

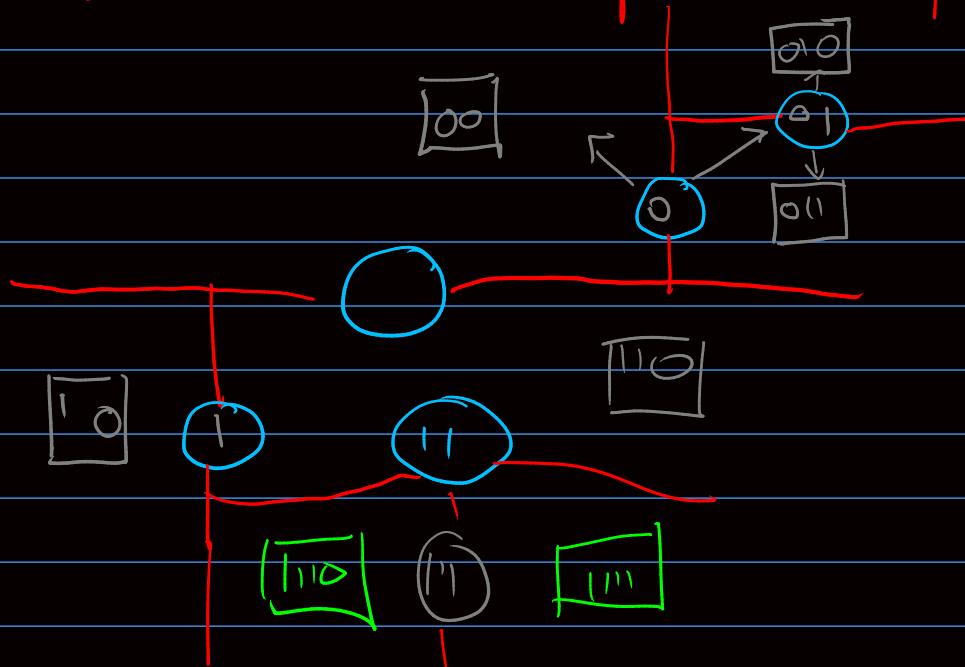


→ Considering a full blown tree, no overfitting

2000, <sup>training</sup> examples → 2000 leaves

to look at the decision boundaries

to divide the 2-D space into 2 parts



- In this case, the max features a decision tree can take is based on the number of training examples.
- In reality there will be exp. nof leaves
  - Decision tree doesn't work
  - You lose features

Curse of Dimensionality:

Imagine attr taking 10 values

You see the samples for each 10

So now you need atleast 10 samples

If there are 2 attr.  $\Rightarrow 10 \times 10 = 100$  samples

3 attr = 1000 samples

30 attr  $\Rightarrow 10^{30}$  samples

$\Rightarrow$  That's why local constancy assumption

for a good number of classical ML

the curse of dimensionality holds

# L7

---

## Universal approximation theorem

---

$$f : R^d \rightarrow R$$

- Can be represented as a network with only one hidden layer
- There will exist an activation function
  - The problem is that this theorem states that there exists the function, and doesn't give the neural net itself
- also in any machine learning problem you don't even know the function

## How many number of hidden units must be there to figure out the function?

---

- It is exponential (takes too long to understand)
- The weight matrix becomes too large to compute

$l$  layers

- $O\left(\binom{n}{d}^{d(l-1)} n^d\right)$  decision regions, for  $d$  features, and  $n$  neurons per hidden layer
- if  $l \geq 5$  then it's called as a deep neural network

## Linear Models

---

- Parameters are linear
  - Eg. Support vector machines
  - or use a linear activation also
- Classical ML can do nonlinear shit too btw.

if  $t$  is target

$\phi$  is transform

$\{(\phi(x_1), t_1), (\phi(x_2), t_2), (\phi(x_3), t_3), \dots, (\phi(x_d), t_d)\}$

Phi can be exponential

we can use a linear model on top of this

## How to make the raw features to matured features?

1. the specific basis functions can be applied by the domain expert, to get the matured features, and then apply the linear model
  - Example, BMI
2. Kernels in SVMs, going with some kind of functions to transform. Eg. Gaussian Kernel



$$(x^{(1)}, t_1) (x^{(2)}, t_2) \dots (x^{(n)}, t_n)$$

$$x = (x_1, x_2, \dots, x_d)$$

$$y(x, w)$$

$$\hookrightarrow y = w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_m \phi_m(x) + b$$

$$y = w_1 x_1 + w_2 x_2 + \dots + w_m x_m$$

$y$  is linear in  $\phi$  (matured features)

but not in  $x$

SVM

$$x_1, x_2, \dots, x_d$$

$$\underbrace{\phi_1, \phi_2, \dots, \phi_m}_{\text{You may not know}}$$

$$k(x^i, x^j) = \frac{-\|x^i - x^j\|^2}{\ell}$$

[

Using SVMs with a kernel  $\rightarrow$  each vector  
is a support vector  $\Rightarrow$  overfit

In SVMs, we fix our features

DL  $\rightarrow$  we get a ton of features

$\rightarrow$  we don't know how many are appropriate

Problem is, architecture is mostly trial & error

$x_{11}, x_{12}, \dots, x_{1d}$ , apply some nonlinear functions  $[\phi_1, \phi_2, \dots, \phi_m]$  to the vector

→ What we get is a matured feature vector  
it may be a linear model of matured features  
but non linear model of original features  
if the basis functions  $[\phi_1, \dots, \phi_m]$  are non linear



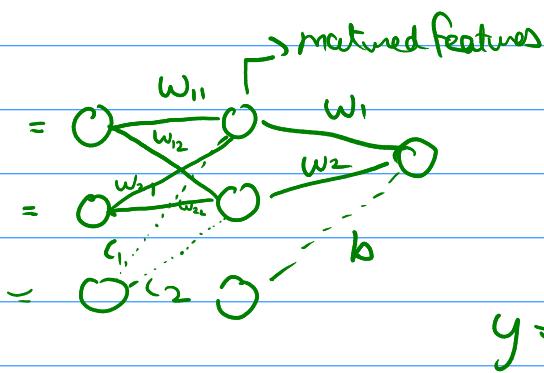
## Feed forward networks

$x_1$	$x_2$	$t$
0	0	0
0	1	1
1	0	1
1	1	0

4 examples

$$\sum_{i=1}^4 ((w_1 x_1^{(i)} + w_2 x_2^{(i)} + b) - t)^2$$

Linear regression  $l = \text{neural network}$   
 $w_1 = 0 \quad w_2 = 0 \neq b = 0.5$



$$h = f^{(1)}(x; W, C)$$

$$y = f^{(2)}(h, w, b) = w^T h + b$$

$$y = f^{(2)} f^{(1)}(x; W, C, w, b)$$

What will be  $f^{(1)}$ ?

L.F.  $\begin{bmatrix} 2 & 1 \end{bmatrix}$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \in \mathbb{R}^2 \text{ to } \mathbb{R}^3 \text{ can be } \begin{bmatrix} f_1(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \\ f_2(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \\ f_3(\begin{bmatrix} 2 \\ 1 \end{bmatrix}) \end{bmatrix}$$

$$\begin{bmatrix} 3 & 15 \\ 4 & 26 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 10 \\ 16 \end{bmatrix}$$

$$\left. \begin{array}{l} f_1 = 3x_1 + 4x_2 \\ f_2 = x_1 + 2x_2 \\ f_3 = 5x_1 + 6x_2 \end{array} \right\} \begin{array}{l} \text{linear combination} \\ \Rightarrow \text{linear transformation} \end{array}$$

## Affine transformation

Affine transform: linear transformation + constant vector added to result

$$f^{(l)}(x; w, c) = \underbrace{w^T x + b}_{\text{Affine Transform}} + \underbrace{c}_{\text{matured space}}$$

For brevity (and to make life easier), let us just consider a linear transformation

$$\begin{aligned} y &= w^T (W^T x) \\ &= (\omega')^T x \text{ where } (\omega')^T = w^T W^T \end{aligned}$$

$$\sum_{i=1}^4 (w^T x - t_i)^2$$

1' SO, The resultant matured features of a linear transformation is a linear model. (which can't fit a non-linear function like XOR)

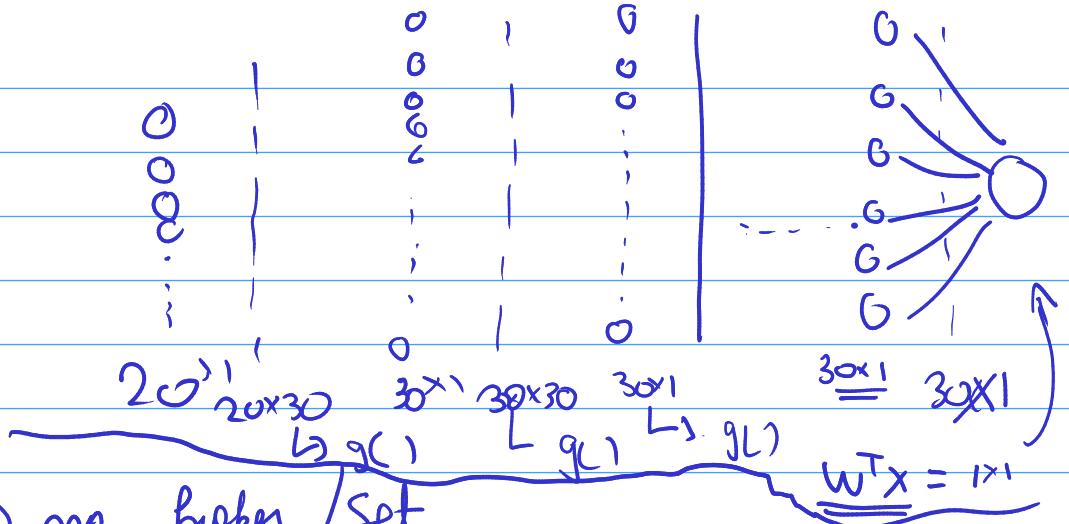
So as a result of this problem, all of deep learning and some machine learning algorithms use something else on top of it

$$\begin{bmatrix} 1 \\ x \end{bmatrix}^T \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} + \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} w_{11} + w_{21} + c_1 \\ w_{12} + w_{22} + c_2 \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

Add a non-linear function  $= \begin{bmatrix} g(w_{11} + w_{21} + c_1) \\ g(w_{12} + w_{22} + c_2) \end{bmatrix}$   
one of them could be

$g(z) = \max(z, 0) \Rightarrow$  This function is an activation

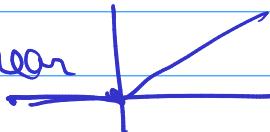
applied to every hidden layer  $\leftarrow$  function



There's no proper / set

"algorithm" to choose an activation function, you have to experiment/make a judgement to choose

RELU: piecewise linear



but not actually linear  
→ nonlinear

- It gives nonlinear properties, and can be used as activation
  - ⇒ but you need to do the math involved.
- It even has some biological motivation.

Similarly, there is not set algo for Number of layers  
we will need insights for the domain/  
experimentation if the above insights are  
there.

First Sec HOME INSERT DRAW VIEW

$x_1, x_2, \dots, x_m$   $\Gamma_{\text{model}} \sim \mathcal{N}(\mu, \sigma^2)$

$$\prod_{i=1}^m p(x_i | \mu, \sigma^2) = \prod_{i=1}^m p(x_i | \mu, \sigma^2)$$

$$\text{arg max}_{\mu, \sigma^2} \prod_{i=1}^m p(x_i | \mu, \sigma^2) = \text{arg max}_{\mu, \sigma^2} \log \left( \prod_{i=1}^m p(x_i | \mu, \sigma^2) \right)$$

$$= \text{arg max}_{\mu, \sigma^2} \sum_{i=1}^m \log(p(x_i | \mu, \sigma^2))$$

$$= \text{arg max}_{\mu, \sigma^2} \frac{1}{m} \sum_{i=1}^m \log(p(x_i | \mu, \sigma^2))$$

to make life easier

$$X = x \quad \text{discrete}$$

$$E(f(x)) = \sum_{i=1}^m f(x_i) P(X=x_i)$$

Normal distribution  $X \sim \hat{P}_{\text{data}}$

$$\frac{1}{m} \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right)$$

$$= \frac{1}{m} m \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{m} \sum_{i=1}^m \left( (-1) \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$= \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{m} \sum_{i=1}^m \frac{(x_i - \mu)^2}{2\sigma^2}$$

$P_{\text{data}}(x)$   $P(x_i | \theta)$   $\hat{P}_{\text{data}}(x)$   
 actual distribution model discrete

$$P_{\text{model}}(x; \theta) \quad \hat{P}_{\text{data}}(x)$$

(discrete)

$D_{KL}(\hat{P}_{\text{data}} \parallel P_{\text{model}})$  shows how dissimilar they are.

$$E_{x \sim \hat{P}_{\text{data}}} [\log \hat{P}_{\text{data}}(x) - \log P_{\text{model}}(x)]$$

$$= E_{x \sim \hat{P}_{\text{data}}} \log (\hat{P}_{\text{data}}(x)) - E_{x \sim \hat{P}_{\text{data}}} (\log P_{\text{model}}(x))$$

("from" follows)

$$= \sum_{i=1}^m \frac{1}{m} \log \left( \frac{1}{m} \right) - \dots$$

Given  
(there are  $m$  things  
you can pick any)

$$= -\log m -$$

$$= \underset{\mu, \sigma^2}{\operatorname{argmin}} \quad E \left[ \frac{-1}{\sqrt{2\pi} \sigma} + \frac{1}{m} \sum_{i=1}^m \left( \frac{x_i - \mu}{2\sigma^2} \right)^2 \right]$$

KL divergence is called cross entropy  
(where the points taken are the same as the data or not)

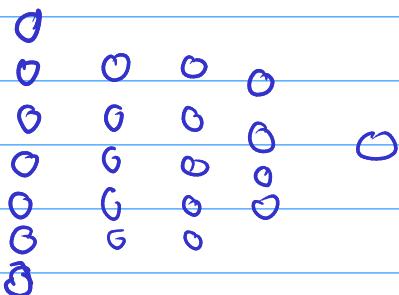
Cross entropy ( $\hat{P}_{\text{data}}, \hat{P}_{\text{model}}(x; \mu, \sigma^2)$ )

$$= - E_{x \sim \hat{P}_{\text{data}}} (\log P_{\text{model}}(x))$$

Whenever you are finding the parameters by MLE, the error function to that is the result of the likelihood function, you call the error function as a "cross entropy" function

This applies to regression as well! MSE if  
so can be got from a likelihood function too,  
it's also a cross entropy error!

Deep learning fixes the family of functions and uses that family



$$p(y_1, y_2, \dots, y_m | x_1, \dots, x_m) = \prod_{i=1}^m p(y_i | x_i)$$

$$\max \log \left( \prod_{i=1}^m p(y_i | x_i) \right)$$

$$\text{min} = \log \left( -\min \sum_{i=1}^m -p \log p(y_i | x_i) \right)$$

$$y(\theta) = -\log p(x|\theta)$$

$$w_1 w_2 \dots = -\log((2\gamma-1)z)$$

$$= \sigma(-(2\gamma-1)z)$$

~~not activation~~ ~~softmax~~ ~~constant~~

$$= z((1-2\gamma)z)$$

@  $\gamma=1$

$$z((1-2\gamma)z) = z(-z)$$

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

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

$$\frac{\partial y(\theta)}{\partial z}$$

$$\sigma((1-2\gamma)z)(1-2\gamma)$$

when  $z \gg 0$

$$z(-z) \rightarrow 0$$

when  $z \ll 0$

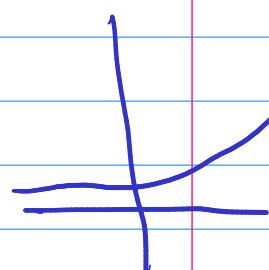
$$z(-z) \rightarrow \text{large}$$

@  $\gamma=0$   $z \gg 0$

$z(z) \Rightarrow \text{large}$  ( $w^{(k+1)}$  changes quite a lot)

$z \ll 0$

$z(z) \rightarrow 0$  ( $w$  won't change)

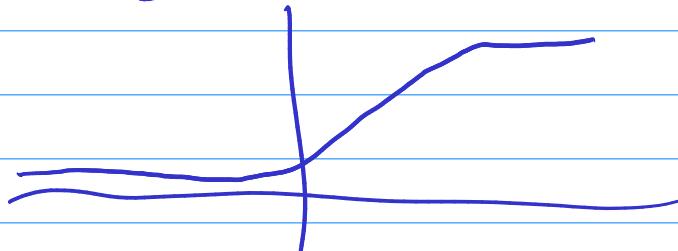
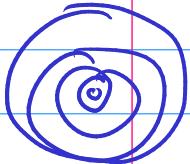


(Take example, treat as SGD)

why is this cost function any better than MSE

PC

$$\begin{aligned} & \max \log P(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_m) \\ &= \max \log \prod_{i=1}^m P(y_i | x_i) \\ &= \max \log \left[ \sigma(z_i)^{y_i} (1 - \sigma(z_i))^{1-y_i} \right] \\ &= \max \sum (y_i \log \sigma(z_i) + (1-y_i) \log (1 - \sigma(z_i))) \\ &= \min \left[ \sum -[y_i \log \sigma(z_i) + (1-y_i) \log (1 - \sigma(z_i))] \right] \end{aligned}$$



$$= \min -[y \log \sigma(z) + (1-y) \log (1 - \sigma(z))]$$

## Activation

Leaky Relu

PReLU

Maxout

RBF

Soft plus

Hard tail

Sigmoid

$$\min_{\omega, b} L(\omega, b, x^{(i)}, t^{(i)}) = \cancel{\text{something}}$$

$$\omega^{(k+1)} = \omega^k - \eta \left( \frac{\partial L}{\partial \omega} \right) \quad \omega_1 = \omega_i^{(k)}$$

$\frac{\partial L}{\partial \omega_1} \quad \frac{\partial L}{\partial \omega_2}$   
 $\frac{\partial L}{\partial \omega_{1000c}}$

$$\frac{\partial z}{\partial t} = f(t)$$

$$\frac{\partial z}{\partial y} = \frac{\partial f(y)}{\partial y}$$

$$\frac{\partial z}{\partial x} = \frac{\partial f(g(x))}{\partial g(x)} \cdot \frac{\partial g(x)}{\partial x}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i}$$

$$\nabla_x z = \frac{\partial z}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}^T \begin{bmatrix} \frac{\partial z}{\partial y_1} \\ \vdots \\ \frac{\partial z}{\partial y_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial y}{\partial x} \end{bmatrix}^T \nabla_y z$$

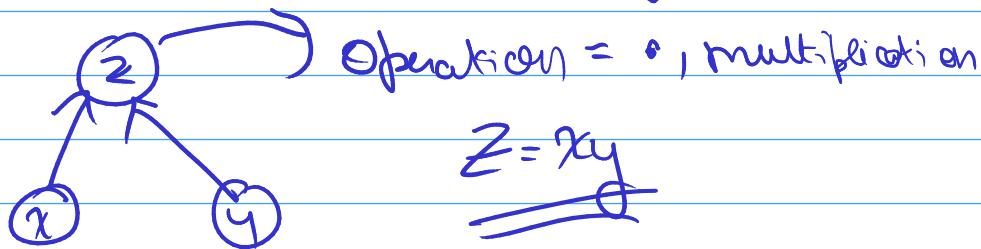
If  $x$  is not a vector, but of <sup>a</sup> higher order (matrix or str.)  
then people see  $x$  as a vector as a whole

→ write the matrix as follows and actually expand this internally

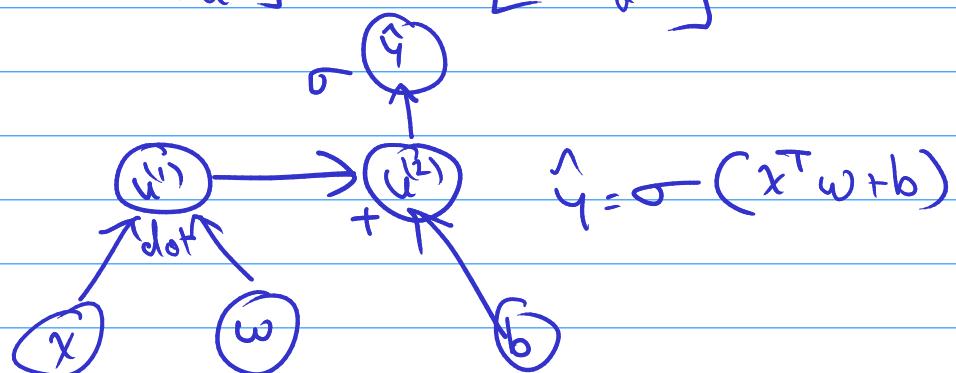
[This is brainfuck on]

To ease us in this process, we use a computation graph, which is a discrete structure

Ex Computational graph of  $xy$

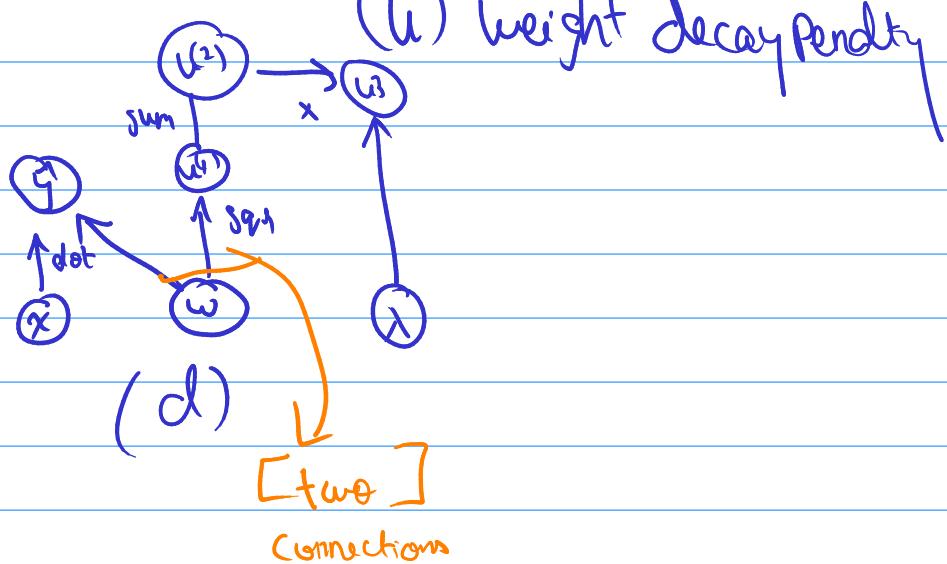


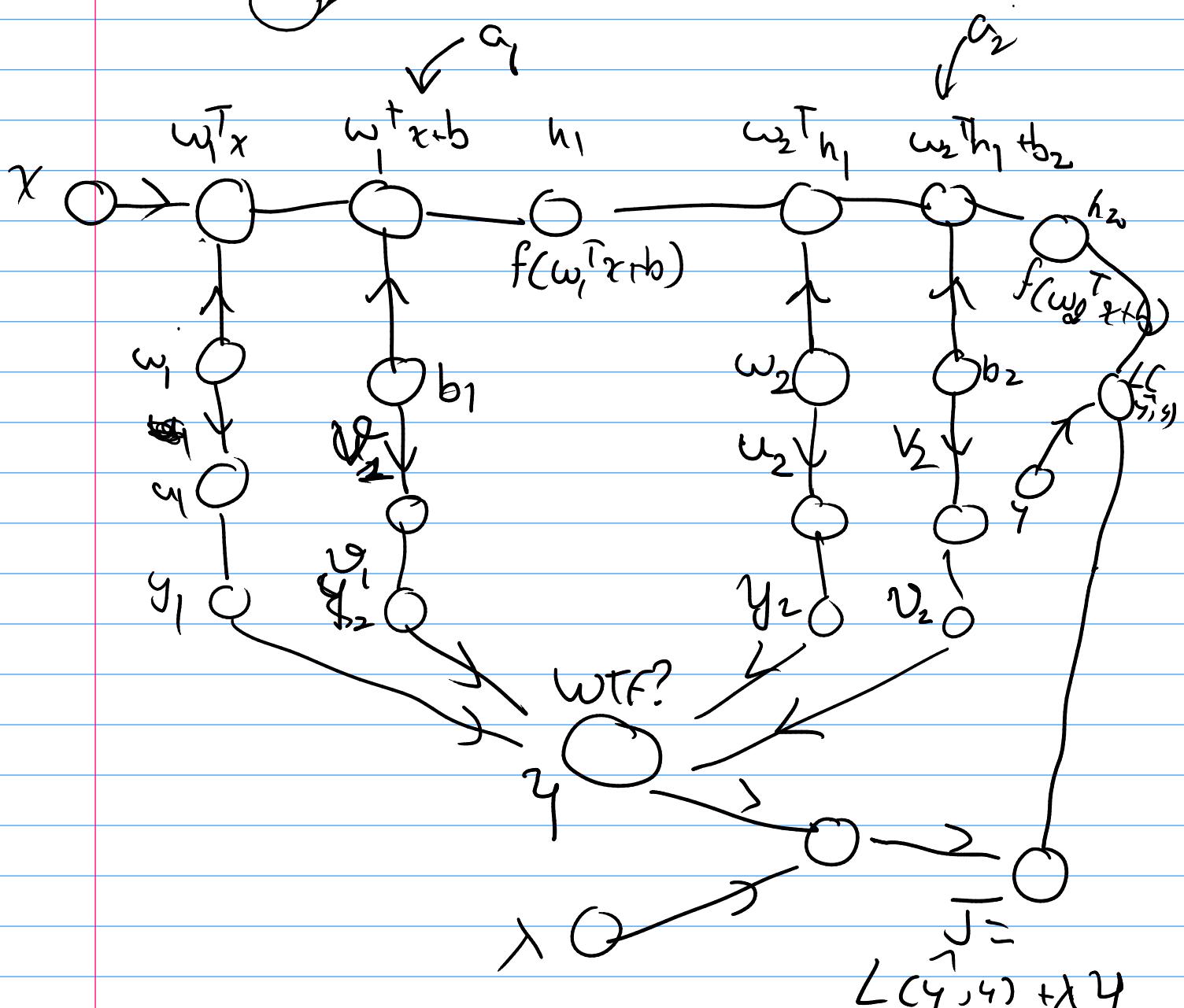
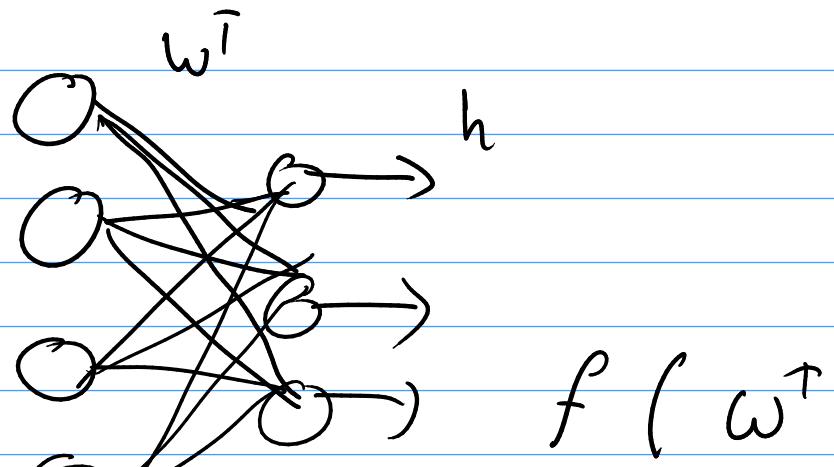
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$



Weight decay (i)  $\hat{y}$  Two operations on 1 variable

(ii) weight decay penalty





$y = \Omega(0)$  (regularized)

$$\min L(\hat{y}, \hat{y}) + \lambda \Omega(\theta)$$

↓  
w, b

$$\Omega(\theta) = \omega_1^2 + \omega_2^2 + \dots + \omega_{10}^2 + \omega_{\theta}^2$$

$$\begin{matrix} \nabla_{\omega_1} J & \nabla_{b_1} J \\ \nabla_{\omega_2} J & \nabla_{b_2} J \\ \nabla_{\omega_3} J & \nabla_{b_3} J \end{matrix}$$

A matrix of  
weights can  
be treated  
as a vector of  
vectors

How can we find the derivative of a scalar  
with a vector

3 layered NN

$$\begin{aligned} \nabla_{a_3}(L(\hat{y}, y)) &= D_{a_3}(h_3) \nabla_{h_3}(L(\hat{y}, y)) \\ &= \nabla_{a_3}(f(a_3)) \odot g = f'(a_3) \odot g \end{aligned}$$

①      ↓ element-wise mul

$$\frac{\partial L}{\partial a_3} = \frac{\partial L}{\partial h_3} \cdot \frac{\partial h_3}{\partial a_3}$$

∴

$$\nabla_{a_3}(J) = \nabla_{a_3}(L(\hat{y}, y) + \lambda \Omega(\theta))$$

$$= \nabla_{w_3}(L(\hat{y}, y)) + \lambda \nabla_{w_3}(\Omega(\theta))$$

$$= \nabla_{a_3}(L(\hat{y}, y)) \nabla_{w_3}(a_3)$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial w_3} + \lambda \nabla_{w_3}(\Omega(\theta))$$

$$\left[ \textcircled{2} \quad \nabla_{a_3}(L(\hat{y}, y)) = f'(a_3) \odot \hat{y} \right] \dots \textcircled{1}$$

$$f'(a_3) \odot \hat{y} \quad \nabla_{w_3} (\omega_3^T b_2 + b_3)$$

$$\frac{\partial J}{\partial w_3} \Big|_{w_3=w_3^{(k)}}$$

$$\underline{\underline{\text{Eq}}} \quad \frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial b_3} = f'(a_3) \odot \nabla_{b_3} (\omega_3^T b_3)$$

keep applying  
Chain rule

$$\nabla \omega_1(\mathbf{J})$$

$$\nabla b_1(\mathbf{J})$$

$$\nabla \omega_3 \mathbf{J} = \nabla_{\mathbf{a}_3} (\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y})) + \lambda \nabla \Omega(\mathbf{G})$$

$$= \nabla_{\mathbf{a}_3} (\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y})) \nabla_{\mathbf{a}_3}(\mathbf{a}_3) + \lambda \nabla_{\mathbf{a}_3}(\Omega(\mathbf{G}))$$

$$= \underbrace{\mathbf{G}(\mathbf{f}'(\mathbf{a}_3)) \mathbf{h}_2^\top}_{\mathbf{g}(\mathbf{f}'(\mathbf{a}_3)) \mathbf{h}_2^\top} + \lambda \nabla_{\mathbf{a}_3}(\Omega(\mathbf{G}))$$

$$\nabla_{b_3} \mathbf{J} =$$

$$\mathbf{g}(\mathbf{f}'(\mathbf{a}_3))^\top + \lambda \nabla_{b_3}(\Omega(\mathbf{G}))$$

Dont worry

I moved the text to the  
next page

$$\begin{aligned}\nabla_{h_2}[L(\hat{y}, y)] &= \nabla_{h_2} a_3 \nabla_{a_3} L(\hat{y}, y) \\ &= \nabla_{h_2} (\omega_3^T h_2 + b_3) (g \odot f'(a_3)) \\ &= \omega_3^T (g \odot f'(a_3))\end{aligned}$$

$$g \leftarrow \nabla_{h_2} (L(\hat{y}, y)) = \omega_3^T (g \odot f'(a_3))$$

$$\begin{aligned}\nabla_{a_2}(L(\hat{y}, y)) &= \nabla_{h_2}(L(\hat{y}, y)) \nabla_{a_2}(h) \\ &= g \nabla_{a_2}(f(a_2)) = g \cdot f'(a_2)\end{aligned}$$

$$\begin{aligned}\nabla_{w_2} J &= \nabla_{w_2}(L(\hat{y}, y)) + \lambda w_2 (-\Omega(\theta)) \text{ (2 layers)} \\ &= \nabla_{a_2}(L(\hat{y}, y)) \nabla_{w_2}(a_2) + \lambda \nabla_{w_2}(-\Omega(\theta)) \\ &= (g \cdot f'(a_2)) (\nabla_{w_2}(\omega_2^T h_1 + b_2)) + \lambda \nabla_{w_2}(-\Omega(\theta))\end{aligned}$$

$$\begin{aligned}\nabla_{b_2}(J) &= \nabla_{b_2}(L(\hat{y}, y) + \lambda -\Omega(\theta)) \\ &= \nabla_{a_2}(L(\hat{y}, y)) \nabla_{b_2}(a_2) + \nabla_{b_2}(\lambda -\Omega(\theta)) \\ &= (g \cdot f'(a_2)) \nabla_{b_2}(\omega_2^T h_1 + b_2) + \nabla_{b_2}(\lambda -\Omega(\theta)) \\ &= (g \cdot f'(a_2)) + \lambda \nabla_{b_2}(-\Omega(\theta))\end{aligned}$$

$$\begin{aligned}\nabla_{h_1}(L(\hat{y}, y)) &= \nabla_{a_2}(L(\hat{y}, y)) \nabla_{h_1}(a_2) \\ &= (g \cdot f'(a_2)) \nabla_{h_1}(\omega_2^T h_1 + b_1) \\ &= (g \cdot f'(a_2)) \omega_2^T \\ g &\leftarrow \nabla_h(L(\hat{y}, y))\end{aligned}$$

$$111^y \quad \nabla_{w_1}(\mathcal{J})$$

can be found

$$\nabla_b(\mathcal{J})$$

from chain rule

in textbook he takes gradient  
of loss function with regularization

$$g \leftarrow \nabla_{\hat{y}}(\mathcal{J}) \quad \text{this is incorrect!}$$

$$\nabla_{\hat{y}}(\mathcal{J}) = \nabla_{\hat{y}}(L(y, \hat{y}) + \lambda \Omega(w, b))$$

$$= \nabla_{\hat{y}}(L(y, \hat{y})) + \lambda \nabla_{\hat{y}}(\Omega(w, b))$$

both  $\Omega$  &  $L$  are functions of  $w_3$ !!

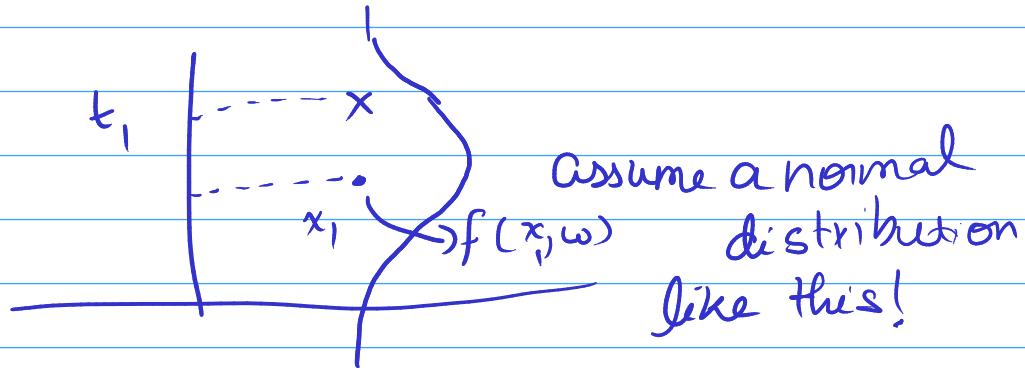
This gradient is not zero!

HOWEVER

You may say that TB considers  
regularization as part of the  
loss function  $\Leftrightarrow$  i.e.  $L(\hat{y}, y)$  has  
the regularized term inside

$$\text{e.g.: } L(\hat{y}|y) = \text{MSE} + \lambda \Omega(\theta)$$

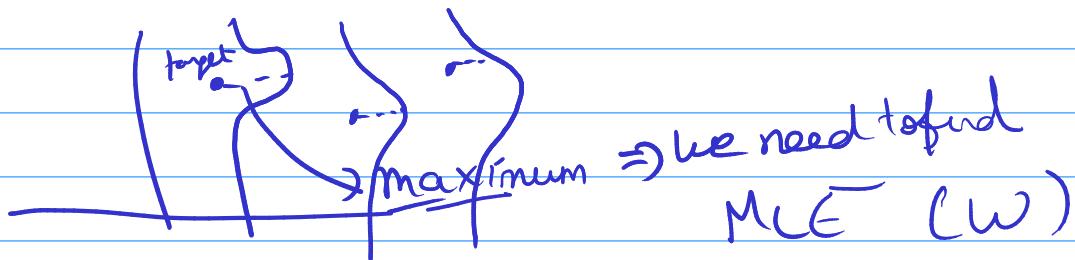
$$y/x \sim N(\bar{y} / f(x_w), \bar{\sigma})$$



Consider covariance matrix to be  $\Sigma$ . For now

now  $f(x)$  is quite away from target

now let's have another  $w_2$  such that



Covariance matrix can also be determined  
 $\Leftrightarrow$  positive, semidefinite  
you need a more complex neural network

What cost functions are good??

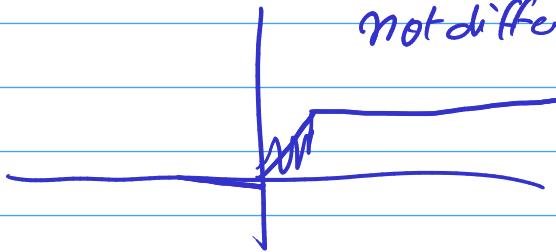
$$\rightarrow p(y=1) = \sigma(w_0 + w_1 \bar{w}_1 + w_2 \bar{w}_2 + \dots + w_n \bar{w}_n)$$

what if

$$p(y=1) = \max \{0, \min \{1, w^T h + b\}\}$$

gradients are zero outside

not differentiable, also saturates,



we need a better cost function !!

$$-\log(p(\dots)) \dots$$

So we need a cost function that saturates

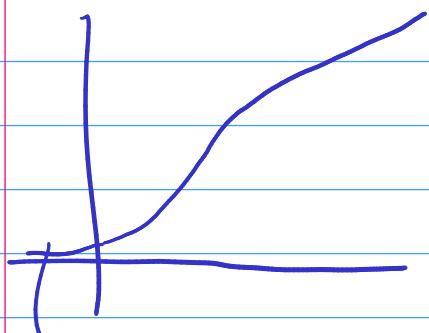
at correct classification

& doesn't saturate @ incorrect classifications

$y x$	1	0
$p(y x)$	$p$	$1-p$

$p$  is a bernoulli

Cost functions aren't forced, it's a consequence  
of an assumption  
of a normal distribution



} if predicted output is low  
(-ve)  
but target's positive } Saturation

$\rightarrow (Z)$ , if there is small change  $Z \ll 0$   
then function won't change

Ex.  $Z \rightarrow \infty$  or  $Z \rightarrow -\infty$

and in center interval, it changes a lot!

My error function should change when I'm not near global or local minima

in SGD, since you're doing it example by example

you want to improve vs drastically  
 $\Rightarrow$  Function shouldn't saturate

## Deriving Cost

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

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

assume  $\ell_p(z) = \log(1+e^z)$  → softplus function

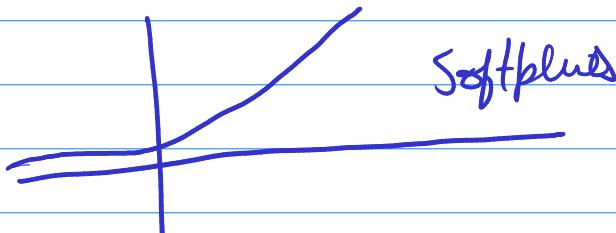
$$z \rightarrow +\infty \quad \log(1+e^z) \rightarrow \infty$$

$$z \rightarrow -\infty \quad \log(1+e^{-z}) \rightarrow 0$$

$$x^{+\infty} = \max\{0, x\}$$

Softplus version  
of max  
function

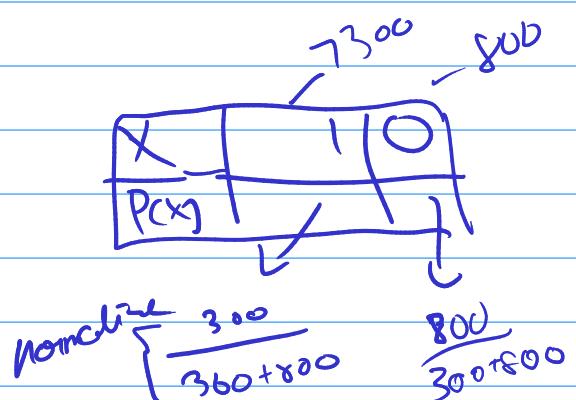
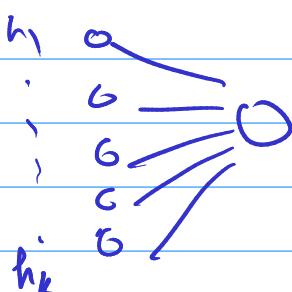
Very  
noice



$$\ell_p(1+e^{-z}) = \log\left(\frac{e^z+1}{e^z}\right) = \log\left(\frac{1}{e^{-z}}\right) \\ = -\log(\sigma(z))$$

$$\log(\sigma(z)) = -\ell_p(-z)$$

$$z = w_0 + w_1 h_1 + \dots + w_k h_k$$



I want to combine  $z$  with target to arrive @ cost  
function that  
doesn't saturate

So they choose  $y_2$  as that combination such that

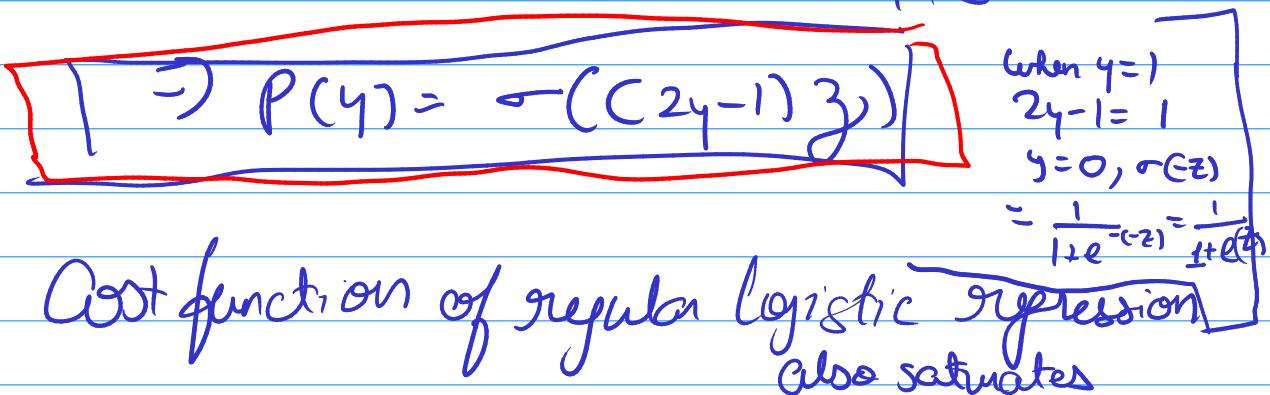
$$\log(\hat{p}(y)) = y_2 = \hat{p}(y) = e^{y_2}$$

$$P(y) = \frac{e^{y_2}}{\sum_i e^{y_i}} = \frac{e^{y_2}}{1 + e^z} \quad (y \text{ takes values } 0, 1)$$

(Normalized)

A "kind of" probability distribution

$$P(y=1) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}} P(y=0) = \frac{1}{1 + e^z}$$



$$P(y_1, \dots, y_m | x_1, \dots, x_m) = \prod_{i=1}^m P(y_i | x_i)$$

$$\max \log \left( \prod_{i=1}^m P(y_i | x_i) \right)$$

$$\max \sum_{i=1}^m \log(P(y_i | x_i))$$

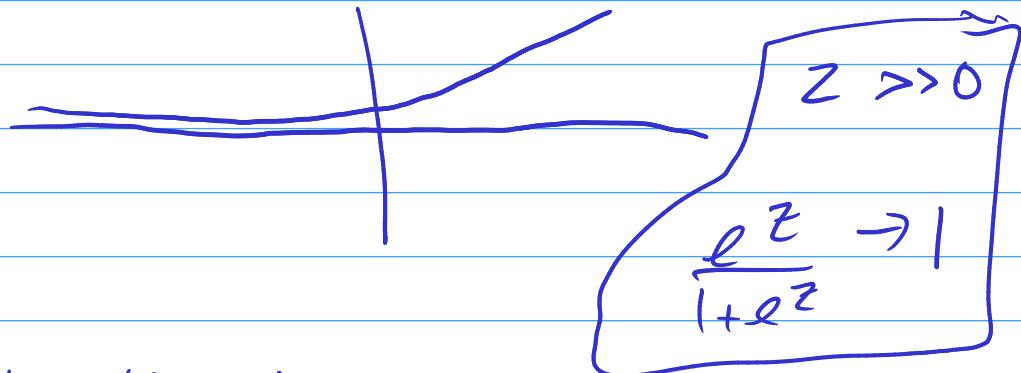
$$\min \sigma \sum_{i=1}^m -\log(P(y_i | x_i))$$

$$\text{Cost func: } J(\theta) = -\log P(y|x)$$

$$= -\log(\sigma((2y-1)z)) = \ell(-((2y-1)z))$$

$$= \varphi((1-\gamma)z)$$

remember  $\varphi$  is the softplus function



$$\omega^{(k+1)} = \omega^{(k)} - \eta \underbrace{\frac{d\ell}{d\omega}}$$

gradient should be small, //  $z \gg 0$ , & prediction is correct

This  $\ell$  function achieves that as, if  $y=1$

$$\varphi((1-\gamma)z)$$

$$\ell(-z) \approx 0 \text{ since } z \gg 0$$

$z \gg 0, y=0$ , you want high gradients

then  $\varphi'(z) = \rightarrow \infty$   
you get

for  $z \ll 0 \quad y=0 \quad$  then you get  $\varphi'(z) \rightarrow 0$

$z \ll 0 \quad y=1 \quad \ell(-z) \rightarrow -z$

Why tho do I function  
WTF is wrong with what I have?

$$P(Y_i=1/x_i) = \sigma(z_i)$$

$$P(Y_i=0/x_i) = 1 - \sigma(z_i)$$

$$P(Y_i) = (\sigma(z_i))^{y_i} (1 - \sigma(z_i))^{1-y_i}$$

$$p(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_m) = \prod_{i=1}^n (p(y_i/x_i))$$

$$\stackrel{\text{prob}}{=} \prod_{i=1}^n (\sigma(z_i))^{y_i} (1 - \sigma(z_i))^{1-y_i}$$

$$= \log(p(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_m))$$

$$\max \sum_{i=1}^m \left[ y_i \log(\sigma(z_i)) + (1-y_i) \log(1-\sigma(z_i)) \right]$$

$$= \min \sum_{i=1}^n - (y_i \log \sigma(z_i) + (1-y_i) \log(1-\sigma(z_i)))$$

$Z \gg 0 \quad y=1 \rightarrow \text{saturates}$

$Z \gg 0 \quad y=0 \quad \text{also saturates};$

$\sigma(z)$  saturates

$\Rightarrow \log(1-\sigma(z))$  saturates (ultimately)

when  $\sigma(z)$   
doesn't change  
much at all)

$\Rightarrow$  problem //

Up next: softmax, onehot

& how to deal with them!

Alright, let's say  $m$  outputs are there

Then how do we predict

$(\mathbf{w}^T \mathbf{h} + b)$  is input

$$(\mathbf{w}^T \mathbf{h} + b) = z_i := \log \hat{p}(y=i/x)$$

unnormalized  
log probabilities

$$\Rightarrow e^{w^T h + b} = \hat{p}(y=i/x)$$

$$p(y=i/x) = \frac{e^{z_i}}{\sum_{i=1}^m e^{z_i}} \text{ for } m \text{ classes}$$

Why distribution? This is a multinomial dist.



→ so we can "turn"  $z$  into a probability this way.

→ we don't normally maximize THE probabilities  
so we ~~not~~ maximize the (unnormalized) log probabilities

then we normalize

Softmax = activation

$$\text{Softmax} = \frac{e^{z_i}}{\sum_{i=1}^m e^{z_i}}$$

one hot encoding

$$\min_m -\log \prod_{i=1}^n \left( p(y=2|x_i) \right)^{y_i} \left( p(y \neq 2|x_i) \right)^{1-y_i} \left( p(y=2/x) \right)$$

This however not for sure that you'll hit local minima often  
some iterations

We might need early stopping

Other activation funcs

ReLU: differentiability problem

@ '0' take either LHS, or RHS derivative

Leaky ReLU

$$h_i = g(z_i) = \max(0, z_i) + \alpha \min(0, z_i)$$

$\alpha = -1 \Rightarrow$  absolute value activation function

$$\begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$$

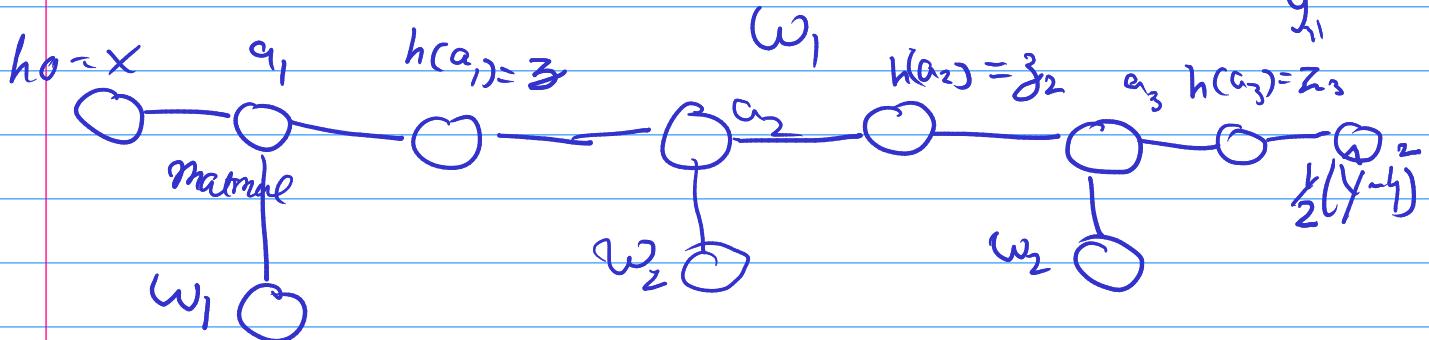
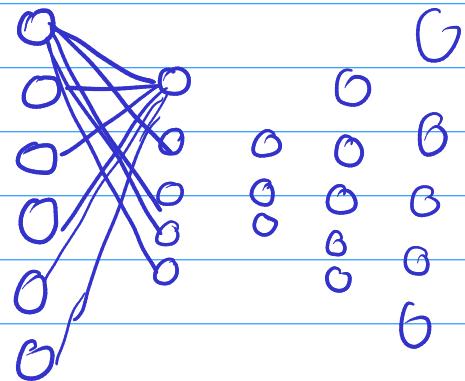
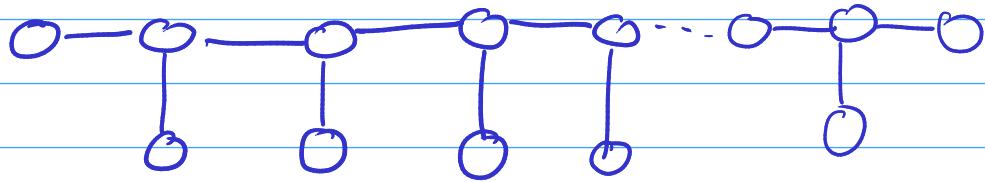
take max of each  $w^T x$  th in each division  
and all neurons in that division  
have max of division in that (maxout)  
output

This is maxout

Dont use Sigmoid,  $\tanh(z)$  saturates lesser, use that instead  
(make sure you shift origin of the graph)  
softmax  $\log(t + e^z)$

hardtanh  $\max(-1, \min(1, a))$

$$E(\omega) = \boxed{\text{ }}$$



$$\omega^{(k+1)} \leftarrow \omega^{(k)} - \eta \frac{\partial E}{\partial \omega} \Big|_{\omega=\omega^k}$$

$$\hookrightarrow \nabla E(\omega) \Big|_{\omega=\omega_k}$$

$$\nabla_A(z) = \nabla_C(z) \quad \nabla_A(c) = G_B^T$$

$$\nabla_C z_{20 \times 1} = \boxed{\quad}_{20 \times 1}$$

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{10} \\ a_{21} & a_{22} & \cdots & a_{2,10} \\ \vdots & \vdots & \ddots & \vdots \\ & & & a_{10,10} \end{bmatrix}_{20 \times 10} \rightarrow \nabla_A(\beta)$$

So they flatten it  $\rightarrow 200 \times 1$

find derivative

and then update

Treat matrix as a vector (vectorize)

$$\frac{\partial}{\partial \alpha} \sigma(\alpha) = \sigma(\alpha) \sigma'(1)$$

$(x_1, t_1) \quad (x_2, t_2) \quad \dots \quad (x_n, t_n)$

$$y = w_0 + w_1 x^1 + w_2 x^2 + \dots + w_q x^q$$

$$E(w) = \sum_{n=1}^N (t_n - y(x_n, w))^2$$

Squared sum of all weights  $\left\| w \right\|^2 \leq \eta \Rightarrow$

$$J(\theta) = \sum_{n=1}^N (t_n - y(x_n, w))^2 + \frac{\alpha}{2} \|w\|^2$$

0  
0 0  
0 0  
0 0  
0 0  
0 0  
0 0

$$\frac{1}{2} (y_n - y)^2 + \frac{\alpha}{2} \|w\|^2$$

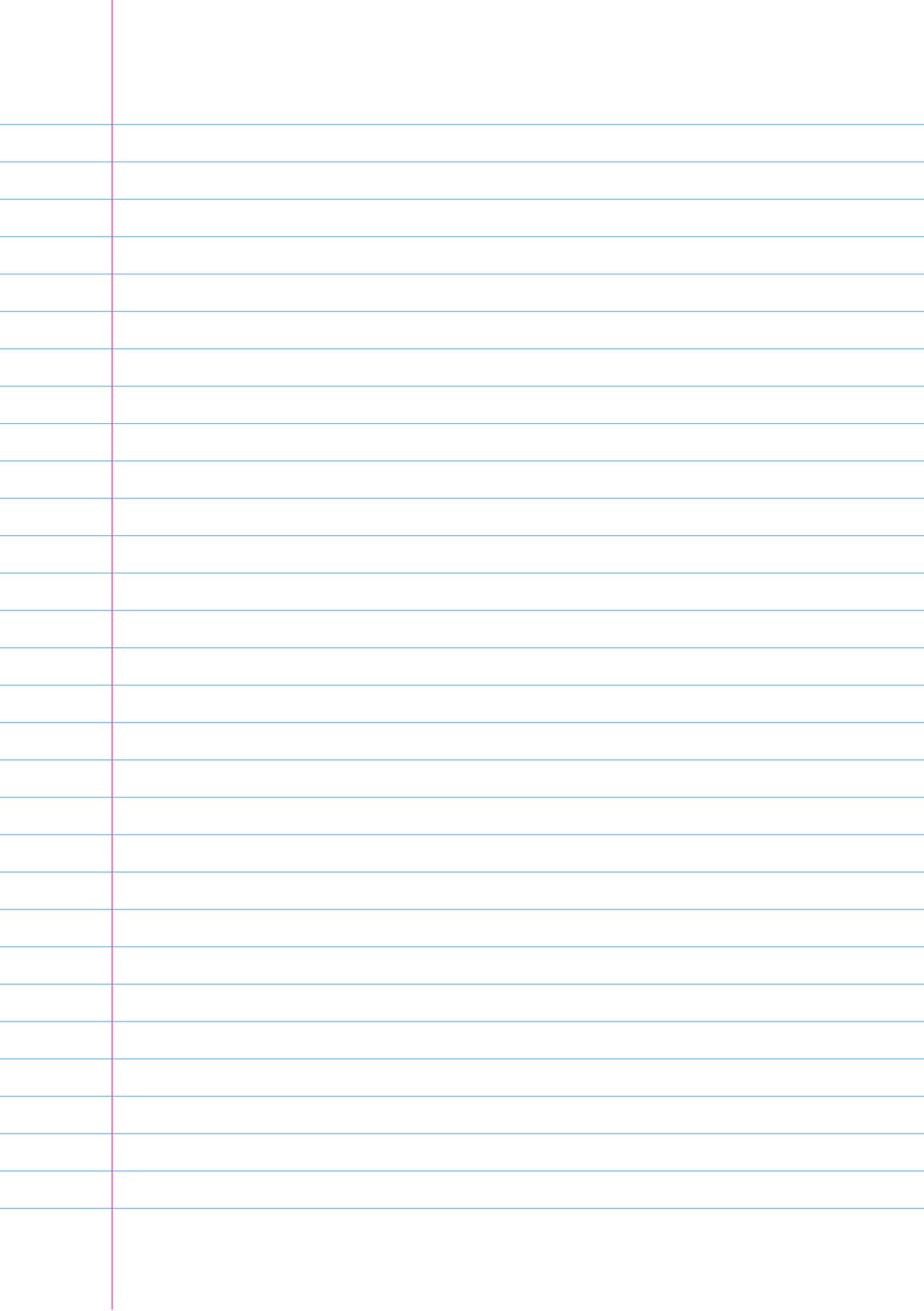
↓  
not biases

bias  $\rightarrow$  extra node with '1' weight

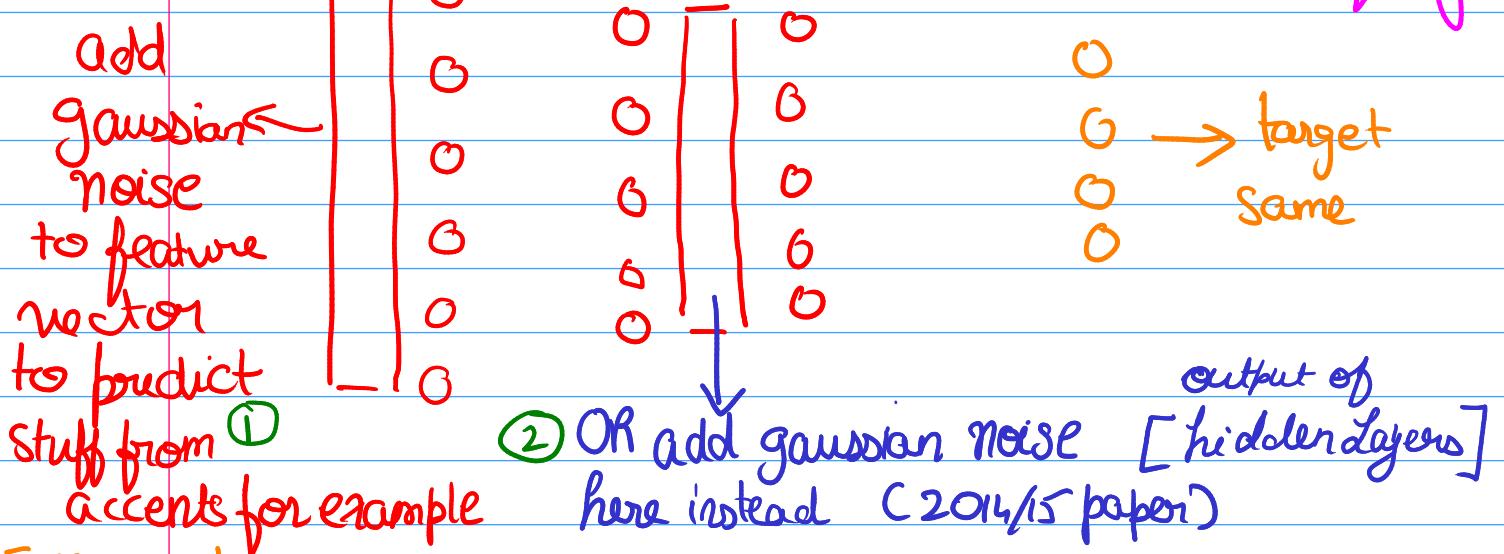
We regularize weights, but not biases  
normally

it is to provide a bias; it's hardwired!

Can be updated  
but rarely is



## Some few other techniques to combat overfitting



[mean at the datapoint]  $\omega^{(k+1)} = \omega^{(k)} - \eta (\nabla E(\omega) |_{\omega=\omega^{(k)}})$  [① & ②]

target can be considered the same

$\Rightarrow$  more datapoints!

$$\omega^{(k+2)} = \omega^{(k+1)} - \begin{bmatrix} n \\ 0 \\ \vdots \\ s_e \end{bmatrix} - \eta (\nabla E(\omega) |_{\omega=\omega^{(k)}}) \quad ③$$

(3) Add noise to the weights themselves

① Seems nice

② Fine

③ Sounds like WTF but works

Another technique

$$(x^{(1)}, t_1) (x^{(2)}, t_2) \dots (x^{(m)}, t_m)$$

$$(0 \quad 1 \quad 0 \ 0 \ 0 \cdots \ 0)$$

Crossentropy  $p^{(0)} p^1 p^2 p^3 p^4 p^5 p^6$

Estimated probabilities functions of w & other params

We want to maximise probability  
⇒ we minimise error

but then if 010000 is wrong  
we can't acc for it

(scope for  
error  
in target  
attribute)

if we acc for the probability of error to be  $\epsilon$

also hyperparameter

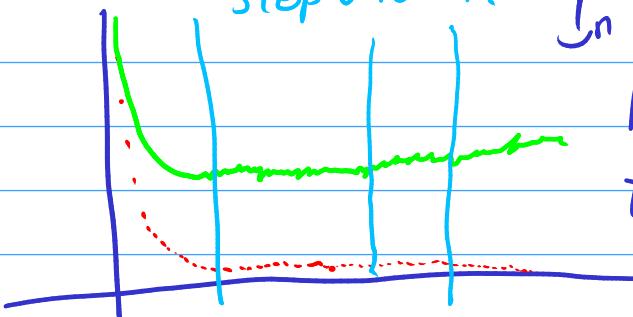
$p_1 \frac{\epsilon}{5} p_2 1-\epsilon p_3 \frac{\epsilon}{5} p_4 \frac{\epsilon}{5} p_5 \frac{\epsilon}{5} p_6 \frac{\epsilon}{5}$  → Softens the error function

This is our new crossentropy

from this we get an error

## Early stopping

Stop where??



In a way where computational power is reduced and fit is optimal

$$\text{SGD} \quad \underbrace{\theta_{k+1}}_{\substack{\text{find training} \\ \text{validation error}}} = \theta_k - \eta (\nabla E(\theta))|_{\theta=\theta_k}$$

$$\theta^* : \theta_{100}$$

$\theta_{200}$  has lower validation error

$$\Rightarrow \theta^* : \theta_{200}$$

$\theta_{300}$  has higher validation error

$$\theta^* : \theta_{200} \text{ (Don't replace)}$$

$\theta_{400}$  also has higher than  $\theta_{200}$

→ 100 epochs  
This is 100

is called

"patience" Dark red  
refer to observe worsening of validation error

Sliding window over epochs

keep a counter everytime

reset if error ↓ otherwise ↑ 1

if counter = patience

R↑

not necessarily increase in worsening

worsening → worse than  $\theta^*$

Let  $n$  be the number of steps between evaluations.

Let  $p$  be the "patience," the number of times to observe worsening validation set error before giving up.

Let  $\theta_0$  be the initial parameters.

```
 $\theta \leftarrow \theta_0$ 
 $i \leftarrow 0$ 
 $j \leftarrow 0$ 
 $v \leftarrow \infty$ 
 $\theta^* \leftarrow \theta$ 
 $i^* \leftarrow i$ 
while  $j < p$  do
    Update  $\theta$  by running the training algorithm for  $n$  steps.
     $i \leftarrow i + n$  → Update  $n$  times (Iterations)
     $v' \leftarrow \text{ValidationSetError}(\theta)$ 
    if  $v' < v$  then
         $j \leftarrow 0$ 
         $\theta^* \leftarrow \theta$ 
         $i^* \leftarrow i$ 
         $v \leftarrow v'$ 
    else
         $j \leftarrow j + 1$ 
    end if
end while
```

Best parameters are  $\theta^*$ , best number of training steps is  $i^*$

Algorithm determines the best amount of time to train. The meta algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.

Let  $n$  be the number of steps between evaluations.

Let  $p$  be the "patience," the number of times to observe worsening validation set error before giving up.

Let  $\theta_0$  be the initial parameters.

$$\theta \leftarrow \theta_0$$

$$i \leftarrow 0$$

$$j \leftarrow 0$$

$$v \leftarrow \infty$$

$$\theta^* \leftarrow \theta$$

$$i^* \leftarrow i$$

while  $j < p$  do

    Update  $\theta$  by running the training algorithm for  $n$  steps.

$$i \leftarrow i + n$$

$$v' \leftarrow \text{ValidationSetError}(\theta)$$

    if  $v' < v$  then

$$j \leftarrow 0$$

$$\theta^* \leftarrow \theta$$

$$i^* \leftarrow i$$

$$v \leftarrow v'$$

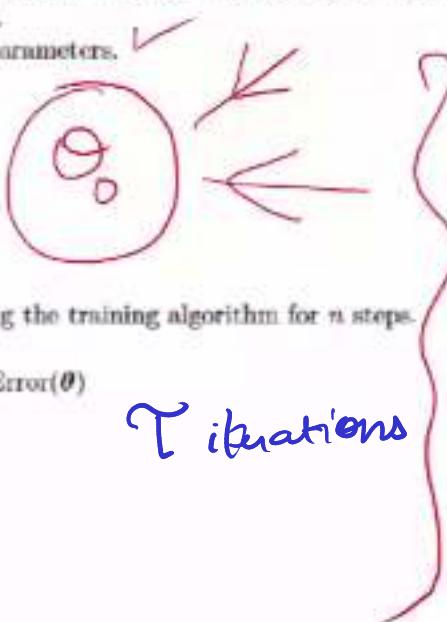
    else

$$j \leftarrow j + 1$$

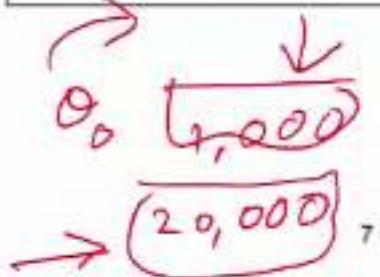
    end if

end while

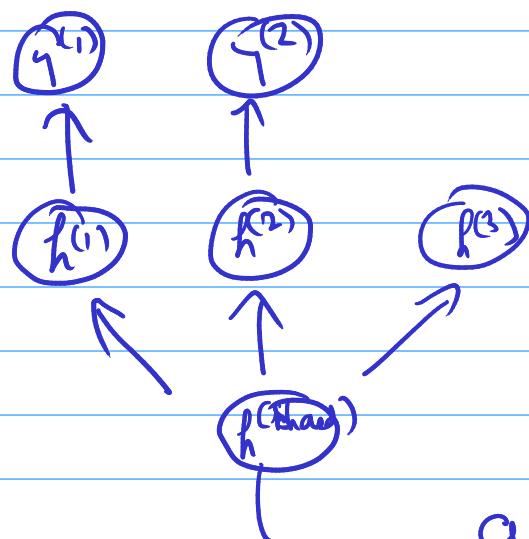
Best parameters are  $\theta^*$ , best number of training steps is  $i^*$



Algorithm determines the best amount of time to train. The meta algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.



Note: SGD does implicit regularization  
(monides)



no. of steps  
for early  
stopping is  
a hyperparameter

Find this by training  
a small subset, then  
use the number to train the whole model  
(Shut up & TRAINNN)

## Flavor #2

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.  
Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (algorithm 7.1) starting from random  $\theta$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This updates  $\theta$ .

$\epsilon \leftarrow J(\theta, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$

while  $J(\theta, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$  do

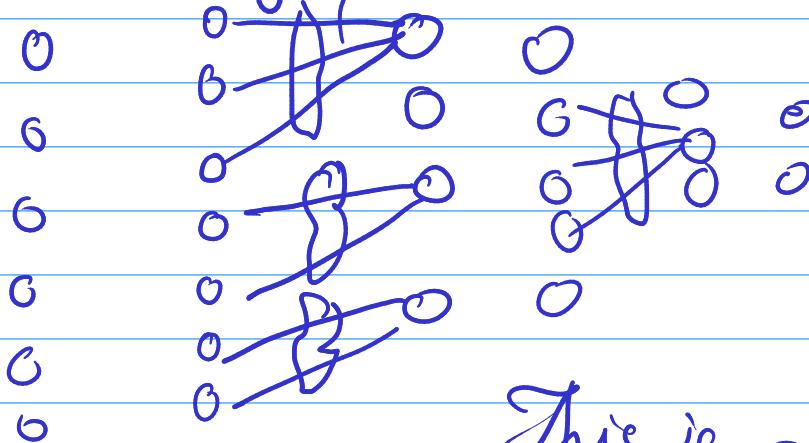
    Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $n$  steps.

end while

$\theta^*$  for this  
don't throw  $\theta^*$  away

70%

Weight sharing - same



(adversarial)

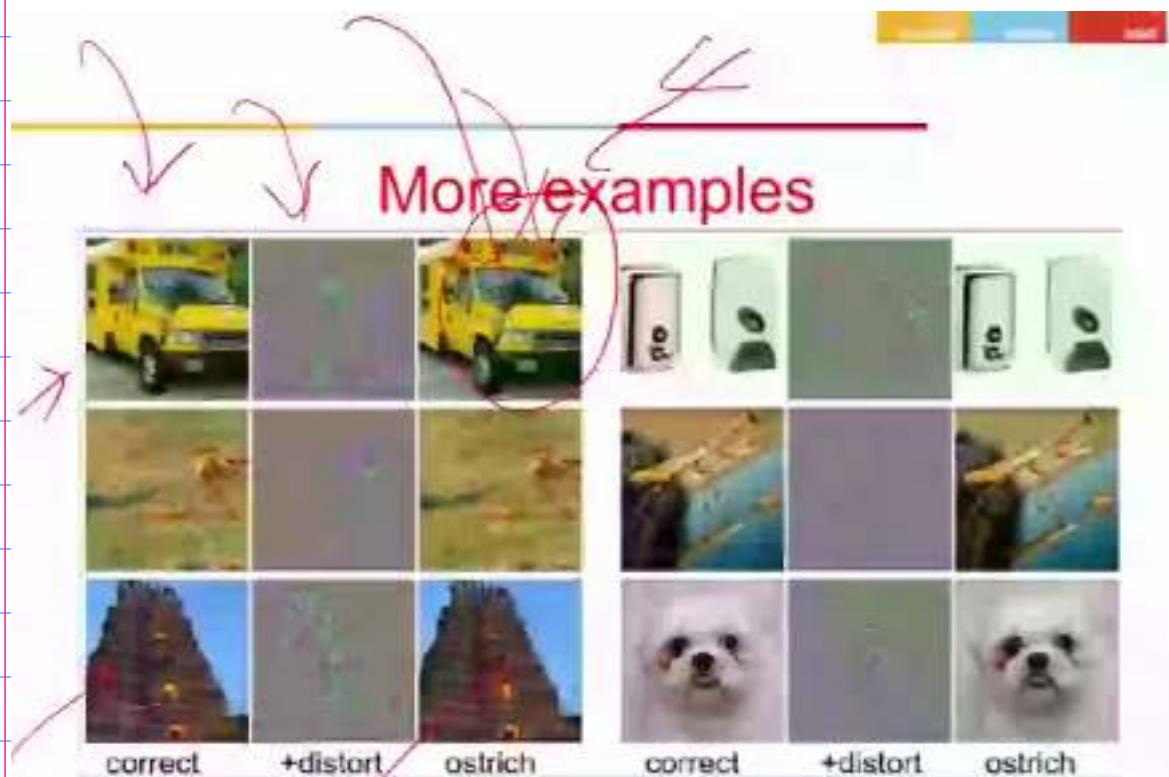
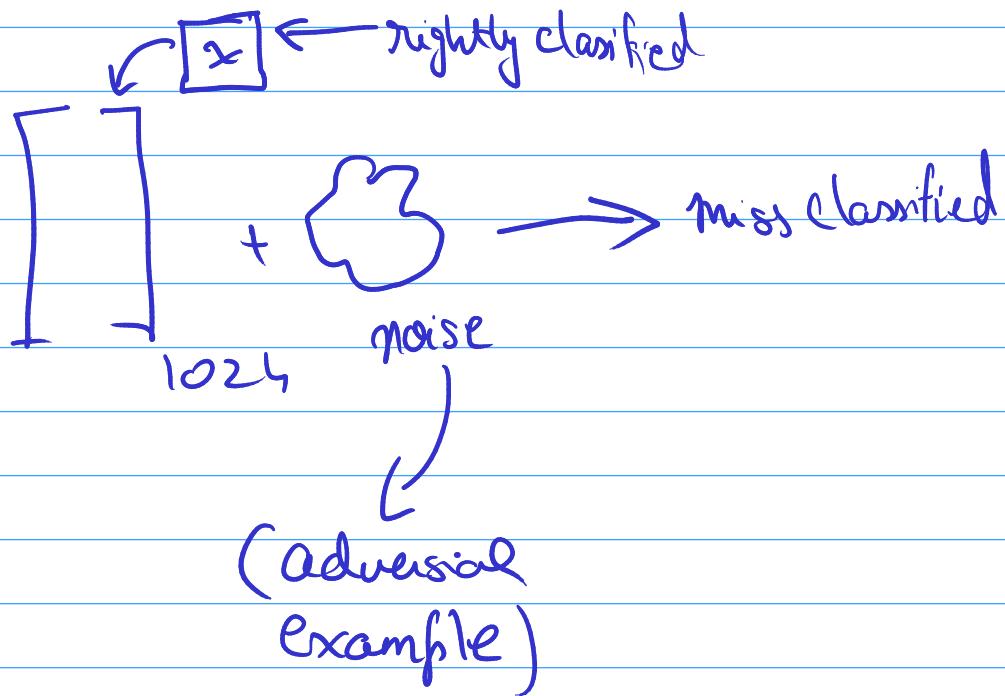
This is also a regularizer,  
overfitting doesn't happen

Adversarial Training

Szegedy

2014

$\times$  ↪ slightly classified



goodfellow 2014, explaining adversarial -- a composition of functions.  
 But relu type functions composed can be incredibly linear most of the time

a heavily linear function will look like a lin. comb. of weights

a close neighbourhood can be  $W \cdot (x + \epsilon)$  ( $x, w$  and  $\epsilon$  are vectors)

$$f(x + \epsilon) \approx f(x) + \frac{\partial f}{\partial x}(\epsilon)$$

$$\frac{1}{x + \epsilon} \approx \frac{1}{x}$$

if  $f$  is quadratic,  $\epsilon$  will become  $\epsilon^2$  etc.  
 and vanish

Every layer

$$\left[ \begin{matrix} w_1 \\ w_2 \\ \vdots \\ w_{30000} \end{matrix} \right] \left( \left[ \begin{matrix} x_1 \\ \vdots \\ x_{30000} \end{matrix} \right] + \left[ \begin{matrix} \epsilon_1 \\ \vdots \\ \epsilon_{30000} \end{matrix} \right] \right)$$

$\epsilon$  is non linear

$\Rightarrow x$  becomes highly non linear

$$f \left( \left[ \begin{matrix} x_{11} \\ x_{12} \\ \vdots \\ x_{100} \end{matrix} \right] + \left[ \begin{matrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_{100} \end{matrix} \right] \right) \text{ changes quite a lot}$$

$$f(x) - f(x+\epsilon)$$

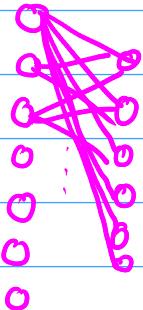
apply 'f' on this

$$g = f_0 \circ f_1 \circ f_2 \dots f_n(x)$$

$g(y) - g(y+\epsilon)$  will be way more than

$$f(x) - f(x+\epsilon)$$

$$f = \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_{100} x_{100}$$



$$f \approx \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_{100} x_{100}$$

Adding noise

$$\omega^T x + \omega^T \delta$$

$$\omega^T x + \omega^T \begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_{100} \end{bmatrix}$$

the quantity

$\det$   
 $\delta_i$  be  $+\epsilon$   
 if  $\omega \geq 0$   
 $\delta_i$  be  $-\epsilon$   
 if  $\omega < 0$

Whitebox attack

$$x = x + \delta$$

$$\delta = \epsilon \begin{bmatrix} \text{sign}(\omega_1) \\ \vdots \\ \text{sign}(\omega_{100}) \end{bmatrix}$$

$$\omega^T x + \epsilon (\|\omega\|_1) = \omega^T x + \epsilon (|\omega_1| + |\omega_2| + \dots + |\omega_{100}|)$$

hyperplane

$$x \rightarrow (x + \delta)$$

$$f(x) \rightarrow f(x + \delta)$$

even

Since  $f$  is linear, if  $x$  is close to  $x_0$

$f(x')$  can be completely different

BUT Why?? (This seems contradictory)  
 You can't work this lol

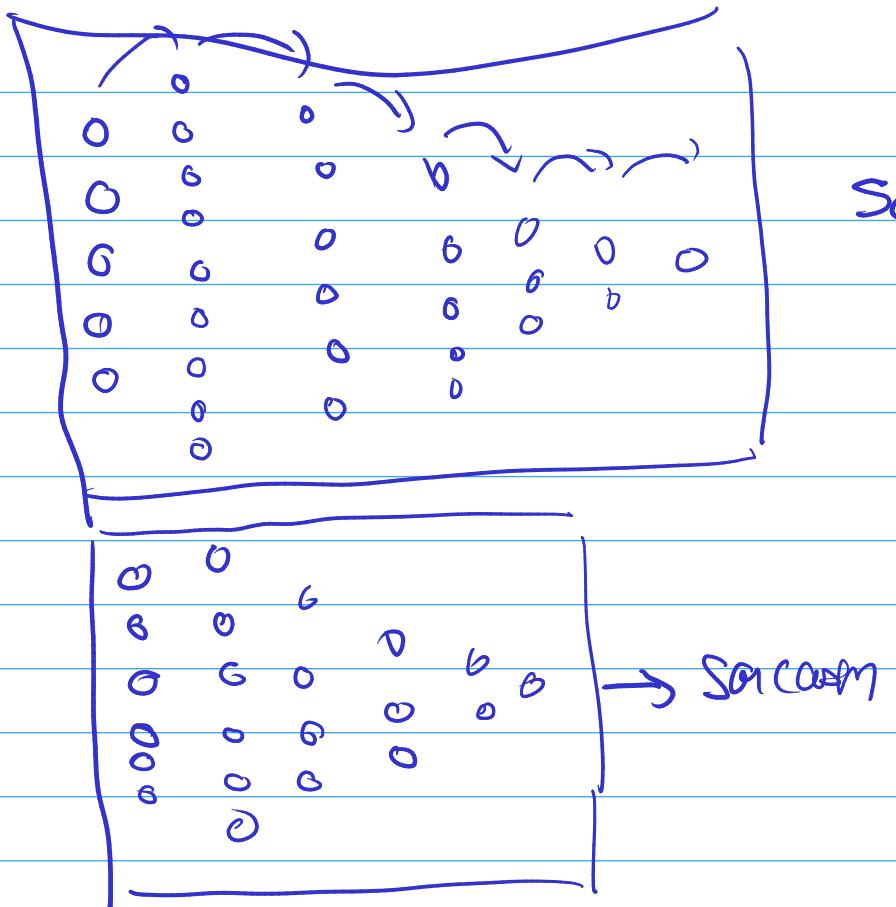
Basically, this is an explanation, not a proof  
 for adversarial examples

It works for not all examples,  $f$  is normally  
piecewise linear

(we are also NOT comparing with polynomial & exponential; we are just saying  
linearity itself makes values grow!  
Compared to sigmoid etc.)

SVMs are also linear but have large decision regions, so they don't face this problem.

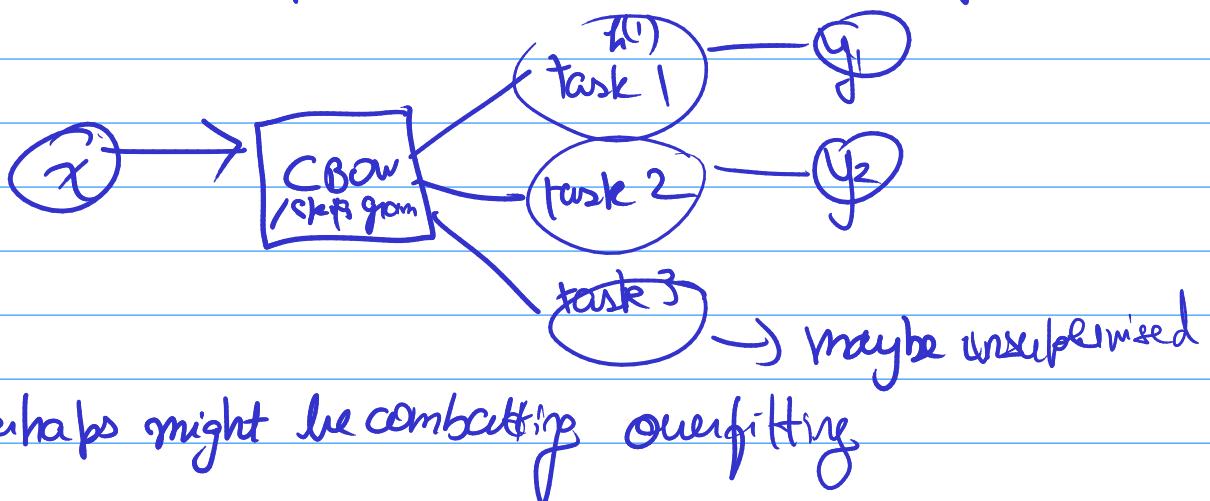
⇒ Training on adversarial examples makes model more robust



Whether it is a sentiment analysis or sarcasm predictor, to get a very nice numerical equivalent that is some

Since otherwise, you might overfit, since the word vector representations may be prejudiced!

DO make a wordembedding Commonly  
and then split



$$\min J(\theta; x, y) + \lambda g(\theta)$$

$$\underbrace{\quad}_{\theta = (\theta_1, \theta_2, \dots, \theta_n)}$$

$$g(\theta) = \|\theta\|^2$$

$$g(\theta) = \|\theta\|_1$$

Now rather than directly using weights, we use some other function

$$+ \lambda \Omega(h)$$

Some functions have <sup>been</sup> proposed

$$c_1 \ c_2 \ - \ - \ - \ c_{20}$$

20 classifiers

$$r_1 \ r_2 \ - \ - \ - \ r_{20}$$

Let  $\epsilon_1 \ \epsilon_2 \ \dots \ \epsilon_k$   
errors

Let them follow a normal distribution  $N(\mu, \sigma^2)$

↗ classifications

Ensembler       $C_1 \ C_2 \ \dots \ C_k$   
                   \underbrace{\hspace{10em}}\_{\text{majority vote}}

→ improves your testing error since majority of models aren't wrong.

→ Bagging ensembler

↳ Model is just the same eg: Decision tree

Consider a sample set       $T_1 \ T_2 \ \dots \ T_{1000}$

$C_1 \rightarrow$  sample with replacement       $T_1 \ T_{500} \ T_1 \ T_{10} \ T_{625} \ T_{323} \ \dots$

$C_2$  similarly

:

$C_k$  (it's also bagging ensembler)

Or

develop  $k$  regression models (varying them)

$M_1 \ M_2 \ \dots \ M_k$

Errors  $\rightarrow E_1 \ E_2 \ \dots \ E_k$

$$\frac{\sum E_i}{k}$$

Let these  $E_i$ 's follows  $N((0, 0, 0, \dots, 0), \text{Cov}(E_1, E_2, E_3, \dots, E_k))$

$$E(E_i) = 0$$

$$\begin{aligned} \text{Var}(E_i) &= E(E_i - E(E_i))^2 \\ &= E(E_i^2) = \sigma^2 \end{aligned}$$

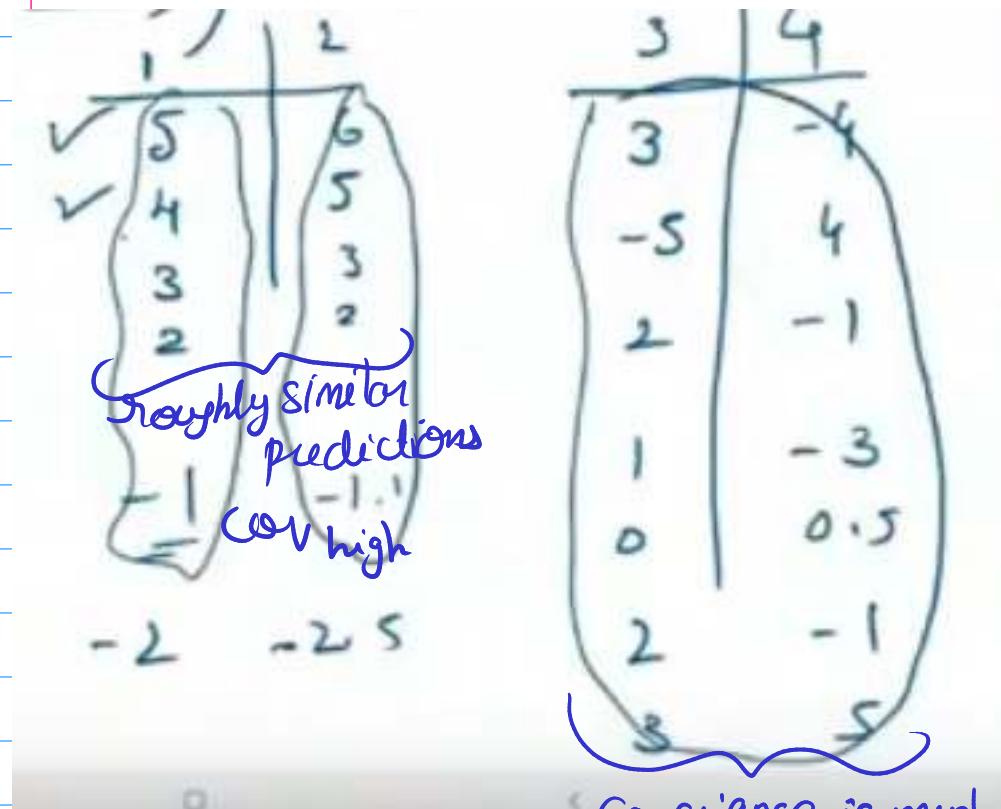
An assumption

$$\text{Var}(E_1) = \text{Var}(E_2) = \dots = \text{Var}(E_K) \\ = \sigma^2$$

$$\text{Cov}(x, y) = E((x - E(x))(y - E(y)))$$

$$\text{Cov}(E_1, E_2)$$

Now the variation of  $x$  from  $E(x)$   
is varying wrt to variability  
variation of  $y$  from  $E(y)$



$$\begin{aligned} \text{Cov}(E_1, E_2) &= E(E_1 - E(E_1))(E_2 - E(E_2)) \\ &= E(E_1 E_2) \\ &= C \end{aligned}$$

We assume all covariances are same

$$\epsilon_1 \quad \epsilon_2 \quad \dots \quad \epsilon_k$$

$$\begin{matrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_k \end{matrix} \left[ \begin{matrix} v & c & c & c & c & c \\ c & v & c & c & c & c \\ c & c & v & c & c & c \\ c & c & c & v & c & c \\ c & c & c & c & v & c \\ - & - & - & - & - & - \end{matrix} \right]$$

Errors  $\rightarrow \epsilon_1, \epsilon_2, \dots, \epsilon_k$

$$\frac{1}{k} \sum_{i=1}^k \epsilon_i$$

$$\frac{1}{k} \sum_{i=1}^k \epsilon_i \quad \dots \quad \frac{1}{k} \sum_{i=1}^k \epsilon_i$$

(1st sample)

Errors: Can you average these?

NO  $E\left(\frac{1}{k} \sum_{i=1}^k \epsilon_i\right)$

Concancels

$$E\left(\left(\frac{1}{k} \sum_{i=1}^k \epsilon_i\right)^2\right)$$

$$= \frac{1}{k^2} E\left(\left(\sum_{i=1}^k \epsilon_i\right)^2\right)$$

$$= \frac{1}{k^2} E\left(\sum_{i=1}^k \epsilon_i^2 + \sum_i \sum_j \epsilon_i \epsilon_j\right)$$

$$= \frac{1}{k^2} E\left(\sum_{i=1}^k \epsilon_i^2\right) + \frac{2}{k^2} \sum_i E(\epsilon_i \epsilon_j)$$

But then isn't  
error squared  
already  
(for argument's  
sake, let's assume  
they're not)

$$= \frac{1}{k^2} \left[ E \left( \sum_{i=1}^k \epsilon_i^2 \right) + \sum_i \sum_j c \right]$$

$$= \frac{1}{k^2} \left[ \sum_{i=1}^k E(\epsilon_i^2) + \sum_i \sum_j c \right]$$

$$= \frac{1}{k^2} [k \sigma^2 + k(k-1)c]$$

$$= \frac{1}{k} \sigma^2 + \frac{k-1}{k} c$$

$$c = 0$$

if  $M_1, M_2, \dots, M_k$  all behave  
differently

$$\text{then } E = \frac{1}{k} \sigma^2$$

$M_1$  every error has variance  $\sigma^2$  (which is  $E(\epsilon_i^2)$ )  
 $M_2$   
 $\vdots$   
 $M_k$

But the ensembles together have a variance of  $\frac{1}{k}$

But  $V$  is a function of  $E$   
 $\Rightarrow$  even  $V$  is reduced

$\Rightarrow$  If all models are independent,  $V$  comes down by  
a factor of  $k$   
 $\Rightarrow$  impressive

If  $c \neq 0$

$$C = E(\epsilon_i \epsilon_j) \approx E(\epsilon_i^2) = \sigma^2$$

$\overbrace{\quad\quad\quad}^{\text{if } \epsilon_i = \epsilon_j}$

If all ensembles are same then you get

Your common variance to be:  $\frac{1}{k} \sigma^2 + \frac{k-1}{k} \sigma^2$

$$\Rightarrow \sigma^2 \geq \sigma^2 \geq \frac{\sigma^2}{k}$$

$\Rightarrow$  Worst case Ensembles perform as good as individual models

Similarly  $C_1 \leq C_2 \dots \leq C_k$

$E_{\text{Ensemble}}(C_1, C_2, \dots, C_k)$

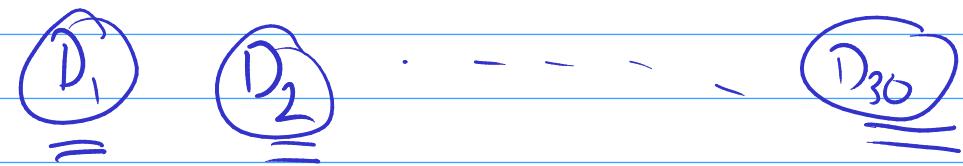
will perform better, baber mat nikal bc

\* It enhances the performance of the model

→ acts as a "kind of" regularizer

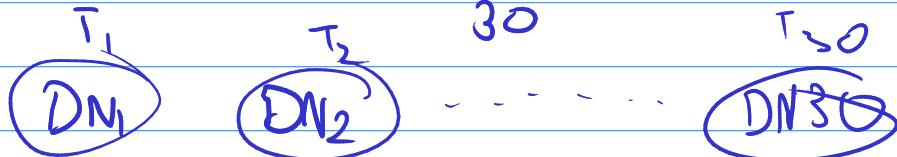
without being a regularizer in the strict sense.

30 models



Test example  
 $P_{30}$

$$\underbrace{P_1 + P_2 + \dots + P_{30}}_{30}$$

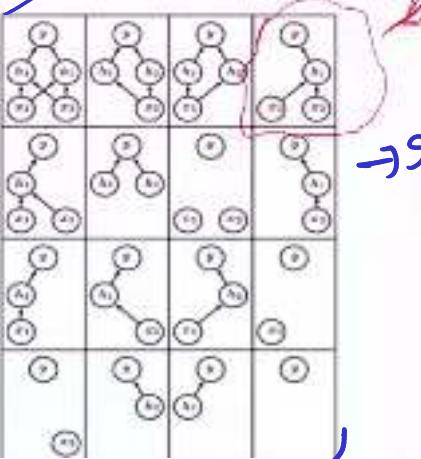
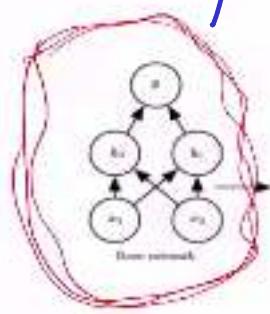


6 classifiers  
train

Szegedy  
ILSVRC

An Ensembler built  
Bagging

exponential  
number of models



Srivastava → Alter the structure  
and build an ensembler  
out of it! #gebrekt

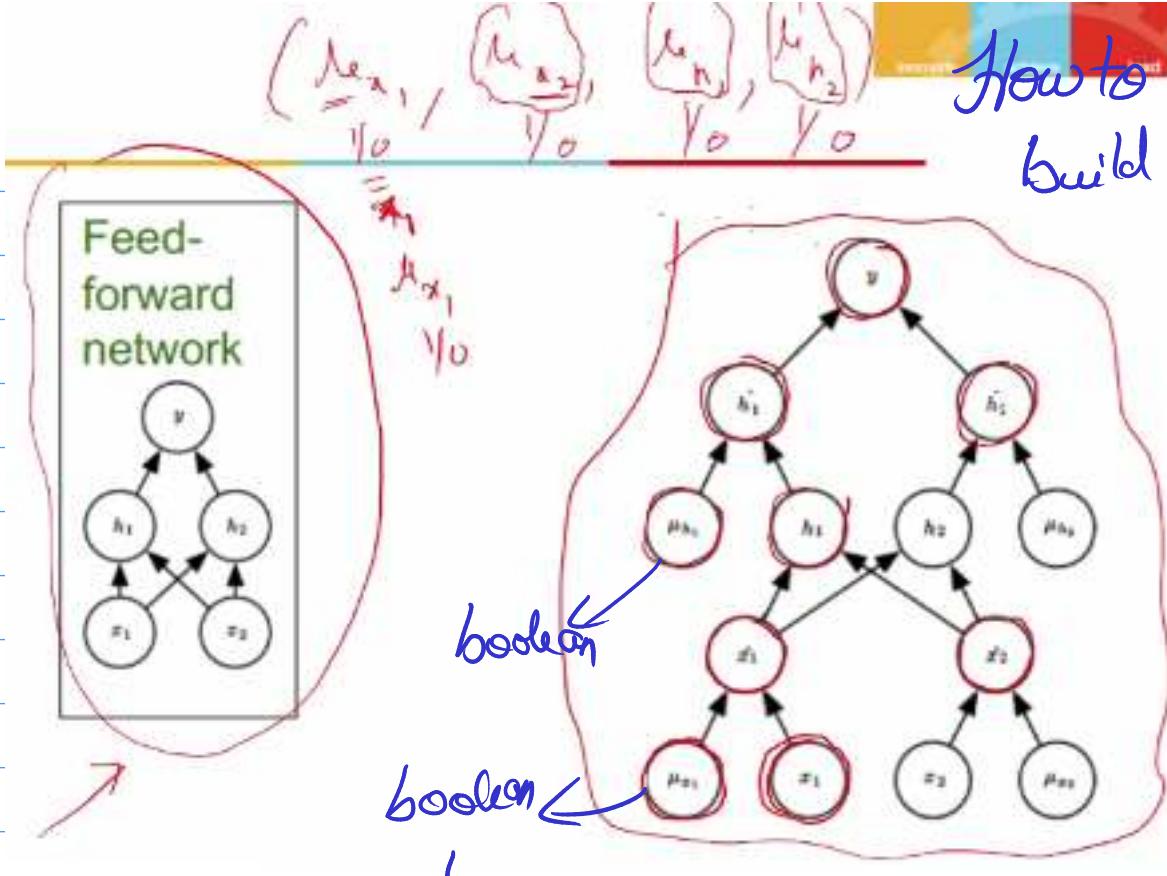
→ remove nodes and/or connections

instead of same model different  
splits, diff models (with same  
backbone) and some split

But in this train all these  
one in parallel, and use their  
aggregate → ensembler

Apparently this is Dropout!

How to  
build dropout?



→ vary this vector.  
it acts as a regularizer  
and we want look into it next class by

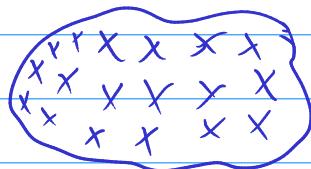
# Finding the gaps in ML/DL

(min 0-1 loss) < discrete function  
integer optimisation  
D-1 (not even integer)  
optimisation problem

very difficult to solve

so they want the surrogate optimisation problems!!

$$\left[ E_{(x,y)} \left( L(f(x;\theta), y) \right) \right] \\ (x,y) \sim P(\text{data})$$



\* Let's assume that the dataset, and testing one coming from the same probability distribution



→ We are never minimizing expectation before

$$\iint (p(x,y)) L(f(x;\theta), y) dx dy$$

) We are never doing this ( we're never talking about gradient descent )

Instead, we are doing

$$\min \frac{1}{n} \sum_{i=1}^n L(f(x_i; \theta), y)$$

$$\min_{(x,y) \sim P_{\text{data}}} E(L(f(x; \theta), y)) \quad (\text{not a KL divergence})$$

$$\frac{\partial E}{\partial \theta_k} = \frac{1}{n} \sum_{i=1}^n \frac{\partial E_i}{\partial \theta_k} \quad n = 10,000,000,000$$

Zumeist

→ SGD / minibatch



{we estimate  $\frac{\partial E}{\partial \theta_k}$  by  $\frac{\partial E_{100}}{\partial \theta_k}$ }

or

$$\text{Estimate } \frac{\partial}{\partial \theta} (E_2 + E_5 + \dots + E_{75} + E_{1000})$$

minibatch gradient descent

batch gradient descent  
typically means classical

is bigger batch size, better estimate? (if yes)

minibatch, stochastic, (online)

↓  
used "most recent" stuff (samples)  
in timestamp or sth for example

## Algorithm: SGD update at training iteration $k$

Require: Learning rate  $\epsilon$

Require: Initial parameter  $\theta$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute gradient estimate:  $\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

    Apply update:  $\theta \leftarrow \theta - \epsilon \hat{g}$

end while

$$\theta \leftarrow \theta - \epsilon (\hat{g}) \quad , \text{million } m$$

A weight with larger derivative will have larger updates  
 Actually we see the opposite

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \epsilon (\nabla_{\theta_i} (\sum_i L(f(x_i^{(i)}; \theta), y)))$$

The learning parameter should not be constant for all iterations  
 (reason is out of scope)  
 Also, from parameter to parameter

## Adagrad

(Let's do this before Adam)

$\epsilon \rightarrow$  global

Initial  $\theta$

$\delta = 10^{-7}$

$r = 0$  (accumulation of gradient)

while

minibatch  $m$

$$g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i, \theta), y^{(i)})$$

$$r \leftarrow r + g \odot g \quad \text{element wise product}$$

$$\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot g \quad \text{but } g \text{ is a vector}$$

$$\theta \leftarrow \theta + \Delta \theta$$

$$\frac{\epsilon}{\delta + \sqrt{r}}$$

$L_1, L_2$   
Dropout  
Early stopping  
Data augment  
Multitasking  
Adversarial training

$$\frac{\epsilon}{\delta + \sqrt{r}} = \begin{bmatrix} \frac{\epsilon}{\delta + \sqrt{r_1}} \\ \frac{\epsilon}{\delta + \sqrt{r_2}} \\ \frac{\epsilon}{\delta + \sqrt{r_3}} \\ \vdots \end{bmatrix}$$

where  $r =$

$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ \vdots \\ r_n \end{bmatrix}$$

# RMSProp

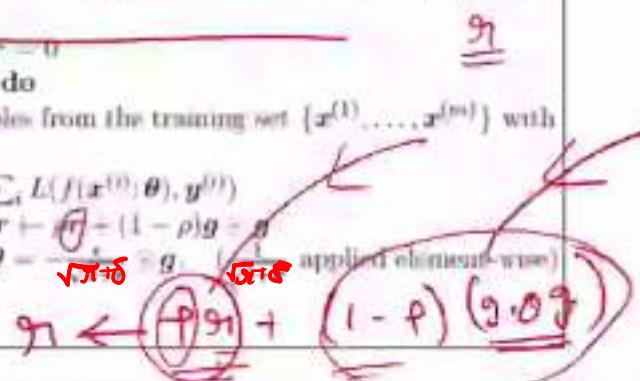
- Modifies AdaGrad for a nonconvex setting  
– Change gradient accumulation into exponentially weighted moving average  
– Converges rapidly when applied to convex function

Craig  
Hinton

## The RMSProp Algorithm

```
Require: Global learning rate  $\epsilon$ , decay rate  $\rho$ .  
Require: Initial parameter  $\theta$   
Require: Small constant  $\delta$  usually  $10^{-6}$ , used to stabilize division by small numbers.  
Initialize accumulation variables  $r = 0$   
while stopping criterion not met do  
    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .  
    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$   
    Accumulate squared gradient:  $r \leftarrow \rho r + (1 - \rho)g^2$   
    Compute parameter update:  $\Delta\theta = \frac{\sqrt{r}}{\epsilon} \cdot g$  (applied element-wise)  
    Apply update:  $\theta \leftarrow \theta + \Delta\theta$   
end while
```

0.6  
=



MIT Press Heimann et al.

## Algorithm: SGD with momentum

```
Require: Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .  
Require: Initial parameter  $\theta$ , initial velocity  $v$ .
```

while stopping criterion not met do

Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

Compute gradient estimate:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

Compute velocity update:  $v \leftarrow \alpha v - \epsilon g$

Apply update:  $\theta \leftarrow \theta + v$

end while

→ We aren't comparing Adagrad with momentum

SGD+

→ SGD + momentum may perform better than adagrad at times

But it is better than SGD!

## Algorithm: SGD with Nesterov momentum

Require: Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

Require: Initial parameter  $\theta$ , initial velocity  $v$ .

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding labels  $y^{(i)}$ .

    Apply interim update:  $\bar{\theta} \leftarrow \theta + \alpha v$  This line is added from plain momentum

    Compute gradient (at interim point):  $g \leftarrow \frac{1}{m} \nabla_{\bar{\theta}} \sum_i L(f(x^{(i)}; \bar{\theta}), y^{(i)})$

    Compute velocity update:  $v \leftarrow \alpha v - \epsilon g$  This line is added from plain momentum

    Apply update:  $\theta \leftarrow \theta + v$

end while

## Algorithm: RMSProp with Nesterov momentum

Require: Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$ .

Require: Initial parameter  $\theta$ , initial velocity  $v$ .

Initialize accumulation variable  $r = 0$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

    Compute interim update:  $\bar{\theta} \leftarrow \theta + \alpha v$

    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\bar{\theta}} \sum_i L(f(x^{(i)}; \bar{\theta}), y^{(i)})$

    Accumulate gradient:  $r \leftarrow \rho r + (1 - \rho)g \odot g$

    Compute velocity update:  $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$ . ( $\frac{1}{\sqrt{r}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + v$

end while

Das has left the meeting

$$v \leftarrow \alpha v - \left( \frac{\epsilon}{\sqrt{r}} \odot g \right)$$

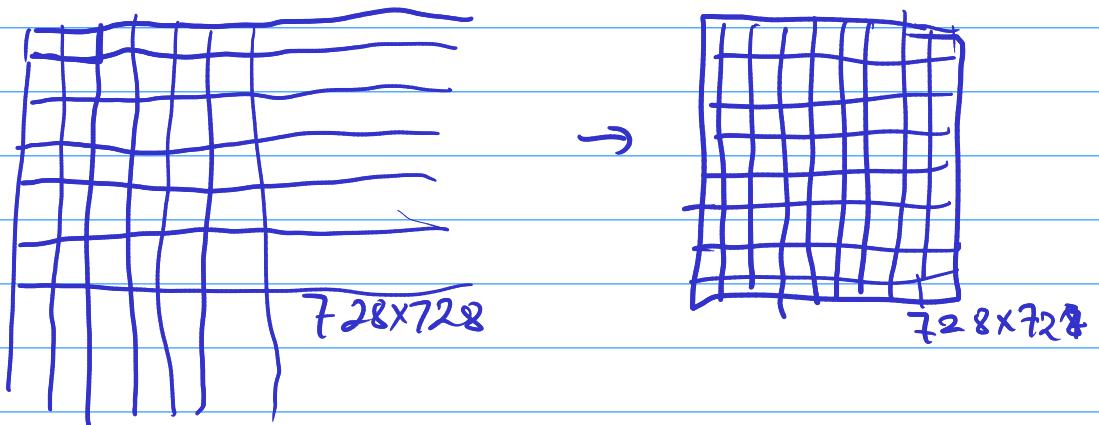
# The Adam Algorithm

Require: Step size  $\epsilon$  (Suggested default: 0.001)  
Require: Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1]$ .  
(Suggested defaults: 0.9 and 0.999 respectively)  
Require: Small constant  $\delta$  used for numerical stabilization. (Suggested default:  
 $10^{-8}$ )  
Require: Initial parameters  $\theta$   
Initialize 1st and 2nd moment variables  $s = 0, r = 0$   
Initialize time step  $t = 0$   
while stopping criterion not met do  
    Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with  
    corresponding targets  $y^{(i)}$ .  
    Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$   
     $t \leftarrow t + 1$   
    Update biased first moment estimate:  $s \leftarrow \rho_1 s + (1 - \rho_1)g$   
    Update biased second moment estimate:  $r \leftarrow \rho_2 r + (1 - \rho_2)g \odot g$  (CS705)  
    Correct bias in first moment:  $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$   
    Correct bias in second moment:  $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$  } to bring the moments as  
    expected  
    Compute update:  $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$  (operations applied element-wise)  
    Apply update:  $\theta \leftarrow \theta + \Delta\theta$  (They are lower)  
end while

\* Influential  
but sparse  
features  
give lower  
updates

specific way of updating  
since  $s=t=0$

256



in an ANN

if we have an i/p matrix of <sup>large</sup>  $200 \times 200$   
and a hidden layer of  
 $200 \times 200$

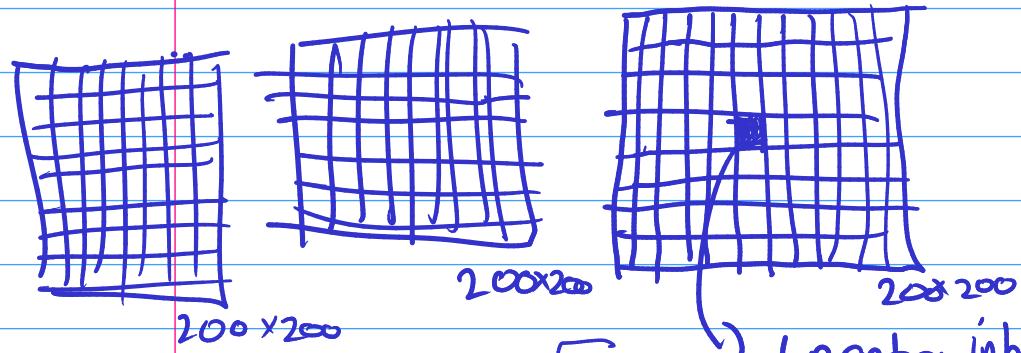
Weight matrix:  $(200 \times 200) \times (200 \times 200)$

for a fully connected network.

Since we're accounting for the dependency of every pixel

But

When we're discussing something like edge detection  
we don't need anything aside from the neighbouring pixels



→ we get an input here  
for that, we need to multiply  
the  $200 \times 200 \Rightarrow 40000 \times 1$ , with  
let's say the column of the said pixel  
So we just need the neighbouring points to  
have some weight & rest all zero

Since there are a LOT of weights that are zero, we  
don't need so much  
So we can just

## ID case

NOTE:

This is not predicting future.  
we are trying to reduce noise  
enhance

### Convolution

Let's say we want to know where a spaceship is given that we know some info that is noisy

(9.28.37)      (9.28.38)      (9.28.39)  
we have these readings too  
(lets say microseconds)

(9.29) → to find lets consider that the only measurement device is noisy

To find the true current position

Now lets find a weighted average of all the previous 60 seconds

9.29

So we have  $x(a)$ ,  $w(t-a)$

@  $t = 9.29$   
 $a = 9.28.37$

$$D(t) = \int x(a) w(t-a) da$$

→ time is 't'

→  $t-a$  is different  
for all 'a'.

Your weight  $w(t-1)$  will be the longest and then  $w(t-2)$ ,  $w(t-3)$  etc...

so you can treat  $w$  as a probability distribution

in order to get an "average"

$$S(t) = (x * w) t$$

$\curvearrowleft$   
x convolved

$$(x * w)(t) = \int_{t_1}^{t_2} x(a) w(t-a) da$$

this is very similar to a convolution formula

$x(a)$  will be some kind of function which we don't know if continuous

So since measurements are discrete

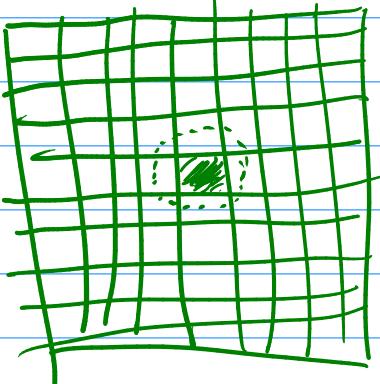
we get  $\sum_{a=6_1}^{t_2} x(a) w(t-a)$

$$= \sum_{a=9, 28, 51}^{9, 29} x(a) w(t-a)$$

this way we can reduce the effect of noise

We are not predicting future!

### 2D case



→ for edge detection

Some take a linear combination of nearby pixels

→ here, the linear combination expression will be the same for all pixels

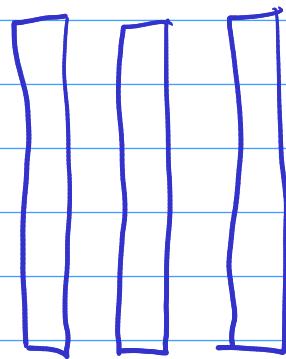
$$S(t) = (x * \omega)(t) = \sum_a x(a) w(t-a)$$

$$S(i,j) = (I * K)(i,j) = \sum_m \sum_n I(m,n) k(i-m, j-n)$$

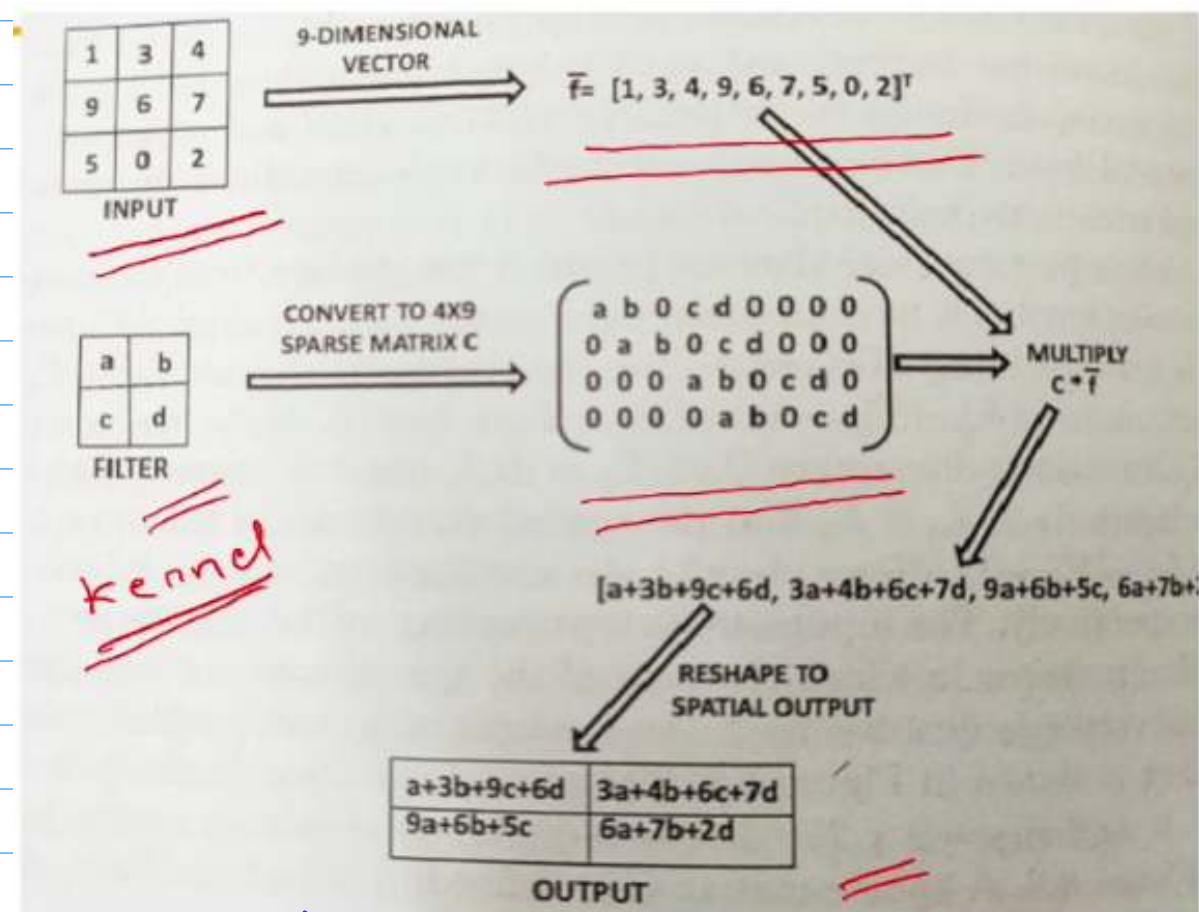
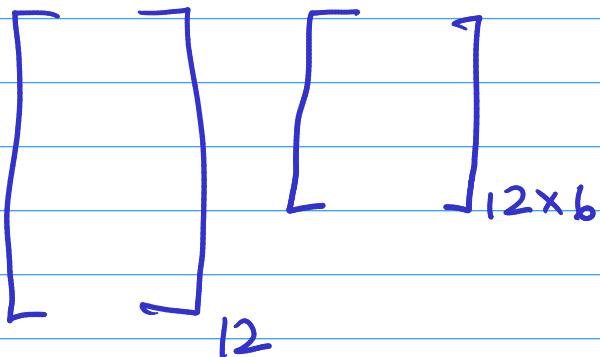
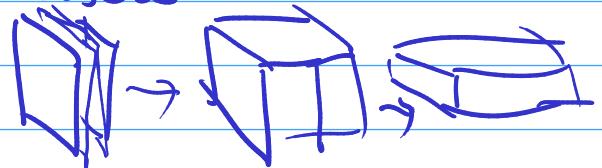
$$S(i,j) = (I * K)(i,j) = \sum_m \sum_n I(i+m, j+n) k(m,n)$$



This is more a convolution w.r.t &



You can give higher dimensional input as tensors.



during backprop, mark the weights that are zero

Some weights have to be the same

So we minimize loss function

$$\min(L) \text{ subject to}$$

some weights being  
same

This is a very complex optimization problem (?)

So let's say  $\omega_{11} = \omega_{22} = \omega_{33} = \omega_{44} = \omega_1$

$$\frac{\partial L}{\partial \omega_1} = \sum \left( \frac{\partial L}{\partial \omega_{11}} \cdot \frac{d\omega_{11}}{d\omega_1} \right)$$

$$= \frac{\partial L}{\partial \omega_{11}} \cdot \frac{\partial \omega_{11}}{\partial \omega_1} + \frac{\partial L}{\partial \omega_{22}} \cdot \frac{\partial \omega_{22}}{\partial \omega_1}$$

$$+ \frac{\partial L}{\partial \omega_{33}} \cdot \frac{\partial \omega_{33}}{\partial \omega_1} + \frac{\partial L}{\partial \omega_{44}} \cdot \frac{\partial \omega_{44}}{\partial \omega_1}$$

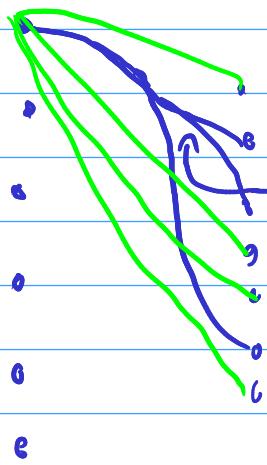
$$\text{but } \omega_{11} = \omega_1 \Rightarrow \frac{\partial \omega_{11}}{\partial \omega_1} = 1$$

$$\Rightarrow \frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial \omega_{11}} + \frac{\partial L}{\partial \omega_{22}} + \frac{\partial L}{\partial \omega_{33}} + \frac{\partial L}{\partial \omega_{44}}$$

$$\omega_1^{(2)} = \omega_1^{(1)} - \eta \frac{\partial L}{\partial \omega_1} \Big|_{\omega_2 = \omega_3 = \omega_4}$$

This way all of them will be updated the same way!  $\hookrightarrow$  MAO

This is the concept of weight sharing



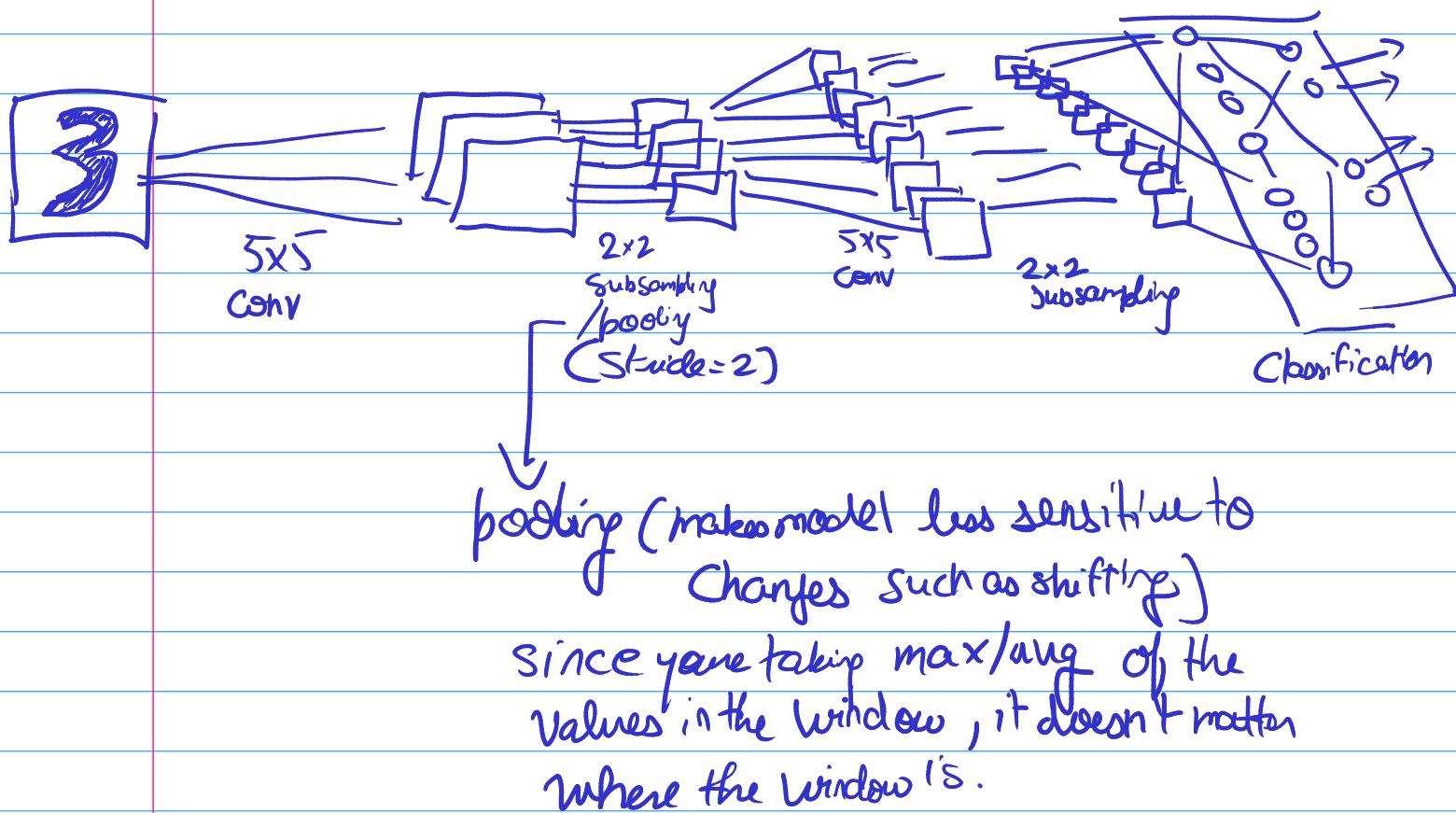
Shared weight  $s$ ,  
they get updated same  
way

→ In this case, by expanding it to a sparse matrix we increase space complexity and reduce time complexity

→ Another way is to just aggregate over a sliding window, which is time intensive but less space  
But internally we don't even actually store a sparse matrix, so most implementations choose sparse conversion

A computational graph does the trick as well

weights as nodes, and multiple edges from weights  $\Rightarrow$  weight sharing



Convolution  $\rightarrow$  Strong prior? low entropy

weights of one hidden unit  $\rightarrow$  close to the neighbour, identical but shifted

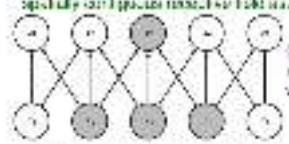
low variance

$\Rightarrow$  local interactions

Convolutional net is similar to a fully connected net but with an initially strong prior over its weights

- says that the weights for one hidden unit must be identical to the weights in its neighbors, but shifted in space

- Prior also says that the weights must be zero, except for 1. the small spatial and frequency field around the first hidden unit.



Convolution with a kernel of width 3 in a hidden unit. It has 9 weights which are the same for all.

- Convolution introduces an infinity norm prior probability distribution over the parameters in a layer

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$(w_0, w_1, w_2, w_3)$$

Pooling also has strong prior

$N$  convolutions  $\rightarrow$  (pooling)

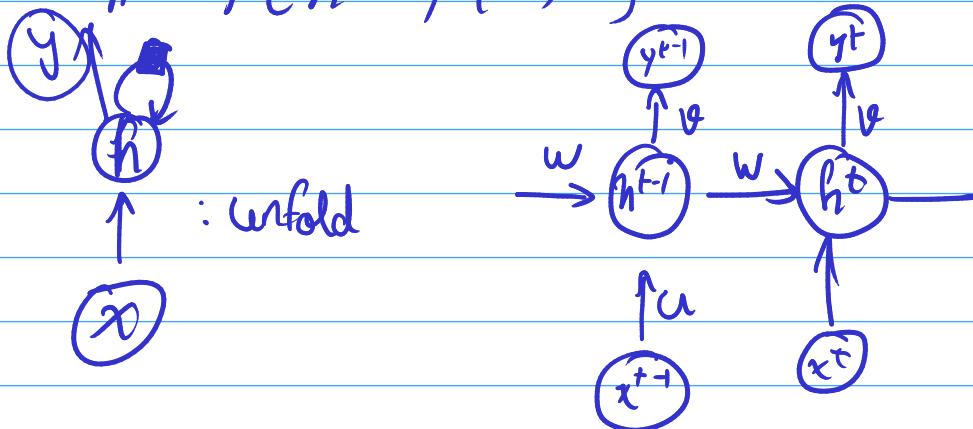
(needn't be after [conv | pooling])

wedid till

q-4 : 9.1D  $\rightarrow$  Neuroscientific basis  
(read when free)

## Recursive Function

$$f^{(t)} = f(h^{(t-1)}, x^{(t)}; \Theta)$$



$$\begin{cases} S^{(1)} = g^{(1)}(x^{(1)}) \\ S^{(2)} = g^{(2)}(x^{(1)}, x^{(2)}) \\ S^{(3)} = g^{(3)}(x^{(1)}, x^{(2)}, x^{(3)}) \\ \vdots \\ S^{(k)} = g^{(k)}(x^{(1)}, x^{(2)}, \dots, x^{(k)}) \end{cases}$$

today's infosys stock price  $\rightarrow$  today's condition  
 yesterday's MSE  
 dependent on Day before's ...

$$\begin{matrix} \textcircled{1} & \textcircled{2} & \{3\} & \{4\} & \{5\} \\ 1 & 1 & 2 & 3 & 5 \end{matrix}$$

$$f_n = f_{n-1} + f_{n-2}$$

$$\begin{matrix} \{0 & 1 & 2 & 3 & 4 & 5 & 6 \\ , & 1 & 2 & 3 & 5 & 8 & 13 \end{matrix}$$

Let the data be  $s_1, s_2, s_3, \dots, s_{t-1}, s_t$  where  $\theta$

$$s_t = f(s_{t-1}, \theta) \neq d_t = 3s_{t-1} + 5$$

↳ how do we discover this

Day	1	2	3	...	900	901	902	$\}$ Stock market
	$d_1, d_2, d_3$				$d_{900}$	$d_{901}$	$d_{902}$	

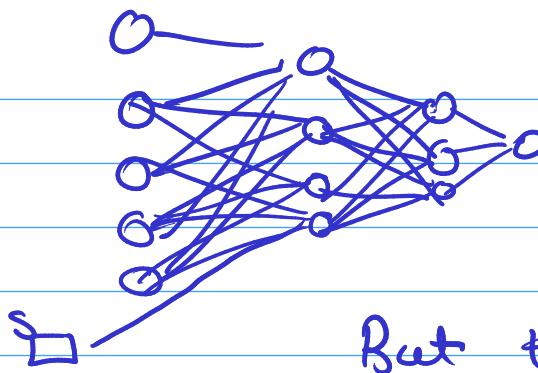
now lets say that  $f(x_1, x_2, \dots, x_{50})$  be the feature vector for 1 day

$$s_{903} = f(x_1, x_2, \dots, x_{50}) \text{ is wrong}$$

$$s_{903}: f(s_{902}, x_1, x_2, \dots, x_{50})$$

↳ consider this also as a feature

Develop a model



But that's stupid

$s_{902}$  depends on the feature vector  
of the day 902

So we get

$$g(x_1^{(902)}, x_2^{(902)}, x_3^{(902)}, \dots, x_{50}^{(902)}; x_1^{(903)}, x_2^{(903)}, \dots, x_{50}^{(903)})$$

so we do

$x_1^{\text{Cur}}$	0		
$x_2^{\text{Cur}}$	0	6	
:	0	6	
$x_{50}^{\text{Cur}}$	6	0	
0, Prev	0	0	0
$x_2^{\text{Prev}}$	0	0	
:	0	0	
$x_{50}^{\text{Prev}}$	0	0	

and build a model

assuming that today's condition  
does not depend on day before's  
state

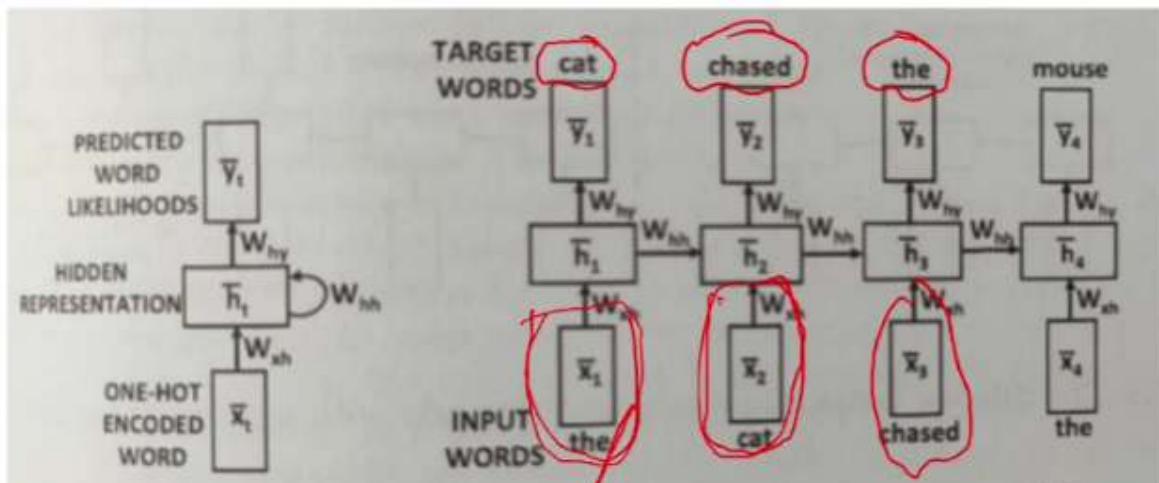
But as per financial  
experts, diff number of  
days are required!

$\Rightarrow$  each model for # of days? TF?

Another problem

The Cat chased the . . . . .

Now the issue is , the later words can influence the sentence



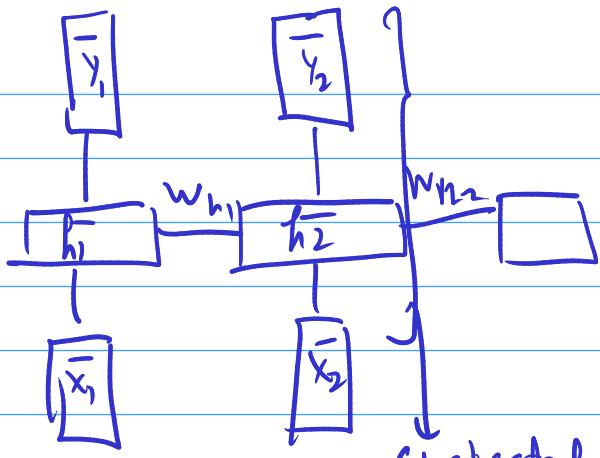
} + chased = the

The essence of the previous  
input & the current input } that  
hinders  
remains  
the same  
other the  
sequence

$f$  is the same }  
$$\begin{cases} f(h_0, x_1) \\ f(h_1, x_2) \\ f(h_2, x_3) \end{cases}$$
       $f(h_3, x_4)$

if you believe that this  $f$  must be the same

almost... }  
f needn't depend on position }  
So the output is independent of the position of input in the sequence  
by assuming that the function is same



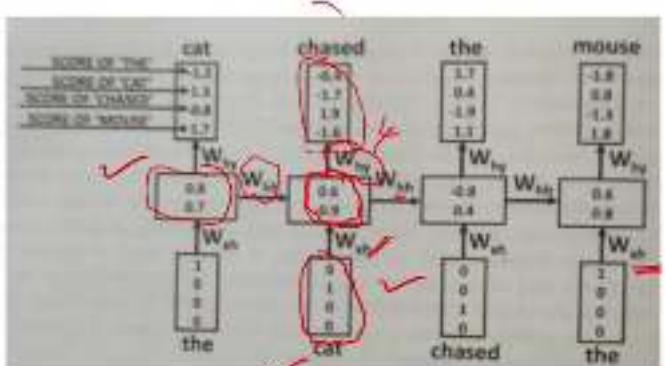
Abstracted info  
is sent to the next guy  
(Same dimension as  $X$  by the way)

so

our abstraction  $h_2 = \sigma(\underbrace{w_1 h_1 + w_2 h_2}_{\text{encode}})$

'h' is as follows:

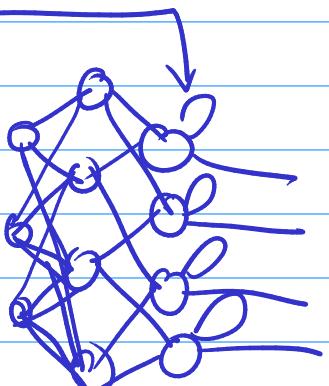
Why  $h_2 \rightarrow \text{word}$   
decode, get score/word vector



$$\bar{h}_t = \tanh(w_{xh} \bar{x}_t + w_{hh} \bar{h}_{t-1})$$

$$\bar{y}_t = w_{yh} \bar{h}_t$$

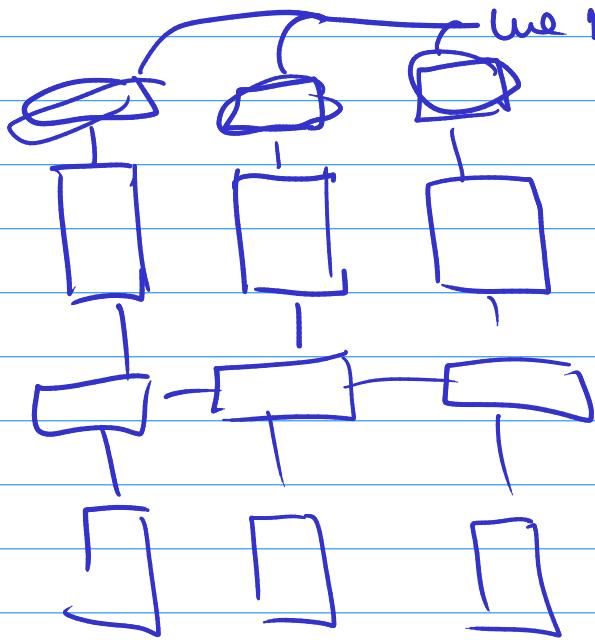
This is loop unrolling



$$\bar{h}_t = \tanh(w_{xh} \bar{x}_t + w_{hh} \bar{h}_{t-1})$$

NOW WE know the ground truth whilst training  
if ground truth

is  $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$  then '1' has to be maximised.



we minimise the loss for entire sequence

The sequence length can vary but that's not a problem since we built the model with that in mind as seen below

$$L = -\log(p_i^{y_i})$$

for gradient descent

$$\frac{\partial L}{\partial w_{2n}} \quad \frac{\partial L}{\partial w_{nn}} \quad \frac{\partial L}{\partial w_{hy}}$$

$$\begin{bmatrix} 1.8 \\ 1.3 \\ 0.8 \\ 1.7 \end{bmatrix} \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \hat{y}_4 \end{bmatrix}$$

$\frac{\partial L}{\partial \hat{y}_i}$  we can find

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = W_{hy} h_2$$

$$\frac{\partial L}{\partial w_{x_h}} : g \underset{\text{accumulated error}}{\overbrace{x}}$$

also  $g^{(2)} x^{(2)}$   
 $g^{(3)} x^{(3)}$  (So we get multiple gradients)

$$L = -\log(w_{xy} h)$$

how to maintain  $w_{x_h}$  to be the same?

Similar to CNNs  
 (Refer notes)  $\frac{\partial L}{\partial w_{x_h}^{(1)}} \frac{\partial L}{\partial w_{x_h}^{(2)}} \frac{\partial L}{\partial w_{x_h}^{(3)}} \frac{\partial L}{\partial w_{x_h}^{(4)}}$

$$\frac{\partial L}{\partial w_{x_h}} = \frac{\partial L}{\partial w_{x_h}} * \left( \frac{\partial w_{x_h}}{\partial w_{x_h}^{(1)}} + \frac{\partial w_{x_h}}{\partial w_{x_h}^{(2)}} + \frac{\partial w_{x_h}}{\partial w_{x_h}^{(3)}} \right)$$

$$w_{x_h} = w_{x_h}^{(1)} = w_{x_h}^{(2)} = w_{x_h}^{(3)} \Rightarrow \frac{\partial w_{x_h}}{\partial w_{x_h}^{(4)}} \Rightarrow 1$$

RNN  $\rightarrow$  Recurrence

$$h_t = f(h_{t-1}, x_t)$$

$$= f(f(h_{t-2}, x_{t-1}), x_t)$$

$$f = \tanh(b + \omega h_{t-1})$$

### Problem

(1) Any big sentence, answer question  
forget stuff  
(carry stuff)

$$h_{t+1} = \tanh(b + \omega h_t + ux^t)$$

$$\frac{\partial L}{\partial h_{t+1}} \underset{?}{=} c$$

$$\left[ \because \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} \right]$$

$$= \frac{\partial L}{\partial h_{t+1}} \underset{?}{=} \boxed{c}$$

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial h^2} \underset{?}{=} \boxed{c}$$
$$\hookrightarrow \frac{\partial L}{\partial h_3} \underset{?}{=} \boxed{c}$$
$$= \frac{\partial L}{\partial h_3} \underset{?}{=} \boxed{c^2}$$

$$\omega = Q A Q^T \Rightarrow \omega^k = Q A^k Q^T$$

(positive semidefinite)

$$A = \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ \vdots & & & & \\ 0 & 0 & 0 & \cdots & \lambda_k \end{bmatrix}$$

Orthogonal matrices

$$\begin{aligned} \therefore \omega^k &= Q A^k Q^T \\ &= Q \Lambda^k Q^T \end{aligned}$$

$$\frac{\partial L}{\partial h_i} = \frac{\partial L}{\partial h_k}, Q A^k Q$$

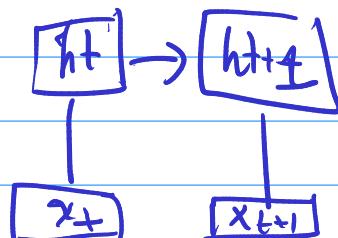
$$\downarrow \quad \lambda_1 \quad \lambda_2 \quad \dots \quad \lambda_k$$

$$\lambda^{40} \geq 0.000\dots01$$

$$\begin{aligned} \lambda &= 2 \\ \lambda^{40} &\Rightarrow 2^{40} \end{aligned}$$

→ So gradients become very small (vanish)  
Very large (explode)

$$h_{t+1} = \tanh(\omega h_t + u x_{t+1})$$



Now we want to forget some stuff  
as well

But how?

$$f: \sigma(b^f + U^f x_t + \omega^f h_{t-1})$$

multiply it by  $h^{t-1}$

Now we get how much we  
forget

Now we need a remember function

(But why? Can't we just add

$$b + Ux^t + Wh^{t-1})$$

$$g = b^g + U^g x + \omega^g h^{t-1}$$

Since it's always good to get an idea of it apparently

we take this and multiply it with  $(b + Ux^t + Wh^{t-1})$

$$f^{(t)} h^{(t-1)} + g^t (b + Ux^t + Wh^{t-1})$$

But thoda wait karle, there's a  
change now

$$f_i^{(t)} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right)$$



$$\delta_i^{(t)} = f_i^{(t)} \delta_i^{(t-1)} + g_i^{(t)} \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right)$$

$$g_i^{(t)} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right)$$

$$h_i^{(t)} = \tanh(\delta_i^{(t)}) q_i^{(t)}$$

$$q_i^{(t)} = \sigma \left( b_i^q + \sum_j U_{i,j}^q x_j^{(t)} + \sum_j W_{i,j}^q h_j^{(t-1)} \right)$$

} But here, we don't do  
 $\tanh(b + (Ux_t + Wh_{t-1}))$   
 here, unlike RNN

(we do

fraction of new person also taken

$$h^{(t+1)} = \tanh(\delta^{(t)}) \sigma(b + Ux_t + Wh_{t-1})$$

now  $s^{(t)} =$

$$f^{(t)} \delta^{(t-1)} + g^{(t)} \sigma(b + Ux_t + Wh_{t-1})$$

So now, we need to find out

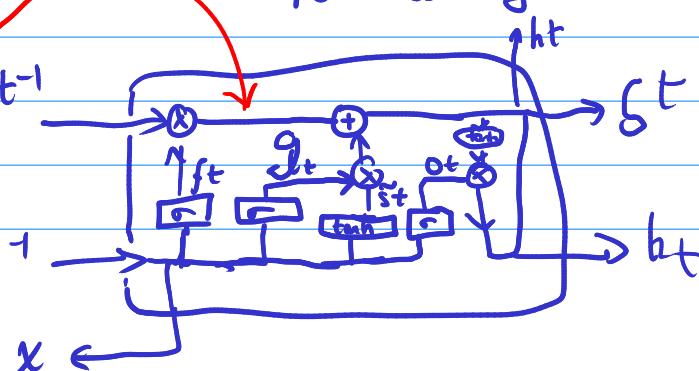
$$b^f \ U^f \ w^f \quad b, w, U, v$$

$$b^g \ U^g \ w^g$$

fraction for remembering

Our backward  
 is now  
 on 's'  
 and is  
 very simple!!

for decoding



the derivative of the product  
will give a sum though  
instead of products

maybe that's why gradients don't  
explode

Weighted average

$$h_i^{(t)} = u_i^{(t)} h_i^{(t-1)} + (1 - u_i^{(t)}) \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} r_j^{(t-1)} h_j^{(t-1)} \right)$$

- Where  $u$  stands for the update gate and  $r$  for reset gate. Their value is defined as usual:

$$u_i^{(t)} = \sigma \left( b_i^u + \sum_j U_{i,j}^u x_j^{(t)} + \sum_j W_{i,j}^u h_j^{(t)} \right) \quad \text{and} \quad r_i^{(t)} = \sigma \left( b_i^r + \sum_j U_{i,j}^r x_j^{(t)} + \sum_j W_{i,j}^r h_j^{(t)} \right)$$

reset  
prev input

GRU

