Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
References

# Machine Learning

## In Robotics

Vinh Ngu

Universität Hamburg
Fakultät für Mathematik, Informatik und Naturwissenschaften
Fachbereich Informatik, Arbeitsbereich TAMS
Masterprojekt SoSe 17

1. Juni 2017

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
References

Agenda

## Agenda

- Introduction to Machine Learning
- Object Detection
- Convolutional Neural Networks
- Grasping with the help of CNN's

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

**What is machine learning?**
How does it work?

# What is machine learning?

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

**What is machine learning?**
How does it work?

# What is machine learning?

- 1950's
- Arthur Samuel as pioneer
- World's first self-learning program - "checkers"
- Detects hidden patterns
- "Cognitive"functions that humans associate with other human minds
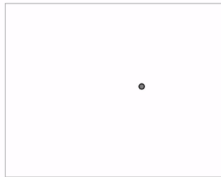- **Learning** and **Problem**-solving

# How do machines learn?

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

**Training Data**

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

Training Data

Feature Extraction

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

Training Data

Feature Extraction

Machine Learning Model

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

Introduction
**Machine Learning**
Object Detection
Convolutional Neural Network
Grasping
References

What is machine learning?
**How does it work?**

# Object Detection

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

**Selective Search**
Learning vs Rules

## Selective Search

- Published by University of Trento, Italy and University of Amsterdam, the Netherlands, 2012 [1]

- Combines the strength of both an exhaustive search and segmentation.

- Uses use bag-of-words for object recognition.

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

Selective Search
Learning vs Rules

# What we want is...

## Object Recognition

**Goal:**



## Problem
Where to look at?

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

Selective Search
Learning vs Rules

## One solution

### Idea
Let's check everywhere in the image
for a possible object!
(Exhaustive-Search in combination
for instance "Lampert"[3])

### Problem
Extremely slow, must process tens of
thousands of candidate objects.

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

Selective Search
Learning vs Rules

# One solution

### Idea
Let's use an object detector first!

### Problem
What about oddly-shaped objects?
Will we need to scan with windows
of many different shapes?

Not objects

Might be objects

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

**Selective Search**
Learning vs Rules

## Final solution

### Idea

Let's perform a segmentation first, then run the object recognition.
Using the segmentations as candidate for possible objects.



### Advantages

Can be efficient, makes no assumptions about object sizes or
shapes.

Introduction
Machine Learning
**Object Detection**
Convolutional Neural Network
Grasping
References

Selective Search
**Learning vs Rules**

## 2012

- 2012: Alex Krizhevsky has won the "ILSVRC" (ImageNet Large-Scale Visual Recognition Challenge)
- 2012 first year where a CNN was used to achieve a top 5 test error rate of 15.4%. (2nd 26.2%)
- CNN's grew prominence.

## Learning vs Rules

- At the beginning classification uses predefined rules
- The definition of rules becomes impossible by complex images
- Artificial intelligence are used to extract the most relevant characteristics
- Still, modern systems do not learn directly from pixel level

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
R-CNN workflow
R-CNN details

# Convolutional Neural Network

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
R-CNN workflow
R-CNN details

## What are CNN's?

- A feed-forward artificial neural network

- Inspired by the organization of the animal visual cortex

- Grew prominence in 2012. (Alex Krizhevsky, Classification error: 26% –> 15%.

- Mainly used for image/video recognition and natural language processing.

- Facebook, Google, Amazon.

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
**Region-based CNN**
R-CNN workflow
R-CNN details

# Region Based Convolutional Neural Networks (R-CNN)

### R-CNN's

- Published by the University of California - 2014 [3]
- Combination of CNN and a domain-specific fine-tuning-method.

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
**Region-based CNN**
R-CNN workflow
R-CNN details

## Intresting facts about R-CNN's

| Method | Dataset | maP |
|---|---|---|
| Selective Search [1] | PASCAL VOC 2010 | 35.1% |
| R-CNN | PASCAL VOC 2010 | 53.7% |
| OverFeat [2] | ILSVRC2013 | 24.3% |
| R-CNN | ILSVRC2013 | 31.4% |

## Definitions

- maP := mean average precision

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
**R-CNN workflow**
R-CNN details

# Workflow of R-CNN



R-CNN workflow

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
R-CNN workflow
**R-CNN details**

## Module design

1. Generates category-independent region proposals.
2. CNN that extracts a fixed-length feature vector from each region.
3. A set of class-specific linear SVMs.

## SVM

- **S**upport **V**ector **M**achines.
- Are supervised learning models with associated learning algorithms with the goal of classification and regression analysis.

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
R-CNN workflow
**R-CNN details**

## Module design

1. Generates category-independent region proposals.

   - Makes use of "Selective Search"[1]

2. CNN that extracts a fixed-length feature vector from each region.

   - Convolutional Neural Network predict the object classes.
   - Using the deep-learning framework "Caffe"[1].

3. A set of class-specific linear SVMs.

   - For object recognition.

Introduction
Machine Learning
Object Detection
**Convolutional Neural Network**
Grasping
References

Introduction
Region-based CNN
R-CNN workflow
**R-CNN details**

Further methods based on R-CNN's

- 2013 : R-CNN [3]
- 2015 : Fast-R-CNN [4]
- 2016 : Faster-R-CNN [2]
- 2017 : Mask R-CNN [3]

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
**Grasping**
References

Problems
Solutions
Conclusion

# Learning to Grasp from 50K Tries and 700 Robot Hours

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
**Grasping**
References

**Problems**
Solutions
Conclusion

## General problems in robot's grasping

- Using grasp-data set with human labeling can be quite challenging.
  1. Object can be grasped in multiple ways.
  2. Human notions of grasping are biased by semantics
- Biggest vision-based grasping dataset is only about 1k images. [1]
  1. Objects in isolation.
  2. Could lead to a bad performance under other environments.
  3.

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
**Grasping**
References

Problems
**Solutions**
Conclusion

## Approaches to deal with the problems.

- Using unlabeled grasping dataset.
  1. Self-supervising algorithm that learns to predict grasp locations via trial and error.
- Created their own dataset for grasping. [1]
  1. 50k items has been collected in 700h of trial and error.

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
**Grasping**
References

Problems
Solutions
**Conclusion**

## Results

- Test the grasp model both on novel objects and training objects under different pose conditions.
- Still failures even by 700h of "practice".

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
References

Problems
Solutions
Conclusion

Thanks for listening!

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
**References**

**Page 1**
Page 2
Page 3
Page 4

# References

📄 Lerrel Pinto and Abhinav Gupta. *The Supersizing Self-supervision: Learning to Grasp from 50K Tries and 700 Robot Hours*. The Robotics Institute, Carnegie Mellon University.

📄 Marco Pedersoli. *Hierarchical Multiresolution Models for fast Object Detection*. Universitat Autònoma de Barcelona, Bellaterra, April 2012.

📄 Ross Girshick, Jeff Donahue, Trevor Darrell and Jitendra Malik. *Rich feature hierarchies for accurate object detection and semantic segmentation*. University of California, Berkeley, 2014.

# References

📄 Rodrigo Benenson, Mohamed Omran, Jan Hosang, Bernt Schiele. *Ten Years of Pedestrian Detection, What Have We Learned?*. Max Planck Institut for Informatics Saarbrücken, Germany.

📄 Pierre Sermanet, David Eigen, Xiang Zhang, Michael Mathieu, Rob Fergus and Yann LeCun
*Integrated Recognition, Localization and Detection using Convolutional Networks*. Courant Institute of Mathematical Sciences, New York University

📄 Ross Girshick, Jeff Donahue, Trevor Darrell and Jitendra Malik. *Rich feature hierarchies for accurate object detection and semantic segmentation*. UC Berkeley, California, 2014

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
References

Page 1
Page 2
**Page 3**
Page 4

References

J. Uijlings, K. van de Sande, T. Gevers, and A. Smeulders. *Selective search for object recognition.*. IJCV, 2013.

B. Alexe, T. Deselaers, and V. Ferrari. *Measuring the objectness of image windows.* TPAMI, 2012.

C. H. Lampert, M. B. Blaschko, and T. Hofmann. *Efficient subwindow search: A branch and bound framework for object localization.* TPAMI, 2009.

Ross Girshick. *Fast R-CNN.* ICCV, 2015.

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
**References**

Page 1
Page 2
Page 3
**Page 4**

# References

📄 Jia, Yangqing and Shelhamer, Evan and Donahue, Jeff and Karayev, Sergey and Long, Jonathan and Girshick, Ross and Guadarrama, Sergio and Darrell, Trevor. *Caffe: Convolutional Architecture for Fast Feature Embedding.* 2014.

📄 Shaoqing Ren, Kaiming He, Ross Girshick and Jian Sun. *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks.* 2016.

📄 Kaiming He, Georgia Gkioxari, Piotr Dollár and Ross Girshick. *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks.* Facebook AI Research (FAIR), 2017.

Introduction
Machine Learning
Object Detection
Convolutional Neural Network
Grasping
**References**

Page 1
Page 2
Page 3
**Page 4**

# References

📑 Yun Jiang, Stephen Moseson, and Ashutosh Saxena. *Efficient grasping from rgbd images: Learning using a new rectangle representation.* 2011.