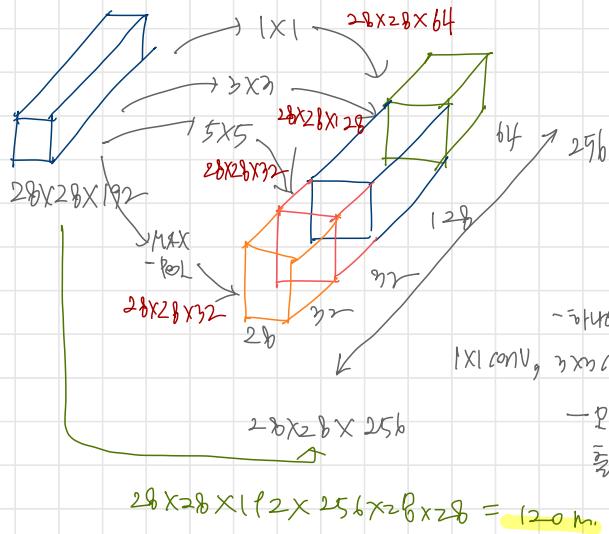



1. Motivation for Inception Network



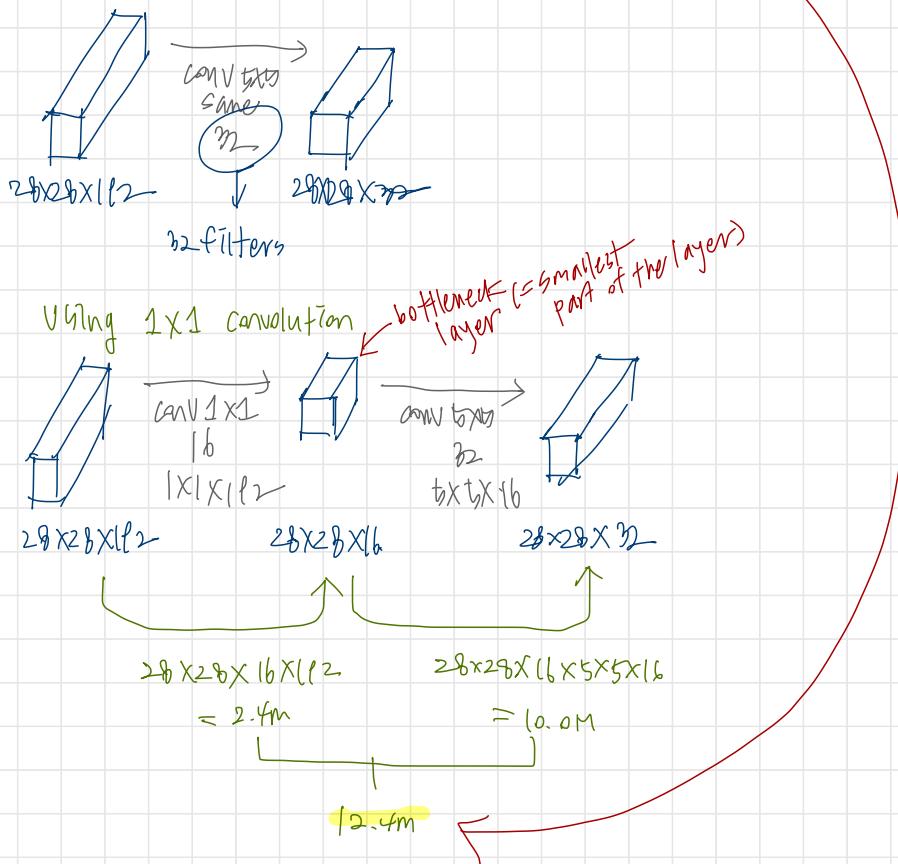
- 각 브랜치마다 별도로 계산하고 그 결과를 합침

- 모든 output을 stack/concatenate
할 때에 차원이 맞지 않음.

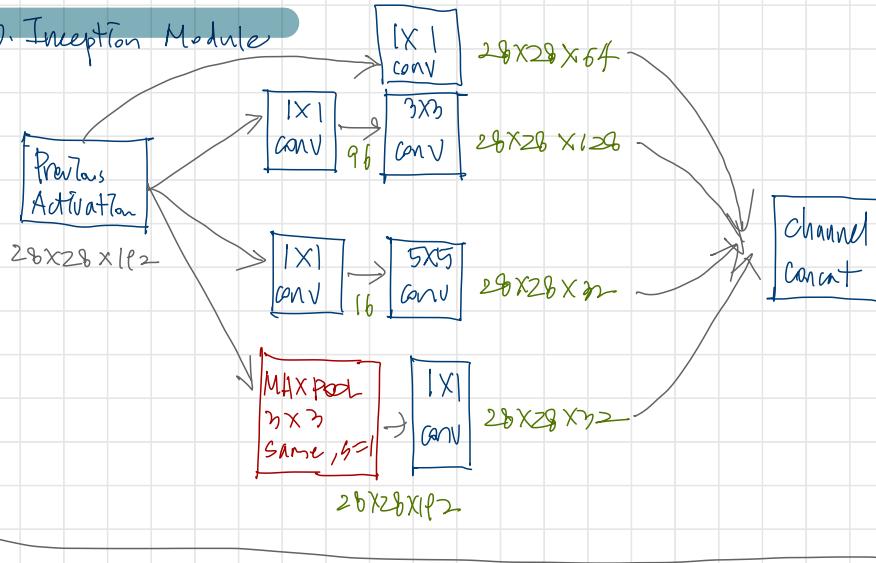
- 모든 output은 stack/concatenate

할 때에 차원이 맞지 않음.

2. The problem of Convolutional Cost



IMPROVED
COMPUTATIONAL
COSTS



2) 다른 Inception Module

Inception Network

- Inception Network은 전부 노드별로 흐름이 되며, 마지막 출력층으로 Fully Connected Layer + Softmax가 존재한다. 그걸에 특이한 점은 중간에 몇 가지로 나뉘어 Inception Module에서 예측을 받았을 때 FCLayer + Softmax가 중간층을 처리한다.

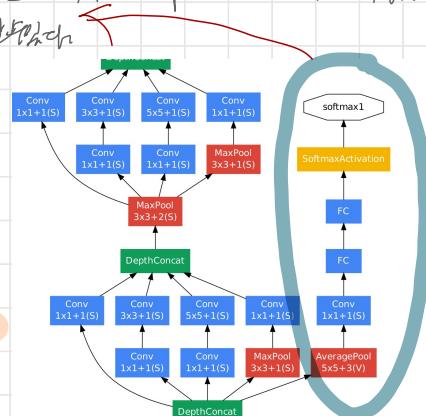
이렇게 궁금해 사용한 예술을 드린

이제가의 결과는 예측하는지 알기

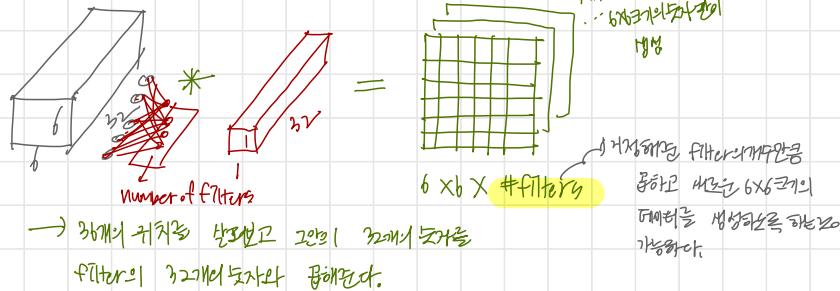
아주나쁘지 않다는 것은 행복한다.

↳ REGULARIZATION (26/31:24)

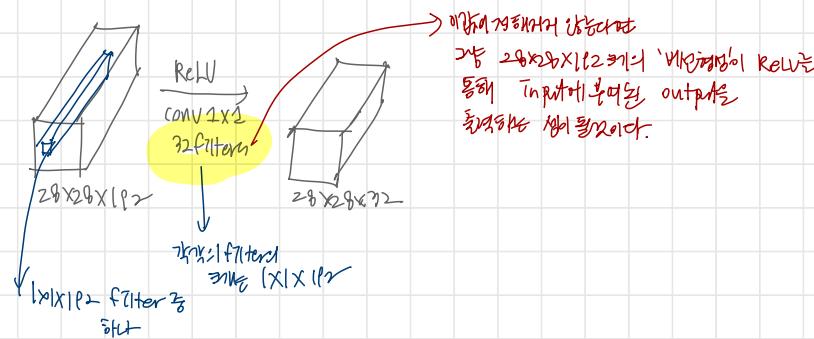
의 문제는 σ_{true} overfitting
방법은 같은 match.



Cf. What does a 1×1 convolution do?



∴ Fully Connected NN applied to the $32 \text{ numbers} \times \text{number of filters}$



Abstract

We propose a deep convolutional neural network architecture codenamed Inception, which was responsible for setting the new state-of-the-art for classification and detection in the ImageNet Large-Scale Visual Recognition Challenge 2014, (ILSVRC14). The main hallmark of this architecture is the improved utilization of the computing resources inside the network. This was achieved by a carefully crafted design that allows for increasing the depth and width of the network while keeping the computational budget constant. To optimize quality, the architectural decisions were based on the Hebbian principle and the intuition of multi-scale processing. One particular incarnation used in our submission for ILSVRC14 is called GoogLeNet, a 22 layers deep network, the quality of which is assessed in the context of classification and detection.

6) 예상대로 'Inception'이라는 이름의 deep convolution NN architectures 제작된다.

ILSVRC 2014 대회에서 우승한 모델은 GoogleNet (GoogLeNet)이 Inception라는 이름으로 구현된다.

cf. Hebbian Principle = "neurons that fire together, wire together"
= 동시에 활성화된 노드들 간에 연결성이 있다

1 Introduction

In the last three years, mainly due to the advances of deep learning, more concretely convolutional networks [10], the quality of image recognition and object detection has been progressing at a dramatic pace. One encouraging news is that most of this progress is not just the result of more powerful hardware, larger datasets and bigger models, but mainly a consequence of new ideas, algorithms and improved network architectures. No new data sources were used, for example, by the top entries in the ILSVRC 2014 competition besides the classification dataset of the same competition for detection purposes. Our GoogLeNet submission to ILSVRC 2014 actually uses $12\times$ fewer parameters than the winning architecture of Krizhevsky et al [9] from two years ago, while being significantly more accurate. The biggest gains in object-detection have not come from the utilization of deep networks alone or bigger models, but from the synergy of deep architectures and classical computer vision, like the R-CNN algorithm by Girshick et al [6].

Another notable factor is that with the ongoing fraction of mobile and embedded computing, the efficiency of our algorithms – especially their power and memory use – gains importance. It is noteworthy that the considerations leading to the design of the deep architecture presented in this paper included this factor rather than having a sheer fixation on accuracy numbers. For most of the experiments, the models were designed to keep a computational budget of 1.5 billion multiply-adds at inference time, so that they do not end up to be a purely academic curiosity, but could be put to real world use, even on large datasets, at a reasonable cost.

— 망상 CNN과 함께 가로, 세로 detection 및 classification을 동시에 처리하기 위해 디자인되었는데, 이를 대비해 양의 경쟁에서 기여하는 것이 아닌 혁신적인 네트워크 설계가 이루어졌다.

특히 GoogleNet은 수동적 처리방법 대비하여 동적 프로세싱이라는 parameter를 확장하여 높은 정확도를 보인다.

— 또한 유연성과 확장성을 위해 각각의 계층을 별개로 처리하는 모듈화된 구조로 "ModularNet"이 된다

① 네트워크가 깊어짐

② Inception Module라는 핵심 패턴

2 Related Work

Starting with LeNet-5 [10], convolutional neural networks (CNN) have typically had a standard structure – stacked convolutional layers (optionally followed by contrast normalization and max-pooling) are followed by one or more fully-connected layers. Variants of this basic design are prevalent in the image classification literature and have yielded the best results to-date on MNIST, CIFAR and most notably on the ImageNet classification challenge [9, 21]. For larger datasets such as Imagenet, the recent trend has been to increase the number of layers [12] and layer size [21, 14], while using dropout [7] to address the problem of overfitting.

Despite concerns that max-pooling layers result in loss of accurate spatial information, the same convolutional network architecture as [9] has also been successfully employed for localization [9, 14], object detection [6, 14, 18, 5] and human pose estimation [19]. Inspired by a neuroscience model of the primate visual cortex, Serre et al. [15] use a series of fixed Gabor filters of different sizes in order to handle multiple scales, similarly to the Inception model. However, contrary to the fixed 2-layer deep model of [15], all filters in the Inception model are learned. Furthermore, Inception layers are repeated many times, leading to a 22-layer deep model in the case of the GoogLeNet model.

Network-in-Network is an approach proposed by Lin et al. [12] in order to increase the representational power of neural networks. When applied to convolutional layers, the method could be viewed as additional 1×1 convolutional layers followed typically by the rectified linear activation [9]. This enables it to be easily integrated in the current CNN pipelines. We use this approach heavily in our architecture. However, in our setting, 1×1 convolutions have dual purpose: most critically, they are used mainly as dimension reduction modules to remove computational bottlenecks, that would otherwise limit the size of our networks. This allows for not just increasing the depth, but also the width of our networks without significant performance penalty.

The current leading approach for object detection is the Regions with Convolutional Neural Networks (R-CNN) proposed by Girshick et al. [6]. R-CNN decomposes the overall detection problem into two subproblems: to first utilize low-level cues such as color and superpixel consistency for potential object proposals in a category-agnostic fashion and then use CNN classifiers to identify object categories at those locations. Such a two stage approach leverages the accuracy of bounding box segmentation with low-level cues, as well as the highly powerful classification power of state-of-the-art CNNs. We adopted a similar pipeline in our detection submissions, but have explored enhancements in both stages, such as multi-box [5] prediction for higher object bounding box recall, and ensemble approaches for better categorization of bounding box proposals.

① LeNet-5 + AlexNet

Stacked Conv Layer

+

Contrast Normalization / Max Pooling Layer

+

Fully Connected Layer

- ImageNet dataset을 다음 막수,

Layer의 수가 적어 dropout 적용하기 어렵기 때문에

Max-pooling layers는 확장하여 accurate spatial info를 얻는다.

우리에게도 AlexNet은 놀라워라를 원한다.

② R-CNN

① Multi-Box Prediction

② Ensemble Approach

↳ Multi-Box Prediction
 - low-level features를 기반으로 selective search 기법으로 사용하는데, 이는 exhaustive search와 segmentation을 결합한 것이다.
 - 실제로 Time 단위로 Time을 대상으로 한 이미지를 Time을 대상으로 한 이미지를 통해 각각의 이미지에서 동일한 input에 대해서 multiple instances를 detecting

LeNet

AlexNet

Serre et al.

NIN

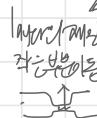
R-CNN

NIN

- 다양한 표면적을 가진 물체

- 1×1 Conv Layer는 멀티스케일 특징을 갖도록 증가시킨다.

↳ 이유: ① computational bottleneck을 차복하기 위해서 차원축소 사용



② 경계의 거리와 함께 depth를 유지하지 않음

3 Motivation and High Level Considerations

The most straightforward way of improving the performance of deep neural networks is by increasing their size. This includes both increasing the **depth** – the number of levels – of the network and its width: the number of units at each level. This is as an easy and safe way of training higher quality models, especially given the availability of a large amount of labeled training data. However this simple solution comes with two major drawbacks.

Bigger size typically means a **larger number of parameters**, which makes the enlarged network more prone to **overfitting**, especially if the number of labeled examples in the training set is limited. This can become a major bottleneck, since the creation of high quality training sets can be tricky

and expensive, especially if expert human raters are necessary to distinguish between fine-grained visual categories like those in ImageNet (even in the 1000-class ILSVRC subset) as demonstrated by Figure 1.

Another drawback of uniformly increased network size is the **dramatically increased use of computational resources**. For example, in a deep vision network, if two convolutional layers are chained, any uniform increase in the number of their filters results in a quadratic increase of computation. If the added capacity is used inefficiently (for example, if most weights end up to be close to zero), then a lot of computation is wasted. Since in practice the computational budget is always finite, an efficient distribution of computing resources is preferred to an indiscriminate increase of size, even when the main objective is to increase the quality of results.

The fundamental way of solving both issues would be by ultimately moving from fully connected to sparsely connected architectures, even inside the convolutions. Besides mimicking biological systems, this would also have the advantage of firmer theoretical underpinnings due to the groundbreaking work of Arora et al. [2]. Their main result states that if the probability distribution of the data-set is representable by a large, very sparse deep neural network, then the optimal network topology can be constructed layer by layer by analyzing the correlation statistics of the activations of the last layer and clustering neurons with highly correlated outputs. Although the strict mathematical proof requires very strong conditions, the fact that this statement resonates with the well known Hebbian principle – neurons that fire together, wire together – suggests that the underlying idea is applicable even under less strict conditions, in practice.

On the downside, todays computing infrastructures are very inefficient when it comes to numerical calculation on non-uniform sparse data structures. Even if the number of arithmetic operations is reduced by 100x, the overhead of lookups and cache misses is so dominant that switching to sparse matrices would not pay off. The gap is widened even further by the use of steadily improving, highly tuned, numerical libraries that allow for extremely fast dense matrix multiplication, exploiting the minute details of the underlying CPU or GPU hardware [16, 9]. Also, non-uniform sparse models require more sophisticated engineering and computing infrastructure. Most current vision oriented machine learning systems utilize sparsity in the spatial domain just by the virtue of employing convolutions. However, convolutions are implemented as collections of dense connections to the patches in the earlier layer. ConvNets have traditionally used random and sparse connection tables in the feature dimensions since [11] in order to break the symmetry and improve learning, the trend changed back to full connections with [9] in order to better optimize parallel computing. The uniformity of the structure and a large number of filters and greater batch size allow for utilizing efficient dense computation.

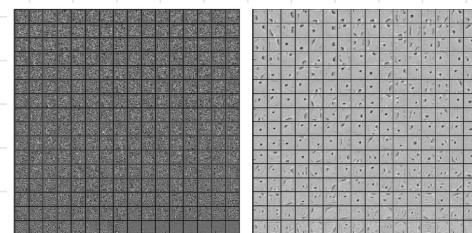
This raises the question whether there is any hope for a next, intermediate step: an architecture that makes use of the extra sparsity, even at filter level, as suggested by the theory, but exploits our current hardware by utilizing computations on dense matrices. The vast literature on sparse matrix computations (e.g. [3]) suggests that converting sparse matrices into relatively dense submatrices tends to give rise to an practical performance for sparse matrix multiplication. It does not seem far-fetched that similar methods would be utilized for the automated construction of non-uniform deep-learning architectures in the near future.

The Inception architectures started out as a case study of the first author for assessing the hypothetical output of a sophisticated network topology construction algorithm that tries to approximate a sparse structure implied by [2] for vision networks and covering the hypothesized outcome by dense, readily available components. Despite being a highly speculative undertaking, only after two iterations on the exact choice of topology, we could already see modest gains against the reference architecture based on [12]. After further tuning of learning rate, hyperparameters and improved training methodology, we established that the resulting Inception architecture was especially useful in the context of localization and object detection as the base network for [6] and [5]. Interestingly, while most of the original architectural choices have been questioned and tested thoroughly, they turned out to be at least locally optimal.

→ DNN의 깊이와 넓이의 학습
깊이(depth)는 학습에 필요한 계산량이 증가하는 경향
+ 넓이(= layer의 노드 수)는
① Parameter ↑ → overfitting
② computational resource ↑

strategy → Fully Connected / Conv Layer의 mix
sparse하게 만들면서 sparsity를 살피기

: sparse 시장 대체로 gradient를 계산하는데 있어 계산량이 줄어들면 좋지만, 이런
strict condition은 어떤 계산량을 줄여야 하는지에 대한 node는
계산량을 줄여야 하는지에 대한 조건은 계산량을 줄여야 하는지에 대한 조건은



(a) Without dropout (b) Dropout with $p = 0.5$.

이전에 dropout은 확률적 sparsity를 추구한 데에서 대비해
확률에 의해 가능성을 좁힌 느낌을 빼고 있다.

즉, dropout은 확률적 sparsity를 추구하는 대신
확률에 의해 가능성을 좁힌 느낌을 빼고 있다.

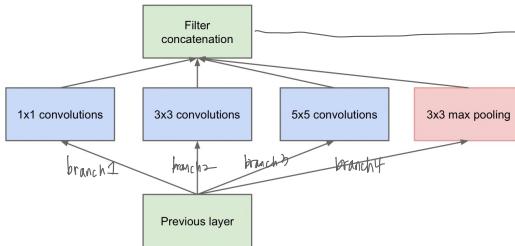
다른 한 쪽은 확률적 sparsity를 추구하는 대신
확률에 의해 가능성을 좁힌 느낌을 빼고 있다.

4 Architectural Details

The main idea of the Inception architecture is based on finding out how an optimal local sparse structure in a convolutional vision network can be approximated and covered by readily available dense components. Note that assuming translation invariance means that our network will be built from convolutional building blocks. All we need to do is to find the optimal local construction and to repeat it spatially. Arora et al. [2] suggests a layer-by-layer construction in which one should analyze the correlation statistics of the last layer and cluster them into groups of units with high correlation. These clusters form the units of the next layer and are connected to the units in the previous layer. We assume that each unit from the earlier layer corresponds to some region of the input image and these units are grouped into filter banks. In the lower layers (the ones close to the input) correlated units would concentrate in local regions. This means, we would end up with a lot of clusters concentrated in a single region and they can be covered by a layer of 1×1 convolutions in the next layer, as suggested in [12]. However, one can also expect that there will be a smaller number of more spatially spread out clusters that can be covered by convolutions over larger patches, and there will be a decreasing number of patches over larger and larger regions. In order to avoid patch alignment issues, current incarnations of the Inception architecture are restricted to filter sizes 1×1 , 3×3 and 5×5 . However this decision was based more on convenience rather than necessity. It also means that the suggested architecture is a combination of all those layers with their output filter banks concatenated into a single output vector forming the input of the next stage. Additionally, since pooling operations have been essential for the success in current state of the art convolutional networks, it suggests that adding an alternative parallel pooling path in each such stage should have additional beneficial effect, too (see Figure 2(a)).

As these "Inception modules" are stacked on top of each other, their output correlation statistics are bound to vary: as features of higher abstraction are captured by higher layers, their spatial concentration is expected to decrease suggesting that the ratio of 3×3 and 5×5 convolutions should increase as we move to higher layers.

④ 017년 layer별 각 unit의 간접성이 커지면서, 이를 unit별 filterbank로 고려하는 방
→ 각 filter size를 통해 각 unit별 cluster로 고려하는 방 (Patch Alignment mechanism)



(a) Inception module, naïve version

Patch-alignment Trace

Filtered image patch patch size 10x10

Patch size 8x8

Patch size 4x4

Patch size 2x2

Patch size 1x1

→ result of naive version requires computation,

the clusters of the filters will be
a single output vector converted to the
next layer.

① 각 filter size별 convolution layer는 해당하는 해당 cluster로 대응되는 데에 맞춰, 이를 통해 filter 사이즈별로 대응된다.

problem

optimal space structure는 어떠한지를 찾기 위해 다양한 방법으로 고려하는 방향.

→ expand module은 inception module을 확장하는 방

This leads to the second idea of the proposed architecture: judiciously applying dimension reductions and projections wherever the computational requirements would increase too much otherwise. This is based on the success of embeddings: even low dimensional embeddings might contain a lot of information about a relatively large image patch. However, embeddings represent information in a dense, compressed form and compressed information is harder to model. We would like to keep our representation sparse at most places (as required by the conditions of [2]) and compress the signals only whenever they have to be aggregated en masse. That is, 1×1 convolutions are used to compute reductions before the expensive 3×3 and 5×5 convolutions. Besides being used as reductions, they also include the use of rectified linear activation which makes them dual-purpose. The final result is depicted in Figure 2(b).

In general, an Inception network is a network consisting of modules of the above type stacked upon each other, with occasional max-pooling layers with stride 2 to halve the resolution of the grid. For technical reasons (memory efficiency during training), it seemed beneficial to start using Inception modules only at higher layers while keeping the lower layers in traditional convolutional fashion. This is not strictly necessary, simply reflecting some infrastructural inefficiencies in our current implementation.

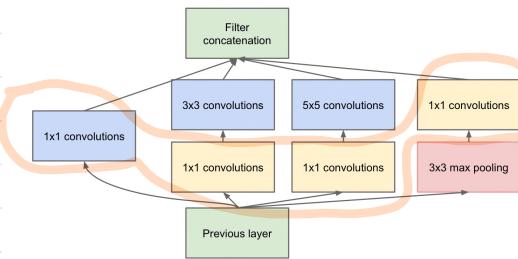
One of the main beneficial aspects of this architecture is that it allows for increasing the number of units at each stage significantly without an uncontrolled blow-up in computational complexity. The ubiquitous use of dimension reduction allows for shielding the large number of input filters of the last stage to the next layer, first reducing their dimension before convolving over them with a large patch size. Another practically useful aspect of this design is that it aligns with the intuition that visual information should be processed at various scales and then aggregated so that the next stage can abstract features from different scales simultaneously.

The improved use of computational resources allows for increasing both the width of each stage as well as the number of stages without getting into computational difficulties. Another way to utilize the inception architecture is to create slightly inferior, but computationally cheaper versions of it. We have found that all the included knobs and levers allow for a controlled balancing of computational resources that can result in networks that are 2–3× faster than similarly performing networks with non-Inception architecture, however this requires careful manual design at this point.

이해도는 잘 되었지만 예상보다 더 어렵다.
정리해 주시면 좋을 것 같다.

Embedding = turns positive integers into dense vectors of a fixed size
Input = (batchsize, inputlength)
Output = (batchsize, inputlength, outputdim)

1x1 patch size는 the expensive convolutional dimensionality reduction을 가능케하는 중요한 관리 요인이다.
Patch size는 연산 효율성을 크게 증가시킨다.
Patch size를 다양한 scale로 처리해야 하기 때문에 각각의 stage에서 서로 다른 scale을 블록화하여 feature를 쉽게 추출하게 된다.



(b) Inception module with dimension reductions

1x1 convolution filters
3x3, 5x5를 거친 경우
Necessary processing power를 대체할 수 있다.
Model depth
부피적이고 깊은 layer를 처리할 때 이를
줄여준다.

embedding을 적용한 후에 정복을 완료하기 시작할 때,

이전 단계에 대한 정보가 필요하다. 그래서 각각의 1x1의 convolution layer는 활성화
함수를 갖지만, 비활성화된 activation function은 의미가 없어.

type	patch size/ stride	output size	depth	#1x1	#3x3 reduce	#3x3	#5x5 reduce	#5x5	pool proj	params	ops
convolution	7x7/2	112x112x64	1							2.7K	34M
max pool	3x3/2	56x56x64	0								
convolution	3x3/1	56x56x192	2		64	192				112K	360M
max pool	3x3/2	28x28x192	0								
inception (3a)		28x28x256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28x28x480	2	128	128	192	32	96	64	380K	304M
max pool	3x3/2	14x14x480	0								
inception (4a)		14x14x512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14x14x512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14x14x512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14x14x528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14x14x832	2	256	160	320	32	128	128	840K	170M
max pool	3x3/2	7x7x832	0								
inception (5a)		7x7x832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7x7x1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7x7/1	1x1x1024	0								
dropout (40%)		1x1x1024	0								
linear		1x1x1000	1							1000K	1M
softmax		1x1x1000	0								

All the convolutions, including those inside the Inception modules, use rectified linear activation. The size of the receptive field in our network is 224×224 taking RGB color channels with mean subtraction. “#3x3 reduce” and “#5x5 reduce” stands for the number of 1×1 filters in the reduction layer used before the 3×3 and 5×5 convolutions. One can see the number of 1×1 filters in the projection layer after the built-in max-pooling in the pool proj column. All these reduction/projection layers use rectified linear activation as well.

ReLU activation function

ReLU (ReLU)

The network was designed with computational efficiency and practicality in mind, so that inference can be run on individual devices including even those with limited computational resources, especially with low-memory footprint. The network is 22 layers deep when counting only layers with parameters (or 27 layers if we also count pooling). The overall number of layers (independent building blocks) used for the construction of the network is about 100. However this number depends on the machine learning infrastructure system used. The use of average pooling before the classifier is based on [12], although our implementation differs in that we use an extra linear layer. This enables adapting and fine-tuning our networks for other label sets easily, but it is mostly convenience and we do not expect it to have a major effect. It was found that a move from fully connected layers to average pooling improved the top-1 accuracy by about 0.6%, however the use of dropout remained essential even after removing the fully connected layers.

Layer-wise loss 71.4%

Layer-wise 21 layers of tree
71.4%

→ Networked layers where layers share gradient

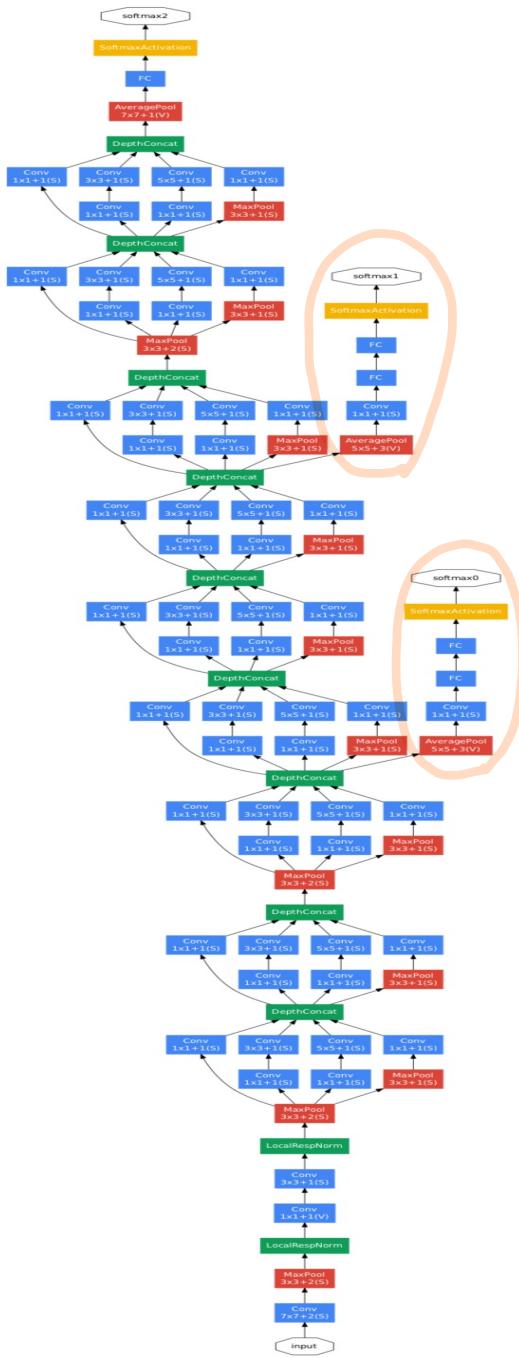
Given the relatively large depth of the network, the ability to propagate gradients back through all the layers in an effective manner was a concern. One interesting insight is that the strong performance of relatively shallower networks on this task suggests that the features produced by the layers in the middle of the network should be very discriminative. By adding auxiliary classifiers connected to these intermediate layers, we would expect to encourage discrimination in the lower stages in the classifier, increase the gradient signal that gets propagated back, and provide additional regularization. These classifiers take the form of smaller convolutional networks put on top of the output of the Inception (4a) and (4d) modules. During training, their loss gets added to the total loss of the network with a discount weight (the losses of the auxiliary classifiers were weighted by 0.3). At inference time, these auxiliary networks are discarded.

→ Layer-wise loss 71.4% (Layer-wise loss 71.4%)

→ Layer-wise loss 71.4% (Layer-wise loss 71.4%)

→ Layer-wise loss 71.4% (Layer-wise loss 71.4%)

Networked layers
layered feature map
shallow module
deep module
auxiliary classifier



* Auxiliary Network

- An average pooling layer with 5×5 filter size and stride 3, resulting in an $4 \times 4 \times 512$ output for the (4a), and $4 \times 4 \times 528$ for the (4d) stage.
- A 1×1 convolution with 128 filters for dimension reduction and rectified linear activation.
- A fully connected layer with 1024 units and rectified linear activation.
- A dropout layer with 70% ratio of dropped outputs.
- A linear layer with softmax loss as the classifier (predicting the same 1000 classes as the main classifier, but removed at inference time).

6 Training Methodology

Our networks were trained using the DistBelief [4] distributed machine learning system using modest amount of model and data-parallelism. Although we used CPU based implementation only, a rough estimate suggests that the GoogLeNet network could be trained to convergence using few high-end GPUs within a week, the main limitation being the memory usage. Our training used asynchronous stochastic gradient descent with 0.9 momentum [17], fixed learning rate schedule (decreasing the learning rate by 4% every 8 epochs). Polyak averaging [13] was used to create the final model used at inference time.

Our image sampling methods have changed substantially over the months leading to the competition, and already converged models were trained on with other options, sometimes in conjunction with changed hyperparameters, like dropout and learning rate, so it is hard to give a definitive guidance to the most effective single way to train these networks. To complicate matters further, some of the models were mainly trained on smaller relative crops, others on larger ones, inspired by [8]. Still, one prescription that was verified to work very well after the competition includes sampling of various sized patches of the image whose size is distributed evenly between 8% and 100% of the image area and whose aspect ratio is chosen randomly between 3/4 and 4/3. Also, we found that the photometric distortions by Andrew Howard [8] were useful to combat overfitting to some extent. In addition, we started to use random interpolation methods (bilinear, area, nearest neighbor and cubic, with equal probability) for resizing relatively late and in conjunction with other hyperparameter changes, so we could not tell definitely whether the final results were affected positively by their use.

- 0.9 momentum \rightarrow asynchronous SGD method
- Fixed Learning Rate Schedule (8% ~ 100% learning rate \rightarrow $f(\text{epoch}) = \text{lr} \cdot e^{-\frac{\text{epoch}}{8}}$ ($\text{lr} = 0.96$))
- Polyak Averaging
 - instead of using the actual parameter vector,
keep a moving average of the parameter vector and use that at test time.
- 다이버레이티브 Random Interpolation Method \rightarrow resize \rightarrow patch sampling \rightarrow SPCA gradient descent
 - 8% - 10% \rightarrow patch sampling \rightarrow SPCA gradient descent
 - photometric distortion \rightarrow image augmentation

9 Conclusions

Our results seem to yield a solid evidence that approximating the expected optimal sparse structure by readily available dense building blocks is a viable method for improving neural networks for computer vision. The main advantage of this method is a significant quality gain at a modest increase of computational requirements compared to shallower and less wide networks. Also note that our detection work was competitive despite of neither utilizing context nor performing bounding box

regression and this fact provides further evidence of the strength of the Inception architecture. Although it is expected that similar quality of result can be achieved by much more expensive networks of similar depth and width, our approach yields solid evidence that moving to sparser architectures is feasible and useful idea in general. This suggest promising future work towards creating sparser and more refined structures in automated ways on the basis of [2].

- 딥러닝 모델의 성능은 각각의 convolution layers를 통한 optimal sparse structure로 구현되는,
단계별로 계산량이 줄어든다.
- dense building blocks를 활용하는 Inception computational cost는 훨씬 적고
상당히 효율적인 성능을 보여준다.
- object detection과 context의 상관성, bounding box regression에 대한 성능은 Inception
이든 Inception Architecture의 경우 크게 차이가 없다.