

Master's Thesis
Deep Learning for Visual Recognition

Remi Cadene
Supervised by Nicolas Thome and Matthieu Cord

Wednesday 7th September, 2016

Contents

Introduction	1
Context	1
Contributions	3
1 Convolutional Neural Networks	5
1.1 Layers	5
1.1.1 Linear or Fully Connected	5
1.1.2 Activation or Non Linearity	5
1.1.3 Spatial Convolution	7
1.1.4 Spatial Pooling	8
1.1.5 Batch Normalization	8
1.2 Convolutional Architectures	9
1.2.1 CNNs (LeNet)	9
1.2.2 Deep CNNs	10
1.2.3 Very Deep CNNs	10
1.2.4 Residual CNNs	12
1.3 Training Methods	12
1.3.1 From Scratch	12
1.3.2 Transfer Learning	13
1.3.3 Loss functions	14
1.3.4 Optimization algorithms	15
1.3.5 Regularization Approaches	15
1.4 Interpretability	16
1.4.1 Definitions	16
1.4.2 Simple Visualization Techniques	17
1.4.3 Advanced Visualization Techniques	18
1.5 Invariance	19
1.5.1 Translation invariance	20
1.5.2 Rotation invariance	21
1.5.3 Scale invariance	21

2	Transfer Learning for Deep CNNs	23
2.1	Medium dataset of food (UPMC Food101)	23
2.1.1	Context	23
2.1.2	Previous work	23
2.1.3	Experiments	25
2.2	Small dataset of objects (VOC2007)	28
2.2.1	Context	28
2.2.2	Previous work	28
2.2.3	Experiments	29
2.3	Small dataset of roofs (DSG2016 online)	30
2.3.1	Context	30
2.3.2	Experiments	31
2.4	Conclusion	32
3	Weakly Supervised Learning	34
3.1	Introduction	34
3.1.1	Definition	34
3.1.2	Multi Instance Learning	34
3.1.3	Spatial Transformer Network	35
3.2	Applying fine tuning to Weldon	38
3.2.1	Context	38
3.2.2	Experiments	39
3.3	Study of Spatial Transformer Network	39
3.3.1	Context	39
3.3.2	Previous work	39
3.3.3	Experiments	40
3.4	Conclusion	42
	Conclusion	44
	Appendices	45
	A Overfeat	46
	B Vgg16	47
	C InceptionV3	48

Abstract

The goal of our research is to develop methods advancing automatic visual recognition. In order to predict the unique or multiple labels associated to an image, we study different kind of Deep Neural Networks architectures and methods for supervised features learning. We first draw up a state-of-the-art review of the Convolutional Neural Networks aiming to understand the history behind this family of statistical models, the limit of modern architectures and the novel techniques currently used to train deep CNNs. The originality of our work lies in our approach focusing on tasks with a low amount of data. We introduce different models and techniques to achieve the best accuracy on several kind of datasets, such as a medium dataset of food recipes (100k images) for building a web API, or a small dataset of satellite images (6,000) for the DSG online challenge that we've won. We also draw up the state-of-the-art in Weakly Supervised Learning, introducing different kind of CNNs able to localize regions of interest. Our last contribution is a framework, build on top of Torch7, for training and testing deep models on any visual recognition tasks and on datasets of any scale.

Acknowledgements

I would specifically like to thank Prof. Matthieu Cord and Assoc. Prof. Nicolas Thome for supervising my research on this project and providing resources for the experiments. Additionally, I thank all the people at LIP6 for the perfect working atmosphere. Lastly, I thank my family and friends for their love and support.

Introduction

Context

Since the beginning of the Web 2.0, the amount of visual data has grown exponentially. As an example, the director of Facebook AI Research Yann LeCun has said that almost 1 billion new photos were uploaded each day on Facebook in 2016 ¹. Thus, computer vision has become ubiquitous in our society, with many applications such as search engine, image understanding, medicine and self-driving car. Core to many of these applications are visual recognition tasks namely image classification, localization and detection. While this seems natural to humans, those tasks are difficult due to the large number of objects in the world, the continuous set of viewpoints from which they can be viewed, the lighting in scene, color variations, background clutter, or occlusion. Sometimes those visual tasks can even be challenging for untrained humans when several classes look very similar such as in fine grained recognition, or very different such as when age, gender, version, etc. are present.

It has long been the goal of computer vision researchers to have a flexible representation of the visual world in order to recognize objects in complex scenes. During the 2000's, the best accuracy was obtained using a hand-crafted model called the Bag of visual Words (BoW). In a first step, robust features extractors (e.g. SIFT [25]) were applied to the dataset for extracting local descriptors from the images. Then, a clustering algorithm (e.g. K-Means) was used to obtain a visual descriptor codebook. Finally, each image were assigned to their own representation in a lower space thanks to a pooling step aggregating all the descriptors. Later, a classifier (e.g. Support Vector Machine and kernel methods) was trained on top of the vectorial representation obtained using BoW. Recent developments in Deep Learning have greatly advanced the performance of these state-of-the-art visual recognition systems to the extent of sweeping aside the hand crafted models such as BoW. Nowadays a lot of products in the industry have benefited from the past years of research in Deep Learning. We can cite Google Photos, Flickr and Facebook, three of the world's

¹<https://youtu.be/vlQomVlaNFg>

largest photo sharing services, that use Deep Learning technologies to efficiently order and sort out piles of pictures ², to better target advertising, or to find people associated to faces [45]. Startups also are using Deep Learning to build better recognition products and to revolutionize the market providing new services ³. Furthermore, it is used to build physical products such as self-driving cars, drones and any kind of robots equipped with cameras ⁴.

Deep Learning can be summed up as a sub field of Machine Learning studying statical models called deep neural networks. The latter are able to learn complex and hierarchical representations from raw data, unlike hand crafted models which are made of an essential features engineering step. This scientific field has been known under a variety of names and has seen a long history of research, experiencing alternatively waves of excitement and periods of oblivion [37]. Early works on Deep Learning, or rather on Cybernetics, as it used to be called back then, have been made in 1940-1960s, and describe biologically inspired models such as the Perceptron, Adaline, or Multi Layer Perceptron [35, 12]. Then, a second wave called Connectionism came in the 1960s-1980s with the invention of backpropagation [36]. This algorithm persists to the present day and is currently the algorithm of choice to optimize Deep Neural Networks. A notable contribution is the Convolutional Neural Networks (CNNs) designed, at this time, to recognize relatively simple visual patterns, such as handwritten characters [21]. Finally, the modern era of Deep Learning has started in 2006 with the creation of more complex architectures [13, 2, 34, 9]. Since a breakthrough in speech and natural language processing in 2011, and also in image classification during the scientific competition ILSVRC in 2012 [18], Deep Learning has conquered many Machine Learning communities, such as Reddit, and won challenges beyond their conventional applications area ⁵.

Especially during the last four years, Deep Learning has made a tremendous impact in computer vision reaching previously unattainable performance on many tasks such as image classification, objects detection, object localization, object tracking, pose estimation, image segmentation or image captioning [10]. This progress have been made possible by the increase in computational resources, thanks to frameworks such as Torch7, modern GPUs implementations such as Cudnn, the increase in available annotated data, and the community-based involvement to open source codes and to share models. These facts allowed for a much larger audience to acquire the expertise needed to train modern convolutional networks. Thus, larger and deeper architectures are trained on bigger datasets to achieve better accuracy each year. Also, already trained models have shown astonishing results when transfered on smaller datasets and evaluated on different visual tasks. Hence, a lot of pretrained models are available on the web. However, CNNs still possess inherent

²<http://googleresearch.blogspot.fr/2013/06/improving-photo-search-step-across.html>

³<http://blog.ventureradar.com/2016/01/19/18-deep-learning-startups-you-should-know>

⁴<http://fortune.com/2015/12/21/elon-musk-interview>

⁵<http://blog.kaggle.com/2014/04/18/winning-the-galaxy-challenge-with-convnets>

limitations. From a theoretical perspective, Deep Neural Networks are not well understood due to their non convex property. Despite numerous efforts, a proof of convergence to good local minima has never been found. Thus, most of the research made in this field are experimentally driven and empirical [48]. From a practical perspective, their need for large amounts of training samples does not provide them the ability to generalize when trained on small and medium datasets. In this context, weakly supervised learning methods, that we describe in the next chapter, can be applied to overcome this limitations. Nevertheless, Deep Neural Networks seems to be the most promising kind of models for solving visual recognition.

The progress, needs and expectations of Deep Learning are undoubtedly signs of the Big Data era, where images of any kind and computable capabilities are more available than ever before. In this context, the Convolutional Neural Networks are the most efficient statistical model for visual recognition. The goal of this work is to produce an analysis of the state-of-the-art methods to train such models, to explain their limitations and to propose original idea to overcome the latter.

Contributions

This master’s thesis introduces a number of contributions to different aspects of visual recognition. However our work is focused on classifying images and recognizing objects using global labels (e.g. one label to indicates the presence or absence of the object).

- In the first chapter, we draw up the state of the art explaining how Convolutional Neural Networks (CNNs) achieve such good accuracy, describing different architectures and clarifying their limits. In the last chapter, we also draw up the state of the art in Weakly Supervised Learning which gather methods to improve the accuracy of CNNs trained on a few amount of images.
- In the second chapter, we apply Deep Neural Networks on a medium dataset of food recipes. We notably show that training CNNs From Scratch with large amount of parameters is possible on this kind of datasets reaching way better accuracy than the BoW model. We also show that Fine Tuning of pretrain models is essential, especially on small dataset. We illustrate this last point by presenting our winning solution of the Data Science Game Online Selection, an image classification challenge made for master and PhD students from all around the world.
- In the last chapter, we study techniques to overcome the limited spatial invariance capacity of CNNs without the use of rich annotations such as bounding boxes. The first technique consists in providing such invariance directly by feeding the networks augmented images. The second consists in using a novel approach developed by a

PhD student at LIP6 which gives the network the ability to localize regions of interest. The third one consists in using a first network to localize precisely the object and a second network to classify the resulting image. Our main results are synthesized at the end of this chapter.

- Overall, this study has helped us to develop *Torchnet-vision*, a framework build on top of Torch7 that serves as a plugin for Torchnet (a high level deep learning framework) for training easily the last architectures and pretrain models on several datasets. Reproducibility of a large amount of our experiments is ensured by the fact that we provide links to our code in footnotes.

Chapter 1

Convolutional Neural Networks

1.1 Layers

1.1.1 Linear or Fully Connected

Mathematically, we can think of a linear layer as a function which applies a linear transformation on a vectorial input of dimension I and output a vector of dimension O . Usually the layer has a bias parameter.

$$y = A \bullet x + b$$
$$y_i = \sum_{j=1}^I (A_{i,j} x_j) + b_i$$

The linear layer is motivated by the basic computational unit of the brain called neuron. Approximately 86 billion neurons can be found in the human nervous system and they are connected with approximately 10^{14} - 10^{15} synapses. Each neuron receives input signals from its dendrites and produces output signal along its axon. The linear layer is a simplification of a group of neuron having their dendrites connected to the same inputs. Usually an activation function, such as sigmoid, is used to mimic the 1-0 impulse carried away from the cell body and also to add non linearity. However we consider here that the activation function is the identity function that output real values.

1.1.2 Activation or Non Linearity

The capacity of the neural networks to approximate any functions, especially non-convex, is directly the result of the non-linear activation functions. Every kind of activation function

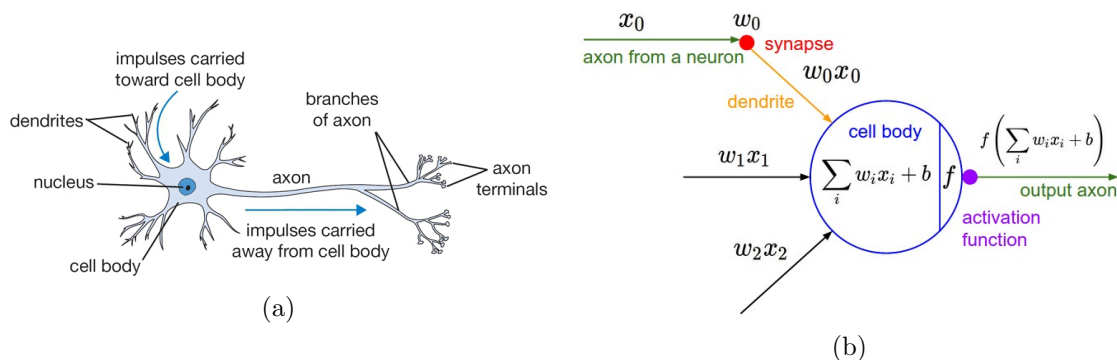


Figure 1.1: A cartoon drawing of a biological neuron (a) and its mathematical model (b).

takes a vector and performs a certain fixed point-wise operation on it. There are three main activation functions.

Sigmoid The Sigmoid non-linearity has the following mathematical form

$$y = \sigma(x) = 1/(1 + \exp^{-x})$$

It takes a real value and squashes it between 0 and 1. However, when the neuron's activation saturates at either tail of 0 or 1, the gradient at these regions is almost zero. Thus, the backpropagation algorithm fail at modifying its parameters and the parameters of the preceding neural layers.

Hyperbolic Tangent The TanH non-linearity has the following mathematical form

$$y = 2\sigma(2x) - 1$$

It squashes a real-valued number between -1 and 1. However it has the same drawback than the sigmoid.

Rectified Linear Unit The ReLU has the following mathematical form

$$y = \max(0, x)$$

The ReLU has become very popular in the last few years, because it was found to greatly accelerate the convergence of stochastic gradient descent compared to the sigmoid/tanh functions due to its linear non-saturating form (e.g. a factor of 6 in [18]). In fact, it does not suffer from the vanishing or exploding gradient. An other advantage is that it involves cheap operations compared to the expensive exponentials. However, the ReLU removes all the negative informations and thus appears not suited for all datasets and architectures.

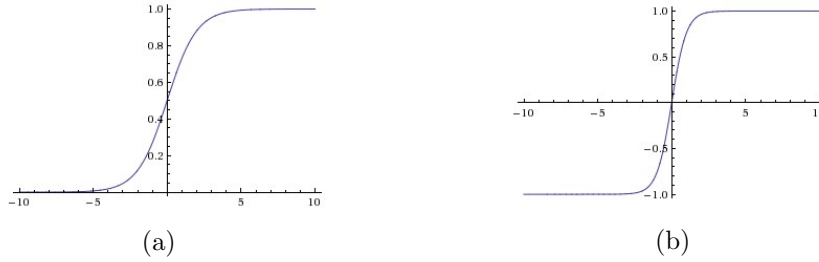


Figure 1.2: A sigmoid (a) and a tanh (b).

1.1.3 Spatial Convolution

Regular Neural Networks, only made of linear and activation layers, do not scale well to full images. For instance, images of size $3 \times 224 \times 224$ (3 color channels, 224 wide, 224 high) would necessitate a first linear layer having $3 * 224 * 224 + 1 = 150,129$ parameters for a single neurone (e.g. output). Spatial convolution layers take advantage of the fact that their input (e.g. images or feature maps) exhibits many spatial relationships. In fact, neighboring pixels should not be affected by their location within image. Thus, a convolutional layer learns a set of N_k filters $F = f_1, \dots, f_{N_k}$, which are convolved spatially with input image x , to produce a set of N_k 2D features maps z :

$$z_k = f_k * x$$

where $*$ is the convolution operator. When the filter correlates well with a region of the input image, the response in the corresponding feature map location is strong. Unlike conventional linear layer, weights are shared over the entire image reducing the number of parameters per response and equivariance is learned (i.e. an object shifted in the input image will simply shift the corresponding responses in a similar way). Also, a fully connected layer can be seen as a convolutional layer with filter of sizes $1 \times 1 \times inputSize$.

It is important to highlight that a spatial convolution is not defined by the spatial size of the input feature maps (e.g. wide and high), neither by the size of the output feature maps, but by the number of filters (e.g. number of output channels), the properties of its filters (e.g. number of input channels, wide, high) and the properties of the convolution (e.g. padding, stride). Animations showing different kind of convolution can be viewed on line ¹.

¹https://github.com/vdumoulin/conv_arithmetic

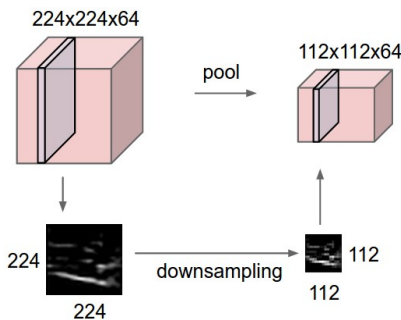


Figure 1.3: The illustration of a spatial pooling operation in 2×2 regions by a stride of 2 in the high direction, and 2 in the width direction, without padding.

1.1.4 Spatial Pooling

In Convolutional Neural Networks, a pooling layer is typically present to provide invariance to slightly different input images and to reduce the dimension of the feature maps (e.g. wide, high):

$$p_R = P_{i \in R}(z_i)$$

where P is a pooling function over the region of pixels R . Max pooling is preferred as it avoids cancellation of negative elements and prevents blurring of the activations and gradients throughout the network since the gradient is placed in a single location during backpropagation.

The spatial pooling layer is defined by its aggregation function, the high and width dimensions of the area where it is applied, and the properties of the convolution (e.g. padding, stride).

1.1.5 Batch Normalization

This layer quickly became very popular mostly because it helps to converge faster [14]. It adds a normalization step (shifting inputs to zero-mean and unit variance) to make the inputs of each trainable layers comparable across features. By doing this it ensures a high learning rate while keeping the network learning.

Also it allows activations functions such as TanH and Sigmoid to not get stuck in the saturation mode (e.g. gradient equal to 0).

1.2 Convolutional Architectures

A lot of convolutional architectures have been developed from the 1990's. In this section, we make an inventory of the most known architectures². Each one represent a step further for more advanced visual recognition.

1.2.1 CNNs (LeNet)

LeNet-5 This kind of architecture is one of the first successful applications of CNNs. It was developed by Yann LeCun in the 1990's and was used to read zip codes and digits. This architecture, with regard to the modern ones, differs on many points. Thus, we will limit ourselves on the most known, LeNet-5 [22], and we will not delve into the details. In overall this network was the origin of much of the recent architectures, and a true inspiration for many people in the field.

LeNet-5 features can be summarized as:

- sequence of 3 layers: convolution, pooling, non-linearity,
- inputs are normalized using mean and standard deviation to accelerate training [20],
- sparse connection matrix between layers to avoid large computational cost
- hyperbolic tangent or sigmoid as non-linearity function,
- trainable average pooling as pooling function,
- fully connected layers as final classifier,
- mean squared error as loss function.

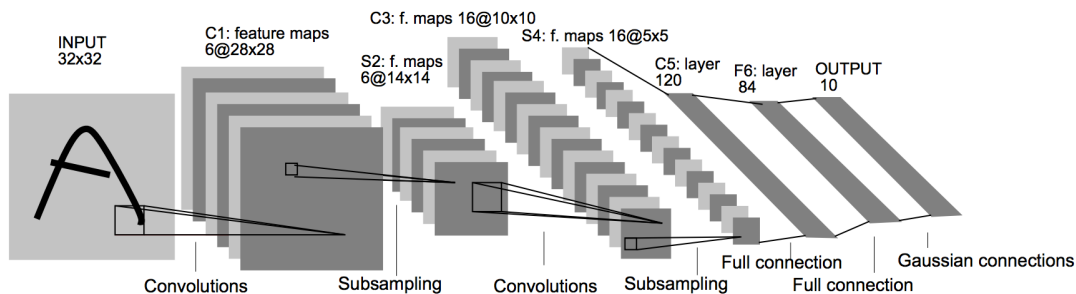


Figure 1.4: Architecture of LeNet-5, an old convolutional neural network for digits recognition.

²<http://culurciello.github.io/tech/2016/06/04/nets.html>

1.2.2 Deep CNNs

AlexNet It is one of the first work that popularized convolutional networks in computer vision. AlexNet [19] was submitted to the ImageNet ILSVRC challenge of 2012 and significantly outperformed the other hand crafted models (accuracy top5 of 84% compared to the second runner-up with 74%). This network, compared to LeNet, was deeper (60 millions of parameters) and bigger (5 convolutional layers, 3 max pooling and 3 fully-connected layers). At this time, the authors provided a multi-GPUs implementation in CUDA to bypass the memory needs. It popularized:

- the ReLU as non-linearity function of choice,
- the method of stacking convolutional layers plus non-linearity on top of each other without being immediately followed by a pooling layer,
- the method of overlapping Max Pooling, avoiding the averaging effects of Average Pooling.

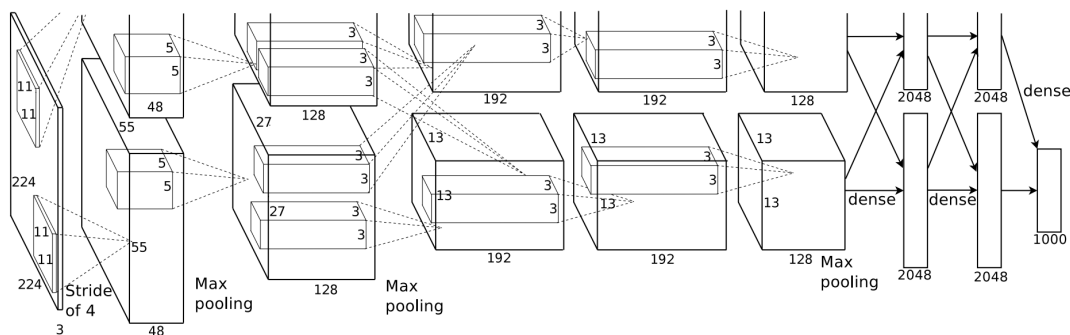


Figure 1.5: An illustration of the architecture of AlexNet. One GPU runs the layer-parts at the top of the figure while the others runs the layer-parts at the bottom..

Overfeat or ZFNet It was the winner architecture of ILSVRC2013 [38] with almost 140 millions of parameters. Based on AlexNet, the size of its middle convolutional layers have been expanded. Also, the stride and filter size on its first layer have been made smaller.

1.2.3 Very Deep CNNs

VeryDeep or VggNet It was the runner-up architecture of ILSVRC2014 [40] with almost 140 millions of parameters. Its main contributions were to show that depth is a critical

component for good performance, to use much smaller 3×3 filters in each convolutional layers, and also to combine them as a sequence of convolutions. The great advantage of VggNet was the insight that multiple 3×3 convolution in sequence can emulate the effect of larger receptive fields, for examples 5×5 and 7×7 . These ideas will be also used in more recent network architectures as Inception and ResNet.

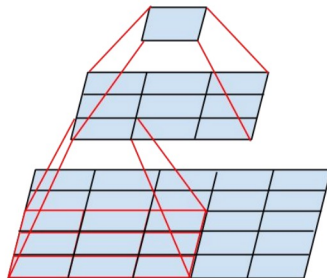


Figure 1.6: Filter of 5×5 or more can be decomposed with multiple 3×3 convolutions such as in VGG.

GoogLeNet or Inception It was the winner architecture of ILSVRC2014 [43]. Its main contribution was the development of an Inception Module that dramatically reduced the number of parameters (40 millions) [32]. Also, it eliminated a large amount of parameters by using average pooling instead of fully connected layers at the top of the convolutional layers. Further versions of the GoogLeNet has been released. The most recent architecture available is InceptionV3 [44]. Notably, it uses batch normalization.

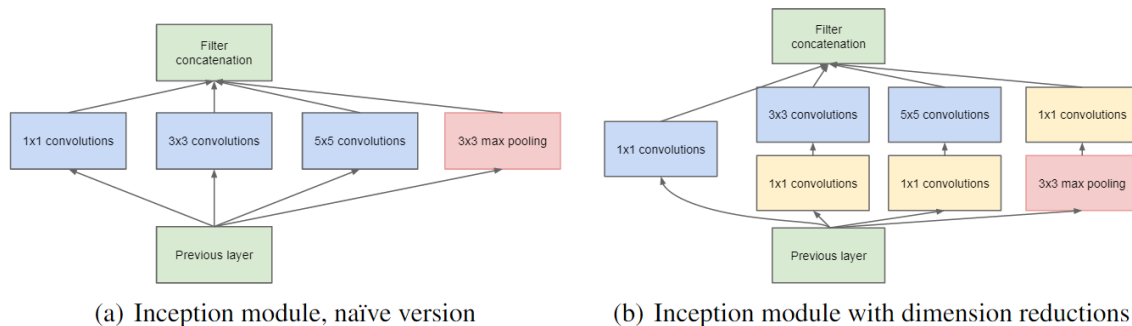


Figure 1.7: 1×1 convolutions are used to decrease the input size before 3×3 convolutions in order to provide more combinational power such as in GoogLeNet.

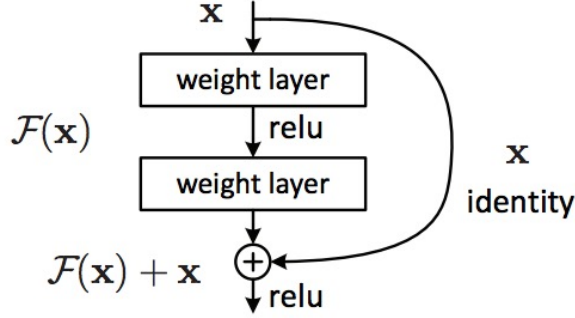


Figure 1.8: A skip connection is used to bypass the input to the next layers such as in ResNet.

1.2.4 Residual CNNs

ResNet It was the winner architecture of ILSVRC2015 [?] with 152 layers. Its main contribution was to use batch normalization and special skip connections for training deeper architectures. ResNet with 1000 layers can be trained with those techniques. However, it has been empirically found that ResNet usually operates on blocks of relatively low depth ($\sim 20 - 30$ layers), which act in parallel, rather than serially flow the entire length of the network [46].

1.3 Training Methods

1.3.1 From Scratch

Initialization All the network parameters are generally initialized with Layer-sequential unit-variance (LSUV) (e.g. each parameters as Gaussian random variables with mean 0 and standard deviation $\frac{1}{n_{inputs}}$ and biases are initialized to zero). Since the LSUV initialization works under assumption of preserving unit variance of the input, pixel intensities are given after subtracting the mean and dividing by the standard deviation. More information can be found in the chapter 3 of Michael Nielsen’s book [30]. In case of pretrain networks, the mean and std of the original dataset are kept.

Loss function To quantify the capacity of the network to approximate the ground truth labels for all training inputs, we define a loss function which takes as inputs the weights, biases, and examples from the training set. For instance, the loss could be the number

of images correctly classified. However, the most efficient way to find the weights and biases, regarding the number of parameters, is to use an algorithm similar to the Stochastic Gradient Descent (SGD). In order to do so, if our chosen loss function is not smooth, we have to choose a surrogate loss (e.g. derivable function) such as Mean Square Error or Cross Entropy.

Backpropagation For each examples, we compute the prediction and its associated loss. We sum up all the loss to compute the final error. Then we use the backpropagation algorithm to propagate the error in order to compute the partial derivatives $\frac{\delta E}{\delta w}$ and $\frac{\delta E}{\delta b}$ of the cost function E for all weights w and bias b . In this work, our goal is not to explain in details how works the backpropagation algorithm. We advise the curious reader to read the chapter 2 of Michael Nielsen’s book [30].

Optimization Once all the derivatives are computed, we update our parameters using a chosen optimization technique such as SGD. We then iterate the predication (e.g. forward pass), the backpropagation of errors (e.g. backward pass) and the optimization until convergence hoping to find a local minimum low enough to ensure good predictions. Even if the chosen surrogate loss function of a neural network is non convex, SGD works well in practice.

Grid search It is common to explore manually the space of hyperparameters such as learning rate, weight decay, learning rate decay, amount of dropout, not to mention the architectures hyperparameters, in order to obtain the best performance in terms of both accuracy and training time.

[28] made an evaluation on the influence of architecture choices and optimization hyperparameters on ImageNet. While there are very few theoretical studies, technical studies of this kind can help the reader to reduce the space of hyperparameters to explore.

1.3.2 Transfer Learning

Features Extraction It consists in extracting features from the network by forwarding examples. Transformations to the examples are possible such as horizontal flip. Then the associated features to the example are aggregated whether by averaging or stacking them. Finally, a classifier is trained and tested on the features. Typically, the later is a Support Vector Machine with a linear kernel.

Fine Tuning It consists in training a pretrained network on a smaller dataset. Typically, the last fully connected layers, which can be viewed as classification layers, are reset and a smaller learning rate is applied to the pretrained layers. By doing so, the goal is to adapt the features to the new dataset. More different is the latter from the original dataset, more parameters/layers must be reset.

1.3.3 Loss functions

In this subsection, we present the three most used loss function to train deep neural networks for classification.

Mean Square Error (MSE) It is a multi class loss formerly used to train neural networks.

$$Loss(x, y) = \frac{1}{n} \sum_i |x_i - y_i|^2$$

with x a vector of n predictions, and y a binary vector full of 0 besides a 1 in the corresponding class dimension .

Cross Entropy It is a multi class loss which is nearly a better choice than MSE.

$$Loss(x, y) = - \sum_i y_i * \log\left(\frac{\exp(x_i)}{(\sum_j \exp(x_j))}\right)$$

with x a vector of n predictions, and y a binary vector full of 0 besides a 1 in the corresponding class dimension .

In fact, it may happen that the initialization of the parameters result in the network being decisively wrong for some training input (an output neuron will have saturated near 1, when it should be 0, or vice versa). The MSE loss will usually slow down learning, but the Cross Entropy loss won't. More information can be found in the chapter 3 of [30].

Loss Multi Label It is the adaptation of the Cross Entropy loss for multi-label classification. It is a multi-label one-versus-all loss based on max-entropy.

$$Loss(x, y) = - \sum_i (y_i * \log\left(\frac{\exp(x_i)}{1 + \exp(x_i)}\right) + (1 - y_i) * \log\left(\frac{1}{1 + \exp(x_i)}\right))$$

1.3.4 Optimization algorithms

The loss function of a CNN is highly non convex. Hopefully the latter is also fully derivable, so that gradient based optimization algorithms can be applied. However, CNNs are usually made of tens of millions of parameters. Thus, only the first order derivatives are used in practice. In fact, the second derivatives are costly in term of memory and computational effort.

Stochastic Gradient Descent (SGD) It is the main optimization algorithm. It consists in using a few examples to compute the gradient of the parameters with respect to the loss function :

$$\theta_{t+1} = \theta_t - \lambda \cdot \nabla_{\theta_t} L(f_{\theta_t}(x_i), y_i)$$

There is no proof of good convergence. However, this algorithm reaches good local minima in practice, even when the parameters are randomly initialized. One of the reason could be the stochastic property of this algorithm, allowing the latter to optimize different loss functions and thus to get out of bad minima. The other reason could be that a lot of local minima are almost as accurate than the global minima. Answers to this question are still under active research.

Approximation of Second Order Derivatives Other optimization algorithms rely on more advance techniques such as momentum, second order approximation and adaptive learning rates [42, 17]. They are known to converge faster and their parameters are sometimes easier to tune by grid search. However, they take a bit more processing time to compute, but also much more memory (2 to 3 more).

Distributed SGD It is the kind of optimization used in parallel computing environments. Different computers train the same architecture with almost the same parameters values. It allows more exploration of the parameters space, which can lead to improved performance [50].

1.3.5 Regularization Approaches

Deep and large enough neural networks can memorize any data. During training, their accuracy on the trainset typically converges towards perfection while it degrades on the testset. This phenomenon is called overfitting.

Regularization L2 The first main approach to overcome overfitting is the classical weight decay, which adds a term to the cost function to penalize the parameters in each dimension, preventing the network from exactly modeling the training data and therefore help generalize to new examples:

$$Err(x, y) = Loss(x, y) + \sum_i \theta_i^2$$

with θ a vector containing all the network parameters.

Data augmentation It is a method of boosting the size of the training set so that the model cannot memorize all of it. This can take several forms depending of the dataset. For instance, if the objects are supposed to be invariant to rotation such as galaxies or planktons, it is well suited to apply different kind of rotations to the original images.

Dropout Finally, a recent success has been shown with a regularization technique called Dropout [41]. The idea is to randomly set a certain percentage of the activations in each layers to 0. During the training, neurons must learn better representations without co-adapting to each other being active. During the testing, all the neurons are used to compute the prediction and Dropout acts like a form of model averaging over all possible instantiations of the model.

Early stopping It consists in stopping the training before the model begins to overfit the training set. In practice, it is used a lot during the training of neural networks.

1.4 Interpretability

1.4.1 Definitions

The desire for interpretation presupposes that predictions alone do not suffice [23]. Typically, we train models to achieve strong predictive power. However, this objective can be a weak surrogate for the real-world goals of machine learning practitioners.

Intelligibility There are two definitions of interpretability. The first one is linked with understandability or intelligibility, i.e., that we can grasp how the model works. Multiple criteria are used to evaluate if a model is interpretable or not. For instance: *will it converge? do we understand what each parameters represents? is it simple enough to be examined all at once by a human?* Understandable models are sometimes called transparent, while incomprehensible models are called black boxes.

Post-hoc interpretability Discussions of interpretability sometimes suggest that human decision-makers, despite being black boxes, are themselves interpretable because they can explain their actions. Deep learning models are also often considered as black boxes. However visualization techniques can help to generate post-hoc interpretations in order to explain their actions.

1.4.2 Simple Visualization Techniques

The first technique consists in visualizing the images directly. For instance, computing the loss over all the testing set allows to visualize the easiest or hardest examples for the network. Also, computing the activations of a certain layer to identify the k-nearest neighbors based on the proximity in the space learned by the model can be a good way to understand what the chosen layer has learned.

A similar approach consists to visualize high-dimensional distributed representations with t-SNE [26], a technique that renders 2D visualizations in which nearby data points are likely to appear close together.

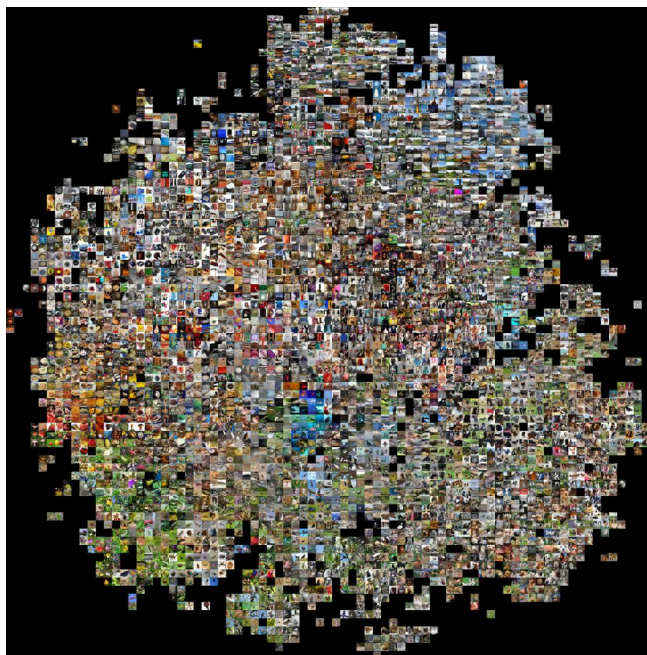


Figure 1.9: Embedded images from ImageNet in 2D space using t-SNE and features extracted from a CNN.

An other technique consists in visualizing the features map generated by the network.



Figure 1.10: Original images with their associated gradient based images.

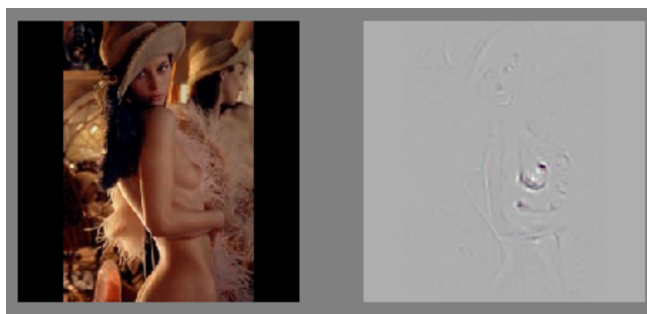


Figure 1.11: The original Lena photo and the resulting deconvolutional image which added to the original could improved Lena to look more like pornography (The goal is to understand what a CNN has learned).

However, this method is not suited when the network is large.

1.4.3 Advanced Visualization Techniques

Gradient based A popular approach is to render the gradient as an image [39]. While this does not say precisely how a model works, it conveys which image regions the current classification depends upon most heavily.

[29] attempt to explain what a network has learned by altering the input through gradient descent to enhance the activations of certain nodes selected from the hidden layers. An inspection of the perturbed inputs can give clues to what the models has learned.

An other technique consists in training deconvolutional neural networks [48, 49]. One example of using this powerful technique is to understand what CNNs are looking at when they see nudity ³.

³<http://blog.clarifai.com/what-convolutional-neural-networks-see-at-when-they-see-nudity>

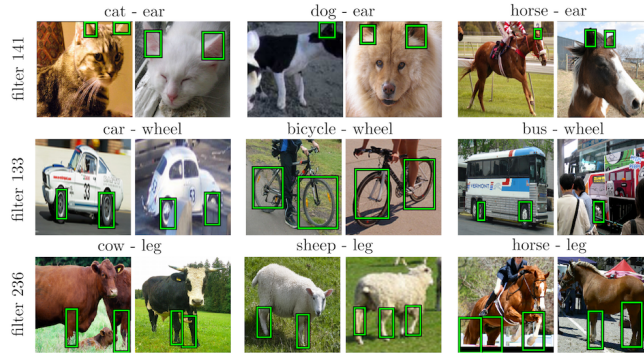


Figure 1.12: Detections performed by filters of AlexNet. The filters are specific to a part and they work well on several object classes containing it.

Semantic parts localization Finally other approaches study whether CNNs learn semantic parts of object classes in their internal representation. They investigate the responses of convolutional filters and try to associate their stimuli with semantic parts.

A recent study [8] discovered that the discriminative power of the network can be attributed to a few discriminative filters specialized to each object class. Despite promoting the emergence of filters learn to respond to semantic parts (and not only object class), they found that only 34 out of 123 semantic parts in PASCAL-Part dataset emerge in AlexNet. Also networks trained for image classification produced the same results than those trained for objects detection or localization.

1.5 Invariance

Convolutional Neural Network is a powerful model able to learn the needed invariance directly from the data. However, when the amount of labeled images (e.g. training set) is not big enough, the architectural choices can limit the capacity of the model to generalize to unknown images (e.g. testing set), or rather the opposite, to improve both the accuracy and training time.

In this section, we explain how CNNs achieve translation, rotation and scale invariance, and we propose some methods to achieve better invariance. Note that the pooling regimes make convolution slightly invariant to translation, rotation and shifting. However, it is not sufficient to be invariant to large spatial transformations.

1.5.1 Translation invariance

Pooling, combined with striding, is a common way to achieve a degree of invariance, but even without this technique, the convolution filters and the fully-connected layers are able to learn spatial invariance [5]. We explain our thinking in the next paragraphs.

In figures 1.13 and 1.14, we can see an ideal network with two convolutional layers with pooling, and two fully-connected layers.

- The first layer filters (which generate the green volume) detect eyes, noses and other basic shapes (in real CNNs, first layer filters match lines and very basic textures).
- The second layer filters (which generate the yellow volume) detect faces, legs and other objects that are aggregations of the first layer filters (real life convolution filters may detect objects that have no meaning to humans).

In figure 1.13, a face is at the corner bottom left of the image (represented by two red and a magenta point). In figure 1.14, the same face is at the corner top left of the image. The same number of activations occurs, but they occur in different regions of the green and yellow volumes. Therefore, any activation point at the first slice of the yellow volume means that a face was detected, independently of the face location. Then the fully-connected (FC) layer is responsible to "translate" a face and two arms to an human body. In each examples, the activation path inside the FC layer was different, meaning that a correct learning at the FC layer is essential to ensure the spatial invariance property.

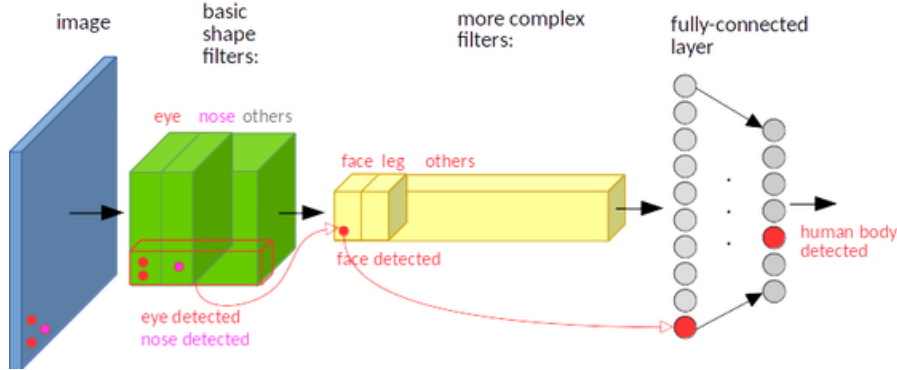


Figure 1.13: Ideal CNN with a human face at the bottom left of the image.

In conclusion, if a CNN is trained showing faces only at one corner, during the learning process, the fully-connected layer may become insensitive to some faces in other corners. Thus, regarding the task and the dataset, it can be suitable to use data augmentation to provide the CNN more examples of the same objects in different image regions. However, for small datasets in which objects can undergo strong translations in the image, some

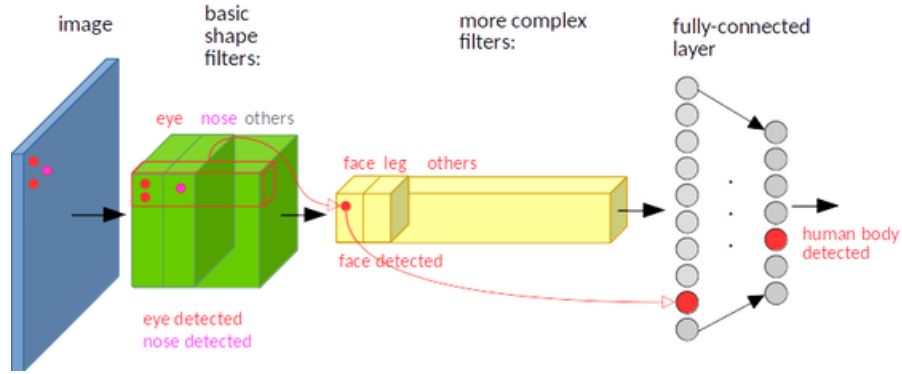


Figure 1.14: Ideal CNN with a human face at the top left of the image.

different approaches are more suitable [7, 3].

1.5.2 Rotation invariance

Different kind of invariance to rotation transformation can be important to learn. The figure 1.15 represents the rotation invariance, where the object keeps the same label regardless its orientation. The figure 1.16 represents the same-equivariance, where the object and its label must keep the same orientation (its is common for segmentation tasks).

The main approach to learn rotation invariance is to use data augmentation to provide the CNN more examples of the same objects with different orientation. However, CNNs will often learn multiple copies of the same filter in different orientations. Some other approaches [6] try to reduce the redundancy, in order to lower the number of parameters, and thus to reduce the risk of overfitting.

1.5.3 Scale invariance

The main approach to learn scale invariance is to use data augmentation to provide the CNN more examples of the same objects with different scales. It is also possible to train several CNNs specialized in each scales and combine their predictions. Finally, other approaches consist to hard code scale invariance into the CNN such as in the Inception architecture or [16].



(a) Plankton



(b) Galaxies

Figure 1.15: Example images for the Plankton and Galaxies datasets, which are rotation invariant.



(a) Satellite image



(b) Building labels

Figure 1.16: Example tile from the Massachusetts buildings dataset, which is same-equivariant to rotation, and corresponding labels

Chapter 2

Transfer Learning for Deep CNNs

2.1 Medium dataset of food (UPMC Food101)

2.1.1 Context

Web API VISIIR Visual Seek for Interactive Image Retrieval (VISIIR) is a project aiming at exploring new methods for semantic image annotation. In this frame, we have contributed to the demo which uses a CNN to recognize food images across 101 categories from the Dataset UPMC_Food101.

UPMC_Food101 It is a large multimodal dataset [47] containing about 100,000 recipes for a total of 101 food categories. Each of them are constituted by around 800 to 950 images from Google Image found using the title of the category. Because of this, this dataset may contain some noise. It is the twin dataset of ETHZ University [4].

2.1.2 Previous work

CNN on UPMC_Food101 [47] compared different visual features on this dataset and found that using VGG19 as features extractor was way more accurate (40.21% top-1 accuracy) than using the BoW model with SIFT features (23.96%). They also compared two architectures and found that a deeper architecture achieved better as features extractor than a shallower (Overfeat with 33.91%).

CNN on ETHZ Food-101 A very recent paper [24] claim the state of the art on two datasets of food images, UEC-256 and ETHZ Food-101. They fine tuned GoogLeNet, a 22-layer network, on these datasets and achieved a 77.4% top-1 accuracy and 93.7% top-5 on ETHZ Food-101.

Image to calories [27] present a system aiming to recognize the contents of a meal from a single image, and then predict its nutritional contents, such as calories. They tested their CNN-based method on a dataset of 75,000 images from 23 different restaurants. Estimating the size of the foods, as well as their labels (2,516), requires solving segmentation and depth / volume estimation from a single image. They first learned a binary classifier between food and non-food class. In order to do so, they modified ETHZ Food-101 and fine tuned GoogLeNet on the new binary dataset called ETHZ Food-101 Background. They used ETHZ Food-201 Segmented to learn segmentation, and Gfood-3d to learn depth and volume. Finally they used USDA NNDB to process the amount of calorie in a certain volume of food. They tested each steps of there method separately and did not provide a end-to-end test.

2.1.3 Experiments

In this subsection, our goal is to find the most accurate supervised learning method on this kind of medium sized dataset and also to understand the capacity of Convolutional Neural Networks. We use the following experimental protocol: we train each models on the same 80% of the dataset and validate the results on the remaining 20%. However, we select the models hyper parameters regarding our score in validation. This is the usual procedure for medium or big sized dataset (e.g. ImageNet). In table 2.2, we summarize our results in order to easily compare the supervised learning methods studied. Notably, our results can be reproduced using our framework ¹.

Forward-backward Benchmark As told earlier, one of our main concern is to find the most efficient techniques to achieve the best accuracy. Since we have multi threaded the data loading and data augmentation procedure, our GPU never need to wait for the latter to finish. Thus, we make a benchmark to compare different architectures and implementations in term of speed. In this study, we use a Nvidia GTX Titan X Maxwell GPU and a Intel Xeon E5-2630 v3 (2.40GHz). We summarize our results in table 2.1. Notably, InceptionV3 combined with cudnn is the fastest network. Still, it has the highest depth and number

¹<https://github.com/Cadene/torchnet-deep6/blob/master/src/main/upmcfood101/inceptionv3.lua>

of modules. However, it has also the lowest amount of parameters. The code for this benchmark can be found online ².

Model	Input size	Implementation	Forward (ms)	Backward (ms)	Total (ms)
Overfeat	221	nn float	13714.15	16277.34	29991.49
Overfeat	221	nn cuda	640.36	1032.78	1673.14
Overfeat	221	cuda	115.07	890.56	1005.63
Vgg16	224	nn float	18532.82	27978.44	46511.26
Vgg16	224	nn cuda	528.00	1407.97	1935.98
Vgg16	224	cuda	458.86	1144.11	1602.96
InceptionV3	399	nn float	31934.56	56328.28	88262.84
InceptionV3	399	nn cuda	787.64	1312.19	2099.83
InceptionV3	399	cuda	239.51	738.50	978.01

Table 2.1: Benchmarks of three architectures used in our study and three implementations: nn float on CPUs (by Torch7), nn cuda on GPUs (by Torch7), cuda R5 on Nvidia GPUs (by Nvidia). The forward and backward pass are averaged over 10 iterations. The measures are made in milliseconds for batches of size 50.

From Scratch better than Hand Crafted In table 2.2, we can see that deep models trained From Scratch, (d) and (g), achieve better accuracy than Hand Crafted models, (a) and (b). In other words, it is possible to train a deep model of 140 millions parameters on a medium dataset made of 80,000 images. This might be possible, because one image is made of $224 * 224 * 3 = 150,528$ pixels. Thus, the trainset is made of 80,000 examples for 150,528 features. Also, we used data augmentation, dropout and early stopping as regularization techniques. We virtually increase the size of the trainset as follow. Each image is randomly rescaled between 224 and 256, then randomly cropped to scale 224, and randomly flipped. We use SGD with Nesterov momentum as our optimization technique. Especially, we use a learning rate decay to slow down the training process after each mini batch updates. We validate the values of our hyper parameters using a unique parameters initialization (e.g. keeping the same seed).

From Scratch better than Features Extraction Deep models trained From Scratch, (d) and (g), achieve far better accuracy than models based on Features Extraction from pretrained deep models on ImageNet. Even if ImageNet contains classes and images which can be found in UPMC_Food101, learning specialized features from scratch is more accurate. This might be possible, because the dataset is big enough and the regularization methods (e.g. data augmentation and dropout) are strong enough.

²<https://github.com/Cadene/torchnet-deep6/blob/master/src/main/cnnbenchmark.lua>

Fine Tuning better than From Scratch Fine tuned models, (e) and (h), achieve better accuracy than From Scratch models, (d) and (g). In the same manner as outlined above, it is efficient to train a deep model on this dataset, but representations learned on ImageNet, even if less accurate, can be also useful. Thus, we experimentally chose to reset all the fully connected layers of our deep networks and keep only the parameters from the convolutional layers (14,714,688 parameters in Vgg16) that we update using a 10 times slower learning rate. In this manner, the networks are able to adapt their representations to UPMC_Food101 benefiting from those learned on ImageNet.

Very Deep better than Deep Overfeat, (c), (d) and (e), has almost the same number of parameters than Vgg16, (f), (g) and (h), but achieve lower accuracy because the later is deeper (16 layers versus 9 in Overfeat). Thus, Vgg16 learn better transferable features than Overfeat. Also, the importance of depth has been experimentally shown on ImageNet.

InceptionV3, the most efficient architecture A fine tuned InceptionV3 reaches the highest accuracy on this dataset. Also, it is 3 times faster to converge than Vgg16, thanks to several factors. Firstly, it has the lowest amount of parameters. Secondly, it has no dropout layers. The latter helps to generalize, but slow down the learning process. Finally, it has batch normalization layers before each non linearity. The latter help to converge faster and also to generalize.

Model	Test top 1 (top 5)	Assoc. train top 1 (top 5)
(a) Bag of visual Words	23.96	–
(b) BossaNova	28.59	–
(c) Overfeat & Extraction	33.91	–
(d) Overfeat & From Scratch	47.46 (69.37)	79.14 (94.49)
(e) Overfeat & Fine Tuning	57.98 (78.86)	89.69 (97.96)
(f) Vgg16 & Extraction	40.21	–
(g) Vgg16 & From Scratch	53.62 (74.67)	88.17 (97.68)
(h) Vgg16 & Fine Tuning	65.71 (82.54)	96.18 (99.39)
(g) InceptionV3 & Fine Tuning	66.83 (84.53)	85.34 (95.91)

Table 2.2: Percentage of good classification on UPMC_Food101. (a), (b), (c) and (f) are reported from [47]

2.2 Small dataset of objects (VOC2007)

2.2.1 Context

PASCAL VOC 2007 dataset The goal of the PASCAL Visual Object Classes Challenge 2007 is to recognize objects from 20 visual object classes in realistic scenes (i.e. not pre-segmented objects). It is fundamentally a multi-label supervised learning problem in that a training set of labelled images is provided (5,000 images compose the training set and 5,000 images compose the testing set). The twenty object classes that have been selected are:

- Person: person,
- Animal: bird, cat, cow, dog, horse, sheep,
- Vehicle: aeroplane, bicycle, boat, bus, car, motorbike, train,
- Indoor: bottle, chair, dining table, potted plant, sofa, tv/monitor

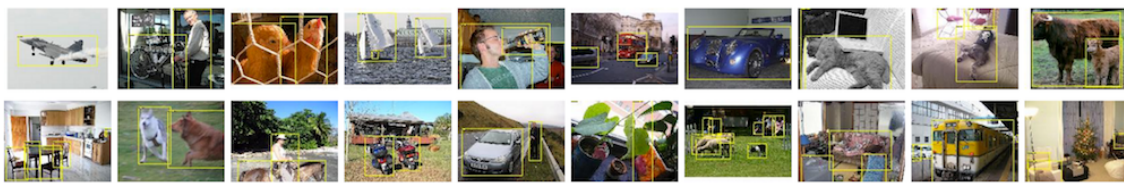


Figure 2.2: Illustration of the 20 object classes from the PASCAL_VOC2007 dataset.

2.2.2 Previous work

Dense Testing [40] evaluated the generalization capacity of their CNNs, namely VGG16 and VGG19, on VOC-2007, VOC-2012, Caltech-101 and Caltech-256. They proposed a method called dense testing. The network is applied densely over the rescaled test images, i.e. the fully-connected layers are first converted to convolutional layers (the first FC layer to a 7×7 convolutional layer, the last two FC layers to 1×1 convolutional layer). The resulting fully-convolutional net is then applied to the whole (uncropped) image (see figure 2.3). The result is a class score map with the number of channels equal to the number of classes. Then, to obtain a fixed-size vector of class scores for the image, the class score map is spatially averaged (Average Pooling). The test set is also augmented by horizontal flipping of the images. Finally, the soft-max class posteriors of the original and flipped images are averaged to obtain the final scores for the image.

They combined three CNNs fine tuned at different scales on ImageNet and achieved an impressive score of 89.7 mean AP on VOC-2007 and 89.3 on VOC-2012. Using the same amount of data, the highest score reported at this time was 82.4 on VOC-2007 and 83.2 on VOC-2012.

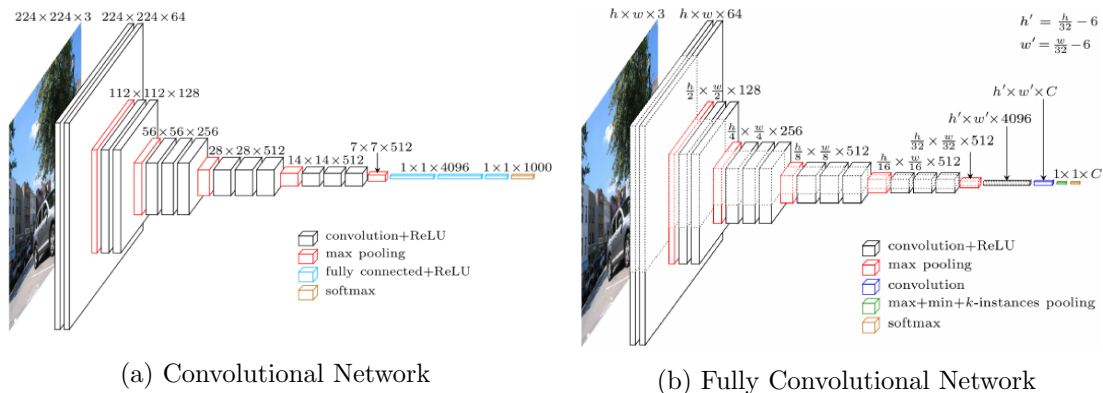


Figure 2.3: Illustration of a possible modification of VGG16 to a fully convolutional architecture in order to process bigger image than of size 3x224x224.

2.2.3 Experiments

Fine tuning outperforms the rest As can be seen in table 2.3, deep features outperform hand-crafted methods and fine tuning Vgg16 pretrained on ImageNet achieves the best accuracy. However the dataset is too small to fit Vgg16 From Scratch. In this experiment, we train our networks on multiple scales varying from 224 to 256 with random horizontal flips, so that the forwarded cropped image is of size 224×224 . Then, we test our models on single scale images (224). In (e), we reset the last two fully connected layers and we train the others with a 10 times smaller learning rate. In (d), we extract deep features from the last ReLU non linearity (before the last fully connected layer) and then we train support vector machines in a one-versus-all strategy cross-validating the regularization parameter.

Model type	Test mAP	Train mAP
(a) BoW	53.2	—
(b) BossaNova and FishersVector	61.6	—
(c) Vgg16 from scratch	39.79	99.73
(d) Vgg16 extraction	83.22	—
(e) Vgg16 fine tuned	85.70	98.81

Table 2.3: Comparison between hand crafted features models (a,b) and deep features models (c,d,e) tested on a single scale (224). (a) and (b) are reported from [1].

2.3 Small dataset of roofs (DSG2016 online)

In this section, we show that fine tuning an InceptionV3 architecture is still the most powerful technique for this kind of dataset. However, we use a bootstrap method to reduce overfitting and thus to achieve the best accuracy on this dataset amongst more than 110 teams.

2.3.1 Context

Data Science Game The DSG ³ is an international student challenge in Machine Learning. The first phase of this challenge is an online selection. It took place on Kaggle from the 14th June 2016 until the 10th July 2016 gathering more than 110 teams from all around the world competing for the first 20 positions. The final phase will took place during the week end of the 10th September in the castle of Cap Gemini.

The team (Jonquille UPMC) that we represented reached the first position of the online selection ⁴. In this section, we will explain our method in details.

Dataset It is a small dataset of satellite images of roofs. The goal is to predict the orientation of the roofs into 4 different categories. Models are evaluated with the multiclass accuracy top1 metric. The public leaderboard is made of 40% of the testing set. The private leaderboard is made of 60% others. The dataset is composed by:

- 8000 images in the training set, 3479 of which are in the category *North-South orientation*, 1856 in the category *East-West orientation*, 859 in the category *Flat roof* and 1806 in the category *Other*,

³<http://www.datasciencegame.com>

⁴<https://inclass.kaggle.com/c/data-science-game-2016-online-selection/leaderboard>

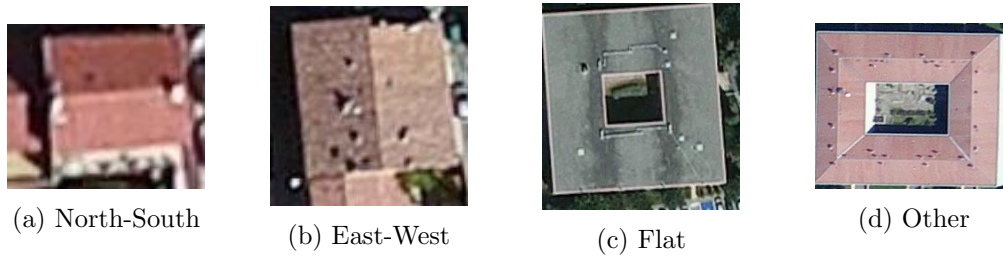


Figure 2.4: Illustrations of the DSG2016 online selection dataset

- 20,760 images in the training set, but without any label,
- 13,999 images in the testing set.

Rules Pretrain models are usable only if they are publicly available and if they are not trained on roof datasets.

2.3.2 Experiments

Winning solution One of the main difficulties of this challenge was the relatively small size of the dataset. Yet, we thought that the latter was big enough to fine-tune a pretrained model on. Our very first submission was a VGG16 finetuned on 80% of the trainset with some data augmentation and validated on the rest. It scored a whopping 82.67% on the public leaderboard (and 81.37% on the private), which put us on the top 3 at the time, and confirmed it was working.

We then tried InceptionV3 because of its good performance in many transfer learning setups. In order to reduce overfitting, we used the following techniques: online data augmentation, the 90 trick (multiplying by 2 the number of images and moving the rotated images from class "north-south" to class "east-west" and vice-versa) and early stopping. We didn't add more features (images height / weight) because it seemed to lead to more overfitting, which was a main concern.

Still, it was difficult to evaluate our models on the little test data that remained. There was a lot of variance between different epochs, folds and hyperparameters. To gain more stability and to use the full training set, we went to the bagging method. We fine tuned some InceptionV3 with different hyperparameters on random bootstraps. The final prediction is the result of a vote among all classifiers.

We also tweaked the predictions of our final submission to counteract the fact that the train set had unbalanced classes, while the test set was balanced. This improved a bit our

score.

Other solutions We fine tuned in a end-to-end manner a Spatial Transformer (STN) architecture with a first InceptionV3 to localize a region of interest and a second InceptionV3 to classify this region. We tried to constraint the kind of spatial transformation the network could achieve, or the number of parameters to tune. However, we think that the localizer learned the identity transformation. Also, it was expensive to train this network (twice the number of parameters compared to a simple InceptionV3) and overfitting was a main concern, thus we stopped to explore this solution.

Second team solution The team ranked second used a semi-supervised technique to augment their training set with unlabeled images. They also used a stacking of 84 models: 12 root models trained on 7 different versions of the images⁵. However they seemed to have overfitted too much the public testing set regarding the gap with their private score.

Fine tuned models	# Models	Data	Public (%)	Private (%)
VGG16	1	80%trainset	82.67	81.37
InceptionV3	1	80%trainset	83.69	82.61
InceptionV3 + STN	1	80%trainset	84.73	84.11
VGG19 (Second Team)	84	100%trainset + no label	87.60	86.38
InceptionV3	91	100%trainset	87.32	86.58
InceptionV3 + Prior	91	100%trainset	87.52	86.76

Table 2.4: Summary of the different models submitted for the DSG online challenge.

2.4 Conclusion

In this chapter, we studied the transfer capacity of the latest convolutional architectures. To do so, we compared three training approaches on several datasets which all have their own particularity : size, semantic distance from ImageNet, geometric variability of their regions of interest.

Recall that our approaches can be synthesize as follow:

- We trained specific deep convolutional networks randomly initializing their parameters (e.g. From Scratch).
- Further, we used pre-trained networks on ImageNet to extract features from the target dataset and trained a linear model (e.g. Features Extraction).

⁵<https://medium.com/@Zelros/how-deep-learning-solved-phase-1-of-the-data-science-game-2712b949963f>

- Finally, we fine tuned previous pre-trained networks to adapt their learned representations to the target dataset (e.g. Fine Tuning).

Each datasets represent a different challenge. In term of accuracy, we illustrated few trends:

- Fine Tuning achieve the best accuracy on medium datasets, and From Scratch achieves better accuracy than Features Extraction.
- Fine Tuning still achieve the best accuracy on small datasets, but From Scratch achieves lower accuracy than Features Extraction. Also, approaches based on a bagging of models are effective to achieve better accuracy.

In term of learning ability, we illustrated few other trends:

- InceptionV3 is the most accurate architecture and its batch normalization layers allow to converge faster.
- Adam is useful to reduce the number of experiments needed to find the optimal set of hyper parameters.
- Early stopping is a good way to control overfitting.

Chapter 3

Weakly Supervised Learning

3.1 Introduction

3.1.1 Definition

We define Weakly Supervised Learning (WSL) as a machine learning framework where the model is trained using examples that are only partially annotated or labeled. For instance, an object detector is typically trained on large collection of images manually annotated with masks or bounding boxes denoting the locations of objects of interest in each image. The reliance on time-consuming human labeling poses a significant limitation to the practical application of these methods. Moreover, manually annotation may not be optimal for the final prediction task. Thus, WSL aims at reducing the amount of human intervention needed.

Despite their excellent performances, current CNN architectures only carry limited invariance properties. Recently, attempts have been made to overcome this limitation using WSL. Many different kind of WSL techniques exists in the literature according to the type of labeled data used, latent variables learned and evaluation methods selected. We also consider that certain papers on multiple instance learning, attention based models or fine grained classification belong to a larger WSL framework. Below, we illustrate our thinking.

3.1.2 Multi Instance Learning

Weakly Supervised Learning (WSL) with Max Pooling In computer vision, the dominant approach for WSL is the Multiple Instance Learning (MIL) paradigm. An image

is considered as a bag of regions, and the model seeks the max scoring instance in each bag. In [31], the authors fixed the weights of a Vgg16 network and learned the last fully connected layer (e.g. 1x1 convolutional layer). They applied Vgg16 to images of bigger size than the original input (224x224). The resulting class score map is spatially reduced using a Max Pooling. They proposed two methods to learn scale invariance. The first was to rescale randomly the input image. The second was to train three specialized CNNs at three different scales and then to average their class scores. The two methods lead to almost the same results. They achieved 86.3 mean AP on VOC-2012. To compare, the same CNN used as a classical features extractor achieved 78.7 mean AP.

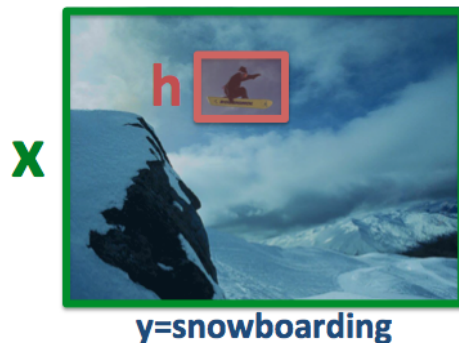


Figure 3.1: Illustration of a Weakly Supervised Learning approach which learn to focus on a region of interest from a global label.

WSL with k Max Min Pooling (WELDON) In [7], the authors proposed to use the same method than in [31], but the class score map is spatially reduced using a sum of k-Min Pooling and k-Max Pooling. As a result, it extends the selection of a single region to multiple high score regions incorporating also the low score regions, i.e. the negative evidences of a class appearance. They also proposed a ranking loss to optimize average precision. Using an average of 7 specialized CNNs at 7 different scales, they achieved 90.2 mean AP on VOC-2007 and 88.5 on VOC-2012. They also claim the state of the art on 5 other datasets. See figure 3.2.

3.1.3 Spatial Transformer Network

Principles The authors of [15] introduce a new learnable module, the Spatial Transformer, which explicitly allows the spatial manipulation of data within the network. This differentiable module can be viewed as a localization network which generate parameters for a grid generator. The grid is then used by a sampler to generate a certain transformation. The latter can be applied on the initial image or on the same feature map used as

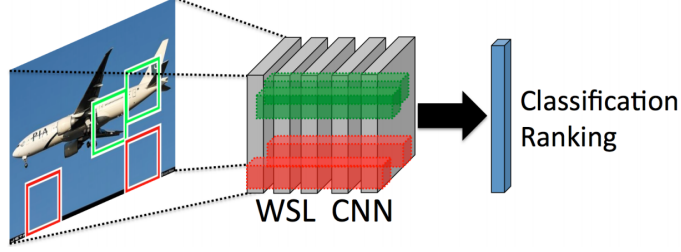


Figure 3.2: Illustration of the WELDON approach, a CNN trained in a weakly supervised manner to perform classification or ranking. It automatically selects multiple positive (green) or negative (red) evidences on several regions in the image.

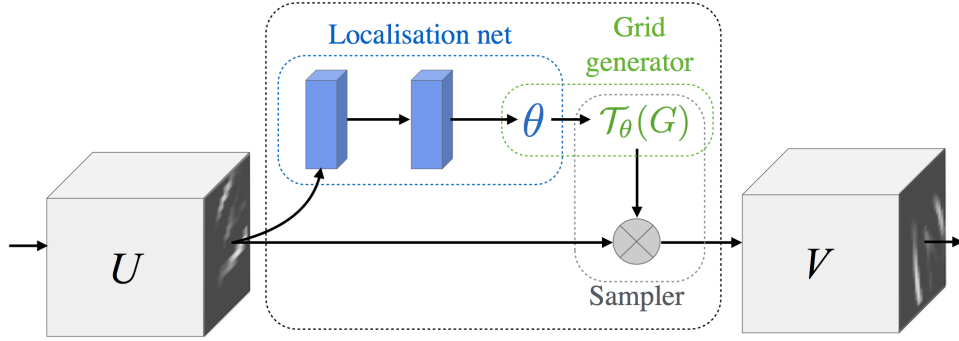


Figure 3.3: Spatial Transformer Module. The input feature map U is passed to a localization network which regresses the transformation parameters θ . The regular spatial grid G over V is transformed to the sampling grid $\mathcal{T}_\theta(G)$, which is applied to U , producing the warped output feature map V .

input to the localization network such as in figure 3.3.

Transformations Multiple transformation can be learned. For instance, $\mathcal{T}_\theta(G)$ can apply a 2D affine transformation A_θ . In this case, the pointwise transformation is

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \mathcal{T}_\theta(G) = A_\theta \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

However the class of transformations \mathcal{T}_θ may be more constrained, such as that used for attention

$$A_\theta = \begin{bmatrix} s & 0 & t_x \\ 0 & s & t_y \end{bmatrix}$$

allowing cropping, translation, and isotropic scaling by varying s , t_x , and t_y . Those parameters can be learned or, to constrain even more, fixed.

Traffic signs classification The first example of application is the supervised classification of a traffic signs dataset (GTSRB ¹). This kind of data is affected by contrast variation, rotational and translational changes. [11] used a spatial transformer network to make classification more robust and accurate. They used a modified version of GoogLeNet with batch normalization and added four different Spatial Transformer module made of a convolutional network as localizer. This method has several advantages over existing state of the art methods in terms of performance, scalability and memory requirement. Also, they reported a lower accuracy for their architecture without spatial transformer modules (99.57% against 99.81%). An implementation is available on the Torch7 blog ².



Figure 3.4: The use of a Spatial Transformer Module on a GTSRB data sample.

Fine grained bird classification The second example is the fine grained classification of a bird dataset (CUB-200-2011). The birds appear at a range of scales and orientations, are not tightly cropped, and require detailed texture and shape analysis to distinguish. This dataset contains 6k training images, 5.8k test images and covers 200 species of birds. They first claimed the state of the art accuracy of 82.3% with an Inception architecture pre-trained on ImageNet and fine tuned on CUB (previous best was 81.0%). Then, they trained a spatial transformer network containing 4 parallel spatial transformer modules parameterised for attention and acting on the input image. They achieved an accuracy of 84.1% outperforming their baseline by 1.8%. The resulting output from the spatial transformers for the classification network is somewhat pose-normalised representation of a bird. It was able to discover and learn part detectors in a data-driven manner without any additional supervision. Also, the use of spatial transformers allowed them to use 448px resolution input images without any impact in performance as the output of the transformed 448px images are downsampled to 224px before being processed.

¹<http://benchmark.ini.rub.de/?section=gtsrb&subsection=news>

²http://torch.ch/blog/2015/09/07/spatial_transformers.html

Model		
Cimpoi '15 [5]	66.7	
Zhang '14 [40]	74.9	
Branson '14 [3]	75.7	
Lin '15 [23]	80.9	
Simon '15 [30]	81.0	
CNN (ours) 224px	82.3	
2×ST-CNN 224px	83.1	
2×ST-CNN 448px	83.9	
4×ST-CNN 448px	84.1	

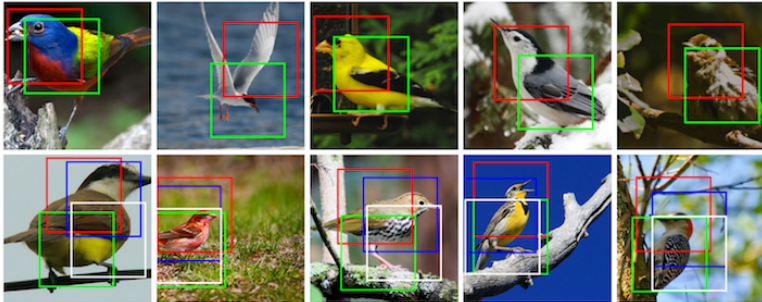


Figure 3.5: Left: the accuracy on CUB-200-2011 bird classification dataset. Spatial transformer networks with two spatial transformers modules (2 x ST-CNN) and four (4 x ST-CNN) in parallel achieve higher accuracy. Right: the transformation predicted by the 2 x ST-CNN (top row) and 4 x ST-CNN (bottom row) on the input image. Notably, for the 2 x ST-CNN, one of the module (shown in red) learned to detect heads, while the other (shown in green) detects the body.

3.2 Applying fine tuning to Weldon

3.2.1 Context

MIT67 It is the database of the Indoor scene recognition challenge ³. It was released during the CVPR 2009 [33]. It contains 67 categories. 5360 images compose the training set. 1340 images compose the testing set. The number of images does not vary across categories. This dataset is quite similar to Pascal Voc in term of size and difficulty, e.g. while some indoor scenes (e.g. corridors) can be well characterized by global spatial properties, others (e.g., bookstores) are better characterized by the objects they contain.



Figure 3.6: Illustration of the 67 indoor categories from the MIT67 dataset.

³<http://web.mit.edu/torralba/www/indoor.html>

3.2.2 Experiments

Fine tuning a Weldon architecture We reproduce the results of [7] and show in table 3.1 that Fine Tuning the Weldon architecture at three different scales improves the overall accuracy. However, the process of Fine Tuning is expensive, thus we do not provide results for bigger scale.

Model	Image size	L6 size	Report (*)	Extraction	Fine Tuning
Vgg16	224×224	1×1	69.93	70.30	70.60
Weldon	249×249	2×2	72.16	72.01	73.4
Weldon	280×280	3×3	72.98	73.80	74.03
Weldon	320×320	4×4	73.40	73.96	74.1

Table 3.1: Accuracy top 1 on MIT67 test set. Without data augmentation and dropout. $k = 1$ aggregation (one region min, one region max). L6 size represents the number of regions (e.g. instances) evaluated during the aggregation. For size 320, 16 regions are evaluated. Results from column Report (*) are reported from [7].

3.3 Study of Spatial Transformer Network

3.3.1 Context

Database In this section, we study the ability of the Spatial Transformer Network to learn large spatial invariance. In order to do so we create a special dataset using the original MNIST dataset padded with 2 black pixels (e.g. all our images are of scale 32×32 instead of the original 28×28). Then, we create a Translated MNIST dataset. All the images are inserted on a 100×100 background full of black pixels, and undergo spatial shifts on the x and y axes. Those shifts are randomly picked up between 0 and 68 (=100-32).

3.3.2 Previous work

Co-localization In the appendix A.2 of [15], the authors explored the use of spatial transformers in a co-localization scenario. Given a set of images that are assumed to contain instances of a common but unknown object class, the model learn from the images only to localize (with a bounding box) the common object. To achieve this, they adopted the supervision that the distance between the image crop corresponding to two correctly localized objects is smaller than to a randomly sampled image crop, in some embedding space. For a dataset $\mathcal{I} = I_n$ of N images, this translates to a triplet loss, where they

minimized the hinge loss

$$\sum_n^N \sum_{m \neq n}^M \max(0, \|e(I_n^T) - e(I_m^T)\|_2^2 - \|e(I_n^T) - e(I_n^{rand})\|_2^2 + \alpha)$$

They used translated (T), and translated and cluttered (TC) MNIST images (28x28) on a 84x84 black background. As features extractor ($e(\cdot)$), they used a pretrain network on MNIST. As localizer, they used a 100k parameters CNN. Also, they used a spatial transformer module parameterized for attention (scale, translation, no rotation).

They measured a digit to be correctly localized if the overlap (are of intersection divided by area of union) between the predicted bounding box and groundtruth bounding box is greater than 0.5. On T they got 100% accuracy. On CT between 75-93%.

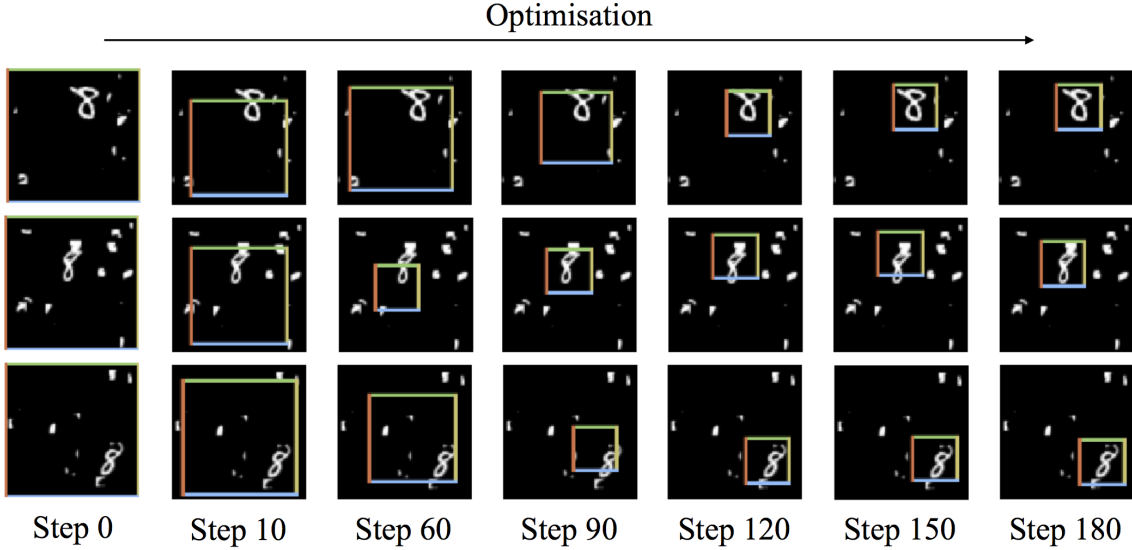


Figure 3.7: Illustration of the dynamics for co-localization. Here are the localization predicted by the spatial transformer for three of the 100 dataset images after the SGD step labelled below. By SGD step 180 the model has process has correctly localized the three digits.

3.3.3 Experiments

Almost all our experiments are made with a batch size of 256, Adam as optimizer, a learning rate of $3e^{-4}$, a learning rate decay of 0 and a weight decay of 0.

On table 3.2, we compare several approaches to a toy problem. The goal for a STN is to localize where the white pixels on the initial images are, then to generate a zoomed-in view of the digit.

ConvNet abilities to learn strong spatial invariance. (a) uses a LeNet-5 as classifier. It is an optimized architecture for the classification of 28x28 digits on a 32×32 background. In this case, LeNet-5 takes as input downsampled images. This process reduces the readability of digits. Thus, LeNet-5 is only able to achieve a 85.07% accuracy.

(b) uses a bigger convolutional network (ConvNet100). We want to measure the cost of the downsampling process. Thus, this network takes as input images of size 100×100 . Finally, it is able to achieve a 96.12%, 11% more than LeNet-5 (a). However, we do not use any data augmentation procedure.

STN abilities (c,d,e,g) uses a LeNet-5 in the same way and apply a transform to the initial image (100×100). However, the latter is also downsampled from 100×100 to 32×32 to fit the input size of the localizer.

Firstly, we show that even when using bad resolution images, the localizer is able to produce good spatial transformations, leading to almost the same accuracy as using the full resolution (d, f, i). Secondly, we show that generating 3, 4 or 6 parameters lead to almost the same accuracy.

Multi Instance Learning (MIL) abilities (f) uses a WSL with Max Pooling architecture as described in subsection 3.1.2. LeNet-5 is transformed to a fully convolutional network in order to take as input images of size 100×100 . The final features map is spatially aggregated by a Max Pooling of the size of the features map (e.g. the output is of size $1 \times 1 \times 10$). This method is the fastest to converge and lead to one of the highest accuracy (99.15%). We have the same results fine tuning a pretrain LeNet-5 on the original MNIST (e.g. a background of size 32×32) with an accuracy of 99.05%.

Model	Input size	Test accuracy top1
(a) LeNet-5	32×32	85.07
(b) ConvNet100	100×100	96.12
(c) STN affine	32×32	99.06
(d) STN translation	32×32	99.10
(e) STN translation+scale	32×32	99.10
(f) MIL MaxPooling	100×100	99.15
(g) STN translation+scale+rotation	32×32	99.18

Table 3.2: Comparison of different methods applied on our MNIST Translated dataset with a background of size 100×100 . We estimate that our results may vary plus or minus 0.10.

3.4 Conclusion

In this chapter, we studied Weakly Supervised Learning (WSL) approaches which can be synthesized as follow:

- Multi Instance Learning (MIL) with Max Pooling which considers an image as a bag of regions and seeks the max scoring region.
- MIL with k Max Min Pooling (e.g. Weldon) which extends the selection of a single region to multiple high score regions and low score regions.
- Spatial Transformer Network (STN) which uses a network (e.g. localizer) that takes as input the original image and generates a transformed image. Thus, the second network (e.g. classifier) takes as input a invariant representation of the object to classify.

In a first section, we studied the Weldon approach on a small and complex dataset (e.g. MIT67).

- We explained that it is well suited for this kind of datasets where regions of interest are multiple and negative evidences of classes are present.
- Especially, we showed that fine tuning is well suited for Weldon architectures.

In a second section, we studied the STN approach.

- Firstly, we explained in which cases STN are used and in which forms.
- Secondly, we compared on a toy dataset (e.g. Translated MNIST) classical Convolutional Neural Networks (CNNs), STN and MIL with Max Pooling approaches in term of spatial invariance capacity. Thus, we showed that STN was able to generate

invariant representations of the digits in order to achieve better accuracy than CNNs and the same accuracy than MIL.

Conclusion

Summary of Contributions In this master’s thesis, we studied deep learning architectures for classifying medium and small datasets of images.

- In a first chapter, we explained how Convolutional Neural Networks can achieve such good accuracy.
- In a second chapter, we showed the efficiency of the Fine Tuning approach on this kind of dataset. We also explained our winning solution to the DSG online challenge based on a bootstrap of fine tuned InceptionV3.
- In a last chapter, we showed the advantages and drawbacks of Weakly Supervised Learning approaches such as Multi Instance Learning (MIL) and Spatial Transformer Networks (STN). Using Fine Tuning, we also improved Weldon, a certain kind of MIL model.

Future Directions In future studies, we would like to adapt the methods developed during this study to multi-modal datasets made of images and texts. Furthermore, we would like to apply Fine Tuning and Weakly Supervised Learning on the last architectures such as Wide Residual Networks. Finally, we would like to explore hybrid weakly supervised architectures counteracting the drawbacks of MIL and STN, and seeking improvements.

Appendices

Appendix A

Overfeat

Layer id	Layer type	Parameters number
(0):	Image (3, 221, 221)	
(1):	Convolution (96, 3x7x7, 2x2, 0x0)	14,208
(3):	MaxPooling (3x3,3x3)	
(4):	Convolution (256, 96x3x3, 7x7)	1,204,480
(6):	MaxPooling (2x2,2x2)	
(7):	Convolution (512, 256x3x3, 1x1, 1x1)	1,180,160
(9):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(11):	Convolution (1024, 512x3x3, 1x1, 1x1)	4,719,616
(13):	Convolution (1024, 1024x3x3, 1x1, 1x1)	9,438,208
(15):	MaxPooling (3x3,3x3)	
(16):	FullyConnected (25600 \rightarrow 4096)	104,861,696
(18):	FullyConnected (4096 \rightarrow 4096)	16,781,312
(20):	FullyConnected (4096 \rightarrow 1000)	4,097,000
(21):	SoftMax	
Total :	9 layers	144,656,488

Table A.1: Deep architecture used in our experiments. A ReLU non-linearity follows each convolutional and fully-connected layers, beside the last one. Convolution (512, 512x3x3, 1x1, 1x1) means 512 filters (e.g. 512 output channels), a kernel size of 256x3x3, 1 step of the convolution to the width and height dimensions, 1 additional zero padded per width to the input, 1 per hight. MaxPooling (2,2,2,2) means a 2D pooling operation on 2x2 pixels neighborhood, by step size of 2x2

Appendix B

Vgg16

Layer id	Layer type	Parameters number
(0):	Image (3, 224, 224)	
(1):	Convolution (64, 3x3x3, 1x1, 1x1)	1,792
(3):	Convolution (64, 64x3x3, 1x1, 1x1)	36,928
(5):	MaxPooling (2x2,2x2)	
(6):	Convolution (128, 64x3x3, 1x1, 1x1)	73,856
(8):	Convolution (128, 128x3x3, 1x1, 1x1)	147,584
(10):	MaxPooling (2x2,2x2)	
(11):	Convolution (256, 128x3x3, 1x1, 1x1)	295,168
(13):	Convolution (256, 256x3x3, 1x1, 1x1)	590,080
(15):	Convolution (256, 256x3x3, 1x1, 1x1)	590,080
(17):	MaxPooling (2x2,2x2)	
(18):	Convolution (512, 256x3x3, 1x1, 1x1)	1,180,160
(20):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(22):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(24):	MaxPooling (2x2,2x2)	
(25):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(27):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(29):	Convolution (512, 512x3x3, 1x1, 1x1)	2,359,808
(31):	MaxPooling (2x2,2x2)	
(32):	FullyConnected (25088 \rightarrow 4096)	102,764,544
(34):	Dropout (50%)	
(35):	FullyConnected (4096 \rightarrow 4096)	16,781,312
(37):	Dropout (50%)	
(38):	FullyConnected (4096 \rightarrow 1000)	4,097,000
(40):	SoftMax	
Total :	16 layers	138,357,544

Table B.1: Deep architecture used in our experiments. A ReLU non-linearity follows each convolutional and fully-connected layers, beside the last one. Convolution (512, 512x3x3, 1x1, 1x1) means 512 filters (e.g. 512 output channels), a kernel size of 256x3x3, 1 step of the convolution to the width and height dimensions, 1 additional zero padded per width to the input, 1 per hight. MaxPooling (2,2,2,2) means a 2D pooling operation on 2x2 pixels neighborhood, by step size of 2x2

Appendix C

InceptionV3

Layer id	Layer type	Parameters number
(0):	Image (3, 299, 299)	
(1):	Convolution (32, 3x3x3, 2x2, 0x0) no bias	864
(2):	Batch Normalization	64*
(4):	Convolution (32, 32x3x3, 1x1, 0x0) no bias	9,216
(5):	Batch Normalization	64*
(7):	Convolution (64, 32x3x3, 1x1, 1x1) no bias	18,432
(8):	Batch Normalization	128*
(10):	MaxPooling (3x3,2x2)	
(11):	Convolution (80, 64x3x3, 1x1, 0x0) no bias	5,120
(12):	Batch Normalization	160*
(14):	Convolution (192, 80x3x3, 1x1, 0x0) no bias	138,240
(15):	Batch Normalization	384*
(17):	MaxPooling (3x3,2x2)	
(18):	4 x Inception1	-
(19):	4 x Inception2	-
(20):	Inception3	-
(20):	2 x Inception4	-
(21):	MaxPooling (8x8,1x1)	
(22):	FullyConnected (2048 \rightarrow 1000)	2,065,392
(23):	SoftMax	
Total :	42 layers	23,816,528

Table C.1: Inception architecture used in our experiments. A ReLU non-linearity follows each convolutional and fully-connected layers, beside the last one. * means that the parameters are not learned by backpropagation. For layers Inception1, Inception2, Inception3 and Inception4, please refer to the original paper [44]. Convolution (512, 512x3x3, 1x1, 1x1) means 512 filters (e.g. 512 output channels), a kernel size of 256x3x3, 1 step of the convolution to the width and height dimensions, 1 additional zero padded per width to the input, 1 per height. MaxPooling (2,2,2,2) means a 2D pooling operation on 2x2 pixels neighborhood, by step size of 2x2.

Bibliography

- [1] Sandra Avila, Nicolas Thome, Matthieu Cord, Eduardo Valle, and Arnaldo De A Araújo. Pooling in image representation: The visual codeword point of view. *Computer Vision and Image Understanding*, 117(5):453–465, 2013.
- [2] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, et al. Greedy layer-wise training of deep networks. *Advances in neural information processing systems*, 19:153, 2007.
- [3] Hakan Bilen and Andrea Vedaldi. Weakly supervised deep detection networks. *CVPR*, abs/1511.02853, 2016.
- [4] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *European Conference on Computer Vision*, pages 446–461. Springer, 2014.
- [5] Jean DaRolt. How is a convolutional neural network able to learn invariant features? <https://www.quora.com/How-is-a-convolutional-neural-network-able-to-learn-invariant-features/answer/Jean-Da-Rolt?srid=NRAY>, year = 2016, note = .
- [6] Sander Dieleman, Jeffrey De Fauw, and Koray Kavukcuoglu. Exploiting cyclic symmetry in convolutional neural networks.
- [7] Thibaut Durand, Nicolas Thome, and Matthieu Cord. Weldon: Weakly supervised learning of deep convolutional neural networks. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2016.
- [8] Abel Gonzalez-Garcia, Davide Modolo, and Vittorio Ferrari. Do semantic parts emerge in convolutional neural networks? *arXiv preprint arXiv:1607.03738*, 2016.
- [9] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. Book in preparation for MIT Press.

- [10] Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, and Gang Wang. Recent advances in convolutional neural networks. *CoRR*, abs/1512.07108, 2015.
- [11] Mrinal Haloi. Traffic sign classification using deep inception based convolutional networks. *arXiv preprint arXiv:1511.02992*, 2015.
- [12] Donald Hebb. 0.(1949) the organization of behavior, 1968.
- [13] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [14] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *JMLR*, abs/1502.03167, 2015.
- [15] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2008–2016, 2015.
- [16] Angjoo Kanazawa, Abhishek Sharma, and David Jacobs. Locally scale-invariant convolutional neural networks. *NIPS Workshop*, 2014.
- [17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, abs/1412.6980, 2015.
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [20] Yann Le Cun, Ido Kanter, and Sara A Solla. Eigenvalues of covariance matrices: Application to neural-network learning. *Physical Review Letters*, 66(18):2396, 1991.
- [21] Yann LeCun and Yoshua Bengio. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- [22] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [23] Zachary C Lipton, David C Kale, Charles Elkan, Randall Wetzell, Sharad Vikram, Julian McAuley, Randall C Wetzell, Zhanglong Ji, Balakrishnan Narayaswamy, Cheng-I Wang, et al. The mythos of model interpretability. *IEEE Spectrum*, 2016.

- [24] Chang Liu, Yu Cao, Yan Luo, Guanling Chen, Vinod Vokkarane, and Yunsheng Ma. Deepfood: Deep learning-based food image recognition for computer-aided dietary assessment. In *International Conference on Smart Homes and Health Telematics*, pages 37–48. Springer, 2016.
- [25] David G Lowe. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. Ieee, 1999.
- [26] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- [27] Austin Meyers, Nick Johnston, Vivek Rathod, Anoop Korattikara, Alex Gorban, Nathan Silberman, Sergio Guadarrama, George Papandreou, Jonathan Huang, and Kevin P Murphy. Im2calories: towards an automated mobile vision food diary. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1233–1241, 2015.
- [28] Dmytro Mishkin, Nikolay Sergievskiy, and Jiri Matas. Systematic evaluation of cnn advances on the imagenet. *arXiv preprint arXiv:1606.02228*, 2016.
- [29] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Inceptionism: Going deeper into neural networks. *Google Research Blog*. Retrieved June, 20, 2015.
- [30] Michael Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015. [Online; accessed 30-July-2016].
- [31] Maxime Oquab, Léon Bottou, Ivan Laptev, and Josef Sivic. Is object localization for free?-weakly-supervised learning with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 685–694, 2015.
- [32] Aaditya Prakash. One by one [1x1] Convolution - Counter-intuitively useful. <http://iamaaditya.github.io/2016/03/one-by-one-convolution>, year = 2016, note =.
- [33] Ariadna Quattoni and Antonio Torralba. Recognizing indoor scenes. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 413–420. IEEE, 2009.
- [34] Marc Aurelio Ranzato, Fu Jie Huang, Y-Lan Boureau, and Yann LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*, pages 1–8. IEEE, 2007.
- [35] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958.

- [36] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- [37] Juergen Schmidhuber. Deep learning in neural networks: An overview. *arXiv*, Apr 2014.
- [38] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. 2013.
- [39] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualizing image classification models and saliency maps. *ICLR Workshop*, 2014.
- [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. 2014.
- [41] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [42] Ilya Sutskever, James Martens, George E Dahl, and Geoffrey E Hinton. On the importance of initialization and momentum in deep learning. *ICML (3)*, 28:1139–1147, 2013.
- [43] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [44] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *CVPR*, 2016.
- [45] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014.
- [46] Andreas Veit, Michael J. Wilber, and Serge J. Belongie. Residual networks are exponential ensembles of relatively shallow networks. *CoRR*, abs/1605.06431, 2016.
- [47] Xin Wang, Devinder Kumar, Nicolas Thome, Matthieu Cord, and Frederic Precioso. Recipe recognition with large multimodal food dataset. In *Multimedia & Expo Workshops (ICMEW), 2015 IEEE International Conference on*, pages 1–6. IEEE, 2015.
- [48] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014.

- [49] Matthew D Zeiler, Dilip Krishnan, Graham W Taylor, and Rob Fergus. Deconvolutional networks. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2528–2535. IEEE, 2010.
- [50] Sixin Zhang, Anna Choromanska, and Yann LeCun. Deep learning with elastic averaging SGD. *NIPS*, abs/1412.6651, 2015.