

01

- ❖ The true challenge proved to be solving tasks that are easy for people to perform but hard for people to describe formally

Application of deep learning :-

① classification

② retrieval

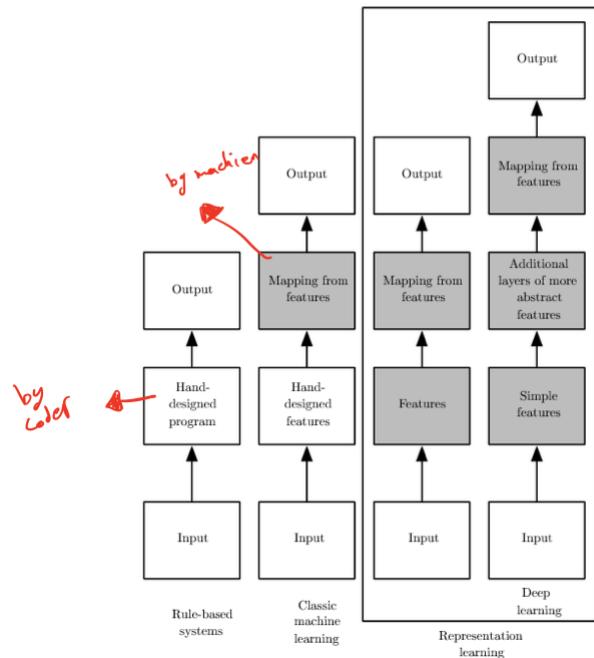
③ detection

④ segmentation

⑤ Image captioning

⑥ Deep Art: Combining Content and Style from Different Images

- ❖ Simple machine learning algorithms depends heavily on the representation of the data they are given



O2 - Classification

Images

Problems with Image classification

- ① view point
- ② background clutter
- ③ illumination
- ④ occlusion
- ⑤ deformation

① binary

② grey

③ color

❖ **Intraclass Variation**
between the samples in
the same class

* **Interclass Variation**
between two classes

Ways to classify Images..

① Knn

- Using distance
- no real training
- complexity \Rightarrow train = $O(1)$, test = $O(n)$
- hyperparameters \Rightarrow K , distance metric (L_1, L_2)
- Set them by: ① train, validation, test

② cross-validation

③ linear classification:

- score function $f(x, w) = wx + b$

03 - Loss function and Optimization

A **loss function** tells how good our current classifier is

- ① SVM loss function
- ② Regularization (add L_2 term)
 → prevent overfitting
- ③ Softmax classifier

What is the softmax loss and the SVM loss if I double the correct class score from 10 → 20?

for SVM will be the same

How to reduce the loss?

- ① trial and error
- ② gradient descent
 - numeric
 - analytic (derivatives)

Multiclass SVM loss function



if the values are different the loss fun. is good

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Loss over full dataset is average:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

$$L = (2.9 + 0 + 12.9) / 3$$

L = 5.27 we try to minimize it

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9 0 12.9		

$$\begin{aligned} &= \max(0, 5.1 - 3.2 + 1) \\ &\quad + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$$

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

$$\begin{aligned} &= \max(0, 2.2 - (-3.1) + 1) \\ &\quad + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 6.3) + \max(0, 6.6) \\ &= 6.3 + 6.6 \\ &= 12.9 \end{aligned}$$

Softmax Classifier (Multinomial Logistic Regression)

- Want to interpret raw classifier scores as **probabilities**
 - Softmax**



cat
car
frog

3.2
5.1
-1.7

$$\exp^{\text{unnormalized probabilities}}$$

24.5
164.0
0.18

Probabilities must be ≥ 0

normalize

0.13
0.87
0.00

Probabilities must sum to 1

$$L_i = -\log (0.13) = 2.04$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cross entropy

$$L_i = -\log P(Y = k|X = x_i)$$

Maximum Likelihood Estimation
Choose weights to maximize the likelihood of the observed data

Softmax vs. SVM

matrix multiply + bias offset

0.01	-0.05	0.1	0.05
0.7	0.2	0.05	0.16
0.0	-0.45	-0.2	0.03

W

-15
22
-44
56

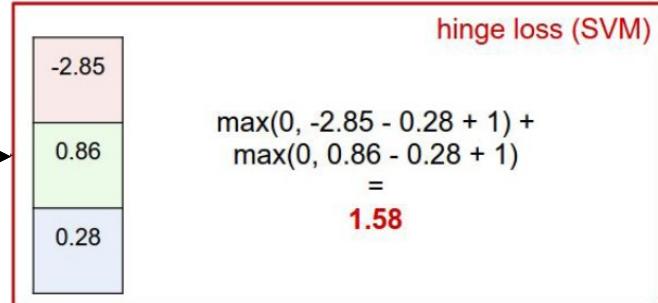
x_i

y_i 2 → inject

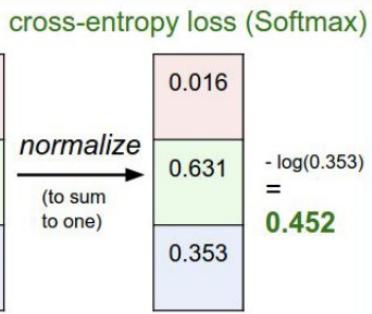
+

0.0
0.2
-0.3

b



most common



Recap

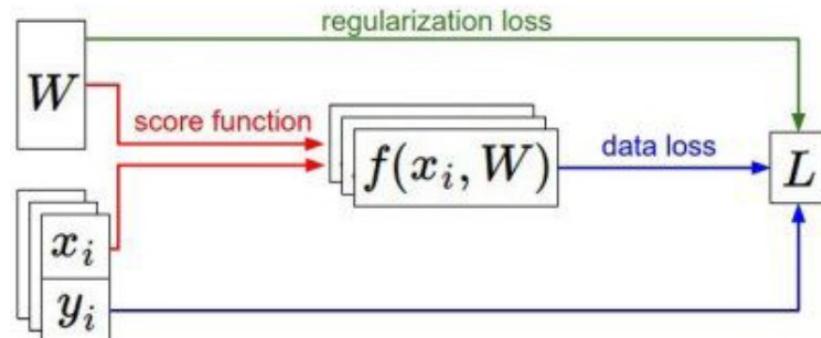
- ❖ We have some dataset of (x, y)
- ❖ We have a **score function**: $s = f(x_i, W)$
- ❖ We have a **loss function**:

Softmax $L_i = -\log \left(\frac{e_k^s}{\sum_j e_j^s} \right)$

SVM $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$

Full loss $L = \frac{1}{N} \sum_{i=1}^N L_i + R(W)$

How do we find the best W ?



OL - NN

* of neurons in the last layer = * of classes

O5 - Multilayer Perceptron

Can't handle non-linear inputs

if $y^i(\sum w_j x_j) > 0$

❖ The perceptron learning algorithm is guaranteed to converge if classes are linearly separable after $(\frac{R}{\gamma})^2$ misclassifications

else if $y^i(\sum w_j x_j) \leq 0$

$$\vec{w} \leftarrow \vec{w} + y^{(i)} x^{(i)}$$

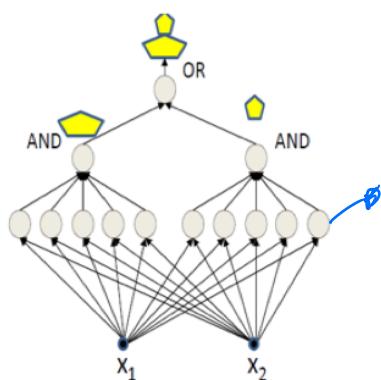
❖ Feed-Forward Neural Network

❖ Has No loops: Neuron outputs do not feed back to their inputs directly or indirectly

repeat until no classification errors

for complex distribution

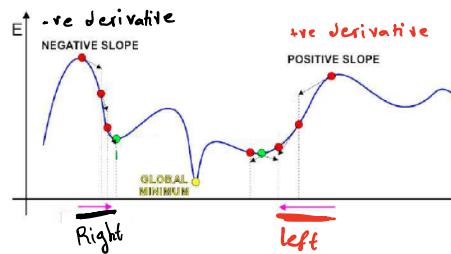
Classification Rule: $\text{Sign}(W^T X_i)$



linear classifier

❖ We don't know the outputs of the individual intermediate neurons in the network for any training input

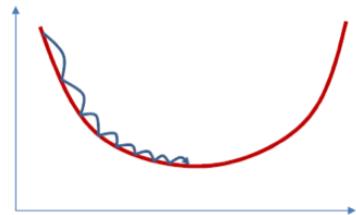
$$\frac{df}{ds} = \sigma(s)(1 - \sigma(s))$$



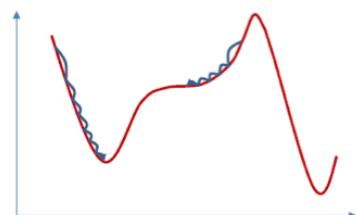
- For multivariate function, to find a **minimum**, move exactly opposite the direction of the gradient

$$w = w - \text{Learning_rate} * dw$$

- For appropriate step size, for **convex** (bowls shaped) functions, gradient descent will always find the minimum



- For **non-convex** functions it will find a local minimum or an inflection point



$$\Sigma w_i x_i + b$$

* the input of neurons is \Rightarrow Weighted sum of input + bias

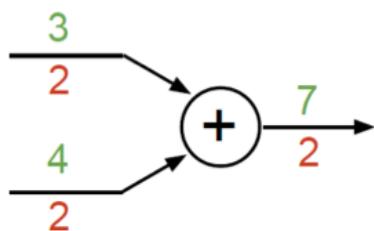
* we can use different activation function in different layers

* we **can't** use different gradient descent algo* in different layers

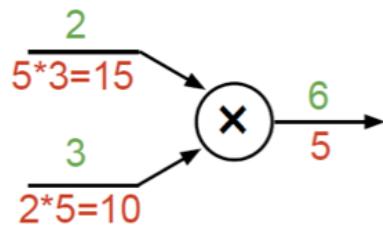
06 - 1 - Backpropagation

Patterns in gradient flow

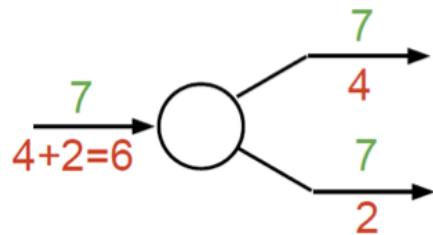
add gate: gradient distributor



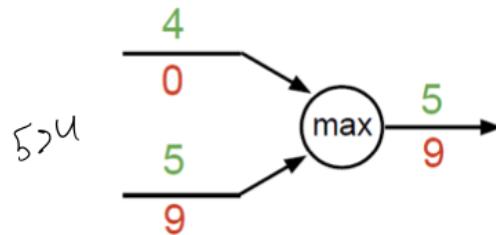
mul gate: "swap multiplier"



copy gate: gradient adder



max gate: gradient route



06 - 2 - Backpropagation

Scalar to Scalar \Rightarrow Regular derivative

Vector to scalar \Rightarrow Gradient (n input, 1 output)

$$\frac{\partial f}{\partial x} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$$

Vector to vector \Rightarrow Jacobian (n input, m output)

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Solve Examples on Vector Backpropagation

07 - Start with PyTorch

* automatic back propagation

* Dynamic Computation graphs

* Tensors are single data type

cannot get GPU tensor as

numpy array directly

\rightarrow you need to convert it

to CPU then to numpy

* .backward() doesn't replace, but accumulates gradients

\rightarrow use
`x.grad.data.zero_()`

losses in PyTorch

- ❖ `nn.CrossEntropyLoss` does both the softmax and the actual cross-entropy.

① forward

② backward

③ optimizing

④ learning (update)

```
x = torch.arange(0., 32)
net = torch.nn.Linear(32, 10)
```

input (features) →
output (classes) →
torch.nn.Linear
to the incoming

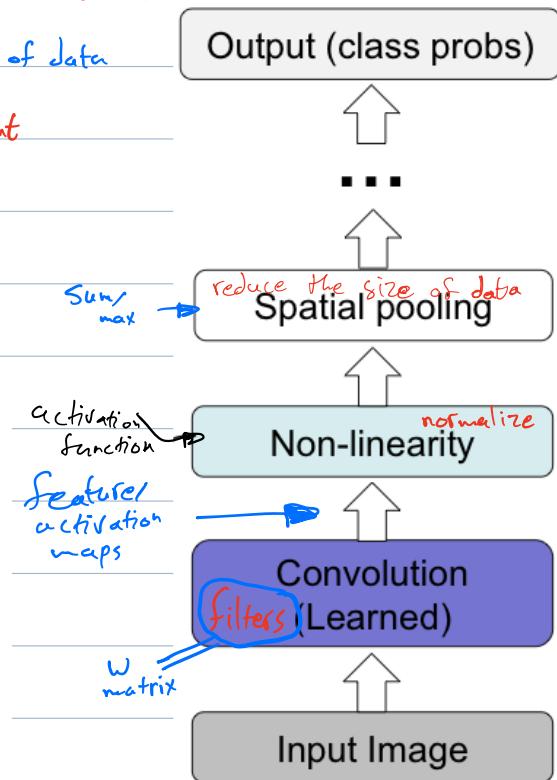
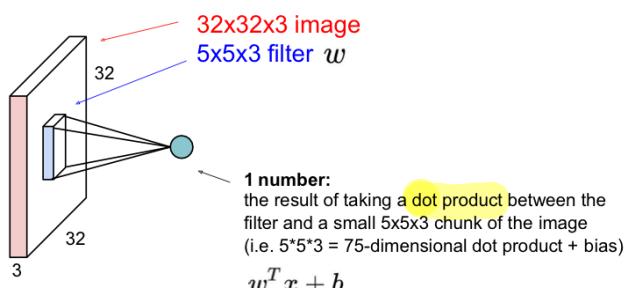
08 - CNN

when we need to keep the locations of input

* Goal \Rightarrow detect a pattern regardless of its location

* normal NN affected by changing the location of data

So we need a new network with Shift invariant

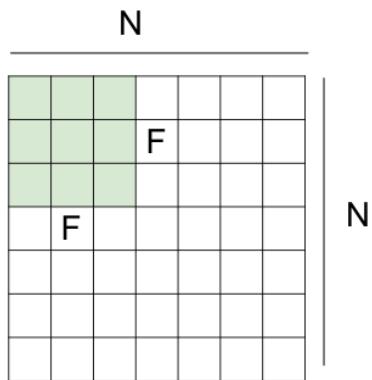


each filter will look for
some feature in the data
(learned automatically)

Stride = Shift

for convolution layers -

Output size:
 $(N - F) / \text{stride} + 1$



to get output = input \Rightarrow Set padding = $(F - 1)/2$
and stride = 1

- Output size:

$(N + 2 * \text{padding} - F) / \text{stride} + 1$

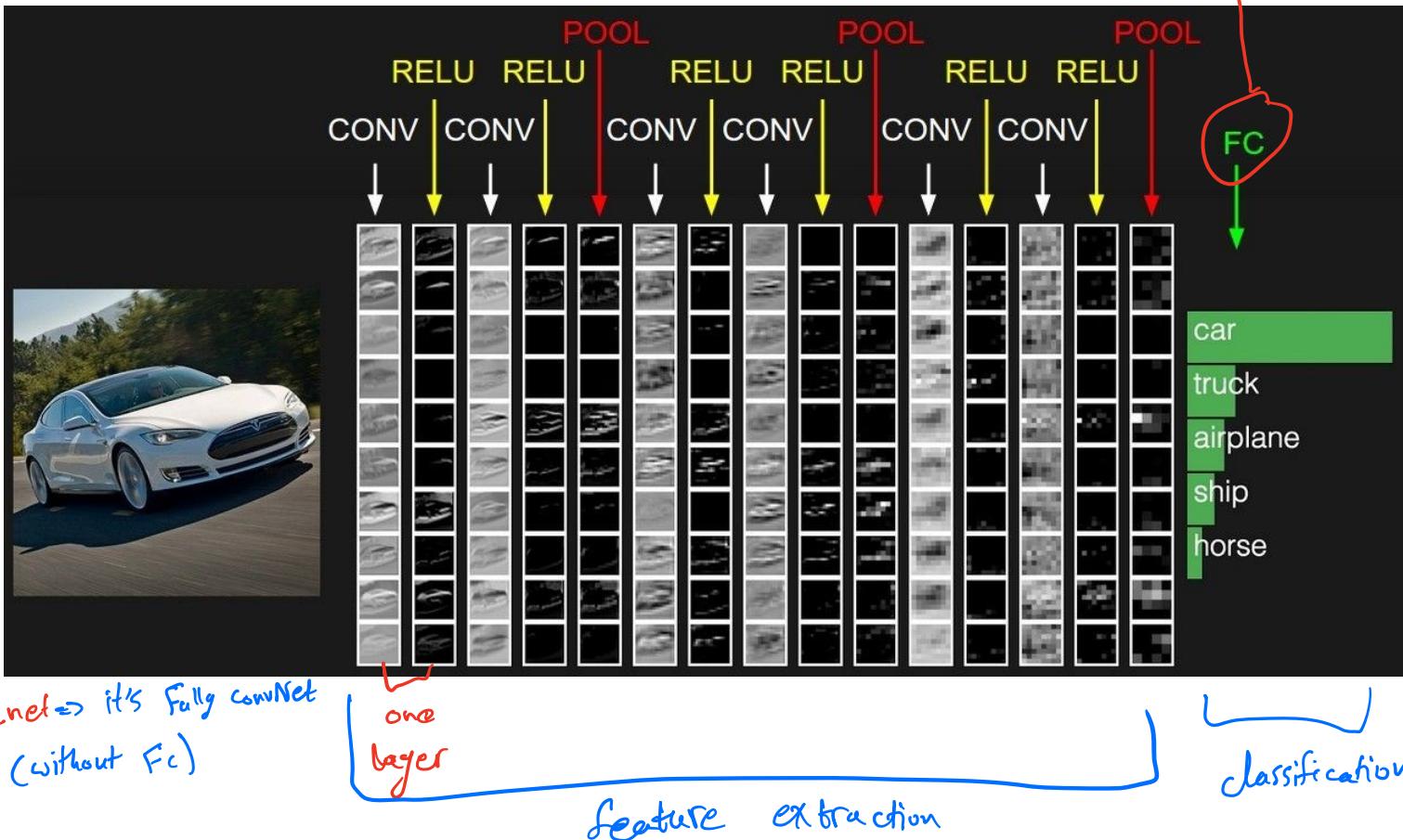
Padding doesn't do anything with resolution of the image

~~Pooling:-~~ helps in reducing overfitting

Max Pooling = take the max of a $n \times n$ filter

there is NO learning parameters

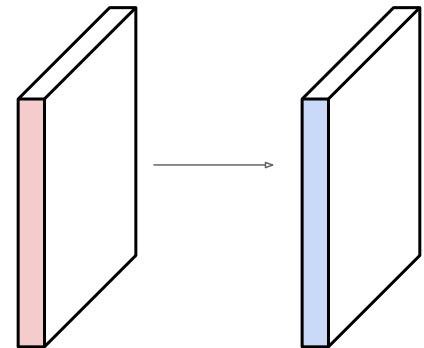
Putting it all together



Examples time:

5 hyperparameters

Input volume: **32x32x3**
① 10 ② 5x5 x3 filters with ③ stride 1, ④ pad 2
⑤ and * of layers

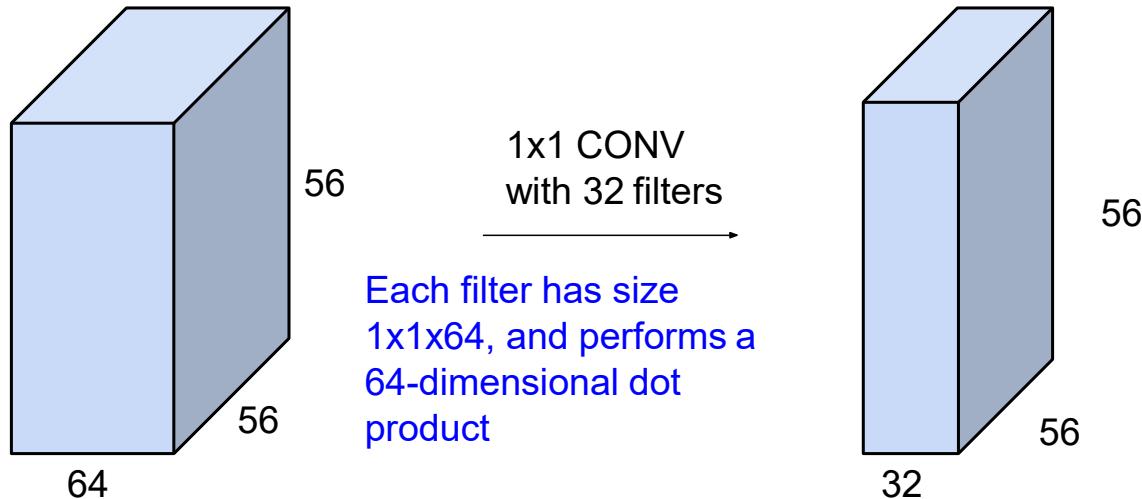


Number of parameters in this layer?

each filter has $5 \times 5 \times 3 + 1$ = 76 params (+1 for bias)
=> $76 \times 10 = 760$

1x1 convolutions

- to reduce the depth of the data (dimensionality reduction)
- used for feature pooling



- because of small kernel size, it doesn't affect by overfitting

O9 - Common CNN

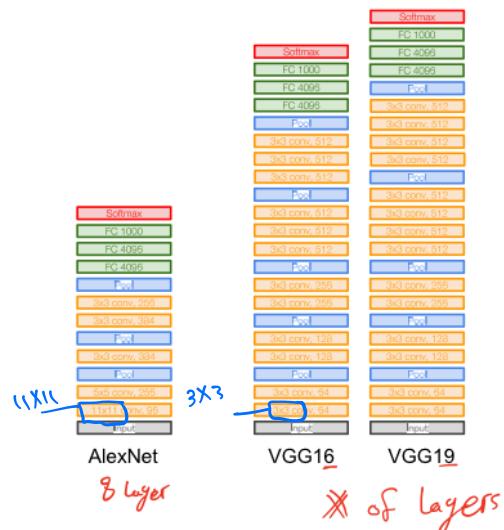
① AlexNet

② VGG16 (138M Parameters)

③ VGG19

Why use smaller filters?

- ① deeper networks
- ② more non-linearities
- ③ fewer parameters

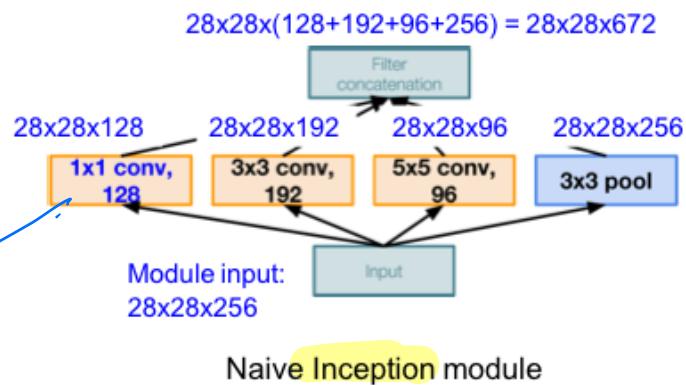


④ GoogleNet

- 22 layers
- No FC layers
- 5M parameters
- 3 outputs

Example:

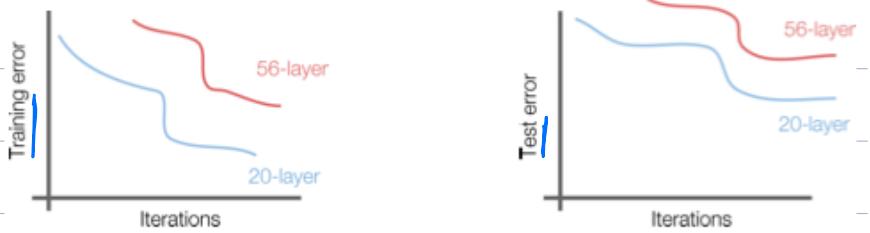
Q3: What is output size after filter concatenation?



"bottleneck"

⑤ ResNet

- 152 layers

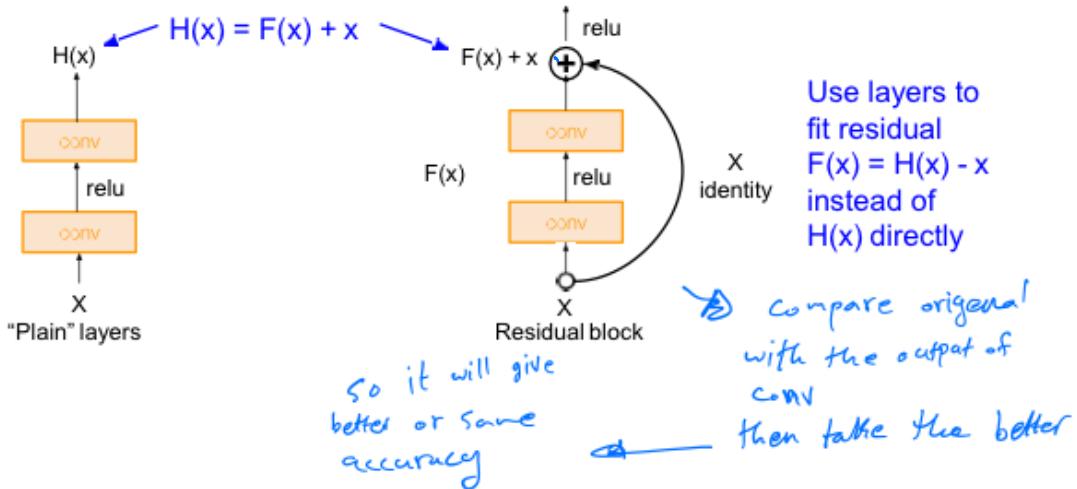


go deeper → overfitting

So
56-layer model performs worse on both training and test error
→ The deeper model performs worse, but it's not caused by overfitting!

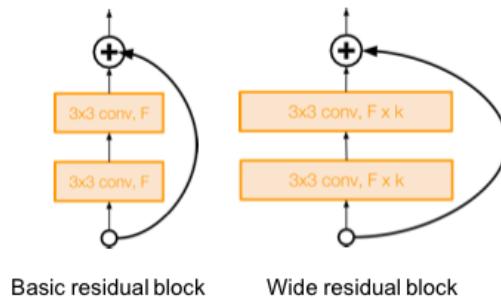
I think because of
vanishing

- Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



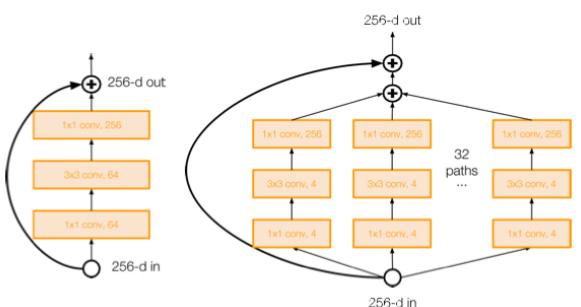
⑥ Wide ResNet

- Increasing width instead of depth
more computationally efficient
(parallelizable)



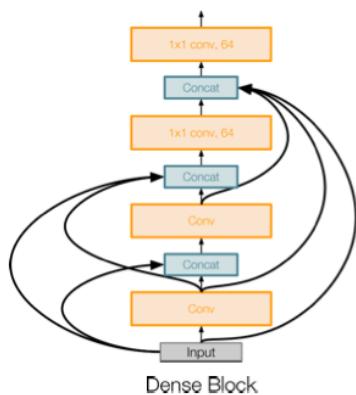
⑦ ResNeXt

apply the Inception module



⑥ DenseNet

- Dense blocks where each layer is connected to every other layer in feedforward fashion



10 - NN Training

1. One time setup

- preprocessing, activation functions, weight initialization, regularization, gradient checking

2. Training dynamics

- babysitting the learning process, parameter updates, hyperparameter optimization

3. Evaluation

- model ensembles

① Data Augmentation (increase # of samples)

translation *shifting and adding padding*

rotation

stretching

shearing,

lens distortions

...*add noise*

--random crops and scales

Should be reasonable
and accepted in the
Problem domain

(2) Preprocessing:-

1. normalize $\frac{x - \mu}{\sigma}$

2. features selection \rightarrow without modify them

3. features reduction \rightarrow return new features sorted by important

(3) Weight Initialization

❖ Q: what happens when all $W=0$?

Every neuron in the network computes the same output, then they will also all compute the same gradients during backpropagation and undergo the exact same parameter updates

Solution
↓

`W = 0.01 * np.random.randn(D, H)`

high $W \Rightarrow$ exploding

low $W \Rightarrow$ vanishing

W and b

Are NOT hyperparameters

Works ~okay for small networks, but problems with deeper networks.

vanishing

Solution
↓

not good with ReLU

weights initialization: "Xavier" initialization $\Rightarrow W = \frac{\text{rand}(\text{fan_in}, \text{fan_out})}{\sqrt{\text{fan_in}}}$

"He init" initialization

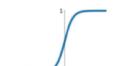
$\Rightarrow W = \frac{\text{rand}(\text{fan_in}, \text{fan_out})}{\sqrt{2/\text{fan_in}}}$

(4) Activation Functions

\rightarrow to decide whether a neuron will fire or not

Sigmoid

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



Leaky ReLU

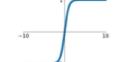
$$\max(0.1x, x)$$



it's not for normalization

tanh

$$\tanh(x)$$

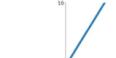


Maxout

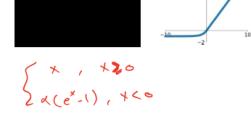
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ReLU

$$\max(0, x)$$



ELU



- Use ReLU. Be careful with your learning rates

- Try out Leaky ReLU / Maxout / ELU / PReLU

- Try out tanh but don't expect much

- Don't use sigmoid

- Sigmoid and Tanh doesn't suffer from vanishing and exploding gradient unlike ReLU

- sum of probabilities of all sigmoid unit $\neq 1$

(5) Batch Normalization

what \rightarrow (Preprocessing at every layer of net)

why

doesn't affect the period of training

- Makes deep networks **much** easier to train!
- Improves gradient flow
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!

How

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad \begin{matrix} \text{to avoid} \\ \text{division by} \\ \text{zero} \end{matrix} \quad // \text{normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

In testing: we get γ and σ
from training

when \rightarrow ① After FC

② After Conv

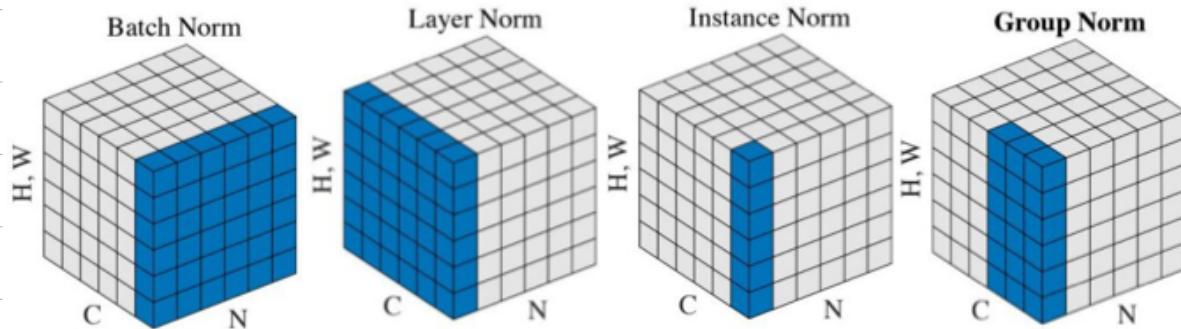
③ Before nonlinearity (activation function)

FC

$$\begin{aligned} \mathbf{x}: N \times D \\ \text{Normalize} \\ \boldsymbol{\mu}, \sigma: 1 \times D \\ \gamma, \beta: 1 \times D \\ \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \sigma + \beta \end{aligned}$$

Conv

$$\begin{aligned} \mathbf{x}: N \times C \times H \times W \\ \text{Normalize} \\ \boldsymbol{\mu}, \sigma: 1 \times C \times 1 \times 1 \\ \gamma, \beta: 1 \times C \times 1 \times 1 \\ \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \sigma + \beta \end{aligned}$$



⑥ Babysitting the Learning Process

Do these
and monitor:
① loss
② accuracy

- ❖ Preprocess data
- ❖ Choose architecture
- ❖ Initialize and check initial loss with no regularization
- ❖ Increase regularization, loss should increase
- ❖ Then train – try small portion of data, check you can overfit *make it overfit*
- ❖ Add regularization, and find learning rate that can make the loss go down
- ❖ Check learning rates in range [1e-3 ... 1e-5]
- ❖ Coarse-to-fine search for hyperparameters (e.g. learning rate, regularization)

if $\text{cost} = \text{nan} \Rightarrow$ learning rate is high

⑦ Hyperparameter Optimization

w and b

are not hyperparameters

* Start with a few epochs then go longer

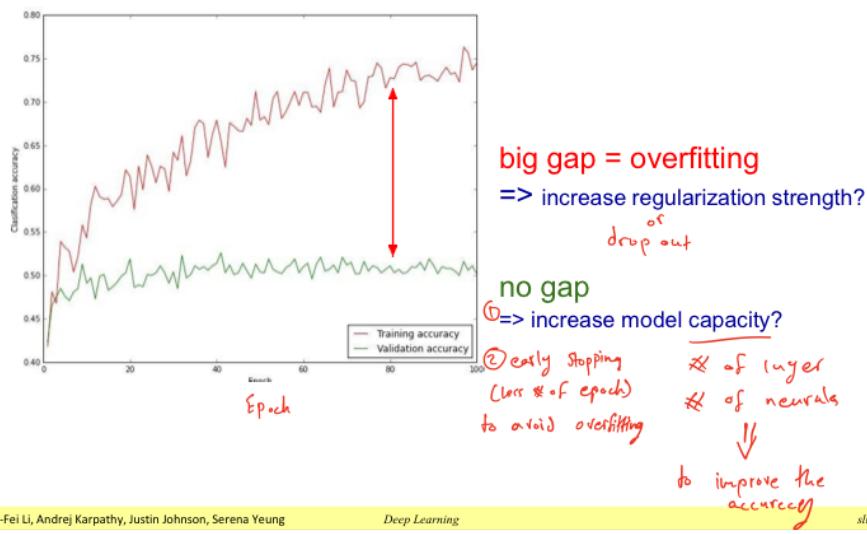
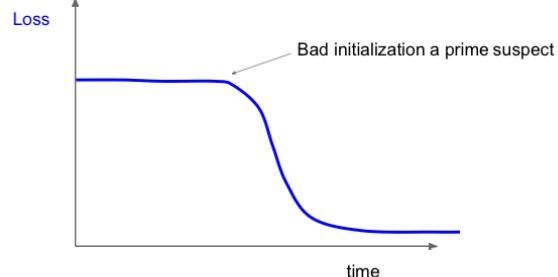
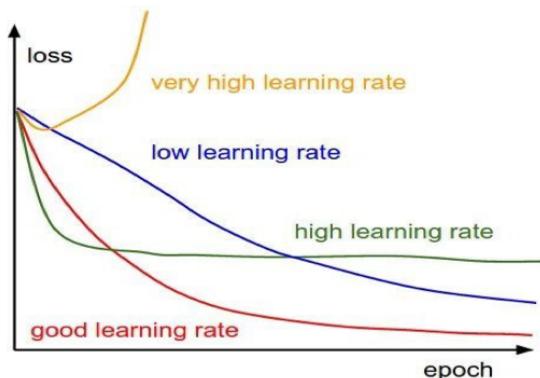
* type of search for best hyperparameters: (to get best accuracy)

① grid set values for every hyperparameter → test the combination

② Random

* Hyperparameters to play with:

- ❖ network architecture
- ❖ learning rate, its decay schedule, update type
- ❖ regularization (L2/Dropout strength)



11 - NN Training

② Optimizers (to optimize gradient)

problem with SGD:

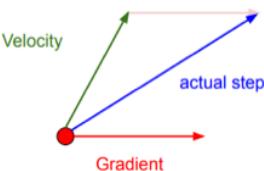
- ① may be very slow
- ② get stuck with local minimum
- ③ can be noisy

like
a ball

① SGD + momentum:-

Parameters: $\rho = 0.9$ OR 0.99

Momentum update:



② AdaGrad

Progress along "steep" directions is damped;
progress along "flat" directions is accelerated

faster

③ RMS Prop

Parameters: β_1 , β_2

to avoid local minima:-

- increase learning rate

- use momentum or adaptive learning

- add some noise while updating weight

④ Adam

AdaGrad

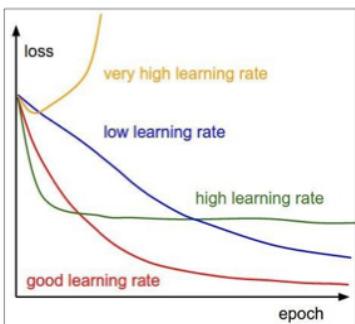
RMSProp

In practice:

- ❖ **Adam** is a good default choice in many cases; it often works ok even with constant learning rate
- ❖ **SGD+Momentum** can outperform Adam but may require more tuning of LR and schedule
 - ✧ Try cosine schedule, very few hyperparameters!

⑨ Learning rate Schedules

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these learning rates is best to use?

A: All of them! Start with large learning rate and decay over time

We decay the learning rate at a few fixed point
e.g. after every 30 epochs

ways to decay:-

① **exponential decay:**

$$\alpha = \alpha_0 e^{-kt}$$

constant

② **1/t decay:**

$$\alpha = \alpha_0 / (1 + kt)$$

epoch

③ **Cosine:**

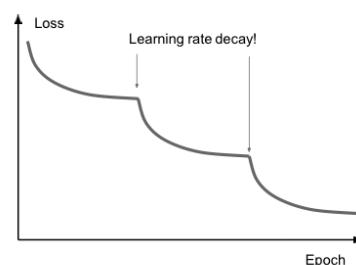
$$\alpha_t = \frac{1}{2} \alpha_0 (1 + \cos(t\pi/T))$$

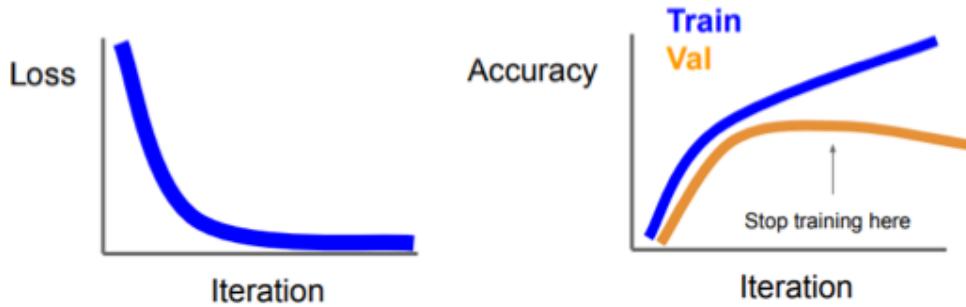
④ **Linear:**

$$\alpha_t = \alpha_0 (1 - t/T)$$

⑤ **Inverse sqrt:**

$$\alpha_t = \alpha_0 / \sqrt{t}$$





- ❖ Stop training the model when accuracy on the validation set **decreases**

(10) Choosing hyperparameters:- (using training Set)

Step 1: Check initial loss → *without decay*

Step 2: Overfit a small sample

Step 3: Find LR that makes loss go down

Step 4: Coarse grid, train for ~1-5 epochs

Step 5: Refine grid, train longer

Step 6: Look at loss curves

Pick best models from Step 4, train them for longer (~10-20 epochs) without learning rate decay

Loss not going down? LR too low, bad initialization

Loss explodes to Inf or NaN? LR too high, bad initialization

Improving single-model performance is easier and more effective than train multiple independent models.

⑪ Improve single-modal performance :-

by Regularization and Dropout

① Regularization \Rightarrow

add term to loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

② Dropout \Rightarrow

What \Rightarrow In each forward pass, randomly set some neurons to zero

Probability of dropping is a hyperparameter; 0.5 is common

Why

- ❖ Dropout forces the neurons to learn "rich" and redundant patterns
- ❖ E.g. without dropout, a noncompressive layer may just "clone" its input to its output
 - ◊ Transferring the task of learning to the rest of the network upstream
- ❖ Dropout forces the neurons to learn denser patterns
 - ◊ With redundancy

each epoch will have different # of neurons

We can apply dropout in hidden and input layers



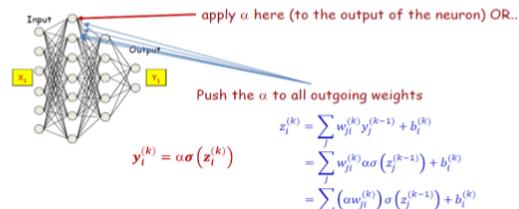
❖ At test time all neurons are active always

\Rightarrow We must scale the activations so that for each neuron:

output at test time = expected output at training time

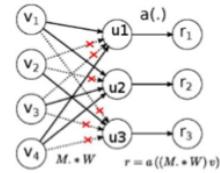
❖ At test time, multiply by dropout probability α

At test



types of Dropout

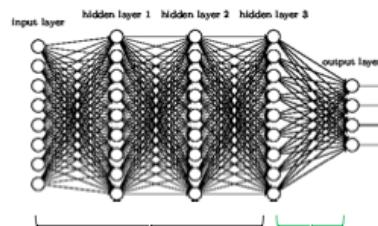
- ❖ Dropconnect drop some weights
 - ◊ Drop individual connections, instead of nodes
- ❖ Shakeout Scale some weights
 - ◊ Scale up the weights of randomly selected weights
 - $|w| \rightarrow \alpha |w| + (1 - \alpha) c$
 - ◊ Fix remaining weights to a negative constant
 - $w \rightarrow -c$
- ❖ Whiteout add some noise
 - ◊ Add or multiply weight-dependent Gaussian noise to the signal on each connection



~~Transfer Learning~~

Your data = target

rich data = source



Set these to the already learned weights from another network

Learn these on your own task
data

freeze them



More generic
More specific

	Very similar dataset	Very different dataset
my data Very little data	Use Linear Classifier on top layer <i>just last layer</i>	You're in trouble... Try linear classifier from different stages
quite a lot of data	relearn Finetune a few layers	Finetune a larger number of layers

→ what model train on it

to prevent overfitting:

- ① Data Augmentation
- ② Dropout
- ③ Early Stopping
- ④ Batch normalization

Example to solve



Forward Pass

$$net_{h1} = w_1 * i_1 + w_2 * i_2 + b_1 * 1$$

$$net_{h1} = 0.15 * 0.05 + 0.2 * 0.1 + 0.35 * 1 = 0.3775$$

$$out_{h1} = \frac{1}{1+e^{-net_{h1}}} = \frac{1}{1+e^{-0.3775}} = 0.593269992$$

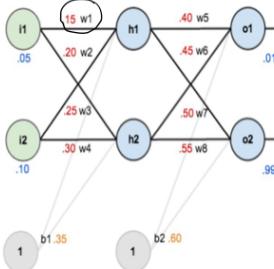
$$out_{h2} = 0.596884378$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$net_{o1} = 0.4 * 0.593269992 + 0.45 * 0.596884378 + 0.6 * 1 = 1.105905967$$

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}} = \frac{1}{1+e^{-1.105905967}} = 0.75136507$$

$$out_{o2} = 0.772928465$$



- Given inputs 0.05 and 0.10, we want the neural network to output 0.01 and 0.99.

Note: E is the loss

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$

$$E_{o2} = 0.023560026$$

$$E_{total} = E_{o1} + E_{o2} = 0.274811083 + 0.023560026 = 0.298371109$$

$$w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$$

$$w_2^+ = 0.19956143 \quad w_3^+ = 0.24975114 \quad w_4^+ = 0.29950229$$

from example 1 find $\frac{\partial J}{\partial w_5}$

$$\frac{\partial E}{\partial w_5} \leftarrow \frac{\partial E}{\partial o_1} \frac{\partial o_1}{\partial y_1} \frac{\partial y_1}{\partial w_5}$$

from example 1 find $\frac{\partial J}{\partial w_1}$

$$\frac{\partial E}{\partial w_1} = \left(\frac{\partial E}{\partial oH_1} \right) \cdot \frac{\partial oH_1}{\partial H_1} \frac{\partial H_1}{\partial w_1}$$

$$\frac{\partial E}{\partial y_1} \cdot \frac{\partial y_1}{\partial oH_1} \frac{\partial oH_1}{\partial H_1} + \frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial oH_1} \frac{\partial oH_1}{\partial H_1}$$

$$\frac{\partial oH_1}{\partial H_1} = \sigma(h_1)(1 - \sigma(h_1)) = 0.593 * (1 - 0.593) = 0.241$$

$$\frac{\partial H_1}{\partial w_1} = i_1 = 0.05$$

$$\frac{\partial E}{\partial y_1} = -(T_1 - o_{y_1}) = -(0.01 - 0.751) \quad 0.7413$$

$$\frac{\partial o_{y_1}}{\partial w_1} = \sigma(y_1)(1 - \sigma(y_1)) = 0.7513 * (1 - 0.7513) \quad 0.1868$$

$$\frac{\partial y_1}{\partial H_1} = w_5 = 0.4$$

$$\frac{\partial E}{\partial y_2} = -(T_2 - o_{y_2}) = -(0.99 - 0.773) = 0.217$$

$$\frac{\partial o_{y_2}}{\partial y_1} = \sigma(y_2)(1 - \sigma(y_2)) = 0.773 * (-0.773) \quad 0.1755$$

$$\frac{\partial y_2}{\partial H_1} = w_6 = 0.55$$

0.0554

+

-0.0209

0.0345

$$\frac{\partial E}{\partial w_1} = 0.0345 * 0.7413 * 0.1868 = 0.000416015$$

