

OD: object detection

التاريخ: / /

Deep Detection

في الصورة الواحدة يمكن تصنيف أكثر من عُرف مثلاً:
يمكن وجود أكثر من عُرف في الصورة (كلب، قطة، بطّة)
وهناك إمكانية لتصنيف كل الأعراض في نفس الصورة.

1. What are the challenges of

التي تواجه (OD) object detection?

1. قد تحتوي الصورة على عدة classes فمثلاً لو كانت صورة حديقة تحتوي على عدة أشياء مختلفة وقد يتواجد عدة أشياء متشابهة وهذا يعني تكرار النماذج.

2. تحديد في تحديد حجم وموقع الكائن الموجود في الصورة.

باستخدام Bounding Box تحديد دقة عملية B.B.

BB: Bounding Box

* object detection evaluation

عند إجراء اختبار "Test" على النموذج المستخدم لتحديد

مواقع الأشياء والكائنات في الصور باستخدام Bounding Box.

يتم التنبؤ بمواقع المربعات وتصنيف الأشياء الموجودة و

تحديد درجة الثقة confidence scores لكل مربع.

ثم سيتم تحديد فيما إذا كان المربع المكتشف هو إيجابي

معيّن True positive أو False positive

ثم يتم رسم منحنى Recall - precision curve لفئة

حساب المسافة تحت المنحنى لتحديد متوسط الدقة Avg precision

كل فئة

يتم حساب معدل متوسط الدقة "mAP" باستخدام هذا

المقياس لتقييم أداء النموذج تحديد المواقع للأشياء وتحديد مدى قدرته على التعرف على الأشياء بدقة وفعالية.

الصيغة:

$$\text{Precision} = \frac{\text{True Positive detections}}{\text{total detections}}$$

$$= \frac{TP}{TP + FP}$$

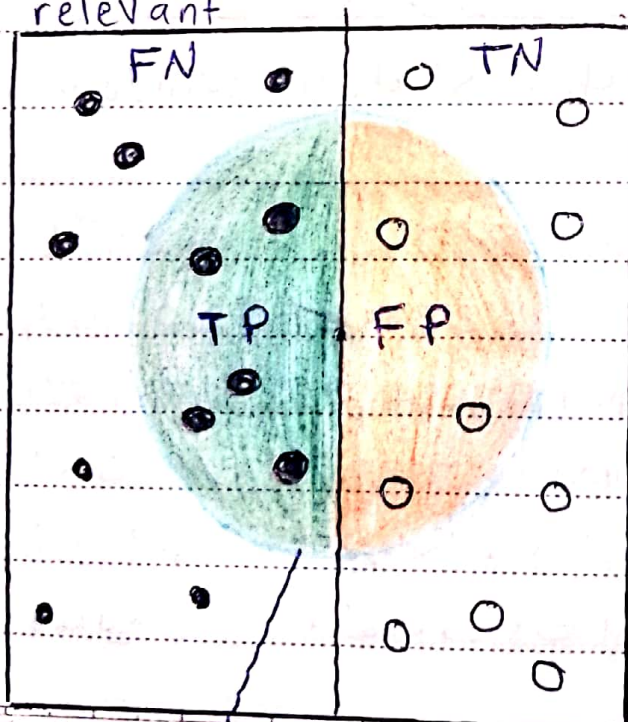
$$\text{Recall} = \frac{\text{True Positive detection}}{\text{total positive test instances}}$$

$$= \frac{TP}{TP + FN}$$

الصيغة Total detections عدد الحالات الكلية التي تم كشفها

TP عدد الحالات المكتشفة الصحيحة

elements relevant



هام FP اكتشاف خاطئ

والتي ليس هو الهدف

الطالب مثلاً صنف الهدف

على أنه دمه وهو ليس دمه

FN : هي ما كان دمه

مثلاً ولم يتم كشفه

TP : ما كان دمه واكتشف

على أنه دمه

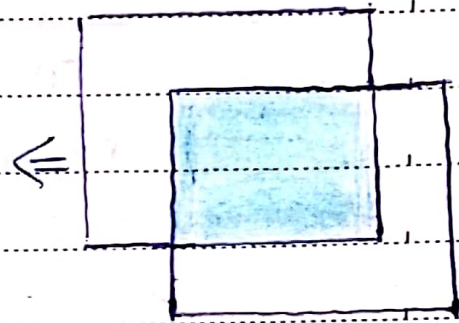
⇒ How many selected items are relevant?

$$\text{Precision} = \frac{TP}{TP+FP}$$

⇒ How many relevant items are selected?

$$\text{Recall} = \frac{TP}{TP+FN}$$

تقاطع
Area (GT ∩ Det) > 0.5
Area (GT ∪ Det)
اتحاد



النقطة التي توي الفرق
كامل المنطقة > 0.5

why Recall increase when Precision decreases? لماذا تزداد Recall عند نقصان Precision

لانخفاض Recall يجب انخفاض المقام و FN لا يجب انخفاض TP لانها موجودة بالسطر والمقام لذلك تنقص FN + TP والقوانين السابقة

conceptual approach: sliding window detection

في الطريقة المستخدمة لتحديد مواقع الأشياء في الصورة مثلاً
تربط window (نافذة = مربع) على كامل الصورة وهذه النافذة
تتوي على نموذج الشفا (detection model) وسيعطى منطقة
المعرف (بمربع) (Bounding Box) وتتاح هذه العملية لك
تقييم النموذج المستخدم لتحديد المواقع ويجب أن يحقق هذا النموذج
قدارية عالية وفقدان معدلات الإيجابيات المأثرة TN للحصول

على نتائج دقيقة، ومن ذلك فإن هذه العملية صعبة التفسير والتطبيق على مجموعة واسعة من المقاييس وفي العوض بالمال الطول (Aspect ratios) وهذه العملية تسهل الكثير من الموارد والدقة لذلك هناك حدود متزايدة لتجنب طوفان تدريب المواقع في الصور مثل استنساخ المقاييس الهرمي "pyramid scale" والتعلم العميق لتجنب دقة التدريب.

* Histogram of oriented gradients (HOG)

يتم تقسيم الصورة إلى كتل (blocks) و كل كتلة تسمى "histogram" لكل كتلة.

Image input → Normalize gamma of color → compute gradients → weights vote into spatial of orientation cells → contrast normalize over overlapping spatial blocks.

* pedestrian detection with HOG

يتم تدريب نموذج لتقدير المشاة في الصور باستخدام آلية دعم

المقاييس الخطية linear svm باستخدام مجموعة من الصور التي

تحتوي على مشاة وغير حاوية الأشياء كمجموعة تدريب.

وعند إجراء الاختبار يتم تطبيق هذا النموذج على Feature map

ويتم البحث عن local maxima في الاستجابة للنموذج ويتم تكرار

هذه العملية على مجموعة من الأبعاد المختلفة لتدريب المشاة في الصور أو

معالجة الفيديو.

R-CNN: Region proposals + CNN Features

R-CNN تعمل وفقاً الخطوات التالية:

Input image → Region proposals → warped image regions → Forward each region through convnet → classify regions with SVMs.

من input image في الصورة التي يتم العمل عليها
و Region proposals في اقتراح مجموعة مناطق ذات حجم مختلفة

من الصورة المدخلة وهذه المناطق قد تحتوي أخطاء مراد كنقطة
و warped image regions تدوير كل منطقة و

تدوير كل حجم ثابت و Forward each region convnet

في تصنيف CNN على هذه المناطق لكشف الأجزاء و

SVMs في القيام بعملية تصنيف لهذه الأجزاء.

إيجابيات R-CNN : Pros

- ① دقة: تحقق نتائج دقيقة في تحديد مواقع الأشياء في الصور
- ② قابلية للتطبيق على الصور من تقاييم البكسل المختلفة

سلبيات R-CNN : Cons

- ① ليست نظام واحد شامل بل يتطلب الأمر تدريب أجزاء مختلفة للتعرف على شكل منفصل
- ② يتطلب التدريب الكثير من الوقت والموارد بحوالي (48) ساعة
- ③ يتطلب التدريب مرور الصور عبر البكسل العصبونية العميقة حوالي "2000" مرة

④ تكون عملية الاستدلال بطيئة (Inference) وتنفرد حوالي

47 ثانية لكل صورة. بإستعمال (VGG-16) وبالتالي هذه الطريقة

غير فعالة في الحالات التي يتطلب فيها الأمر الكشف عن الأشياء بسرعة

* Fast R-CNN

Fast R-CNN تعمل وفقاً الخطوات:

input image \rightarrow Forward whole image through convNet \rightarrow conv5 Feature map of image \rightarrow RoI pooling layer \rightarrow Fully-connected layers \rightarrow Bounding-box regressor

منه convNet لا استخراج المعالم (features) من الصور
conv5 Feature map... في Feature map مجموعة من
والتي تحتوي على معلومات محددة عن المناطق المهمة في الصورة مثل
الحواف والأشكال و RoI pooling Region of Interest
للوصول إلى المعلومات المحددة من المواقع المحتملة التي تم التنبؤ بها
ثم تتم عملية Fully connected ثم Bounding-box
التدريب المواقع والتصنيف.

* Fast R-CNN training

بعد القيام بخطوات Fast R-CNN يتم القيام خطوة
جديدة في multi-task loss أي تستخدم هذه الخطوة لتجنب
دقة التصنيف بالاعتماد على الدالة loss:

$$L(y, p, b, \hat{b}) = \underbrace{-\log p(y)}_{\text{softmax loss}} + \lambda \underbrace{\mathbb{I}[y > 1] L(b, \hat{b})}_{\text{regression}}$$

منه $p(y)$ loss for ground truth class
 \hat{b} predicted class probabilities
predicted box

// Yolo: You only look once //

التاريخ: / /

$$L_{reg}(b, \hat{b}) = \sum_{i=\{x, y, w, h\}} smooth_{L_1}(b_i - \hat{b}_i)$$

$$smooth_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

مقارنة بين Fast R-CNN و R-CNN

؛ Faster R-CNN

system	time	of data	of +12 data
R-CNN	~ 50 s	66.0	-
Fast R-CNN	~ 2 s	66.9	70.0
Faster R-CNN	198 ms	69.6	73.2

objects detection by Yolo

تعمل Yolo على تقسيم صورة الدفق إلى $(K \times K)$ خلية
 (grid = cell) وكل خلية تمثل الغرف المادي فيها ويكون
 هذا الغرف واعتقاداً في هذه الخلية يجب ان يكون مركزها
 (x, y) الخرد للفرص يقع في هذه الخلية
 حيث سيتم تحديد مواقع وجود الكائنات في صورة ما وتحديد

حدودها باستخدام مربعات تسعين boundary boxes

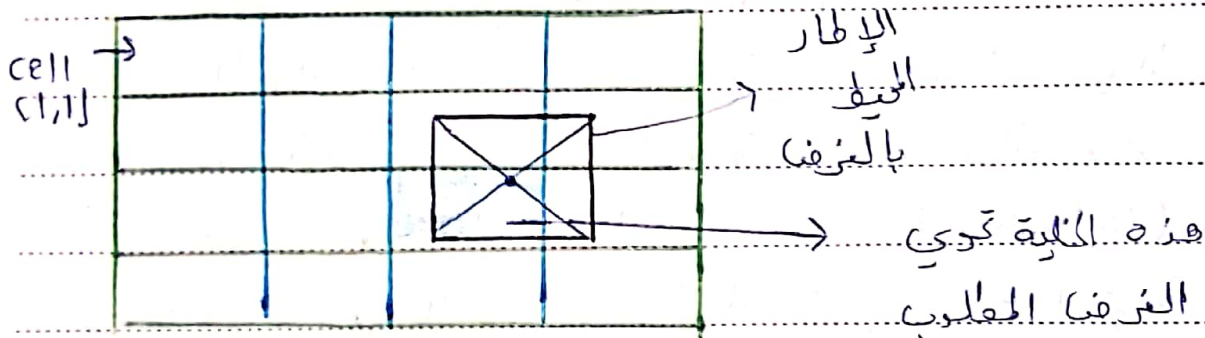
ويحتوي كلا مربع على قيم عناصر (x, y, w, h)

بالإضافة إلى درجة الوثوقية box confidence score

(x, y) هما الإحداثيات الأفقية والعمودية

(w, h) هما العرض والارتفاع

$$0 < x < 1, 0 < y < 1$$



Input image

وتكون خطوات 40.10

1- استخراج مميزات الصورة باستخدام شبكة تحويل الصورة إلى

مربعات ملائمة conv feature maps وبمقاس 7×7

2- يتم إضافة طبقتي Fully connected (FC) إلى map

للتنبؤ بنتائج التصنيف وتحديد عدد المربعات لكل فئة ويتم توقع

نتائج التصنيف على شكل نقاط درجات الثقة score لكل

فئة يمكنه بالإضافة إلى حدود مربعين (bbox) مع الثقة لكل مربع

3- بالنسبة لبيانات PASCAL فإن الحجم المتوقع هو

$$7 \times 7 \times 30 \text{ الـ } 30 \text{ هي ناتج } (4+1) \times 2 + 20$$

الـ 20 هي عدد الفئات يمكنه التصنيف

* الهيكلية الرئيسية لنموذج 10% :

24 convolution layers

Fully connected (FC) طبقتين

7x7 x 1024 منظمات ذات أبعاد

أما لتدريب المربعات المحدودية :

يتم في 10% تدريب عدّة مربعات محدودية bounding boxes

لكل cell في صورة الدفلا بهدف تغطية كلا المرافق

المحتملة لوجود كائنات. ولتدريب المربع الحدودي المؤهل

عن الكائن الحقيقي في الصورة يتم اختيار المربع الذي يحقق

أعلى قيمة لمقياس (IoU) intersection over union

وعند حساب الخسارة loss للتنبؤ بالكائنات الحقيقة يتم

حساب الخسارة فقط للمربع الحدودي الذي تم اختياره

* الحساب loss في نموذج 10% : يتم في 10% استخدام

خسارة مربعات الاختلاف "sum-squared error" بين

التنبؤات و ground truth والتي تتكون من ثلاث

مكونات رئيسية :

1- خسارة التصنيف classification loss

2- خسارة الموقع localization loss : يتم حسابها

باستخدام مقياس (IoU) وهي الخسارة المرتبطة بدقة

تدريب موقع المربع الحدودي

3- خسارة الثقة confidence loss وهي الخسارة المرتبطة

بمدى ثقة وجود الكائن داخل المربع الحدودي

* The localization loss:

$$\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] \\ + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 \\ + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2]$$

حيث λ_{coord} هي

increase the weight for the loss in the boundary box coordinates

* $\mathbb{1}_{ij}^{obj}$: 1 if the j th boundary box in cell i is responsible for detecting the object, otherwise 0.

* confidence loss

$$\sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} (c_i - \hat{c}_i)^2$$

إذا تم الكشف عن كائن في box i

$$\lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} (c_i - \hat{c}_i)^2$$

إذا لم يتم الكشف عن الكائن في box i

\hat{c}_i is the box confidence score of the box j in cell i

$\mathbb{1}_{ij}^{noobj}$ is the complement of $\mathbb{1}_{ij}^{obj}$

λ_{noobj} : weights down the loss when
detecting background

* class Probability

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

ملاحظة: وفقاً لعمل شبكة 10% وتقسيم الصورة إلى خلايا
ونحسب كل خلية كحالة كانت تحوي أو غير خلية أو لا بالتالي
من الممكن أن يكون العرف يغطي أكثر من خلية فهذا ليس
رسم أكثر من إطار على نفس العرف وكلا هذه المشكلة نقوم
بعمل هدف للقيم الدنيا