UCL x DeepMind lecture series

In this lecture series, leading research scientists from leading AI research lab, DeepMind, will give 12 lectures on an exciting selection of topics in Deep Learning, ranging from the fundamentals of training neural networks via advanced ideas around memory, attention, and generative modelling to the important topic of responsible innovation.

Please join us for a deep dive lecture series into Deep Learning!

#UCLxDeepMind

General information



Exits: At the back, the way you came in

Wifi: UCL guest





today's speaker Viorica Patraucean

Viorica is a research scientist at DeepMind, working mainly on Computer Vision related problems, with focus on video processing. She did her PhD in Toulouse, France, on statistical models for image understanding, and then focused on 3D shape- and video-analysis during her postdoctoral work in Paris and Cambridge. Her dream is to contribute to creating a computational model for the human visual system.



TODAY'S LECTURE Vision beyond classification: advanced models for Computer Vision We will cover computer vision tasks beyond image classification (object detection, semantic segmentation, instance segmentation) and associated models for each.

Then we will talk about the benefits of using more than single images as inputs to neural networks (pairs of images, videos) and the tasks that become available (optical flow estimation, action recognition).

In the third part, we will discuss settings that do not require strong supervision, in particular metric learning.

The talk will end with a brief discussion around open questions in computer vision.



Vision beyond classification

Advanced models for Computer Vision

6

Viorica Pătrăucean, Research Scientist

66 A picture is worth a thousand words

Classification models learn only a few

Resnet-50: bicycle, garden



Holy grail a model that achieves human level scene understanding











Sees-it-all model



Beyond supervised image classification

1 Supervised image classification

3 Supervised image classification 2 Supervised image classification

4 Open questions



The deep learning puzzle

By the end of this lecture, you will know how to redefine these building blocks to perform different visual tasks, using different inputs, and different forms of supervision.



DeepMind





Tasks beyond classification

Task definitions	Train and eval	Tricks of the trade
Object detection	Models and losses	Hard negative mining
Semantic segmentation	Metrics and benchmarks	Transfer learning



Important tasks not covered



Generated Caption: two beach chairs under an umbrella on the beach

Pose estimation



Image captioning

Image Captioning: Transforming Objects into Words, Herdade et al, 2019 Towards Accurate Multi-person Pose Estimation in the Wild, Papandreou et al, 2017

Tasks - Increasing granularity

classification



semantic segmentation



object detection



instance segmentation





DeepMind

Task 1 Object detection

Multi-task problem

Classification and localisation



Inputs



Targets



ightarrow Object bounding box (x_c,y_c,h,w)

for all the objects present in the scene



Inputs and targets



Dataset



Dataset

How to learn to predict bbox coordinates?



Recap: Softmax + cross entropy for classification



$$\ell_{\mathrm{CE}}(f_{\mathrm{sm}}(\mathbf{x}), \mathbf{t}) = -\sum_{j=1}^{k} \mathbf{t}_{j} \log[f_{\mathrm{sm}}(\mathbf{x}_{j})] = -\sum_{j=1}^{k} \mathbf{t}_{j} [\mathbf{x}_{j} - \log \sum_{l=1}^{k} e^{\mathbf{x}_{l}}]$$

Assign data points to categories; output is discrete.



Bounding box prediction





Classification

Mistakes are not quantifiable in classification; the data is not ordered.



Bounding box prediction





Classification



In classification, the output is discrete*, in regression the output is continuous.



Quadratic loss for regression



Minimise the mean squared error over samples.

Summary: classification vs regression

Property	Classification	Regression
Basic	map inputs to predefined classes	map inputs to continuous values
Output	discrete values	continuous values
Nature of the data	unordered data	ordered data
Algorithms	logistic regression, decision trees, neural networks	linear regression, neural networks



Quadratic loss for regression

$$\ell_2(\mathbf{x},\mathbf{t}) = \|\mathbf{t}-\mathbf{x}\|^2$$

Ground truth 1



How to deal with multiple targets?



Quadratic loss for regression

$$\ell_2(\mathbf{x},\mathbf{t}) = \|\mathbf{t}-\mathbf{x}\|^2$$

Ground truth 1



Convert regression into classification, by discretising the output values, and then refine through regression.



Classification then regression



Convert regression into classification, by discretising the output values, and then refine through regression.



Classification then regression



one_hot label

Convert regression into classification, by discretising the output values, and then refine through regression.



Two-stage detector

- ➔ Identify good candidate bboxes
- Classify and refine





Identify good candidate bboxes





Identify good candidate bboxes



 \bigcirc scales and ratios for (h, w)



Identify good candidate bboxes

ightarrow Discretise bbox space (x_c, y_c, h, w)

- \bigcirc anchor points for (x_c, y_c)
- \bigcirc scales and ratios for (h, w)
- \Rightarrow *n* candidates per anchor
- predict objectness score for each bbox
- \Rightarrow sort and keep top K



Identify good candidate bboxes

\bigcirc Discretise bbox space (x_c, y_c, h, w)

- \bigcirc anchor points for (x_c, y_c)
- \bigcirc scales and ratios for (h, w)
- \Rightarrow *n* candidates per anchor
- predict objectness score for each bbox
- \Rightarrow sort and keep top K

Refine through regression MLP(4)



Identify good candidate bboxes

\bigcirc Discretise bbox space (x_c, y_c, h, w)

- \bigcirc anchor points for (x_c, y_c)
- \bigcirc scales and ratios for (h, w)
- \Rightarrow *n* candidates per anchor
- predict objectness score for each bbox
- \Rightarrow sort and keep top K

Refine through regression MLP(4)



Figure from Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, Ren et al, 2016
Case study 1: Faster R-CNN

Want to learn more?



Jaderberg et al Spatial Transformer Networks (2015)

Identify good candidate bboxes

\Rightarrow Discretise bbox space (x_c, y_c, h, w)

- \bigcirc anchor points for (x_c, y_c)
- \bigcirc scales and ratios for (h, w)
- \Rightarrow *n* candidates per anchor
- predict objectness score for each bbox
- \Rightarrow sort and keep top K

Refine through regression MLP(4)



Figure from Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, Ren et al, 2016

Case study 2: RetinaNet - one-stage detector





Case study 2: RetinaNet - one-stage detector



Most of the candidate bboxes are easy negatives: poor learning signal.



Issue with one-stage detectors

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$



Issue with one-stage detectors

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$

Faster R-CNN prunes these in stage 1.

One-stage detectors employ hard negative mining heuristics.



Hard negative mining

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$

Faster R-CNN prunes these in stage 1.

One-stage detectors employ hard negative mining heuristics.



Hard negative mining

Want to learn more?



Sung and Poggio Learning and Example Selection for Object and Pattern Detection (1994)

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$

Faster R-CNN prunes these in stage 1.

One-stage detectors employ hard negative mining heuristics.



- 1. Get set of positive examples
- 2. Get a random subset of negative examples (full set is too big)
- 3. Train detector
- 4. Test on unseen images
- 5. Identify false positive examples (a
 person was detected where there was
 none) = hard negatives; add them to
 the training set
- 6. Repeat from 3.



RetinaNet solution

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$

Faster R-CNN prunes these in stage 1.

One-stage detectors employ hard negative mining heuristics.

RetinaNet uses Focal Loss (FL).



RetinaNet solution

Most of the candidate bboxes are background, easy to classify.

The accumulated loss of the many easy examples overwhelms the loss of rare useful examples $\ell_{CE}(p > .5) > 0$

Faster R-CNN prunes these in stage 1.

One-stage detectors employ hard negative mining heuristics.

RetinaNet uses Focal Loss (FL).



Good accuracy @ ~8fps speed



DeepMind

Task 2 Semantic segmentation

Semantic segmentation

Bounding boxes are not good representations for certain types of objects.

We need more refined representations.



Image from COCO dataset - Microsoft COCO: Common Objects in Context, Lin et al, 2014

Semantic segmentation

Inputs





Targets



Class label for every pixel





Semantic segmentation



Dense prediction problem - how to generate an output at the same resolution as the input?



Recap: Pooling



Pooling: compute mean or max over small windows to reduce resolution.



































Other upsampling methods

unpooling with indices SegNet

Deconvolutions DeconvNet

Unpooling: upsample to increase resolution; here 2x2 kernel.

SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, Badrinarayanan et al, 2016 DeconvNet: Learning Deconvolution Network for Semantic Segmentation, Noh et al, 2015



Case study: U-NET





U-Net: Convolutional Networks for Biomedical Image Segmentation, Ronneberger et al, 2015

Case study: U-NET





U-Net: Convolutional Networks for Biomedical Image Segmentation, Ronneberger et al, 2015

Recall RetinaNet - same U shape





Bonus: Instance segmentation

Want to learn more?



Semantic segmentation



Pixel-wise labels can be confusing for overlapping objects in the same category.

Instance segmentation



Object detection + segmentation







DeepMind

Metrics and benchmarks

Evaluation metrics

Want to learn more?



Berman et al. The Lovasz-Softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks (2018)

Classification

Accuracy: percentage of correct predictions

Top-1: top prediction is the correct class

Top-5: correct class is in top-5 predictions

Object detection and segmentation

- → intersection-over-union (IoU)
 - **non-differentiable**: used only for evaluation

$$\mathcal{J}(\mathbf{P},\mathbf{T}) = \frac{\mathbf{P} \bigcap \mathbf{T}}{\mathbf{P} \bigcup \mathbf{T}}$$



Benchmarks

Similar to Imagenet for various tasks

Public platforms for model evaluation

Maintain a leaderboard to track state-of-the-art models





5000 images with high quality annotations · 20000 images with coarse annotations

Dataset Overview



COCO 2019 Object Detection Task



DeepMind

Tricks of the trade

Transfer learning

Let $\mathcal{D} = \{\mathcal{X}, P(X)\}, X = \{x_1, ..., x_n\} \in \mathcal{X}$ be a domain and $\mathcal{T} = \{\mathcal{Y}, f(\cdot)\}, f(x_i) = \hat{y}_i, y_i \in \mathcal{Y}$ a task defined on this domain.

Given a source domain and task $\begin{pmatrix} \mathcal{D}_S \\ \mathcal{T}_S \end{pmatrix}$ and a target domain and task $\begin{pmatrix} \mathcal{D}_T \\ \mathcal{T}_T \end{pmatrix}$, reuse knowledge learnt by $f_S \inf f_T$

$$\begin{pmatrix} \mathcal{D}_S \\ \mathcal{T}_S \end{pmatrix} \longrightarrow \begin{pmatrix} \mathcal{D}_T \\ \mathcal{T}_T \end{pmatrix}$$

Intuition: features are shared across tasks and datasets. Reuse knowledge.



Transfer learning across different tasks

loss



 \bigcirc W₅₋₆ trained from scratch





Want to learn more?



Zamir et al. Taskonomy: Disentangling Task Transfer Learning (2018)

Figure from Taskonomy: Disentangling Task Transfer Learning, Zamir et al, 2018



Transfer learning across different domains

Sim2Real



Train in simulation using RL - \mathcal{D}_S

Use Automatic Domain Randomization: data augmentation + hard negative mining







Beyond supervised image classification











DeepMind



Beyond single image input






























Motion helps object recognition when learning to see.





Motion helps object recognition when learning to see.



The Development of Invariant Object Recognition Requires Visual Experience With Temporally Smooth Objects, Wood and Wood, 2018

Videos

 \rightarrow

 \rightarrow



Natural data augmentation: translation, scale, 3D rotation, camera motion, light changes



CATER: A diagnostic dataset for Compositional Actions and Temporal Reasoning, Girdhar and Ramanan, 2019

Beyond single image input

Inputs	Task definitions	Models	Challenges Obtaining labels	
Pairs of images	Optical flow estimation	Image-based models		
Videos	Action recognition	3D convnets	A note on efficiency	
		Recurrent (not covered)		



DeepMind





Optical flow estimation





Case study - FlowNet

- Encoder-decoder architecture similar to U-NET
 - Supervised training

 \rightarrow

 \rightarrow

 \rightarrow

- Loss: Euclidean distance
- Flying chairs dataset





Case study - FlowNet

- Encoder-decoder architecture similar to U-NET
 - Supervised training

 \leftrightarrow

 \rightarrow

 \rightarrow

 \mapsto

- Loss: Euclidean distance
- Flying chairs dataset
- Sim2Real transfer



FlowNet: Learning optical flow with convolutional networks, Fischer et al, 2015



DeepMind

Video input

Video models from image models

Cityscapes





Improving Semantic Segmentation via Video Propagation and Label Relaxation, Zhu et al, 2019

Video models using 3D convolutions



6

Recap: 2D convolution operation





The **kernel** slides across spatial dimensions.









































Properties of 3D convolutions



Strided, dilated, padded, [...] convolutions apply in 3D as well.



Action recognition

Inputs

- \rightarrow RGB video T x H x W x 3
- \rightarrow (optional) flow map $T \times H \times W \times 2$



Targets



action label one_hot $1 \times N_{classes}$

e.g. cricket shot

Kinetics600 dataset

- 600k training videos, 600 classes
- Curated Youtube videos
- Each video: 250 frames (~ 10 sec.)
- Current accuracy: 81.8 % top-1

Case study: SlowFast



Transfer learning returns

Want to learn more?

Carreira and Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset (2017)

Inflating 2D kernels into 3D



Intuition: a tiled image is a video of a static scene, filmed with a fixed camera.



DeepMind

Challenges

Challenges in video processing

Difficult to obtain labels Large memory requirements High latency High energy consumption

 \rightarrow

 \rightarrow

 \rightarrow

 \rightarrow



Improve efficiency of video models



 \mapsto

Inspiration from biological systems

Maximise parallelism to increase throughput and reduce latency [1, 2]

Exploit redundancies in the visual data to obtain frugal models [3]

1] Massively parallel video networks, Carreira, Patraucean et al, 2018 2] Sideways: depth-parallel training of video models, Malinowski, Swirszcz, Carreira, Patraucean, 2020 3] Blink and you won't miss it: video processing without temporal redundancies, Patraucean et al, 2020



Beyond supervised image classification











DeepMind



Beyond strong supervision

Labelling is tedious - Research topic in itself



Interactive Object Annotation with Polygons

NOTE: If inference is slow due to heavy traffic (benchmark is 0.3 seconds per interaction), please consider trying our demo locally using our available code For sponsorship/donation to help develop this web tool, please contact polyrnn@cs.toronto.edu





Self-supervision - Metric learning



Standard losses (e.g. cross-entropy, mean square error)

Hearn mapping between input(s) and output distribution / value(s)

Metric learning

Iearn to predict distances between inputs given some similarity measure (e.g. same person or not)

Images from VGGFace2: A dataset for recognising faces across pose and age, Cao et al, 2018

Self-supervision - Metric learning

Metric learning

 \rightarrow

 \rightarrow

- Contrastive loss
- Triplet loss
- \rightarrow
- State-of-the-art on representation learning

Applications

- (Multimodal) self-supervised
 representations, e.g. image+sound [1]
- \rightarrow
- Information retrieval [2]
- \rightarrow
- Low-shot face recognition [3]



Metric learning

Contrastive loss (margin loss)

Dataset:

 (r_0, r_1, y) $\begin{cases} y = 1, & \text{if } (r_0, r_1) \text{ same-person} \\ y = 0, & \text{otherwise} \end{cases}$

$$\label{eq:linear} \begin{split} \ell(r_0,r_1,y) &= y \mathrm{d}(r_0,r_1)^2 + (1-y)(\max(0,m-\mathrm{d}(r_0,r_1)))^2 \\ \mathrm{d}(\cdot,\cdot) \ \textbf{-Euclidean distance} \end{split}$$





Metric learning

Contrastive loss (margin loss)

Dataset:

 (r_0, r_1, y) $\begin{cases} y = 1, & \text{if } (r_0, r_1) \text{ same-person} \\ y = 0, & \text{otherwise} \end{cases}$

$$\begin{split} \ell(r_0,r_1,y) &= y \mathrm{d}(r_0,r_1)^2 + (1-y)(\max(0,m-\mathrm{d}(r_0,r_1)))^2 \\ \mathrm{d}(\cdot,\cdot) \ \textbf{-Euclidean distance} \end{split}$$




Metric learning

Contrastive loss (margin loss)

Dataset:

 (r_0, r_1, y) $\begin{cases} y = 1, & \text{if } (r_0, r_1) \text{ same-person} \\ y = 0, & \text{otherwise} \end{cases}$

$$\label{eq:linear} \begin{split} \ell(r_0,r_1,y) &= y \mathrm{d}(r_0,r_1)^2 + (1-y)(\max(0,m-\mathrm{d}(r_0,r_1)))^2 \\ \mathrm{d}(\cdot,\cdot) \ \textbf{-Euclidean distance} \end{split}$$

Difficult to choose *m*





Metric learning

Triplet loss

Dataset:

 (r_a, r_p, r_n) $\begin{cases} (r_a, r_p) & \text{similar} \\ (r_a, r_n) & \text{dissimilar} \end{cases}$

$$\ell(r_a, r_p, r_n) = \max(0, m + d(r_a, r_p)^2 - d(r_a, r_n)^2)$$

better than contrastive loss
relative distances more
meaningful than a fixed margin





Metric learning

Triplet loss

Dataset:

 (r_a, r_p, r_n) $\begin{cases} (r_a, r_p) & \text{similar} \\ (r_a, r_n) & \text{dissimilar} \end{cases}$

$$\ell(r_a, r_p, r_n) = \max(0, m + d(r_a, r_p)^2 - d(r_a, r_n)^2)$$

 better than contrastive loss
relative distances more meaningful than a fixed margin
hard negative mining to select informative triplets

Want to learn more?



Wu et al. Sampling Matters in Deep Embedding Learning (2018)





New state-of-the-art in representation learning

Same data, different augmentations



(a) Original

(f) Rotate {90°, 180°, 270°}



(b) Crop and resize



(g) Cutout



(c) Crop, resize (and flip) (d) Color distort. (drop) (e) Color distort. (jitter)



(h) Gaussian noise



(i) Gaussian blur





(j) Sobel filtering



New state-of-the-art in representation learning

 \Rightarrow

 \mapsto

Composition of data augmentations

Learnable non-linear transformation

Larger mini-batches and longer training



Beyond supervised image classification





03 Supervised image classification





DeepMind

Open questions



Open questions

 \rightarrow

Is vision solved? What does it mean to solve vision?

human level scene understanding - what benchmarks?

How to scale systems up?

model parallelism, better hardware, less supervision - more common sense

What are good visual representations for action?

Unsupervised Learning of Object Keypoints for Perception and Control, Kulkarni, Gupta et al, 2019



Learning to see from static images might make things harder than they should be.

Rethink vision models design and training from the perspective of moving pictures and with the end-goal in mind: intelligent agents that interact with the real world in real time.



Thank you

Questions



DeepMind

Useful resources





Al2 Allen Institute for Al

Computer Vision Explorer

~

×

~

Recognition

Classification

Segmentation

Vision and Language

Pose Estimation

Surface Normals

★ Scene Geometry

Depth

About

Human Centric Vision **^**

Detection

O

Classification

Image classification is the task of assigning an input image, a single label drawn from a fixed set of categories. Image classification models are trained and evaluated on large classification datasets such as ImageNet that has 1000 image categories.

TRY IT FOR YOURSELF

1. Choose an Image



2. Run a model

Л

Synthetic datasets for Computer Vision



SceneNet RGB-D: 5M Photorealistic Images of Synthetic Indoor Trajectories with Ground Truth

John McCormac Ankur Handa Stefan Leutenegger Andrew J. Davison

Dyson Robotics Lab at Imperial College, Department of Computing, Imperial Collge London

https://www.datasetlist.com/



Transporter architecture

