

Deep learning

CNN and supervised learning

March 2019 @ TELECOM Nancy
T. BAGREL

Introduction

Interesting NN architectures (supervised)

Convolutional Neural Network (CNN)

- ▶ Supervised :(
- ▶ Most common architecture for image processing
- ▶ Position-invariant feature extraction (with shared weights)
- ▶ Very efficient

Introduction

Interesting NN architectures (unsupervised)

Restricted **B**oltzmann **M**achine (RBM)

- ▶ Unsupervised!
- ▶ General feature extraction
- ▶ Basic block of **D**eep **B**elief **N**etworks

Autoencoders (AE)

- ▶ Unsupervised!
- ▶ Can be used for denoising
- ▶ Can be used for general feature extraction

Autoencoders

Structure and principle

Principle

- ▶ **Simple** : reproduce the input at the output
- ▶ **How** : pass the input (image) through the NN, compute error from $\text{diff}(I, O)$, adjust weights with GD

Feature extraction

Bottleneck : “center” of the NN where there is the lowest number of neurons. Activation in the “bottleneck” for a given input (image) represents its compressed form.

Autoencoders

Denoising

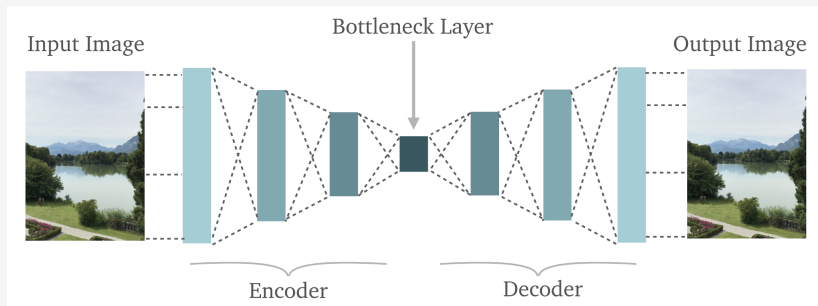


FIGURE – Autoencoder with feature extraction purpose

Denoising

Denoising can be achieved with roughly the same architecture

CNN Components

Gradient Descent

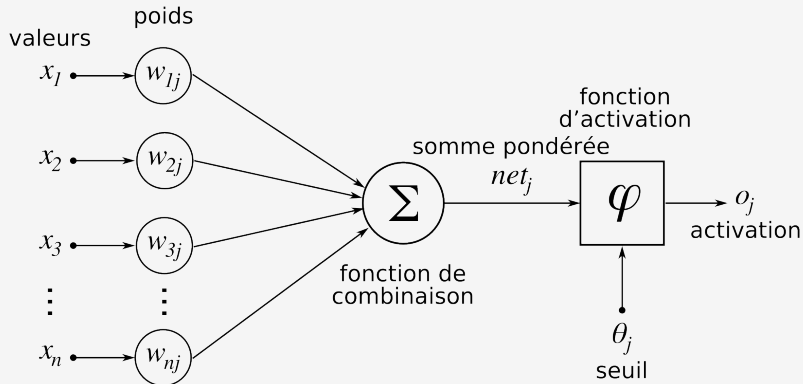


FIGURE – Neuron (perceptron) schema

CNN Components

Gradient Descent

Problem solved by GD

In which direction do we need to update our weights to minimize the error?

Last layer neuron n situation

E : error, L : loss function, e : expected value, o : obtained value, φ : activation function, $x_{n,\cdot}$: inputs for the last layer neuron, $w_{n,\cdot}$: weights for the last layer neuron

$$E = L(e, o) \tag{1}$$

$$= L \left(e, \varphi \left(\sum_k x_{n,k} w_{n,k} \right) \right) \tag{2}$$

CNN Components

Gradient Descent

What we need to know

- ▶ **loss function** $L(e, o)$ and $\frac{\partial L}{\partial o}$
- ▶ **activation function** $\varphi(r)$ and $\frac{d\varphi}{dr}$

What we can compute

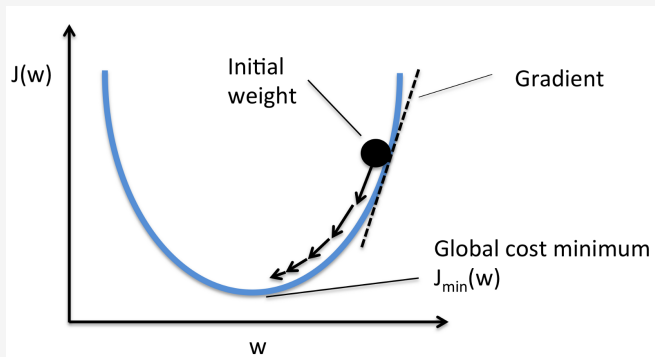
With the **chain rule**, we can compute $\frac{\partial E}{\partial w_{n,k}}(e, o) \quad (\forall k)$

CNN Components

Gradient descent and adjustment made

Weights adjustment (α : learning rate)

$$w_{n,k} \leftarrow w_{n,k} - \alpha \cdot \frac{\partial E}{\partial w_{n,k}}(e, o)$$



CNN Components

Learning rate and advanced adjustment methods

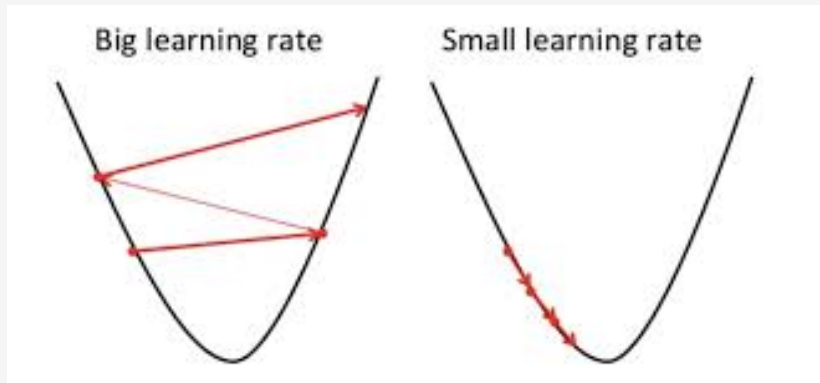


FIGURE – Importance of an appropriate learning rate

CNN Components

Advanced adjustment methods

Called optimizers in **TensorFlow**

Roles

- ▶ Try to prevent staying in a non-global local minima during learning
- ▶ **Examples** : Nesterov Momentum, Adam...

Used in my NNs

I use Nesterov's momentum :
`optimizers.SGD(nesterov=True, momentum=M)`
because it's recommended on most tutorials

CNN Components

GD vs SGD vs BGD

Each one is a different manner to apply Gradient Descent

Differences

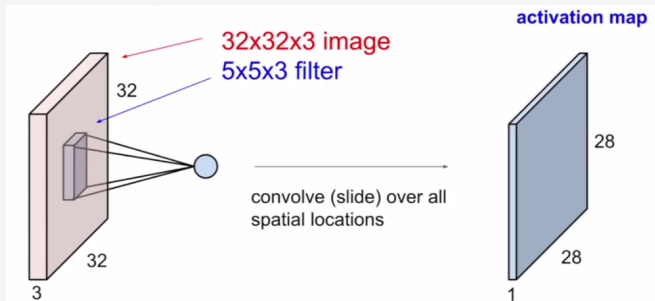
- ▶ **(basic) Gradient Descent** : weights are updated after every epoch (whole dataset pass) → **very costly** but more stable than SGD!
- ▶ **Stochastic Gradient Descent** : weights are updated after **every** sample → **very fast**!
- ▶ **Batch Gradient Descent** : weights are updated after N samples ($10 \leq N \leq 100$: batch size) → **best of both worlds**!

CNN Components

Convolutional layer

The (3D) **dot product** between the **filter weights** and an **input image zone** of the same size/shape gives an activation value. Doing this operation for all possible image zones (with **overlap**) gives the **activation map** for this filter.

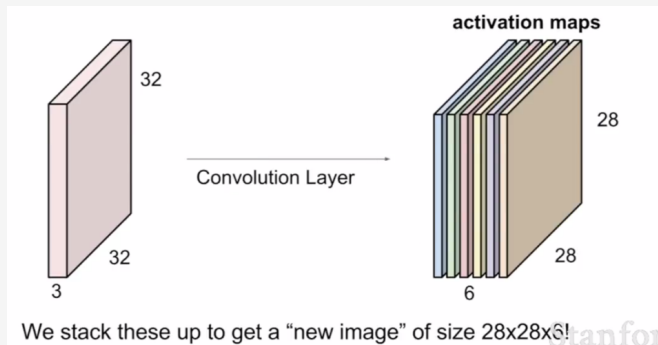
It explains the position invariance of feature extraction for CNNs!



CNN Components

Convolutional layer

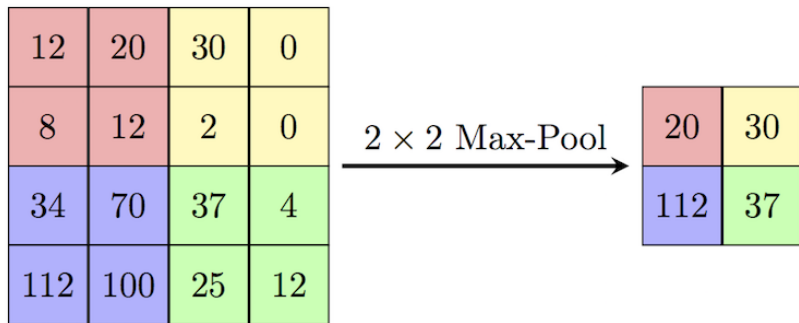
A **convolutional layer** is composed of several filters, each one with its **own weights**. Hence, the output of the convolutional layer is a **stack of activation maps** (one for each filter)



CNN Components

Pooling layer

A **pooling layer** is used to downsize/down-sample its input. **All** the values in a zone of the input volume results in **one** value in the output volume, by application of a simple function. Input zones **must not overlap** here!
A very common example is the *max pooling* 2×2 :



CNN Components

Output layer for classification

The output layer is used to extract classification information from a final dense/convolutional/pooling layer.

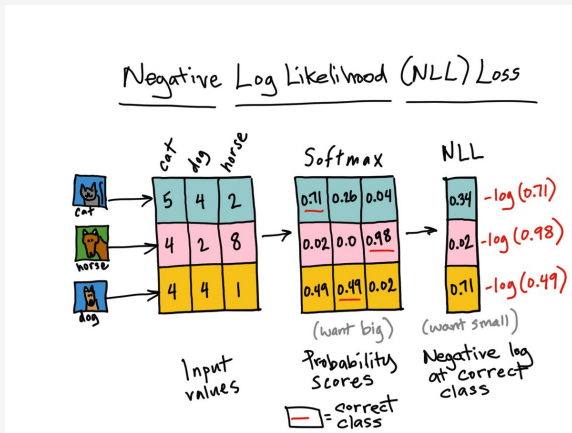
Types

- ▶ For binary classification : Dense layer with **sigmoid** activation function and **1** neuron (1 output value)
- ▶ For multiclass classification : Dense layer with **softmax** activation function and **N_{classes}** neurons, followed by a **$\text{arg_max}(\cdot \cdot \cdot)$** function

CNN Components

Loss function for classification

Cross entropy (aka. Negative Log Likelihood) is the loss function used in classification, which works well with the **sigmoid** or **softmax** activation functions (output layer)



CNN Components

CNN schema for multiclass

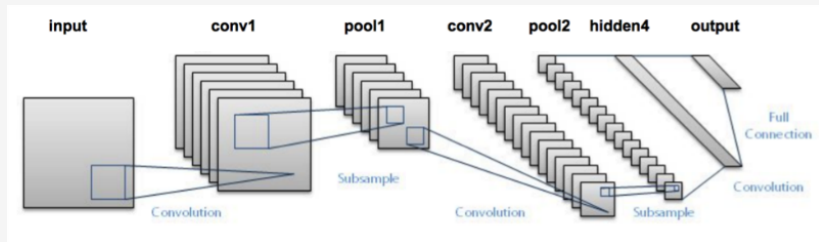


FIGURE – CNN architecture

Medical imaging

Availability

The issue

There are medical image datasets available online, but

- ▶ High resolution
- ▶ Difficult to find the differences or the good class for an untrained eye
- ▶ File format and dataset organisation is not always normalized

Examples

- ▶ List of public datasets : bit.ly/2E0ziVp, bit.ly/2u0Da0Q
- ▶ Datasets : bit.ly/2F34jqb (cancer cells), bit.ly/2u99w9L (neuroimaging), bit.ly/2cGnNQL (Alzheimer)...

Medical imaging

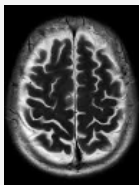
What to do so?

Requirements

- ▶ Large dataset
- ▶ Low-res images (max. 256 x 256)
- ▶ Something which looks like medical images

Idea

Generate **brain 2D images** and add random “tumors” on them



brain_tumor_factory

Let's get our hands dirty!

How to generate?

Home-made C program : `brain_tumor_factory`

- ▶ take a `single` base image
- ▶ add tumor(s) randomly depending on the chosen scenario
- ▶ alter the resulting image, to ensure that each output image is a bit unique

Program size : \simeq 500 LOC

Output : \simeq 1000 images/s

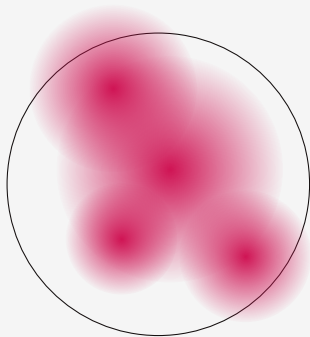
Image size : `96 x 128 x 1` (grayscale)

brain_tumor_factory

Creating tumors : tumor function

Instructions

Take a circle of radius R_{area} , and put $n \in \llbracket n_{\min}, n_{\max} \rrbracket$ “soft discs” of radius $r_i \in [r_{\min}, r_{\max}]$ with their center placed randomly inside the chosen area



brain_tumor_factory

Altering the resulting image : `alter_full` function

Alterations

- ▶ `shrink` : shrink the image on x-axis and/or y-axis
- ▶ `move` : shift the image on x-axis and/or y-axis
- ▶ (`ghost_blur` 0.1%) : simulate motion blur on a shot
- ▶ `alter_contrast` : \pm [brightness range], with custom "brightness center"
- ▶ `alter_shade` : global flat \pm [brightness]
- ▶ `alter_rdc` : per-pixel flat random \pm [brightness]

For each image, alterations parameters are random numbers taken from a `uniform` or `normal` distribution

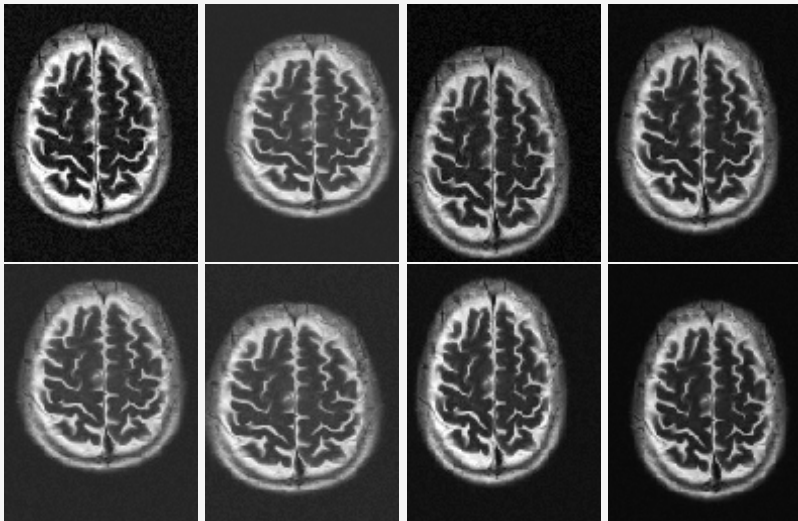
brain_tumor_factory

Base image



brain_tumor_factory

Altered images

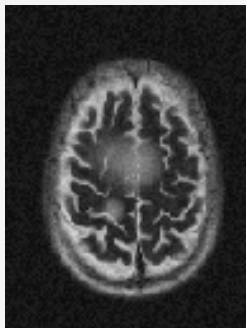


brain_tumor_factory

Syndrome A

Description

- ▶ $\llbracket 1, 2 \rrbracket$ big tumor(s)
- ▶ $\llbracket 0, 2 \rrbracket$ medium tumor(s)

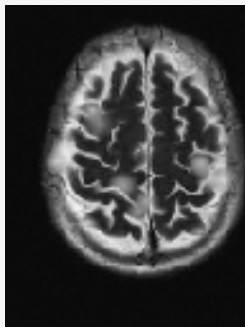


brain_tumor_factory

Syndrome B

Description

- ▶ $\llbracket 4, 8 \rrbracket$ medium tumor(s)
- ▶ $\llbracket 0, 3 \rrbracket$ small/tiny tumor(s)

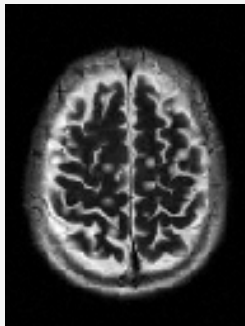


brain_tumor_factory

Syndrome C

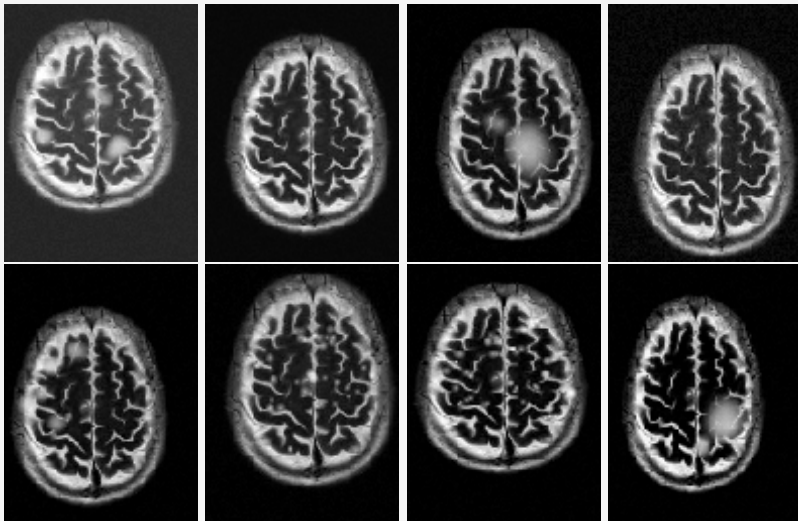
Description

- ▶ only $[[20, 48]]$ small/tiny tumor(s)



brain_tumor_factory

Guess game : EASY mode



brain_tumor_factory

Guess game : HARD mode

Observation

It turns out that it's kinda easy to find the correct syndrome type, isn't it?

Let's try the **HARD** mode then (more realistic) 😈

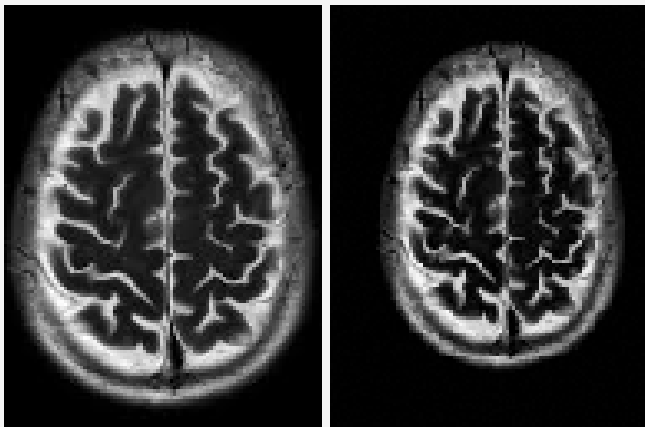
HARD mode specifications

- ▶ Tumor opacity **greatly** reduced
- ▶ Tumor size reduced a bit

The **base image** will always be placed on the **LEFT** of the image whose type must be determined

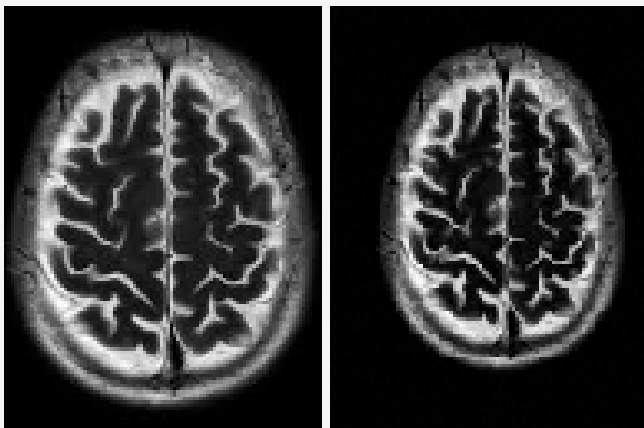
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode

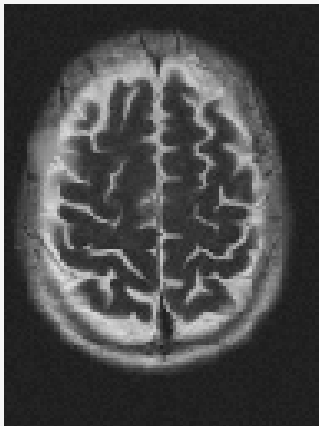
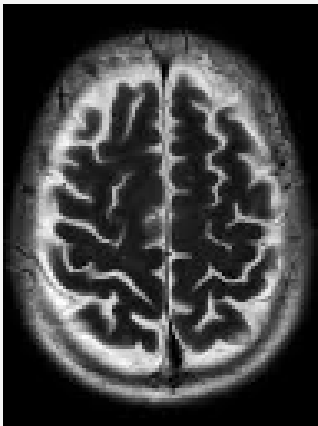


Answer

Syndrome C

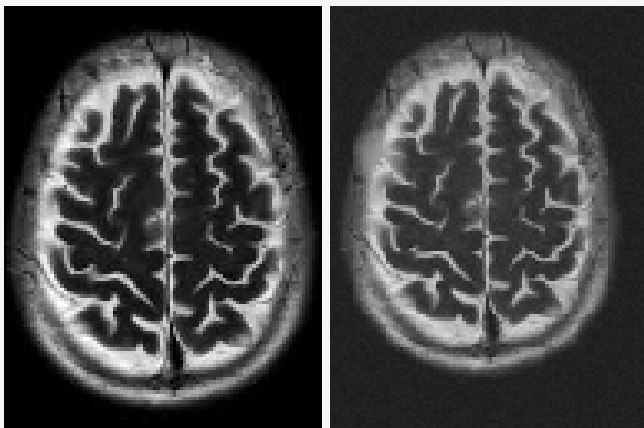
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode

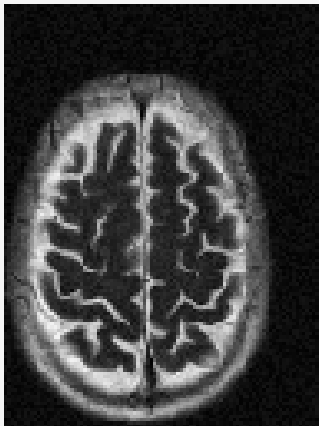


Answer

Syndrome A

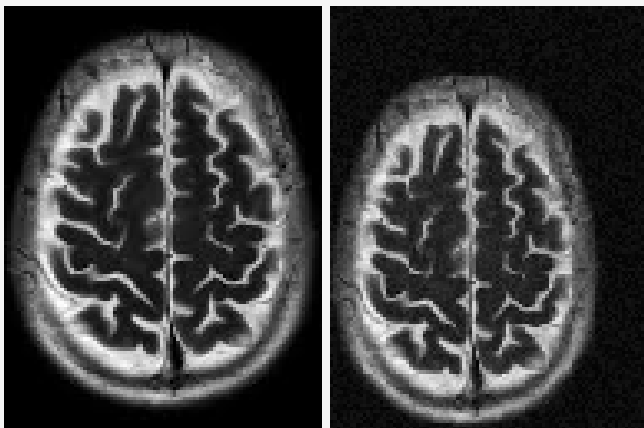
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode

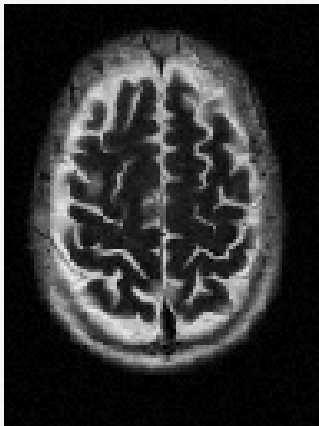


Answer

Z (healthy patient)

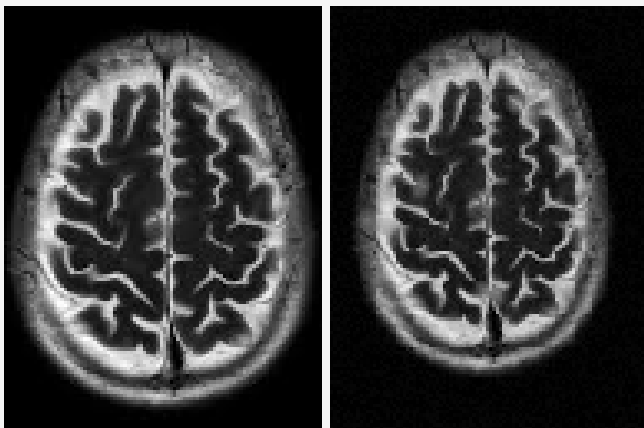
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode

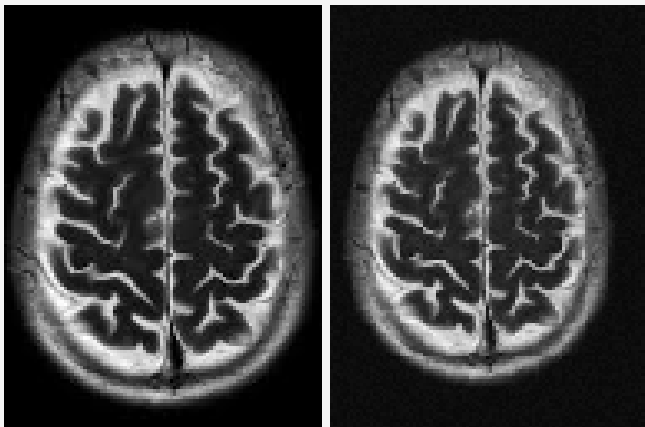


Answer

Syndrome B

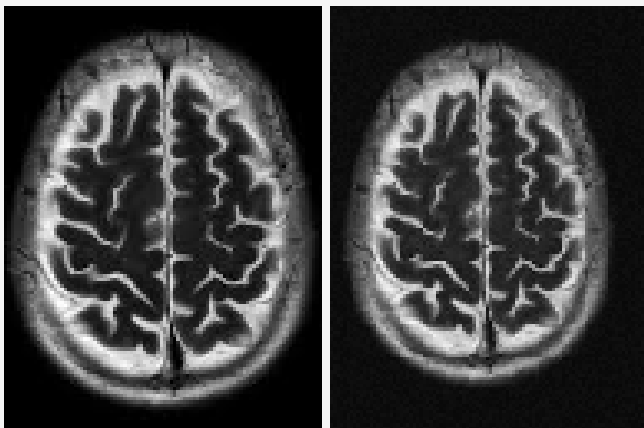
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode

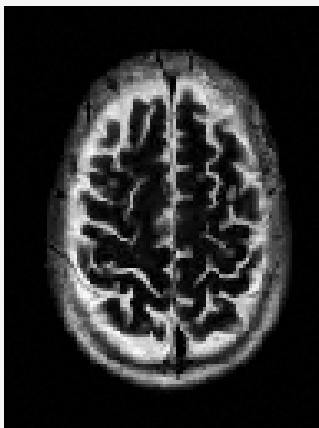
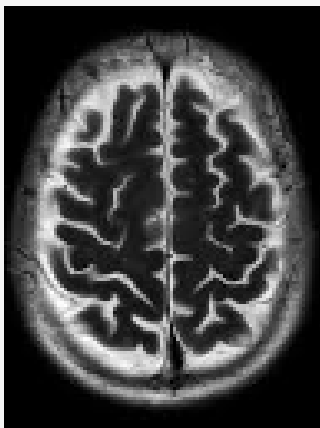


Answer

Z (healthy patient)

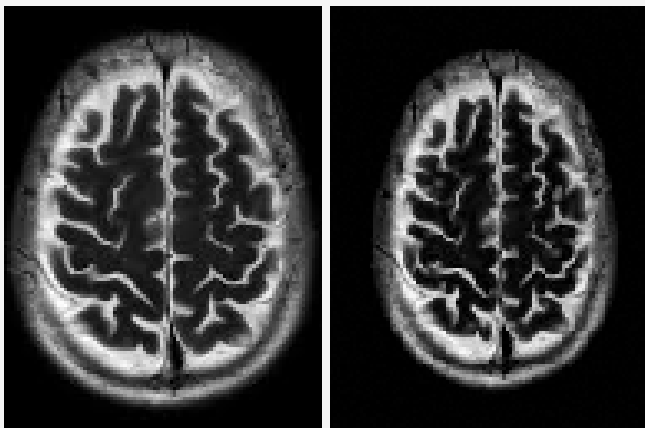
brain_tumor_factory

Guess game : HARD mode



brain_tumor_factory

Guess game : HARD mode



Answer

Syndrome C

brain_tumor_factory

Guess game : HARD mode

Observation

Not so easy this time? 😈

Let's see what deep learning can do

brain_tumor_nn

Description

Framework

- ▶ DL4j / ScalNet (Scala)
- ▶ not really fast (afterwards)
- ▶ Keras-like API

Architecture

- ▶ 3 Conv2D (MaxPooling2D inbetween) : 16 64 256
- ▶ Dense 1024
- ▶ GD (BGD) with Nesterov Momentum
- ▶ (adapted from MNIST NN examples)

brain_tumor_nn

Description

Issue (after 2 hours learning)

- ▶ Accuracy blocked at 0.25
- ▶ loss wasn't decreasing
- ▶ a **single type** predicted everytime

Why (afterwards)

- ▶ I didn't wait long enough (on hard samples)
- ▶ 3 Conv layer NN + slow framework didn't help much neither
- ▶ in addition, learning on mobile CPU (i7) is still slow
- ▶ (but afterwards) architecture wasn't really the issue

brain_tumor_nn2

Description

Framework

- ▶ Tensorflow with Keras
- ▶ Learning on CPU and then on GPU

Architecture

- ▶ Simplified architecture first : only 2 Conv2D layers
- ▶ Learning on **EASY** images first

Dataset : 100 000 L / 10 000 T **EASY** images

Very good results : **97-99 % acc** after 4 epochs (CPU)

brain_tumor_nn2

Training on HARD dataset (v1)

Architecture (v1)

2 Conv2D (64 128) + 1 Dense (1024)

Dataset : 100 000 L / 10 000 T HARD images

Learning : GPU now, about 30x faster (\simeq 1 min/epoch)

Encouraging results : over 80 % acc after an hour, but overfitting at the end : about 100 % acc on train data

Final score on test data : 90 % acc

brain_tumor_nn2

Training on HARD dataset (v2)

Smaller network, for faster convergence, and to reduce overfitting

Architecture (v2)

2 Conv2D (32 64) + 1 Dense (512)

Good results : more than 84 % acc after a few epochs, but still overfitting at the end : about 100 % acc on train data

Final score on test data : 90 % acc

brain_tumor_nn2

Training on HARD dataset (v3)

Even smaller network, in order to reduce overfitting again

Architecture (v3)

2 Conv2D (16 32) + 1 Dense (256)

Good results : more than 86 % acc after a few epochs, but still **overfitting** : about 100 % acc on train data

Final score on test data : 90 % acc

Issues

Overfitting

Description

Overfitting : when the model learns more the specificities of the train dataset samples (ie noise) than the general features of them (which are shared with the test dataset)

Causes

Too many free parameters, too few different training samples

How to counter it?

Regularization is kinda effective

Issues

Overfitting

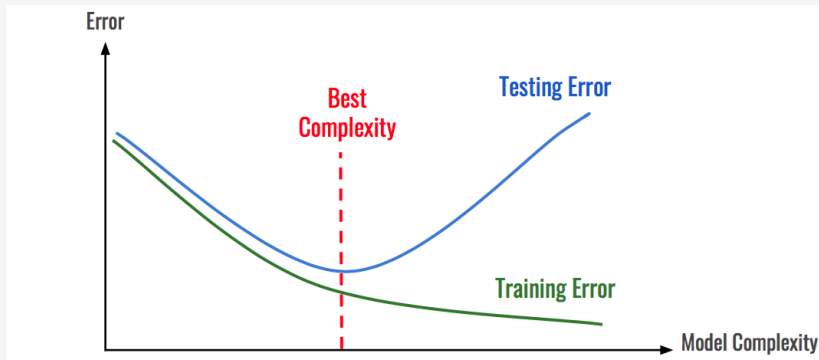


FIGURE – Memorizing is not learning!

Solutions

Regularization

General idea

Regularization : penalize the model when weights are high (decrease “freedom” of the model)

Consequences

- ▶ Counters overfitting (accuracy on test dataset will be closer to the accuracy on train dataset)
- ▶ **BUT** too much regularization can lead to **underfitting**!

Issues

Underfitting

Description

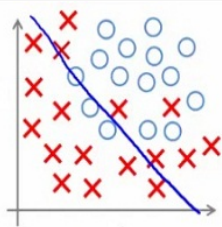
Underfitting : when there's not enough free parameters/freedom to learn the datasets (both training and testing)

How to counter it?

- ▶ Reduce regularization factor
- ▶ Add parameters (extend the model)

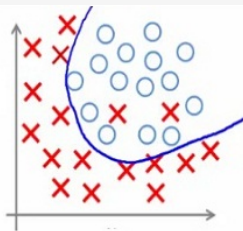
Issues

Underfitting

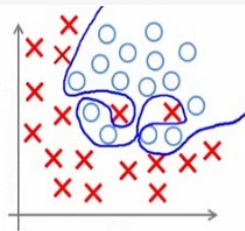


Under-fitting

(too simple to explain the variance)



Appropriate-fitting



Over-fitting

(forcefitting -- too good to be true)

FIGURE – Balance is the key here

brain_tumor_nn2

Training on HARD dataset (v4)

Same architecture as v2, with regularization : L2(0.001)

Architecture (v4)

2 Conv2D (32 64) + 1 Dense (512) L2(0.001)

Good results : over 89 % acc on test data after an hour of training, but underfitting : acc on training dataset no longer improves (blocked at 96 % acc)

Final score on test data : 91 % acc

brain_tumor_nn2

Training on HARD dataset (v5)

Larger model than v4, with less regularization :
L2(0.0005)

Architecture (v5)

2 Conv2D (64 128) + 1 Dense (1024) L2(0.0005)

Very good results : over 91 % acc on test data after an hour of training!

Final score on test data : 93 % acc

Let's try adding more layers now!

brain_tumor_nn2

Training on HARD dataset (v6)

One more conv layer than v5

Architecture (v6)

3 Conv2D (32 64 128) + 1 Dense (1024) L2(0.0005)

Very good results : over 95 % acc on test data after an hour of training!

Final score on test data : 96 % acc

One more layer? Why not!

brain_tumor_nn2

Training on HARD dataset (v7)

One more conv layer than v6

Architecture (v7)

4 Conv2D (16 32 64 128) + 1 Dense (1024) L2(0.0005)

Still very good results, but not improving...

Final score on test data : 95 % acc

What should we do now?

brain_tumor_nn2

Training in HARD dataset (v8 & v9)

Architecture (v8) – the stupid one

4 Conv2D (16 32 64 256) + 1 Dense (2048) L2(0.0001)

Overfitting again, worse than v7 (not a good move...)

But the (final) v9 architecture gave the **best results**!

- ▶ 2x more filters for each conv layer
- ▶ same regularization as v7

Architecture (v9)

4 Conv2D (32 64 128 256) + 1 Dense (1024) L2(0.0005)

Final results on test data : more than **97 % acc**!

brain_tumor_nn2

Conclusions

- ▶ The **simplest** → **most complex** approach paid off
- ▶ It was much easier to reason about the NN when the learning time was very short (with GPU)
- ▶ Each model was trained about **3 hours** (equiv. \simeq 4 days on CPU)

Key points

- ▶ Learning can be disappointing at first (the 3 firsts epochs in my case), with just one class predicted for example
- ▶ Computer can outperform Human on classification with a decent amount of learning

brain_tumor_nn2

Bias towards test data

Problem

Even if the model doesn't learn directly from the test dataset, the `save_best` approach introduces a sort of “indirect learning” from the test dataset

Solution (not implemented here)

Use a 3rd dataset for the final comparison between models, called **validation dataset**

CNN IRL

Real world examples and dataset augmentation

The `brain_tumor_factory` was useful because I had :

- ▶ a normalized dataset
- ▶ as many samples as I wanted

But this is not often the case in real life examples

Solution

Use dataset augmentation, which generates new samples from the real ones, by rotating, resizing...

Credits

Thanks for listening!

Sources

- ▶ <https://oreil.ly/2TLLnUM>
- ▶ <https://bit.ly/2F3exFV>
- ▶ <https://bit.ly/2XSh5Pr>
- ▶ <https://bit.ly/2Ces8to>
- ▶ <https://bit.ly/2NWh0pA>
- ▶ <https://bit.ly/2J8z0mA>
- ▶ <https://bit.ly/2LsZNBM>
- ▶ <https://bit.ly/2VUuK73>
- ▶ <https://bit.ly/2Uv1A0j>