

CSCE
496/896
Lecture 9:
Object
Detection
Stephen Scott
Introduction
Performance
Measures
R-CNN
SPP-net
Fast R-CNN
YOLO

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

sscott@cse.unl.edu



◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ・ つくぐ



CSCE
496/896
Lecture 9:
Object
Detection
Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net

Fast R-CNN

YOLO

- Performance measures
- RCNN
- SPP-net
- Fast RCNN
- YOLO



Performance Measures Mean Average Precision

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN SPP-net Fast R-CNN

YOLO

Mean average precision (mAP) to measure how well

objects are identified

- Recall from Lecture 3
 - **Precision** is fraction of those labeled positive that are positive
 - **Recall** is fraction of the true positives that are labeled positive
 - Precision-recall curve plots precision vs recall





Performance Measures Mean Average Precision (2)

496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net

Fast R-CNN

YOLO

 Given a ranking (by confidence values) of *n* items, average precision at *n* (AP@n) is average of precision values at each position in the ranking:

$$AP = \sum_{k=1}^{n} P(k) \Delta r(k) \quad ,$$

where P(k) is precision at position k and $\Delta r(k)$ is change in recall: r(k) - r(k-1) (= 0 if instance k is negative, = $1/N_p$ if k is one of N_p positives)

- E.g., if ranking = $\langle +, +, -, +, \rangle$, AP@5 = (1)(1/3) + (1)(2/3) + (2/3)(2/3) + (3/4)(1) + (3/5)(1)
- Larger as more positives ranked above negatives
- mAP is mean of average precision across all classes



Performance Measures Intersection Over Union

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN SPP-net Fast R-CNN YOLO **Intersection over union** (IoU) to measure quality of bounding boxes

 Divide the size of the two boxes' intersection by the size of their union





Basic Idea of Object Detection

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net Fast R-CNN YOLO Split input image into $\ensuremath{\textit{regions}}$ and classify each region with

- a CNN and other machinery
 - Region boundary is bounding box
 - Object detected in region is object in BB



Issues:

- Limited to bounding boxes of fixed sizes and locations
- An object could span regions

Region CNN (Girshick et al. 2014)

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net Fast R-CNN YOLO

- R-CNN proposes collection of 2000 regions in image
- Warps each region to match input dimensions (227 × 227 × 3) of CNN to get 4096-dimensional embedded representation
- Classifies each embedded vector with class-specific binary SVMs
- Apply class-specific regressors to fine-tune bounding boxes



R-CNN: Regions with CNN features



Region CNN (Girshick et al. 2014) Example from Girshick (2015)



YOLO





Region CNN (Girshick et al. 2014) Selective Search

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net Fast R-CNN YOLO Popular method to propose Rols: selective search

- Segment the image
- Compute bounding boxes of segments
- Iteratively merge adjacent segments based on similarity
 - Linear combination of similarities of: color, texture, size, shape



◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ・ つくぐ

Goto 2



Region CNN (Girshick et al. 2014) Issues

- 496/896 Lecture 9: Obiect Detection
- Stephen Scott
- Introduction
- Performance Measures
- **R-CNN**
- SPP-net East R-CNN
- YOLO

- Training and detection are slow
- Detection: 13s/image on GPU, 53s/image on CPU
- Due to large number of regions proposed, each run through CNN and classified

◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ・ つくぐ

Nebraska Spatial Pyramid Pooling (He et al. 2015)

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net

Fast R-CNN

- Part of R-CNN's slowdown at test time is running each Rol through ConvNet separately
- To speed up test time, instead put entire image through single ConvNet



- Choose Rols from ConvNet output and run through spatial pyramid pooling (SPP) layer
 - Max/avg pooling with fixed number of bins
 - Produces fixed-length vector regardless of input size

Nebraska Lincoln

Spatial Pyramid Pooling (He et al. 2015) Example from Girshick (2015)



Nebraska Linon Spatial Pyramid Pooling (He et al. 2015) Drawbacks

- CSCE 496/896 Lecture 9: Object Detection
- Stephen Scott
- Introduction
- Performance Measures
- R-CNN
- SPP-net
- Fast R-CNN
- YOLO

- While training is faster then R-CNN, is still slow and disk-intensive
- Cannot efficiently update ConvNet parameters, so kept frozen
 - Each Rol's receptive field covers most of entire image, so forward pass expensive across all images of mini-batch



Fast R-CNN (Girshick 2015) Hierarchical Sampling

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net

Fast R-CNN

YOLO

Similar architecture to SPP-net

 Mini-batches constructed via hierarchical sampling: Sample a similar number of Rols over a smaller number of images







Fast R-CNN (Girshick 2015) Example from Girshick (2015)

CSCE Fast R-CNN (test time) 496/896 Lecture 9: Obiect Detection Linear + Softmax classifier Stephen Scott Bounding-box regressors softmax Introduction Fully-connected layers FCs Performance Measures "Rol Pooling" (single-level SPP) layer **B-CNN** "conv5" feature map of image Regions of SPP-net Interest (Rols) Fast R-CNN from a proposal Forward whole image through ConvNet YOLO method ConvNet Input image



Fast R-CNN (Girshick 2015) Example from Girshick (2015)





You Only Look Once (Redmon et al. 2016)

- CSCE 496/896 Lecture 9: Object Detection
- Stephen Scott
- Introduction
- Performance Measures
- R-CNN
- SPP-net
- Fast R-CNN
- YOLO

- A single, unified network
- Can process 45 frames per second on a GPU (155 fps for Fast YOLO)
- Lower mAP than some R-CNN variants, but much faster

◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ・ つくぐ

● Highest mAP of real-time detectors (≥ 30 fps)



You Only Look Once (Redmon et al. 2016)

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

R-CNN

SPP-net

Fast R-CNN

YOLO

- Divides image into *S* × *S* grid
- Each grid cell predicts *B* bounding boxes, each as (*x*, *y*, *w*, *h*) (coordinates, width, height), and a confidence (five total predictions)



- *x*, *y*, *w*, *h* ∈ [0, 1] (relative to image dimensions and grid cell location)
- Each cell also predicts C class probabilities
- Output is $S \times S \times (5B + C)$ tensor

You Only Look Once (Redmon et al. 2016) Architecture



Leaky ReLU for all layers except output, which is linear

Nebraska

Nebraska Lincol You Only Look Once (Redmon et al. 2016) Training

CSCE 496/896 Lecture 9: Object Detection

Stephen Scott

Introduction

Performance Measures

λ

R-CNN

SPP-net

Fast R-CNN

YOLO

- Pretrained 20 convolutional layers on ImageNet 1000
- Added 4 convolutional layers and 2 connected layers
- Trained to optimize weighted square loss function $\lambda_{coord} = 5$ times more weight on (x, y, w, h) predictions

$$\begin{array}{l} \begin{array}{l} \begin{array}{l} & \sum\limits_{i=0}^{S^{2}}\sum\limits_{j=0}^{B}\mathbb{1}_{ij}^{\text{obj}}\left[(x_{i}-\hat{x}_{i})^{2}+(y_{i}-\hat{y}_{i})^{2}\right] \\ & +\lambda_{\text{coord}}\sum\limits_{i=0}^{S^{2}}\sum\limits_{j=0}^{B}\mathbb{1}_{ij}^{\text{obj}}\left[\left(\sqrt{w_{i}}-\sqrt{\hat{w}_{i}}\right)^{2}+\left(\sqrt{h_{i}}-\sqrt{\hat{h}_{i}}\right)^{2}\right] \\ & \quad +\sum\limits_{i=0}^{S^{2}}\sum\limits_{j=0}^{B}\mathbb{1}_{ij}^{\text{obj}}\left(C_{i}-\hat{C}_{i}\right)^{2} \\ & \quad +\lambda_{\text{noobj}}\sum\limits_{i=0}^{S^{2}}\sum\limits_{j=0}^{B}\mathbb{1}_{ij}^{\text{noobj}}\left(C_{i}-\hat{C}_{i}\right)^{2} \\ & \quad +\sum\limits_{i=0}^{S^{2}}\mathbb{1}_{i0}^{\text{obj}}\sum\limits_{c\in\text{classes}}\left(p_{i}(c)-\hat{p}_{i}(c)\right)^{2} \end{array} \right.$$