

Individualisierte Gesichtsdetektion mit YOLOv2 in Darknet



Img Src: p/reddie.com/darknet & p/reddie.com/darknet/yolo
Darknet ist ein Machine Learning Framework wie Caffe, Tensorflow oder Theano. Darknet ist in C mit Cuda geschrieben.

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

YOLO bedeutet in diesem fall You Only Look Once. YOLO ist ein Singleshot Object Detector. Was das alles bedeutet schauen wir uns als erstes an, zusammen mit den alternativen zu YOLO. Da YOLO der älteste der Single Short Object Detektoren ist, schauen wir uns die funktionsweise davon etwas genauer an. Die meisten anderen funktionieren auf ähnliche weise.

In YOLOv2 wurde eine Hierarchie eingeführt um einfach mehr Klassen für die Detektion zu haben. Um schlussendlich das ganze etwas greifbarer zu machen, gibt es noch ein high-level Beispiel für eine Gesichtsdetektion.

Klassifikation VS Detektion

- **Klassifikation**
 - Erraten des dargestellten Objektes
- **Detektion**
 - Erraten wo Objekte im Bild sind und Klassifikation dieser

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

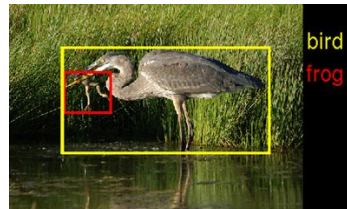
Gesichts-
detektion

Fazit

Demonstrator
Erklärung



Img Src: Wikimedia Commons - Sören T Eriksson



Img Src: Low-Power Image Recognition Challenge (LPIRC 2018)

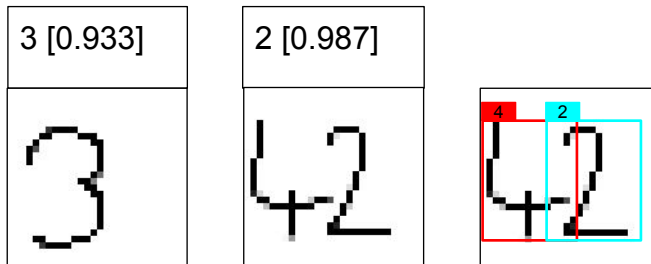
Als erstes: Was ist überhaupt der Unterschied zwischen Klassifikation und Detektion. Klassifikation beschäftigt sich damit eine oder mehrere Bezeichnungen dem Bild zuzuordnen. Heutzutage funktionieren Neuronale Netzwerke sehr gut bei Klassifikationsaufgaben, in denen nur ein Objekt gezeigt wird. Wenn jedoch mehrere bekannte Klassen in einem Bild vorkommen sind die resultate durchzogen. Detektion ist die herausforderung das Objekt auch im Bild zu lokalisieren. Einmal gefunden, müssen die objekte auch noch richtig klassifiziert werden.

Beispiel Klassifikation vs Detektion

■ MNIST Ziffer Klassifikation

- Einzelne Ziffer → klar
- Zwei Ziffern → ?

■ Object Detection evaluiert Regionen



Generiert mit <https://tensorflow-mnist.herokuapp.com/>

Klassifikation Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

Um an den vorhergehenden Sprecher, Dr Bernd Kosch, anzuknüpfen, eine Veranschaulichung mit Zahlen wie vom MNIST Datensatz.

Die Klassifikation einer einzelnen Zahl funktioniert sehr gut.

Jedoch wenn wir 2 unterschiedliche Zahlen haben, was ist die Antwort, wenn wir nur eine Klasse dem Bild zuweisen dürfen?

Mit Detektion schauen wir uns jede mögliche Detektion einzeln an. Das heisst, wir erhalten die Zahlen einzeln als Output.

2-Stage Detektor

- Beispiele
 - R-CNN, Fast R-CNN, Faster R-CNN, Mask R-CNN [Facebook]
 - R-FCN [Microsoft Research]
- Gute Präzision, langsam

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

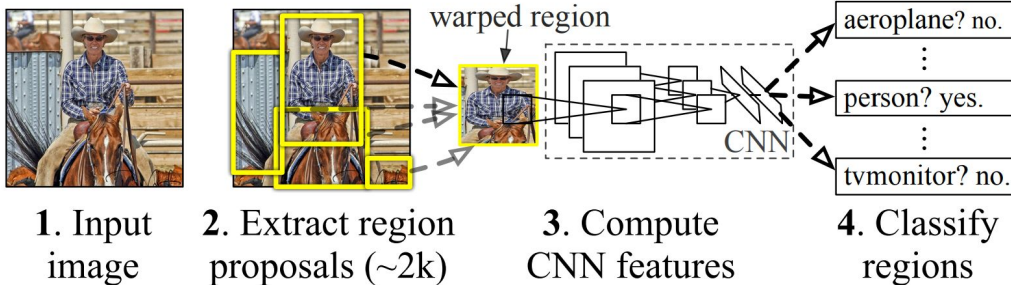
Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

R-CNN: *Regions with CNN features*



Img Src: Girshick, Ross, et al. "Rich feature hierarchies for accurate object detection and semantic segmentation" 2014

4

2-Stage Detektor besteht aus 2 Schritten:

- Algorithmus zum Regionen finden
- Convolutional Neural Network zum Klassifizieren der Region

Der algorithmus um die Region zu finden kann auch ein Neuronales Netz sein.

Für jede möglich Detektion wird die anschliessende Klassifikation einzeln aufgerufen.
Dementsprechend ist auch die geschwindigkeit. Die genauigkeit der Klassifikation ist jedoch gut.

Um kurz auf die Beispiele zu kommen. R-CNN ist der älteste dieser und hat seither viele Verbesserungen genossen.

Die neuste Version, Mask R-CNN ist von Facebook AI Research.

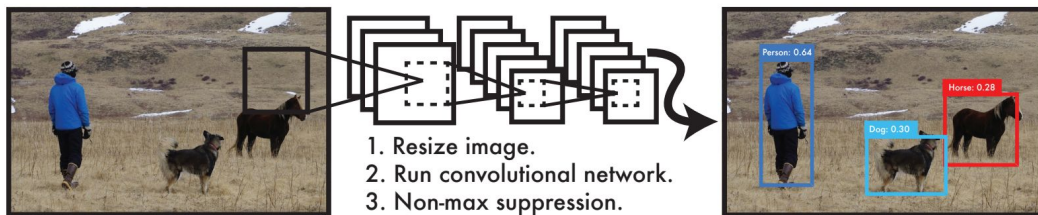
R-FCN hingegen ist von Microsoft Research unterstützt.

Single Shot Detectors

■ Beispiele

- YOLO, YOLOv2, YOLOv3 [University of Washington]
- SSD [Google]
- DSSD [Amazon]
- RetinaNet [Facebook]

■ Schneller als 2-Schritt detektoren, schlechtere Resultate



Img Src: Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." 2016

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

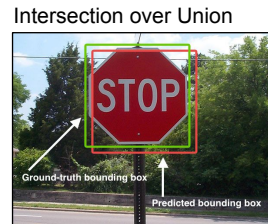
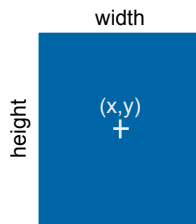
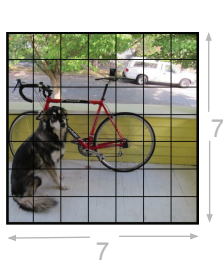
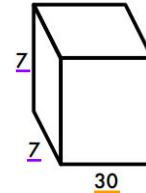
Wie wir gesehen haben, muss das ganze Netz für die Klassifikation mehrfach laufen. Das führt dazu, dass die Zeit um ein Bild zu verarbeiten sehr hoch ist. Single Shot Detectors hingegen, machen die Detektion in einem Schritt. Wie das im Detail aussieht, werden wir uns anschliessend anhand von YOLO anschauen.

Um uns kurz die prominentesten Vertreter zu gemüte führen. Wir sehen, die Tech-Giganten sind gut vertreten. YOLO ist dabei das älteste, und entsprechend das am einfachsten zu verstehende.

"The highest accuracy object detectors to date are based on a two-stage approach popularized by R-CNN, where a classifier is applied to a sparse set of candidate object locations. In contrast, one-stage detectors that are applied over a regular, dense sampling of possible object locations have the potential to be faster and simpler, but have trailed the accuracy of two-stage detectors thus far" - Focal Loss for Dense Object Detection <https://arxiv.org/pdf/1708.02002.pdf>

YOLO - Konzept

- 7 x 7 Grid*
- Pro Zelle 2 Bounding Boxes
 - Positionierung x, y, width, height
 - Wahrscheinlichkeit der Korrektheit der Bounding Box IoU
- 20 Klassen (Pascal VOC)
- $7 \times 7 \times (2 \cdot 5 + 20)**$



src: <http://www.pyimagesearch.com>

* YOLOv2 verwendet die Input größe des Bildes / 32 als Grid
 ** YOLOv2 verwendet Anchor Boxes mit eigenen Klassen -> $7 \times 7 \times 2 \times (5 + 20)$

- Klassifikation
- Detektion
- Neuronale Netze
- YOLO in Detail**
- Hierarchie
- Gesichtsdetektion
- Fazit
- Demonstrator Erklärung

Schauen wir uns diese Box etwas genauer an. Wir sehen die Dimensionen 7x7x30. Das 7x7 ist einfach erklärt. Konzeptuell ist das Bild in ein 7x7 Raster aufgeteilt. Die tiefe von 30 besteht aus 2 Teilen. Bounding boxes und die Klassen, so wie bei einem regulären Netz zur Klassifikation.

Der erste sind 2 bounding boxes. Für eine Bounding box braucht es eine position, an welcher zu Zeichnen begonnen wird. Dazu noch die höhe und breite. Etwas mehr komplexer ist der 5. Wert. Das ist eine Schätzung der Intersection over Union, oder IoU.

Dieser Wert besagt wie gut eine Detektion mit der Realen Position des Objektes im Bild übereinstimmt.

Indem wir die KI diesen Wert Schätzen lassen, erhalten wir einen wert in welchem ausgedrückt ist, wie sicher sich die KI der Schätzung ist.

Noch zu der Anzahl Klassen. In dieser erklärung wird das Pascal VOC: Datenset mit 20 Klassen verwendet.

Entsprechend wenn man das coco datenset mit 80 klassen verwenden würde, hätte man einen 7x7x90 Quader.

YOLOv1 und v2 sind nur wenig unterschiedlich.

Für mehr Details: Andrew Ng "YOLO Algorithm" und Andrew Ng "Anchor Boxes"

YOLO - Cost Function (Mathe Teil)

$$\sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

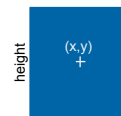
Klassifikation

Klassifikation
Detektion

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i_j}^{\text{obj}} (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2$$

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i_j}^{\text{obj}} (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2$$

Bounding Box



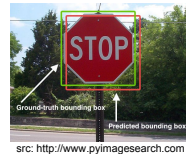
Neuronale
Netze

YOLO in Detail

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i_j}^{\text{obj}} (C_i - \hat{C}_i)^2$$

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{i_j}^{\text{noobj}} (C_i - \hat{C}_i)^2$$

Detektion



Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

Schauen wir kurz in den Mathematischen teil, benötigt für das Training. Hier sehen wir den unterschied zwischen Klassifikation, in Grün, und Detektion, Grün und Blau zusammen, sehr deutlich

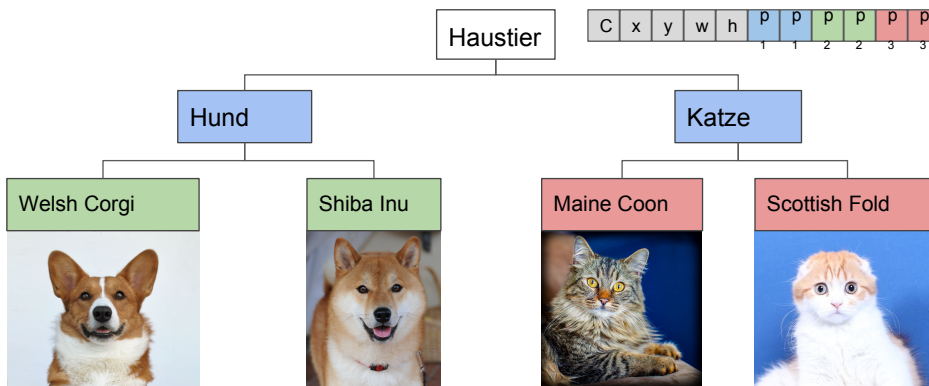
Klassifikation mit squared error

Lambdas sind tuning parameter

$\mathbb{1}^{\text{obj}}$ sind toggles, damit nur falsche Resultate einfluss haben, bei denen überhaupt ein resultat möglich war

Wordtree Hierarchie

- Finden von visuellen Gemeinsamkeiten
- Falls nicht genau Klassifizierbar, wird die Nächste überklasse verwendet
 - Best Effort



Img Src: Wikimedia Commons

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

Die Hierarchie ist die gleiche wie in WordTree. Der das Wurzelement ist der allgemeinste begriff, die Blätter des Baumes sind die genauesten verfügbaren begriffe. Für eine Klassifikation bedeutet das, dass die visuellen gemeinsamkeiten der kinder erlernt werden für einzelne Knoten. Das erlaubt auch eine Teilklassifikation von Objekten, welche sonst unbekannt wären. Somit würde eine andere Hunderasse wie der Valhund nicht als eine beliebige klasse geschätzt werden, sondern als Hund.

Technisch ist das mit einem Softmax per Unterbegriffsgruppe (Co-Hyponyms) gelöst. Das heisst für diese 4 Klassen haben wir 2 zusätzliche zwischen klassen und berechnen das ganze in 3 softmax regionen.

Das heisst die summe für alle wahrscheinlichkeiten einer Farbe ist 1.

Wenn wir einen Schwellenwert nicht erreichen, bleiben wir beim Auswerten einfach bei dieser Klasse stehen.

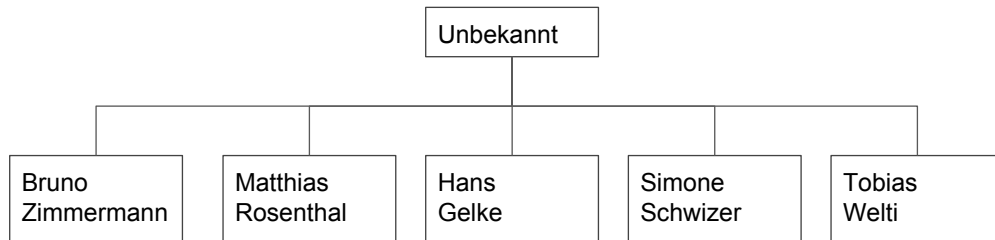
Nehmen wir eine gewöhnliche Hauskatze. Diese hat alle alle merkmale welche sie zu einer Katze machen.

Jedoch kann nicht genau gesagt werden, welche Rasse sie genau ist. Das heisst keine der Katzenrassen ist höher als ein Threshold. Beispielsweise 75%.

YOLO anhand eines Beispiels

■ Ziel

- Gesichtsdetektion
- Klassifikation der Person



Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Um das ganze etwas Greifbarer zu machen ein Beispiel mit meinen Arbeitskollegen. Meine Kollegen möchte ich erkennen, falls es eine andere Person ist, möchte ich diese als "Unbekannt" markieren. Um nochmals auf YOLO zu sprechen zu kommen. Die Implementation dieser Baumstruktur und der Softmax Regionen ist meines Wissens nach bisher nur im Framework "Darknet", für das YOLO ursprünglich geschrieben wurde.

Tree Structure

Tree File		Name File
Designator	Parent Line	
face	-1 ← Root Element	Unbekannt
rosn	0	Matthias_Rosenthal
gelk	0	Hans_Gelke
zing	0	Raphael_Zingg
zimb	0	Bruno_Zimmermann
welo	0	Tobias_Welti
weii	0	Armin_Weiss
sczr	0	Simone_Schwizer
moes	0	Dominik_Moesch
mazl	0	Amin_Mazlounian
kamm	0	Tobias_Kammacher
hubp	0	Philipp_Huber
gana	0	Rico_Ganahl
groo	0	Alexey_Gromov

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

**Gesichts-
detektion**

Fazit

Demonstrator
Erklärung

Diese Struktur muss codiert werden. Grundlegend dafür ist das Tree File.
Die Zeile auf welcher der eintrag steht ist massgebend für die Struktur.
Wir haben ein Label und eine angabe auf welcher Zeile das Elternelement liegt.
Um das ganze etwas schöner anzuzeigen, können wir noch ein Name File erstellen.

Darknet YOLO Components

- .data File für referenz zu Trainingsdaten und Beschriftungen
- .cfg Datei mit der strukturdefinition des Netzes
- .weights Datei mit den erlernten Gewichten
- .tree mit einer Repräsentation der Hierarchie
 - Nur Relevant wenn Hierarchie verwendet wird. Verlinkt über eine Zeile in der .cfg Datei

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

**Gesichts-
detektion**

Fazit

Demonstrator
Erklärung

Darknet unterscheidet sich stark von Tensorflow.

Es hat jedoch einige gemeinsamkeiten mit Caffe.

Wir haben ein .data file um die Trainingsdaten zuzuführen und ein verweis auf die Beschriftung zu haben.

Die .cfg Datei entspricht dem Caffe prototxt. Das ist eine Beschreibung des Models.

Die .weights ist das element das wir im trainings prozess effektiv ändern.

Zu guter letzt brauchen wir noch das .tree file um die Hierarchie bekannt zu geben.

■ Face Detection Data Set and Benchmark (FDDB)

- Bilder von Gesichtern mit Koordinaten



■ Bilder von Arbeitskollegen nach Person sortiert (Klassifikations Daten)



Beispiel für Dateiname
00001_zimb.jpg

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

Für das Beispiel habe wir 2 Datensets.

Das erste ist das FDDB (face detection data set and benchmark).

Dieses ist eine Kollektion von Personenbildern.

Zu jedem Bild gibt es Informationen wo genau die Gesichter in dem Bild sind.

Für uns sind das in der Klassifikation alle "Unbekannt"

Das zweite Datenset ist ein selbsterstelltes.

Es besteht aus den Gesichtern meiner Arbeitskollegen mit den Namen im Dateiname.

Hier sieht man die Schwierigkeit beim Erstellen der Trainingsdaten.

Ein paar Fotoserien machen und in die entsprechenden Ordner kopieren ist schnell gemacht.

Auf jedem Bild das Gesicht finden und genau markieren? Das dauert.

Training

1. Klassifikation trainieren
2. Letzten Layer entfernen - Nur erlernte Filter relevant
3. Detektionstraining mit Fddb Datenset
4. Detektion über Klassifikationsdaten
5. Training mit kombiniertem Datenset

Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

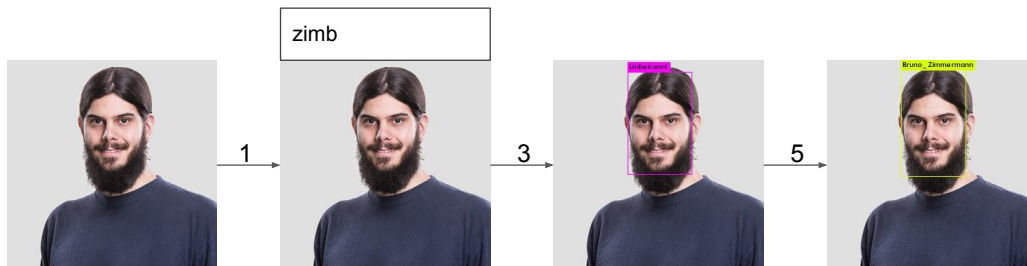
Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

14



Das trainieren des Gesicht detektors lässt sich in 5 Schritte einteilen.

Als erstes wird die Klassifikation trainiert. Dies, weil es wesentlich schneller geht für das Netz zum relevante Filter zu finden in diesem training.

Wenn die qualität ausreichend ist, kann der letzte layer entfernt werden.

Wir wollen ja nicht die Linie für die Klassifikation, sondern dieser Quader von YOLO für die Detektion.

Als nächsten Punkt trainieren wir mit dem Fddb Daten. Nach dem Training kann unsere KI Gesichter finden.

Das Gedächtnis, wie diese zuzuweisen sind haben wir jedoch im 2. Schritt entfernt.

Nun haben wir jedoch Daten, von welchen wir die Klassen kennen und ein Algorithmus um diese Objekte zu finden.

Das lässt sich gut kombinieren um einen Datensatz zu erstellen mit Position und Klasse.

Geht etwas schneller als alles von hand zu markieren.

Nun kann das Netz mit den kombinierten Daten weiter trainiert werden, um die Klassifikation in die einzelnen Personen wieder zu haben.

Fazit Hierarchie

- Fully Convolutional
 - Effizient auf neuromorpher Hardware wie z.B. Movidius
- Out of the Box ist das Resultat teilweise etwas durchgezogen
 - Interessant für Datengenerierung - Viele Daten mit wenig Aufwand
- Statisch → Für neue Klassen muss neu Trainiert werden!

Klassifikation
Detektion

Neuronale
Netze

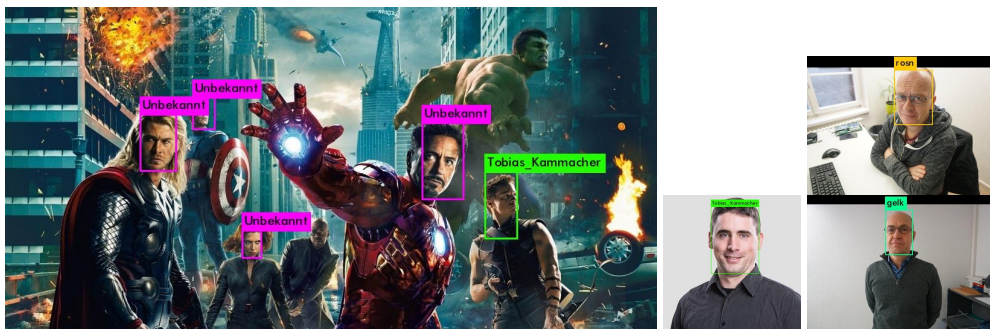
YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung



15

Die Resultate sind nicht wirklich Perfekt. Wir wir bei dem Filmplakat zu sehen, haben wir bereits eine falsch Detektion.

Das liegt an dem Threshold und dem versuch möglichst eine genaue Klasse zuzuweisen.

Das führt mich zur aussage: Es funktioniert recht gut für Detektionen, jedoch sind die Klassifikationen nicht immer gut.

Generell, wenn die Verarbeitungszeit kein Faktor ist, sind 2-Step-Detektoren wie R-CNN besser bisher.

Was von besonderem Interesse ist, ist die möglichkeit einfach Detektionsdaten zu generieren.

Dabei ist die genaue Klassifikation nicht so wichtig, da wir ja die Klasse schon kennen.

Was jedoch alle diese Methoden gemeinsam haben, ist die statische größe für die Klassen. Um mehr Klassen zu erkennen, müsste man diese trainiert haben.

Das funktioniert nicht, für einen Demonstrator wie wir ihn ausgestellt haben. Schauen wir uns das kurz genauer an

Dynamischer Ansatz (Demonstrator)

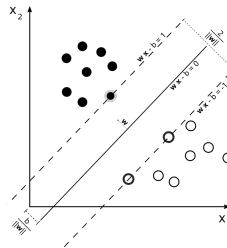
■ Detektions Algorithmus (zB YOLO mit nur Gesicht trainiert)



■ Neuronales Netz zur Klassifikation (VGG16)

- Klassifikation Resultate egal, nur erlernte filter relevant
- 128 Werte, welche ein Gesicht beschreiben

■ Support Vector Machine (SVM)



Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Demonstrator
Erklärung

Wie vielleicht einige von Ihnen gesehen haben, haben wir ein Demonstrator, bei dem neue Gesichter während dem laufen eingescannt werden können. Das widerspricht den letzten paar slides, auf denen wir uns ein statisches Klassifikationsbeispiel angesehen haben.

Wie macht man etwas mehr dynamisches? Wir verwenden eine Pipeline. Ähnlich wie bei den 2-Stage Detektoren.

Wir haben ein Algorithmus um die Gesichter zu finden, also zum Beispiel YOLO ohne die Klassifikation.

Als nächster Schritt schneiden wir das Gesicht aus und füttern es in ein Convolutional Neural Network, von der Struktur her kann es zum Beispiel ein VGG16 sein.

Dabei interessiert uns die Klassifikationswerte selbst nicht, sondern die Werte der Filter. Diese 128 Werte werden dann als Vektor interpretiert verwendet um eine Support Vector Machine zu trainieren.

Diese SVM werden dann bei einer Detektion durchgecheckt. Schlägt es bei keinem Gesicht an, ist es Unbekannt.

Kontakt

Bruno Zimmermann

bruno.zimmermann@zhaw.ch

www.zhaw.ch/ines

<https://blog.zhaw.ch/high-performance/>



Klassifikation
Detektion

Neuronale
Netze

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

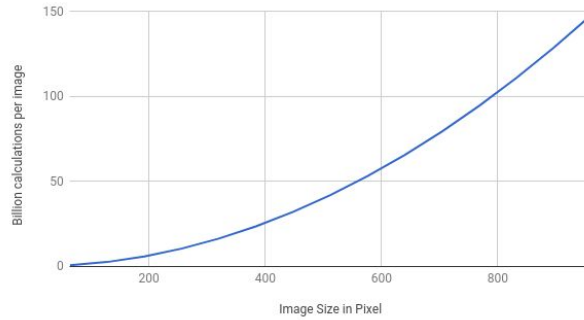
Demonstrator
Erklärung

Appendix

Computation Costs YOLOv2 (with Darknet19)

Factor *32	Size	Billion Calculations
2	64	0.656678912
4	128	2.626715648
6	192	5.910110208
8	256	10.50686259
10	320	16.4169728
12	384	23.64044083
14	448	32.17726669
16	512	42.02745037
18	576	53.19099187
20	640	65.6678912
22	704	79.45814835
24	768	94.56176333
26	832	110.9787361
28	896	128.7090668
30	960	147.7527552

YOLO Architecture Calculation Scaling



Neuronale Netze

Klassifikation Detektion

YOLO in Detail

Hierarchie

Gesichts-detektion

Fazit

Appendix

Configuration File Details

```
[convolutional]
filters=100 = 5*num + 5*classes
size=1
stride=1
pad=1
activation=linear

[region]
anchors = 0.77871, 1.14074, 3.00525, 4.31277, 9.22725, 9.61974
bias_match=1
classes=15
coords=4
num=5
softmax=1
jitter=.2
rescore=1
```

Neuronale Netze

Klassifikation
Detektion

YOLO in Detail

Hierarchie

Gesichts-
detektion

Fazit

Appendix

Mehr Tree Struktur Beispiele

Line Designator Parent Line
Number

```
00 saeugetiere -1
01 canidae 0
02 felidae 0
03 hund 1
04 wolf 1
05 fuchs 1
06 luchs 2
07 hauskatze 2
08 panthera 2
09 corgi 3
10 vallhund 3
11 wildkatze 7
12 sandkatze 7
13 loewe 8
```

Neuronale Netze

Klassifikation Detektion

YOLO in Detail

Hierarchie

Gesichts- detektion

Fazit

Appendix