

---

# Popular Network Architectures (& BatchNorm)

---

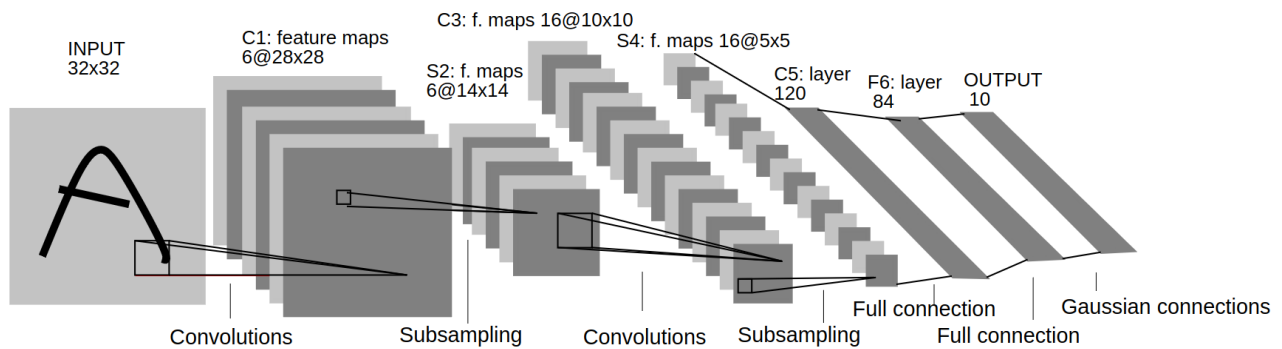
*LeNet-5 [5], MNIST, AlexNet [4], ImageNet, VGG [6], Inception (GoogleLeNet), BatchNorm [3],  
ResNet [2], DenseNet, Squeeze-Excite Net U-Net, Data Augmentation*

*Author: Anton Zhitomirsky*

## Contents

<b>1</b>	<b>LeNet-5</b>	<b>2</b>
1.1	shared weights . . . . .	3
1.2	sub-sampling . . . . .	3
1.3	Loss . . . . .	3
<b>2</b>	<b>AlexNet</b>	<b>4</b>
<b>3</b>	<b>VGG</b>	<b>6</b>
3.1	fewer wide convolutions or more narrow convolutions? . . . . .	7
3.2	The VGG block . . . . .	7
3.3	Performance . . . . .	8
<b>4</b>	<b>Batch Norm</b>	<b>8</b>
4.1	Motivation . . . . .	8
4.2	Solution . . . . .	8
4.3	Advantages . . . . .	9
4.4	Updated Aim . . . . .	9
4.5	Application . . . . .	9
4.6	Randomization during use . . . . .	9

# 1 LeNet-5



**Figure 1:** Architecture of LeNet-5 a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical [5]

```

1 # from https://pytorch.org/tutorials/beginner/blitz/neural_networks_tutorial.html
2 # https://www.analyticsvidhya.com/blog/2023/11/lenet-architectural-insights-and-practical
  ↳ -implementation/
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6
7
8 class Net(nn.Module):
9
10     def __init__(self):
11         super(Net, self).__init__()
12         # 1 input image channel, 6 output channels, 5x5 square convolution
13         # kernel
14         self.conv1 = nn.Conv2d(1, 6, 5)
15         self.pool1 = nn.AvgPool2d(2, stride=2)
16         self.conv2 = nn.Conv2d(6, 16, 5)
17         self.pool2 = nn.AvgPool2d(2, stride=2)
18         # an affine operation: y = Wx + b
19         self.fc1 = nn.Linear(16 * 5 * 5, 120) # 5*5 from image dimension
20         self.fc2 = nn.Linear(120, 84)
21         self.fc3 = nn.Linear(84, 10)
22
23     def forward(self, x):
24         x = self.conv1(x)
25         x = F.relu(x)
26         x = self.pool1(x)
27         x = self.conv2(x)
28         x = F.relu(x)
29         x = self.pool2(x)
30
31         x = torch.flatten(x, 1) # flatten all dimensions except the batch dimension
32
33         x = self.fc1(x)
34         x = F.relu(x)
35         x = self.fc2(x)
36         x = F.relu(x)
37         x = self.fc3(x)

```

code/LetNet.py

LeNet [5] was initially designed for low-resolution, black and white image recognition, specifically for digits. It demonstrated that CNNs could reliably perform both tasks of object localization and recognition (on low-resolution black and white images).

The last steps of the LeNet-5 architecture employ fully connected layers to convert features into final predictions. The scenario of MNIST data didn't matter, because of its small number of output classes. This *complicates the architecture when scaling up, influencing both computational resources and architecture*.

The Convolutional Neural Network Architecture ensures some degree of shift, scale and distortion invariance: *local receptive fields, shared weights* (or weight replication), and spatial or temporal *sub-sampling*.

Once a feature has been detected, its exact location becomes less important. Only its approximate position relative to other features is relevant.

## 1.1 shared weights

This algorithm is particularly useful for shared-weight networks because the weight sharing creates ill-conditioning of the error surface. Because of the sharing, one single parameter in the first few layers can have an enormous influence on the output. Consequently, the second derivative of the error with respect to this parameter may be very large, while it can be quite small for other parameters elsewhere in the network

## 1.2 sub-sampling

A simple way to reduce the precision with which the position of distinctive features are encoded in a feature map is to reduce the spatial resolution of the feature map. This can be achieved with a so-called sub-sampling layers which performs a local averaging and a sub-sampling, reducing the resolution of the feature map, and reducing the sensitivity of the output to shifts and distortions.

## 1.3 Loss

Maximum Likelihood Estimation Criterion (MLE), which is equivalent to the Mean Squared Error (MSE)

$$E(W) = \frac{1}{P} \sum_{p=1}^P y_{D^P}(Z^P, W)$$

Where  $y_{D^P}$  is the output of the  $D_p$ -th RBF unit, i.e. the one that corresponds to the correct class of the input pattern  $Z^p$ .

It lacks three important properties:

1. **Trivial Solution with Adaptation of RBF Parameters:** If the parameters of the Radial Basis Function (RBF) are allowed to adapt, the MLE criterion has a trivial and unacceptable solution. In this solution, all RBF parameter vectors become equal, leading to a constant and unchanging state of the network. The network effectively ignores the input, and all RBF outputs become equal to zero.
2. **Lack of Competition Between Classes:** The MLE criterion lacks competition between different classes in the training process. Introducing a more discriminative training criterion, referred to as the Maximum A Posteriori (MAP) criterion, could address this issue. The MAP criterion aims to maximize the posterior probability of the correct class or minimize the logarithm of the probability of the correct class given the input. Unlike MLE, MAP introduces a competitive element by pushing up the penalties of incorrect classes in addition to pushing down the penalty of the correct class.

3. **Potential Collapsing Phenomenon:** When RBF weights are allowed to adapt with the MLE criterion, there is a risk of a collapsing phenomenon where all RBF centers become equal. The discriminative criterion of MAP prevents this collapsing effect by keeping the RBF centers apart from each other. The discriminative criterion ensures that the posterior probabilities of different classes remain distinct, preventing the network from ignoring the input.

## 2 AlexNet

The AlexNet [4] was introduced in 2012 and introduced more layers for larger inputs and larger filters. Originally, the architecture was split into two columns/streams, but this was a workaround for the hardware limitations at the time.

It used the ImageNet dataset which contained color images with nature objects of larger size compared to MNIST ( $469 \times 387$  vs  $28 \times 28$ ) and 1.2 million images vs 60k images.

Its key improvements on the LeNet were

- Add a dropout layer after two hidden dense layers (better robustness /regularization)  
Dropout allowed for much deeper networks by introducing regularization not just at the input layer but throughout multiple layers of the network. This made it possible to control the complexity of the model more effectively.
- Change activation function from sigmoid to ReLu (no more vanishing gradient)  
enabling training of deeper networks more efficiently.
- MaxPooling instead of Average Pooling  
This made the learned features more shift-invariant, which is important for object recognition. Max pooling generally retains the most salient features and discards less useful information, making the model more robust.
- Heavy data Augmentation like cropping, shifting, and rotation.
- Model ensembling  
Ensemble methods is a machine learning technique that combines several base models in order to produce one optimal predictive model
- Using softmax function for classification
- Increase in kernel size and stride  
design choice to accommodate the higher resolution of images in the ImageNet dataset compared to MNIST.

AlexNet is substantially more complex, being about 250 times more computationally expensive. However, it is only ten times larger parametrically than LeNet-5. This way AlexNet is notorious for its high memory usage.

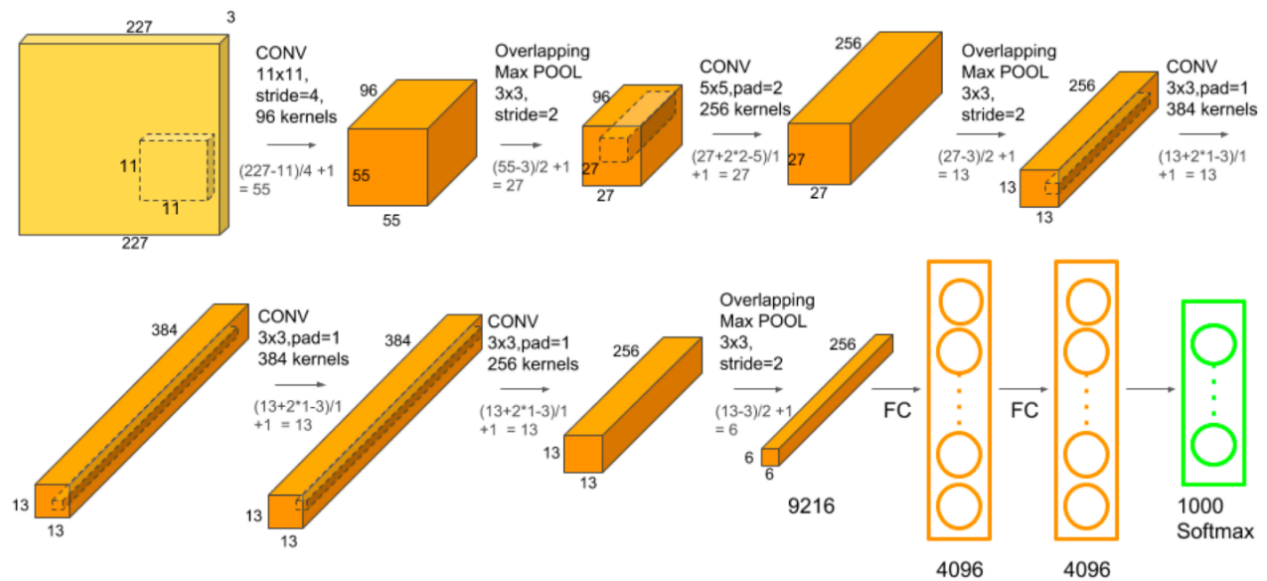


Figure 2: Architecture of AlexNet

```

1 import torch
2 import torch.nn as nn
3
4 class AlexNet(nn.Module):
5     def __init__(self, num_classes: int = 1000, dropout: float = 0.5) -> None:
6         super(AlexNet, self).__init__()
7         self.features = nn.Sequential(
8             nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
9             nn.ReLU(inplace=True),
10            nn.MaxPool2d(kernel_size=3, stride=2),
11            nn.Conv2d(64, 128, kernel_size=5, padding=2),
12            nn.ReLU(inplace=True),
13            nn.MaxPool2d(kernel_size=3, stride=2),
14            nn.Conv2d(128, 256, kernel_size=3, padding=1),
15            nn.ReLU(inplace=True),
16            nn.Conv2d(256, 256, kernel_size=3, padding=1),
17            nn.ReLU(inplace=True),
18            nn.MaxPool2d(kernel_size=3, stride=2),
19        )
20        self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
21        self.classifier = nn.Sequential(
22            nn.Dropout(p=dropout),
23            nn.Linear(256 * 6 * 6, 4096),
24            nn.ReLU(inplace=True),
25            nn.Dropout(p=dropout),
26            nn.Linear(4096, 4096),
27            nn.ReLU(inplace=True),
28            nn.Linear(4096, num_classes),
29        )
30
31    def forward(self, x: torch.Tensor) -> torch.Tensor:
32        x = self.features(x)
33        x = self.avgpool(x)
34        x = torch.flatten(x, 1)
35        x = self.classifier(x)
36        return x

```

code/AlexNet.py

### 3 VGG

The Visual Geometry Group [6] uses the ‘bigger means better’ philosophy. VGG introduced the notion of repeated blocks.

- **NOT Add more dense layers** - computationally expensive
- **NOT Add more convolutional layers** - as the network grows, specifying each convolutional layer individually becomes tedious
- **Group Layers into blocks** - these blocks can easily be parameterized, creating a more organized, modular architecture

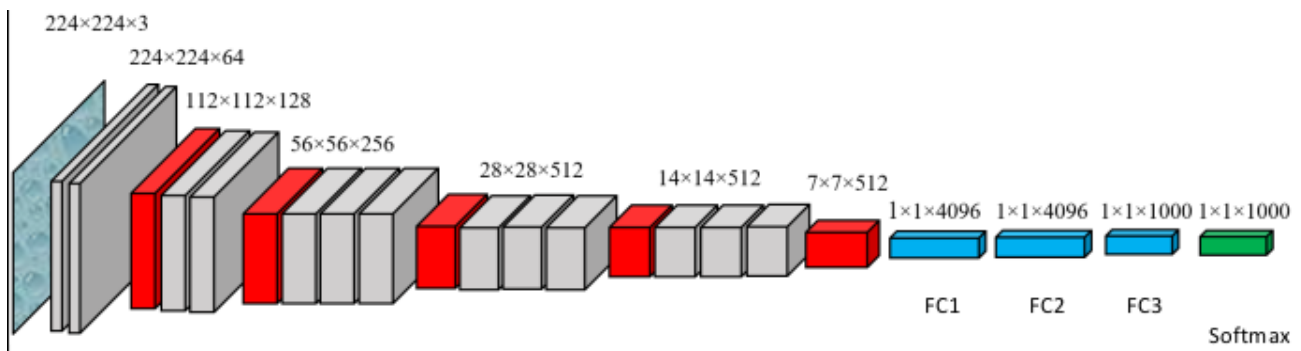


Figure 3: Architecture of VGG

```

1 from functools import partial
2 from typing import Any, cast, Dict, List, Optional, Union
3
4 import torch
5 import torch.nn as nn
6
7 cfgs: Dict[str, List[Union[str, int]]] = {
8     "A": [64, "M", 128, "M", 256, 256, "M", 512, 512, "M", 512, 512, "M"],
9     "B": [64, 64, "M", 128, 128, "M", 256, 256, "M", 512, 512, "M", 512, 512, "M"],
10    "D": [64, 64, "M", 128, 128, "M", 256, 256, 256, "M", 512, 512, 512, "M", 512,
11    ↪ 512, "M"],
12    "E": [64, 64, "M", 128, 128, "M", 256, 256, 256, 256, "M", 512, 512, 512, 512, "M",
13    ↪ 512, 512, 512, 512, "M"],
14 }
15
16 class VGG(nn.Module):
17     def __init__(
18         self, features: nn.Module, num_classes: int = 1000, init_weights: bool = True,
19         ↪ dropout: float = 0.5
20     ) -> None:
21         super(VGG, self).__init__()
22         self.features = features
23         self.avgpool = nn.AdaptiveAvgPool2d((7, 7))
24         self.classifier = nn.Sequential(
25             nn.Linear(512 * 7 * 7, 4096),
26             nn.ReLU(True),
27             nn.Dropout(p=dropout),
28             nn.Linear(4096, 4096),
29             nn.ReLU(True),
30             nn.Dropout(p=dropout),
31             nn.Linear(4096, num_classes),
32         )
33         if init_weights:
34             for m in self.modules():

```

```

32         if isinstance(m, nn.Conv2d):
33             nn.init.kaiming_normal_(m.weight, mode="fan_out", nonlinearity="relu"
    ↪ )
34             if m.bias is not None:
35                 nn.init.constant_(m.bias, 0)
36         elif isinstance(m, nn.BatchNorm2d):
37             nn.init.constant_(m.weight, 1)
38             nn.init.constant_(m.bias, 0)
39         elif isinstance(m, nn.Linear):
40             nn.init.normal_(m.weight, 0, 0.01)
41             nn.init.constant_(m.bias, 0)
42
43     def forward(self, x: torch.Tensor) -> torch.Tensor:
44         x = self.features(x)
45         x = self.avgpool(x)
46         x = torch.flatten(x, 1)
47         x = self.classifier(x)
48         return x
49
50
51 def make_layers(cfg: List[Union[str, int]], batch_norm: bool = False) -> nn.Sequential:
52     layers: List[nn.Module] = []
53     in_channels = 3
54     for v in cfg:
55         if v == "M":
56             layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
57         else:
58             v = cast(int, v)
59             conv2d = nn.Conv2d(in_channels, v, kernel_size=3, padding=1)
60             if batch_norm:
61                 layers += [conv2d, nn.BatchNorm2d(v), nn.ReLU(inplace=True)]
62             else:
63                 layers += [conv2d, nn.ReLU(inplace=True)]
64             in_channels = v
65     return nn.Sequential(*layers)
66
67
68 def _vgg(arch: str, cfg: str, batch_norm: bool, progress: bool, **kwargs: Any) -> VGG:
69     model = VGG(make_layers(cfgs[cfg], batch_norm=batch_norm, **kwargs))
70     return model
71
72 def vgg16(pretrained: bool = False, progress: bool = True, **kwargs: Any) -> VGG:
73     return _vgg('vgg-16', "D", False, pretrained, progress, **kwargs)

```

code/VGG.py

### 3.1 fewer wide convolutions or more narrow convolutions?

Recent comprehensive analysis from papers has shown that using more layers of narrow convolutions outperforms using fewer wide ones. This has been a general trend in network design: having more layers of simpler functions is generally more powerful than fewer layers of more complex functions.

### 3.2 The VGG block

Several  $3 \times 3$  convolutions, padded by one maintains the spatial dimensions from the input to the output layer, and at the end, max-pooling layer of  $2 \times 2$  and stride of 2 halves the resolution.

Combining these blocks with dense layers, creates an entire family of architectures just by varying the number of blocks.

### 3.3 Performance

VGG tends to be a lot slower when compared to the throughput of AlexNet, however, it makes up for in terms of accuracy. While VGG might require more computational resources, it generally provides superior performance.

🔍 28.41

## 4 Batch Norm

The problem deep networks face is with convergence, making training difficult and time consuming. There is a strategy to perform “deep supervision” which involves backpropogating from intermediate stages to help the network learn - however, this has its limitations [TODO].

### 4.1 Motivation

Very deep models involve the composition of several functions or layers. The gradient tells how to update each parameter, under the assumption that the other layers do not change. In practice, we update all of the layers simultaneously; gradients flow from the top layer down to the bottom layers of the network. As a result, the last layers start to adapt first, followed by the layers below them, creating a cascade of adaptations throughout the network.

However, this leads to a problem: As the bottom layers adapt, they change the features that are fed back up to the top layers. This means that the top layers, which had already started to adapt, have to readjust to these new inputs. Training Deep Neural Networks is complicated by the fact that the distribution of each layer’s inputs changes during training, as the parameters of the previous layers change [1]. **Each layer’s learning destabilizes the next, causing a slow convergence process that takes a long time for all layers to adapt properly.**

### 4.2 Solution

Batch normalization can be applied to any input or hidden layer in a network [1].

Batch Norm mitigates this issue by normalizing the features within each mini-batch, thereby stabilizing the training process and speeding up convergence. The idea is to make minor corrections to the layers by adjusting their mean and variance during training.

$$\begin{aligned}\mu_B &= \frac{1}{|B|} \sum_{i \in B} x_i & \sigma_B^2 &= \frac{1}{|B|} \sum_{i \in B} (x_i - \mu_B)^2 + \epsilon \\ x_{i+1} &= \gamma \frac{x_i - \mu_B}{\sigma_B} + \beta\end{aligned}\tag{1}$$

Where  $\gamma$  is variance and  $\beta$  is the mean.  $|B|$  represents the size of the mini-batch, and  $x_i$  are the individual data points in that mini-batch.  $\epsilon$  is a small positive value such as  $10^{-8}$  introduced to avoid the potential problem of undefined gradient.

We compute the mean and variance of a given mini-batch during training. Then, re-normalize each input feature by subtracting its mean and dividing it by its standard deviation. The model also learns two parameters to scale ( $\gamma$ ) and shift ( $\beta$ ) the normalized features.

We then adjust these first and second moments separately from the rest of the network learning, leading to Equation 1.



### 4.3 Advantages

This way, we are able to “normalize” each mini-batch separately, making the training of deep networks more stable and faster (dramatically reducing the number of training epochs required to train deep networks). The normalization can be undone by the network by learning appropriate  $\gamma$  and  $\beta$  parameters if it sees a benefit to it.

### 4.4 Updated Aim

The original motivation of batch normalization is reducing covariate shift. This is not correct. BatchNorm actually worsens covariate shift, yet it still aids in model convergence.

**BatchNorm is effectively acting as a form of regularization by introducing noise into the model.**

Here’s how it works: You calculate the mean and variance of a mini-batch, let’s say, of 64 observations. Since you’re working with a small sample, both the mean and the variance are subject to noise. You then normalize the features using these noisy statistics, introducing a random scale and shift into your model at each batch. This random noise acts as a regularizing factor, which is why you often don’t need additional regularization techniques like dropout when you’re using BatchNorm. However, this property makes BatchNorm sensitive to the size of the mini-batch.

If your mini-batch is too large, you’re not introducing enough noise for effective regularization. If it’s too small, the noise level becomes counterproductive, affecting convergence. This mini-batch size sensitivity becomes especially significant in multi-GPU settings, where batch sizes are often adjusted.

### 4.5 Application

If you’re working with a dense layer, a single normalization is applied to all the activations in that layer.

In the case of a convolutional layer, a separate normalization is performed for each channel.

### 4.6 Randomization during use

You don’t want this randomness when you’re using the model for predictions. So, you fix the gamma and beta parameters that the model has learned during training. Instead of using batch statistics, you use the running average for the mean and variance to normalize the features.

## References

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [2] Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: [1512.03385](https://arxiv.org/abs/1512.03385) [cs.CV]. URL: <https://arxiv.org/abs/1512.03385>.
- [3] Sergey Ioffe and Christian Szegedy. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. 2015. arXiv: [1502.03167](https://arxiv.org/abs/1502.03167) [cs.LG]. URL: <https://arxiv.org/abs/1502.03167>.
- [4] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet classification with deep convolutional neural networks”. In: *Commun. ACM* 60.6 (May 2017), pp. 84–90. ISSN: 0001-0782. DOI: [10.1145/3065386](https://doi.org/10.1145/3065386). URL: <https://papers.nips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a6b1d3d-Paper.pdf>.

- 
- [5] Y. Lecun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324. DOI: [10.1109/5.726791](https://doi.org/10.1109/5.726791). URL: [http://vision.stanford.edu/cs598\\_spring07/papers/Lecun98.pdf](http://vision.stanford.edu/cs598_spring07/papers/Lecun98.pdf).
  - [6] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. 2015. arXiv: [1409.1556](https://arxiv.org/abs/1409.1556) [cs.CV]. URL: <https://arxiv.org/pdf/1409.1556.pdf>.