

# Proseminar „Convolutional Neural Networks - Methoden und Anwendungen“

## Bilderkennung mit Vgg

Justin Schartner

18. Mai 2022

# Struktur

- 1** Einleitung
- 2 Allgemeines
- 3 Architektur
- 4 Training
- 5 Beispiel
- 6 Bewertung
- 7 Ausblick

# Thema

## Problemstellung

Gibt es CNN-Architekturen, welche Bilder noch besser, als schon bekannte Architekturen klassifizieren können?

## Lösung

Durch die Erhöhung der Tiefe eines CNN verspricht man sich genauere Aussagen über Bilder machen zu können.

## Ergebnis

Die VGG-Architektur hat bewiesen, dass die Tiefe eines CNN, eine ausschlaggebende Komponente hinsichtlich der Bilder-Klassifizierung ist.

# Gliederung

- Was ist VGG? **Allgemeines**
- Wie funktioniert VGG, was macht es besonders? **Architektur**
- Wie trainiert man ein VGG-net, was ist wichtig? **Training**
- Wie implementiert man ein VGG-net? **Beispiel**
- Wie gut ist VGG? **Bewertung**

# Struktur

- 1 Einleitung
- 2 Allgemeines**
- 3 Architektur
- 4 Training
- 5 Beispiel
- 6 Bewertung
- 7 Ausblick

# Eckdaten

- **Visual Geometry Group**
- Department of Engineering Science, University of Oxford
- Karen Simonyan und Andrew Zisserman
- **Veröffentlicht:** 4 Sep 2014
- **Letzte Änderung:** 10 Apr 2015

# Idee

## Steigerung der Genauigkeit durch:

- Steigerung der Tiefe, des CNNs
- Schachtelung von Convolutional-Layer-Blöcken
- Einsatz von kleinen 3x3-Filtern und einer Stride von 1

# Bilderkennung



# Struktur

- 1 Einleitung
- 2 Allgemeines
- 3 Architektur**
- 4 Training
- 5 Beispiel
- 6 Bewertung
- 7 Ausblick

# Architekturarten

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 x 224 RGB-image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

# Architektur

## Convolutional Layers

- receptive field:  $3 \times 3$ ,  $1 \times 1$
- activation: ReLU, stride: 1, padding: 1, channels: 64, 128, 512, 512

## Pooling Layer

- Max-Pool
- field:  $2 \times 2$ , stride: 2

## Fully-Connected Layers

- activation: ReLU
- channels: 4096, 4096, 1000

## Softmax Layer



# Convolutional Layers

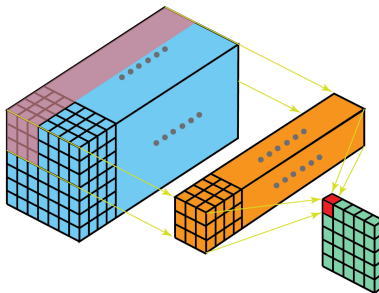


Abbildung: Quelle ?

(Breite x Höhe x Tiefe)  $\xrightarrow{\text{Conv2d}(\text{Filter: } 3 \times 3 \times \text{Tiefe}, \text{ Filter Anzahl: } n)}$  (Breite x Höhe x n)

# ReLU

Rectified Linerar Unit

$$f(x) = \begin{cases} x & x > 0 \\ 0 & \text{sonst} \end{cases}$$

$$f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

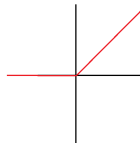


Abbildung:  $f(x)$

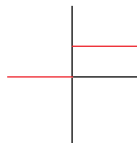


Abbildung:  $f'(x)$

# Max-Pooling

- Stride: 2x2

- Kernel: 2x2

⇒ Die Breite und Höhe wird halbiert

⇒ Daten werden auf die ausschlaggebenden Informationen reduziert

# Fully Connected Layers

- Input: (7x7x512)  
⇒ (7x7x512) Input-Neuronen
- Aktivierungsfunktion: ReLU
- Zwei versteckte Layer mit jeweils 4096 Neuronen
- ImageNet-Klassifizierung von 1000 Klassen  
⇒ 1000 Output-Neuronen

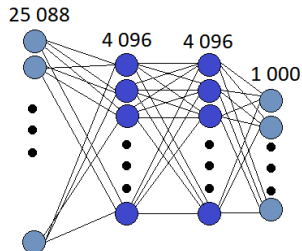


Abbildung: full-yconnected layer



# Softmax

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

- normalisierte Exponentialfunktion
- kategoriale Verteilung
- Transformation in den Wertebereich [0,1]

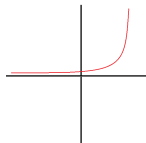


Abbildung:  $e^x$

$$\begin{pmatrix} -0.5 \\ 0.8 \\ 1.3 \end{pmatrix} \xrightarrow{\text{Softmax}} \begin{pmatrix} 0.093 \\ 0.342 \\ 0.564 \end{pmatrix}$$

# Struktur

- 1 Einleitung
- 2 Allgemeines
- 3 Architektur
- 4 Training**
- 5 Beispiel
- 6 Bewertung
- 7 Ausblick



# Training

## Details

- Learning-rate:  $10e-2$  -  $10e-4$
- Momentum: 0.9
- Weight-Decay:  $5e-4$
- Biases wurden mit 0 initialisiert
- VGG16 wurde mit VGG11-Gewichten initialisiert

## **Training** **Bild-Processing**

- Bilder wurden zufällig aus anderen Bildern ausgeschnitten
- Bilder wurden zufällig horizontal gedreht
- Bilder wurden zufällig skaliert

⇒ Filter werden trainiert Features auf verschieden Arten zu erkennen

# Struktur

- 1 Einleitung
- 2 Allgemeines
- 3 Architektur
- 4 Training
- 5 Beispiel**
- 6 Bewertung
- 7 Ausblick

# Vgg initialisieren

1. Variante: Die Komponenten einzeln erstellen und aneinander reihen.

Anpassbarkeit an das Problem

Das Trainieren des Netzwerkes nimmt mehr Zeit in Anspruch

2. Variante: Benutzen von schon bestehenden (und trainierten) Netzwerken.

Sind sofort einsatzbereit

Aufwand ist kleiner

Die vortrainierten Gewichte können die Trainingszeit minimieren

Konfigurierung kann aufwendig sein

# Vgg initialisieren

## 1. Variante

```

30
31
32 def create_fully_connected_layers(self, num_classes):
33     fully_connected_layers = nn.Sequential(
34         nn.Linear(512*7*7, 4096),
35         nn.ReLU(),
36         nn.Dropout(p=0.5),
37         nn.Linear(4096, 4096),
38         nn.ReLU(),
39         nn.Dropout(p=0.5),
40         nn.Linear(4096, num_classes),
41         nn.Softmax(dim=0)
42     )
43     return fully_connected_layers
44

```

```

def create_conv_layers(self, in_channels, architecture):
    layers = []

    # [[64, 64], [128, 128], [256, 256, 256], [512, 512, 512], [512, 512, 512]]
    for block in architecture:
        # [64, 64]
        for layer in block:
            out_channels = layer

            # conv-layer + relu
            layers += [
                nn.Conv2d(in_channels=in_channels, out_channels=out_channels,
                           kernel_size=(3,3), stride=(1,1), padding=(1,1)),
                nn.ReLU()
            ]

            in_channels = out_channels

        # max-pooling
        layers += [nn.MaxPool2d(kernel_size=(2,2), stride=(2,2))]
    return nn.Sequential(*layers)

```



# Vgg initialisieren

## 1. Variante

```

2 import torch.nn as nn
3
4 vgg16 = [[64, 64], [128, 128], [256, 256, 256], [512, 512, 512], [512, 512, 512]]
5 vgg19 = [[64, 64], [128, 128], [256, 256, 256, 256], [512, 512, 512, 512], [512, 512, 512, 512]]
6
7 class VGG_net(nn.Module):
8     def __init__(self, in_channels=3, num_classes=1000, architecture=vgg16):
9         super(VGG_net, self).__init__()
10
11         #Convolutional Layers
12         self.conv_layers = self.create_conv_layers(in_channels, architecture)
13
14         #Fully Connected Layers
15         self.fully_connected_layers = self.create_fully_connected_layers(num_classes)
16
17     def forward(self, x):
18         # x = [2, 3, 224, 224]
19
20         x = self.conv_layers(x)
21         # x = [2, 512, 7, 7]
22
23         x = x.reshape(x.shape[0], -1)
24         # x = [[1 .. (512x7x7)], [1 .. (512x7x7)]]
25
26         x = self.fully_connected_layers(x)
27         # x = [[1 .. 1000], [1 .. 1000]]
28
29         return x

```

# Vgg initialisieren

## 1. Variante

```

VGG_net(
(conv_layers): Sequential(
  (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (1): ReLU()
  (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (3): ReLU()
  (4): MaxPool2d(kernel_size=(2, 2), stride=(2, 2), padding=0, dilation=1, ceil_mode=False)
  (5): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (6): ReLU()
  (7): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (8): ReLU()
  (9): MaxPool2d(kernel_size=(2, 2), stride=(2, 2), padding=0, dilation=1, ceil_mode=False)
  (10): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (11): ReLU()
  (12): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (13): ReLU()
  (14): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (15): ReLU()
  (16): MaxPool2d(kernel_size=(2, 2), stride=(2, 2), padding=0, dilation=1, ceil_mode=False)
  (17): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (18): ReLU()
  (19): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (20): ReLU()
  (21): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (22): ReLU()
  (23): MaxPool2d(kernel_size=(2, 2), stride=(2, 2), padding=0, dilation=1, ceil_mode=False)
  (24): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (25): ReLU()
  (26): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (27): ReLU()
  (28): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
  (29): ReLU()
  (30): MaxPool2d(kernel_size=(2, 2), stride=(2, 2), padding=0, dilation=1, ceil_mode=False)
)
(fully_connected_layers): Sequential(
  (0): Linear(in_features=25088, out_features=4096, bias=True)
  (1): ReLU()
  (2): Dropout(p=0.5, inplace=False)
  (3): Linear(in_features=4096, out_features=4096, bias=True)
  (4): ReLU()
  (5): Dropout(p=0.5, inplace=False)
  (6): Linear(in_features=4096, out_features=1000, bias=True)
  (7): Softmax(dim=0)
)

```

```

56
57 if __name__ == "__main__":
58     net = VGG_net(in_channels=3, num_classes=1000, architecture='vgg16')
59     print(net)
60

```

# Vgg initialisieren

## 2. Variante

```

VGG(
  (features): Sequential(
    (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (1): ReLU(inplace=True)
    (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (3): ReLU(inplace=True)
    (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (5): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (6): ReLU(inplace=True)
    (7): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (8): ReLU(inplace=True)
    (9): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (10): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (11): ReLU(inplace=True)
    (12): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (13): ReLU(inplace=True)
    (14): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (15): ReLU(inplace=True)
    (16): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (17): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (18): ReLU(inplace=True)
    (19): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (20): ReLU(inplace=True)
    (21): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (22): ReLU(inplace=True)
    (23): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
    (24): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (25): ReLU(inplace=True)
    (26): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (27): ReLU(inplace=True)
    (28): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (29): ReLU(inplace=True)
    (30): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (avgpool): AdaptiveAvgPool2d(output_size=(7, 7))
  (classifier): Sequential(
    (0): Linear(in_features=25088, out_features=4096, bias=True)
    (1): ReLU(inplace=True)
    (2): Dropout(p=0.5, inplace=False)
    (3): Linear(in_features=4096, out_features=4096, bias=True)
    (4): ReLU(inplace=True)
    (5): Dropout(p=0.5, inplace=False)
    (6): Linear(in_features=4096, out_features=1000, bias=True)
  )
)

```

```

import torch

if __name__ == "__main__":
    net = torch.hub.load('pytorch/vision:v0.10.0', 'vgg16', pretrained=True)
    print(net)

```

# Vgg trainieren

```

4 def train_network(net, x_train, y_train, epochs):
5
6     #Training-Parameter
7     learning_rate = 0.001
8     momentum = 0.9
9
10    #Loss function
11    criterion = nn.CrossEntropyLoss()
12    #Optimizer
13    optimizer = optim.SGD(net.parameters(),
14                           lr=learning_rate, momentum=momentum)
15
16    for epoch in range(epochs):
17
18        #reset the gradients in net
19        optimizer.zero_grad()
20
21        #forward
22        outputs = net(x_train)
23
24        #calculate loss
25        loss = criterion(outputs, y_train)
26        #calculate gradients in net
27        loss.backward()
28
29        #backward
30        optimizer.step()
31
32

```

# Vgg benutzen

pug 0.9928255081176758

```

1 #get network
2 net = net = torch.hub.load('pytorch/vision:v0.10.0', 'vgg16', pretrained=True)
3
4 #get image
5 image = Image.open('dog1.jpg')
6
7 #get categories
8 categories = loadCategories('imagenet_classes.txt')
9
10
11 #preprocess image
12 input_image = pre_process(image)
13
14 #forward
15 with torch.no_grad():
16     output = net(input_image)
17
18 #apply softmax
19 probs = torch.nn.functional.softmax(output[0], dim=0)
20
21 #result
22 guess, guess_prob = findHighestProb(probs, categories)
23 print(guess, guess_prob)
24

```

# Struktur

- 1 Einleitung
- 2 Allgemeines
- 3 Architektur
- 4 Training
- 5 Beispiel
- 6 Bewertung**
- 7 Ausblick

## Ergebnis

- Die Anwendung von mehreren 3x3 Filtern ersetzt die Funktionalität von bsp. 7x7 Filtern und erhöht die Diskriminativität
- Die Tiefe eines CNNs ist ausschlaggebend für die Genauigkeit
- VGG-Architekturen haben viele Anwendungsbereiche, erzielen auf verschiedensten Datenbanken erfolgreiche Resultate

# Komplexität

- Anzahl an Parametern: 15.1, 15.3, 20.6, 25.9 in Millionen  
Mehr Parameter führen zu einer längeren Trainingszeit
- Entspricht etwa 528MB  
Hoher Speicher-Verbrauch



# Performance



# Struktur

- 1 Einleitung
- 2 Allgemeines
- 3 Architektur
- 4 Training
- 5 Beispiel
- 6 Bewertung
- 7 Ausblick**

# Ausblick