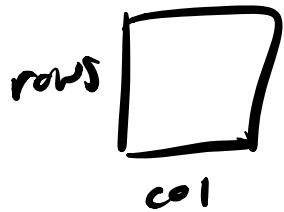


keras - DL framework
built on tensorflow

numpy - height, width, size



Aspect Ratio

CNN

- pixel, not feature extraction

Training

- train, test, validation
- sane aspect ratio
- tune params
- ~1000 - 5000 example images / class
- class imbalances

process

- gather dataset
- preprocess
- split the dataset
- train the classifier ✓
- evaluate

kNN

- easy to implement
 - not actually "learning"
 - requires comparison to ALL data
 - ↳ ACOT of memory \rightarrow not scalable

• ANN

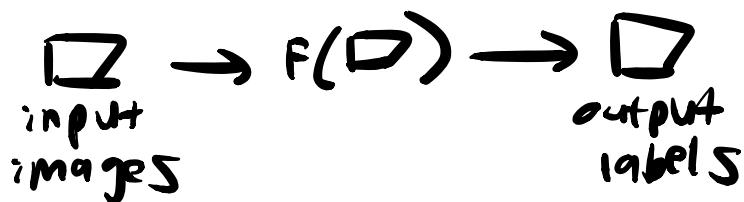
↳ make estimates
→ sacrifice "correctness"

parametrized learning

- learning the necessary parameters of a given model

Linear classifier

- N images in RGB, sized 32x32 pixels



• loss \leftarrow minimize \mathcal{L}

$$\bullet D = 32 \cdot 32 \cdot 3 = 3072 \leftarrow \begin{matrix} \text{image} \\ \text{pixels} \end{matrix}$$

width height channels

$$L = 3 \leftarrow \text{labels}$$

f = $\underset{\text{optimize } L}{\text{W}} x_i + \underset{\text{minimizing loss}}{b}$

A hand-drawn diagram consisting of three horizontal layers. The top layer contains two upward-pointing arrows. The middle layer contains one downward-pointing arrow. The bottom layer contains two upward-pointing arrows. Each layer is enclosed in a bracket-like shape. There are small 'K' symbols at the ends of the top and bottom layers.

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \boxed{\begin{array}{c|cc} 1 & 0 \\ \hline 0 & 0 \end{array}} \quad 3 \times 2 \quad 2 \times 1$$

$3 \times 3072 \quad 3072 \times 1 \quad 3 \times 1$

LOSS

- Hinge loss \leftarrow calculates margin (SVMs)

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

↓
 jth class ↑ score of
 ith data point correct

$$s_j = f(x_i, w), \text{ loss score}$$

- Squared Hinge Loss

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$$

- Avg loss

$$L = \frac{1}{N} \sum_i L_i$$

- Cross-entropy loss (softmax classifiers)

- probabilities for each class take 1

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right)$$

bigger the better
 • correct value guessed

↓
 lower the better

$$\stackrel{i}{=} -\log \left[P(Y=y_i | X=x_i) \right]$$

negative log likelihood

optimization

→ updates EACH iteration

slow

gradient descent

- finds 1 good enough solution

even if it only finds a local minima

- assume it's a convex problem

(aka that we will find the global minima)

→ combine w and b into one matrix

→ find optimal α learning rate

SGD

loop until
- epoch done
- loss \downarrow
- loss not improved

batch = next-training-batch(data, 256)

$w_{gradient} = \text{evaluate.gradient}(loss, data, w)$

batch

$w += -\alpha \cdot w_{gradient}$

→ better to train with MORE EPOCHS

and SMALLER α

more noisy
faster convergence
no negative effects to loss/
accuracy

stochastic gradient descent

- updates weight via small batches of training data

- batchsize:

32

64

128

256

- multiple weight updates per epoch

→ gives model more chances to

→ more updates added to the

learn update's have to do with weight matrix

- Update S

$$w = w - \alpha \nabla_w f(w)$$

- Momentum

$$\vec{v} = \gamma \vec{v}_{t-1} + \alpha \nabla_w f(w)$$

always
use

$$w = w - \vec{v}_t$$

(incl
concepts) \rightarrow increase strength of updates in the same direction
decreases strength when gradients switch directions

$$\gamma = 0.9 \text{ usually}
(or 0.5 until stabilizes)$$

- Acceleration

- corrects momentum
- use for small datasets
quid in large datasets

- Regularization

- helps w/ generalization
prevents overfitting

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(w)$$

$$w = w - \alpha \nabla_w f(w) - \lambda R(w)$$

L2 regularization (weight decay)

$$R(w) = \sum_i \sum_j w_{ij}^2$$

L1 regularization

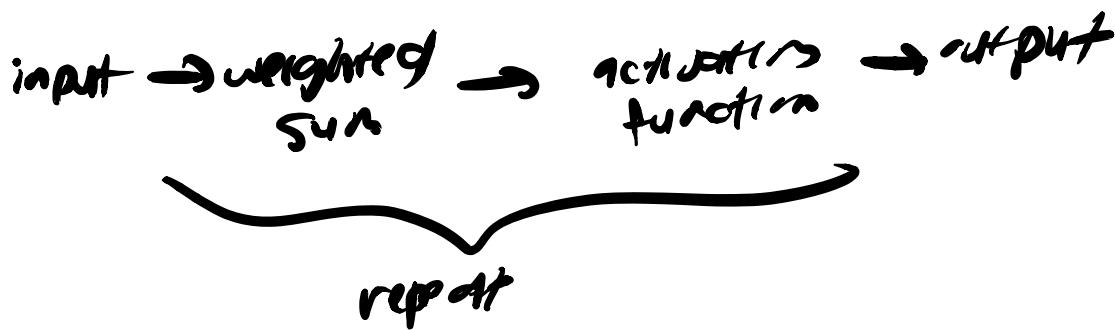
$$R(w) = \sum_i \sum_j |w_{ij}|$$

Elastic Net

$$R(w) = \sum_i \sum_j \beta w_{ij}^2 + \gamma |w_{ij}|$$

Neural Networks

- neurons fire or doesn't fire



- activation functions \leftarrow help learn complex data
linear \rightarrow nonlinear

step

sigmoid

tanh

ReLU

leaky ReLU

ELU

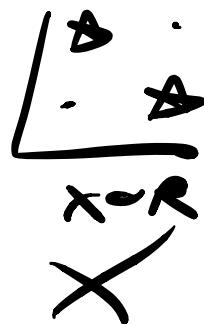
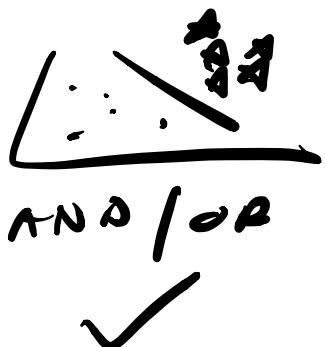
- tips

use ReLU for baseline

$\rightarrow \text{ReLU}$
 $\rightarrow \text{FCU}$

Perception

- linear
- update to until get all correct prediction



Backpropagation

$\xrightarrow{\text{forward pass / propagation phase}}$

$\xleftarrow{\text{backward pass}}$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ij}}$$

w: weight
E: loss
o: output
net: net output

XOR

- can't be solved w/ one layer (perception)
b/c not a linear problem (looping)
- requires multiple layers (NN)
which requires backprop to update weights

(similar to looping in perceptrons)

NN w/ backprop allows us to learn patterns
in datasets that are nonlinearly separable

Backprop

- generalization of gradient descent algorithms
specifically used to train multi-layer
feedforward networks

- requires:

forward pass

backward pass

MNIST vs CIFAR

- basic feedforward network w/ only fully connected
networks not suitable for complex images

\Rightarrow CNN !!!

NN System

- ① dataset
- ② model / architecture
- ③ loss function
- ④ optimization method

creating own NN is
long/prone to bugs

→ use

. Keras

. Theano

. TensorFlow

. PyTorch

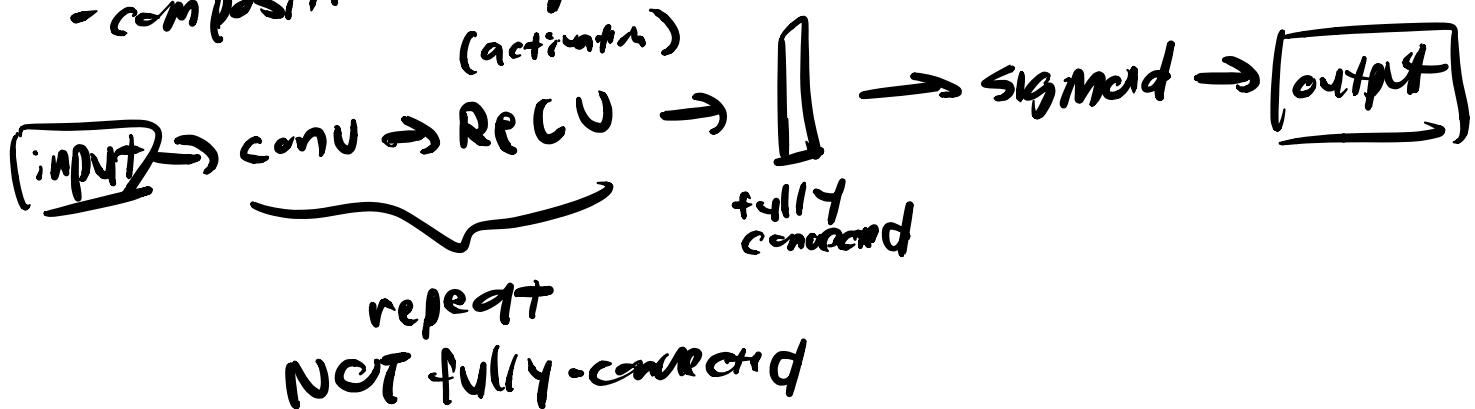
weight initialization

- constant
- uniform / normal distributions
 - LeCun
 - Glorot / Xavier
 - He et al / Kaiming / NSRA

CNN

- local invariance
- compositionality

(activation)



- convolution

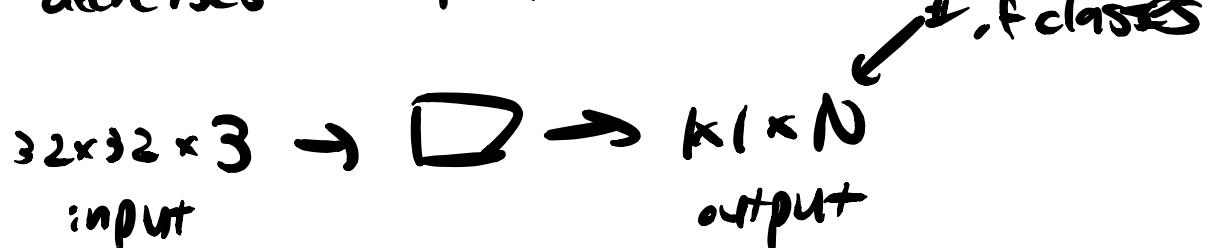


- must use odd kernel size to ensure valid center

- apply rescale to $[0, 255]$ b/c CNN can increase val/size

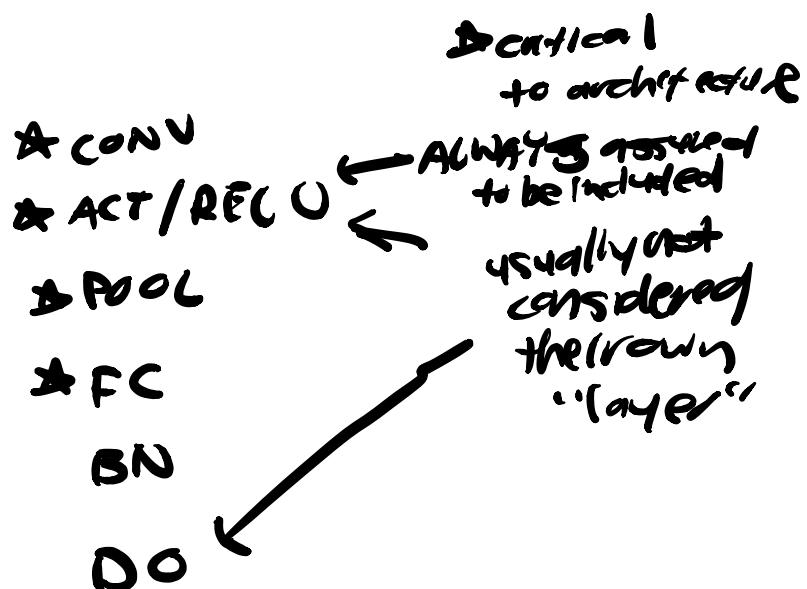
- local connectivity

- CNNs are 3D - width, height, depth/channels
 - subsequent layer neurons only connected to a small region of layer before via weights
 - decreases # of parameters
classes



- layer type S

- conv
 - activation
 - pooling
 - fully-connected
 - batch norm
 - dropout



i.e. INPUT \rightarrow CON \rightarrow RECU \rightarrow FC \rightarrow SOFTMAX

- local connectivity

- connect each neuron only to a local region of the receptive field

i.e. $32 \times 32 \times 3$ input value

3×3 receptive field/kernel size

③ channel depth

$$\Rightarrow 3 \times 3 \times 3 = 27 \text{ weights}$$

i.e. $16 \times 16 \times 94$ in input
where $\xrightarrow{\text{FxF filter}}$
 3×3 receptive field / kernel size

④ channel depth

$$\Rightarrow 3 \times 3 \times 94 = 846 \text{ weights}$$

- size of output volume

- depth $\leftarrow K$ # of filters/channels (K)
depth column: filters belonging to same layer
- stride

- for each step, create a new depth column around local region of image where we convolve each of the K filters

- store output in 3D volume

- small stride \rightarrow larger output volume
- large stride \rightarrow smaller output volume

- reduce spatial dimensions of input

• zero-padding size

- output volume matches input volume

- critical for deep CNN that apply

multiple filters - one to each other

- no extra output size

$$\left(\frac{W - F + 2P}{S} \right) + 1 \quad \xleftarrow{\text{MUST be integer}}$$

W : square input image

F : filter size

S : stride

P : padding

- output of conv layer $W_{out} \times H_{out} \times D_{out}$

$$W_{out} = \frac{W_{in} - F + 2P}{S} + 1$$

$$H_{out} = \frac{H_{in} - F + 2P}{S} + 1$$

$$D_{out} = K$$

- output of activation layer \leftarrow immediately follows conv

$$W_{out} = W_{in}$$

$$H_{out} = H_{in}$$

$$D_{out} = D_{in}$$

- output of pooling layer \leftarrow between layers after activation

- reduce # parameters/computations

- controls overfitting

- max pooling \leftarrow middle of NN

- avg pooling \leftarrow final layer

- usually

$$S=1, S=2$$

poolsize 2×2 or 3×3 (early in NN)

- $H_{out} = \frac{H_{in}-F}{S} + 1$

$$H_{out} = \frac{H_{in}-F}{S} + 1$$

$$D_{out} = D_{in}$$

- common to NCF use pooling in middle of NN,
and only ^{regress w/ FC} avgpooling

Batch Normalization

- reduces # epochs needed to train NN

- "stabilizes" training \leftarrow more stable
loss curve

- lower final loss

- placed AFTER nonlinearity activation

Dropout

- regularization

- for each mini-batch, randomly drop neurons

- reduce overfitting

Architecture

- CONV, FC, BN \leftarrow learn params
- POOL, RELU/ACT \leftarrow perform operation

INPUT \rightarrow [(CONV \rightarrow RELU) \cdot N \rightarrow pool?] \rightarrow [FC \rightarrow RELU] \cdot K \rightarrow FC

$$0 \leq N \leq 3$$

$$0 \leq M$$

$$0 \leq K \leq 2$$

Examples

- INPUT \rightarrow FC
- INPUT \rightarrow [CONV \rightarrow RELU \rightarrow pool] \cdot 2 \rightarrow FC \rightarrow RELU \rightarrow FC
- INPUT \rightarrow [CONV \rightarrow RELU \rightarrow CONV \rightarrow RELU \rightarrow pool] \cdot 3
 \rightarrow [FC \rightarrow RELU] \cdot 2
 \rightarrow FC
- INPUT \rightarrow CONV \rightarrow RELU \rightarrow FC
- AlexNet +
- VGGNet

stacking multiple CONV layers before POOL

- allows CONV layers to develop more complex

features before destructive pooling op

- POOL vs CONV

↑
scatters just NO pool,
and use conv step size to
reduce volume size

rules

- SQUARE input images

- take advantage of linear alg opt libraries
- common sizes:

32, 64, 96, 128, 227, 229

- input layer should be divisible by 2 mult times

- tweak filtersize & stride
- allows spatial inputs to be sampled down via pool operation

- conv layers / filters

• general: $3 \times 3, 5 \times 5$

• only in first conv: $7 \times 7, 11 \times 11$

• too large

- will reduce spatial dim of vol

- too quickly
- stride $s=1$ \leftarrow CONV
learn features
 - use pool to downsample only
 - apply zero-padding
 - to match input/output dim
 - max pooling
 - for downsampling
 - 2x2 receptive field size
 - $s=2$ stride
 - $NO \geq 3$
 - batch normalization
 - expensive op
 - \rightarrow SIGNIFICANT INCREASE IN COMP TIME
 - use off-the-shelf
 - stabilize training to tune hyperparams
 - dropout
 - $p=0.5$
 - r between FC layers

- (sorters) b/w pool & conv,
 $p=0.1 \sim 0.25$
- battles over fitting

- pool \rightarrow conv

- initially explore/learn w/ pool

maxpool \rightarrow conv \rightarrow avg pooling

- rotation / scale

- not exactly rotation/invariant
but filters have learned

- translation

• CNN EXCELS

Keras

- settings

• epsilon: $1e^{-7}$

• float32

• channels_last

• backend: tensorflow

$\xrightarrow{\text{openCV}}$
rows, cols, channels

- C(FAR)

- very easy to overfit due to limited # of low res training samples

- want:

boost accuracy

reduce overfitting

- shallow net

- input \rightarrow conv \rightarrow relu \rightarrow fc \rightarrow output

Model serializations

- saving & loading a trained model /

model.save (filepath)

load_model (filepath)

LeNet + MNIST

hello world!!!

benchmarking
gradient based SGD

- recognizing handwritten digits

INPUT \rightarrow [conv \cdot tanh \cdot pool] \cdot 2

RELU

4 layers!
- 2 conv
- 2 FC

\rightarrow FC \rightarrow tanh \rightarrow FC \rightarrow sigmoid

"hidden"

→ OUTPUT

const. load data
→ 764-d vector
(28x28 grayscale)

- deep

→ CNN layers increase

spatial input dim decrease

Python Imports

① Network Architectural
i.e. shallownet, LeNet

② optimizer

i.e. SGD

③ preprocess / split test train

④ classification Report

VGGNet

① only uses 3×3 filter for CONV layers

② stack MULTIPLE CONV → FC layers
before applying pool

Layer Type	Output Size	Filter Size / Stride
INPUT IMAGE	$32 \times 32 \times 3$	
CONV	$32 \times 32 \times 32$	$3 \times 3, K = 32$
ACT	$32 \times 32 \times 32$	
BN	$32 \times 32 \times 32$	
CONV	$32 \times 32 \times 32$	$3 \times 3, K = 32$
ACT	$32 \times 32 \times 32$	
BN	$32 \times 32 \times 32$	
POOL	$16 \times 16 \times 32$	2×2
DROPOUT	$16 \times 16 \times 32$	
CONV	$16 \times 16 \times 64$	$3 \times 3, K = 64$
ACT	$16 \times 16 \times 64$	
BN	$16 \times 16 \times 64$	
CONV	$16 \times 16 \times 64$	$3 \times 3, K = 64$
ACT	$16 \times 16 \times 64$	
BN	$16 \times 16 \times 64$	
POOL	$8 \times 8 \times 64$	2×2
DROPOUT	$8 \times 8 \times 64$	
FC	512	
ACT	512	
BN	512	
DROPOUT	512	
FC	10	
SOFTMAX	10	

Note TIPS

- smaller α learning rate
 - smaller weight updates
 - reduces overfitting
- USE batch norm!!!

Learning Rate Scheduler

- "learning rate annealing"
- "adaptive learning rates"
- standard weight update formula:
$$w \leftarrow -\alpha \cdot \nabla$$

↑
typically
.1 or 0.01
- beneficial to decrease learning rate over time
- learning rate scheduling
 - decrease gradually per epoch
 - drop at specific epochs

- process

- find good weights early & use high α

~~time-based~~ - tune weights later w/ small α

- learning rate update formula

$$\alpha = \alpha_0 \left(\frac{1}{1 + \text{decay} \cdot \text{iterations}} \right)$$

↑

$$\text{decay} = \frac{\alpha_0}{\# \text{epochs}}$$

~~drop-based~~

- drop learning rate usually by:

$F = 0.15$ or order of magnitude
every fixed # epochs

$$\alpha_E = \alpha_0 F^{\frac{E}{D}}$$

- E: current epoch
- F: learning rate drop factor
- D: specified drop epoch

faster drop (0.25)

→ can stagnate learning of weights
(drop too fast)

BE prepared to spend ALOT of time
training/tuning networks!!!

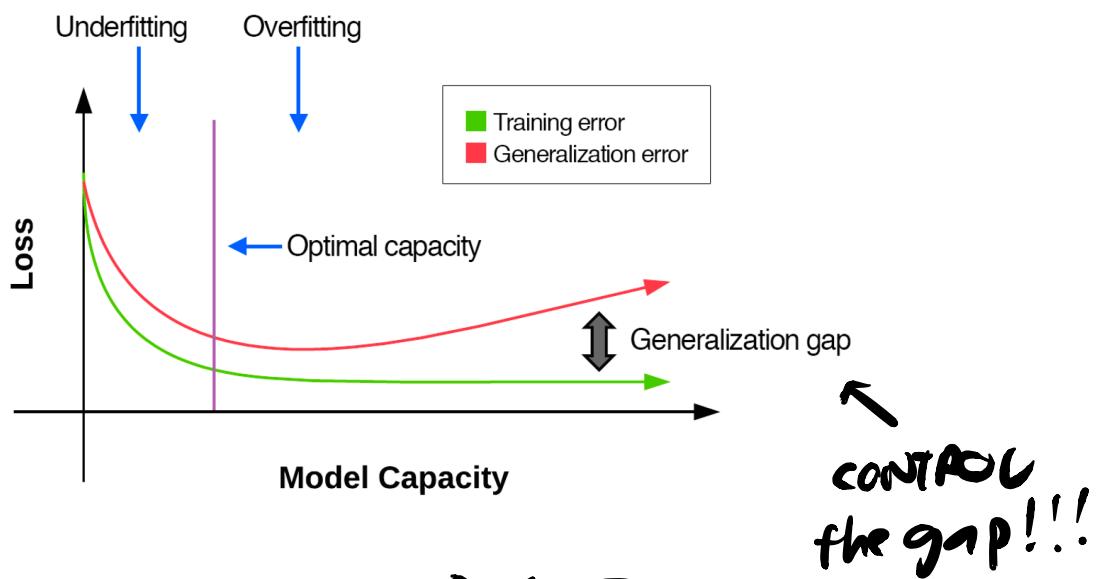
- even simple datasets can take 100s of experiments to get a high accuracy!
- training deep NNs :
 - part science
 - part art ↗ intelligent guess-and-check

spotting Under/over-fitting

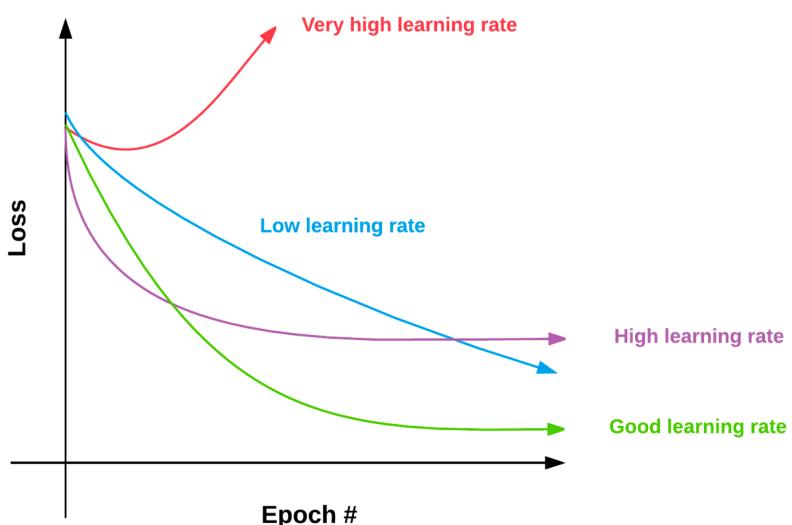
- Training Monitor
 - visualize loss graph after every epoch
- Training goals :
 - reduce training loss
 - ensure gap b/w training/testing is reasonably small

- Adjusting capacity of NN

- increase capacity
 - add more layers + neurons
- decrease capacity
 - fewer layers / neurons
 - apply regularization



- Effects of Learning Rates



- Assessing your model

- are you applying regularization techniques?
 - is the learning rate too high?
 - is your network too deep?

- If validation loss < training loss (AWESOME!)

- data augmentation during training
 - form of regularization enables better generalization
- not training "hard enough"
 - try increasing α learning rate & tweaking regularization strength

- If you think you see signs of underfitting

→ train for another 10-15 epochs to ensure

- If underfitting

→ add more neurons

- If overfitting

→ ^{various} regularizing techniques

checkpointing

- combines :

save/serialize models to disk
monitoring under/overfitting

⇒ Keras: ModelCheckpoint (builtin callback)

- 2 methods:

- checkpoint incremental / improvements

- checkpoint only the best model

Architecture Visualization

- used for:

- debugging ← HIGHLY ENCOURAGED
- publication

- visualizes

(batch_size, height, width, depth)

Goal of ImageNet challenge

- classify input image into 1000 categories

Pretrained Networks

- model size based on # parameters

NOT # training data

- Keras

- VGG16
- VGG19
- ResNet50
- Inception V3
- Xception

} works out
of the box!

VGGNet

- only 3×3 filters
- very slow to train
- large weights ~500MB

ResNet

- deep networks trained via only SGD
- global avg pooling
- smaller weights than ResNet
~100MB

} 224x224 input

Inception v3

- multi-level feature extractor
 - smaller weights than VGGNet and ResNet
 - originally called "GoogLeNet"
- ~90 - 100 MB

299 x 299 input
[]

Xception

- extension of Inception
- smallest pre-trained networks included in the Keras library!

~90 - 100 MB

Squeeze Net

- not included in Keras
- extremely small weights

~4.9 MB

TYPICAL IMAGE SIZES IN ImageNet

- 124
- 227
- 256
- 299

Responsible disclosure

- how to disclose a vulnerability
- contact stakeholders directly
to ensure they know there is an issue!
→ as opposed to immediately
posting vuln. online!

Dataset organization

|root|classname|imagename.jpg

→ |dataset|{1-9}|example.jpg

ROI

"region of interest"

Breaking captchas

- case study

- download set of images
- label / annotate images
for training
- training a CNN on your
custom dataset
- evaluating / testing the
trained CNN

- obtaining and labeling a dataset
can be half (if not more the battle!)

→ try to use traditional
techniques to speed up
labeling process