

Farshad Borjalizade

Object detection with custom
version of yolov2

University of Tehran

Dr.Sajedi

Mordad 1400



فهرست

چکیده	3
مقدمه	4
مجموعه داده	6
پیاده سازی مدل	7
تعریف مدل	9
ساختار شبکه عمیق استخراج ویژگی ها	12
تابع خطا (Loss function)	14
آموزش مدل	15
تفاوت ها و شباهت ها	16
منابع	17

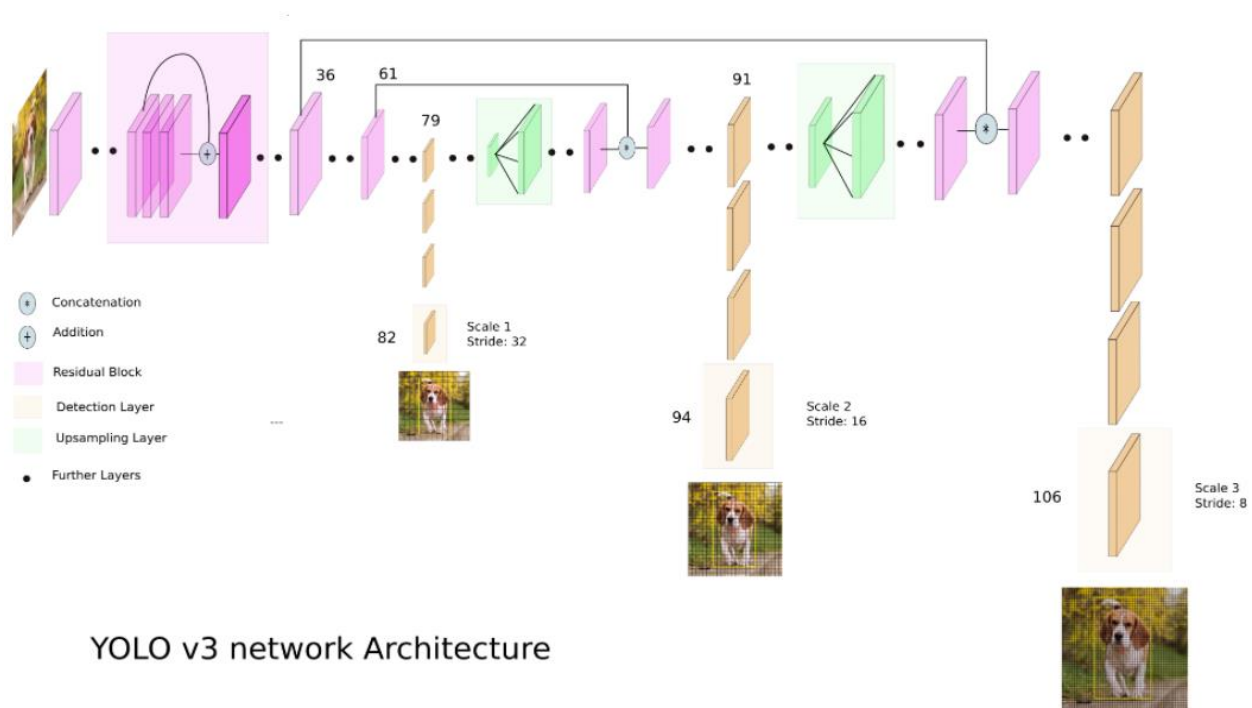


چکیده

در بینایی ماشین تشخیص اشیا (object detection) همواره یکی از چالش برانگیزترین و محبوب ترین موضوعات به دلیل کاربردهای فراوانی که دارد محسوب می شود. روش های متعددی در این زمینه بیان شده است که یکی از بهترین های آن ها استفاده از شبکه های عصبی پیچیده عمیق (DCNNs) است این روش بهترین نتایج را در مقایسه با روش های دیگر گرفته است، این روش برای بهینه سازی پارامترهای مدل و یادگیری از روش های مبتنی بر مشتق (gradient based) استفاده می کند. در این پروژه قصد داریم که با الهام گرفتن از مدل های YOLOv2 و YOLOv3 مدلی را طراحی کنیم که میزان خطای آن کمتر از مدل های گفته شده باشد. شبکه خود را با مجموعه داده Pascal VOC2012 آموزش داده ایم که در مورد آن در بخش مربوطه توضیحات کامل بیان شده است.

در چند سال اخیر یادگیری عمیق (deep learning) به طور گسترده ای در زمینه های مختلف بینایی ماشین (computer vision) از جمله دسته بندی تصاویر (image classification) و تشخیص اشیا و تقسیم بندی تصویر (image segmentation) مورد استفاده قرار گرفته است، یکی از موضوعاتی که بسیار مورد توجه قرار گرفته، تشخیص اشیا است. چالش هایی که در این حوزه مطرح است را می توان به دو قسمت تقسیم کرد اول اینکه مکان شی مورد نظر را در تصویر شناسایی کنیم (location problem) و دوم اینکه که در مکان مشخص شده چه شی ای قرار دارد (category problem). در روش های سنتی تشخیص اشیا از غلتاندن پنجره (sliding windows) با اندازه های متفاوت و کاندید قرار دادن مناطق متفاوت تصویر این کار انجام می شد با ظهور شبکه های عمیق و خصوصاً شبکه های پیچشی (convolutional neural network) ها و استفاده آن در مدل های مختلف R-CNN و Fast R-CNN و Faster R-CNN جهش قابل توجهی در این زمینه رخداد اما یک ایراد بزرگی که به این مدل ها می شد گرفت بار محاسباتی و زمان بسیار آن ها برای یادگیری شبکه بود که منشأ این ایراد را می توان به دو عامل تقسیم کرد اول اینکه این شبکه ها end-to-end نبودند به این معنا که شبکه بدین صورت نیست که تمام کار های تشخیص تصویر به صورت یکپارچه انجام شود و دوم اینکه مدل زمان بسیار زیادی را صرف یادگیری و کاندید قرار دادن نواحی مختلف تصویر برای اینکه مشخص کند آیا در منطقه مورد نظر شی ای وجود دارد یا خیر می کرد.

با معرفی مدل YOLO (you only look once) مشکل کاندید قرار دادن نواحی مختلف تصویر را به عنوان یک مسئله رگرسیون در نظر می گیرد و نواحی محتمل وجود شی را شناسایی می کند، این روش به طور قابل توجهی سرعت یادگیری مدل را افزایش می دهد اما نسبت به مدل Faster R-CNN از میزان خطای بیشتری برای تشخیص مکان اشیا برخوردار است، برای حل این مشکل نسخه های متفاوتی از YOLO منتشر شده است که با استفاده از تکنیک های متنوع سعی بر این داشتند که علاوه بر حفظ سرعت مدل ایراد دقت پایین مدل برای تشخیص مکان را هم رفع کنند. YOLOv3 توانست با استفاده از ساختار شبکه Darknet-53 توازن خوبی بین سرعت و دقت مدل، کار تشخیص اشیا را انجام دهد و به اصطلاح جز مدل های برتر (state-of-the-art) این حوزه قرار گرفت.

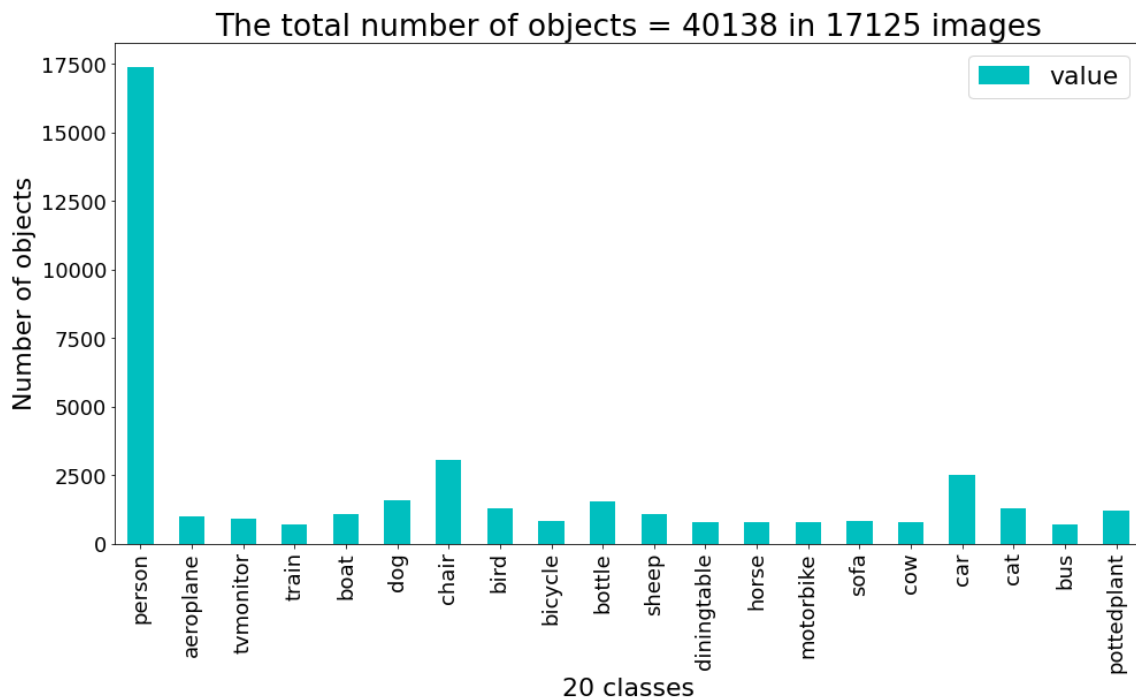


تصویر ۱- نمونه ای از ساختار YOLOv3

ایده مدلی که پیاده سازی کرده ایم ترکیبی از مدل های YOLOv2 و YOLOv3 است که با بسیاری از پارامترها و ساختار مدل های YOLO اشتراک دارد.

مجموعه داده

در این پروژه از مجموعه داده معروف Pascal VOC2012 استفاده می کنیم این مجموعه داده به طور کلی دارای ۲۰ کلاس از ۴ دسته متفاوت که انسان (Person) ۱ کلاس، حیوان (Animal) ۶ کلاس، وسایل نقلیه (Vehicle) ۷ کلاس و وسایل خانه (Indoor) ۶ کلاس می باشد که نتایج را بر روی آن بیان می کنیم در تصویر ۲ تعداد کلاس ها و اشیا موجود در هر کلاس به همراه تعداد تمامی اشیا موجود در مجموعه داده نشان داده شده است.



تصویر ۲- توزیع تعداد اشیا در کلاس ها

```
{'name': 'aeroplane', 'xmin': 86, 'ymin': 115, 'xmax': 312, 'ymax': 270},  
{'name': 'aeroplane', 'xmin': 110, 'ymin': 130, 'xmax': 163, 'ymax': 182},  
{'name': 'person', 'xmin': 162, 'ymin': 266, 'xmax': 177, 'ymax': 339}  
{'name': 'person', 'xmin': 21, 'ymin': 279, 'xmax': 36, 'ymax': 352}
```

تصویر ۳- نمونه ای از داده های برچسب گذاری شده

پیاده سازی مدل

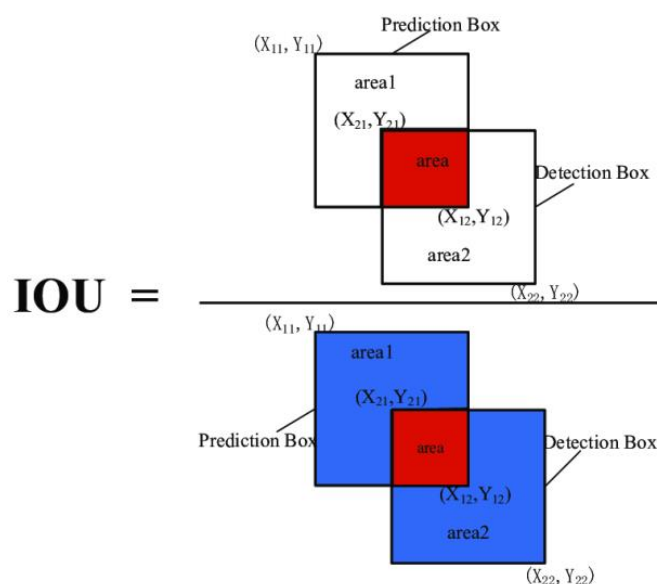
برای پیاده سازی مدل تشخیص شی ابتدا نیاز به یک سری از تعریف جزئیات و مفروضات داریم.

Intersection Over Union – IOU

در بسیاری از مدل های تشخیص اشیا IOU به عنوان معیاری برای ارزیابی در نظر گرفته می شود، در واقع معیاری است برای بررسی میزان احتمال وجود یک شی در یک مکان از تصویر.

فرض کنید که یک تصویر برچسب خورده داریم یعنی تمام اشیا موجود در تصویر را به همراه محل رخداد آن ها (bounding box coordinates) را داریم سپس مدل ما هم یک سری از مکان ها را به عنوان رخداد شی ای در نظر میگیرد معیار IOU کاری انجام می دهد به این صورت است که بررسی می کند چقدر

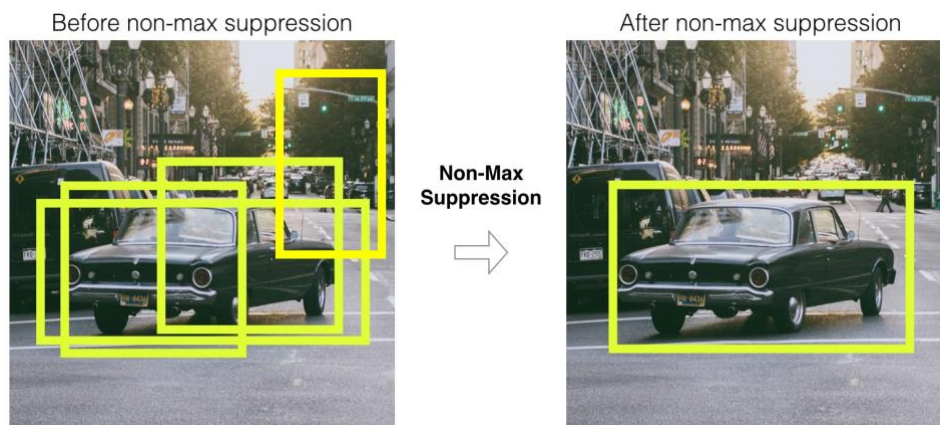
محدوده ای که مدل ما انتخاب کرده است با چیزی که در واقعیت وجود دارد اشتراک و هم پوشانی دارد.



تصویر ۴- نحوه عملکرد الگوریتم IOU

Non Maximum Suppression - NMS

در فرایند آموزش مدل و در بخش تعیین محل رخداد شی در تصویر (bounding box coordinates) تعدادی مکان برای تشخیص شی پیش بینی کرده ایم اما نکته ای که وجود دارد این است که ممکن است این پنجره ها شی های یکسانی را پیش بینی کرده باشند برای مثال اگر در تصویر یک ماشین داشته باشیم قصد داریم که بهترین پنجره ای که مکان رخداد شی را پیش بینی کرده است را داشته باشیم و بقیه پیش بینی ها را حذف کنیم، در واقع هدف حذف کردن پنجره هایی است که یک شی یکسان را نشان می دهند.

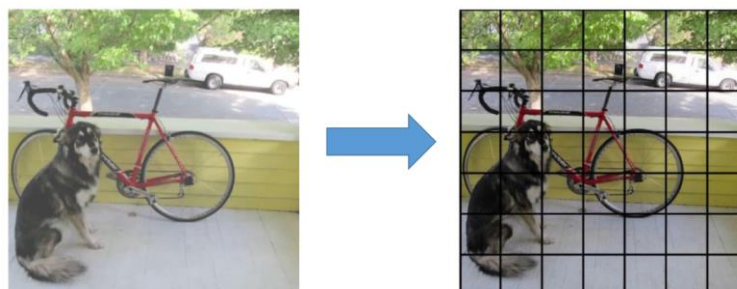


تصویر ۵- مثالی از اعمال الگوریتم NMS

تعریف مدل

ساختار کلی مدل به چهار قسمت تقسیم می شود که عبارت اند از:

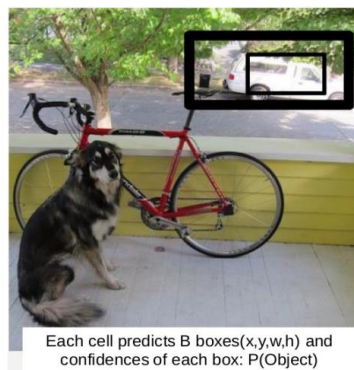
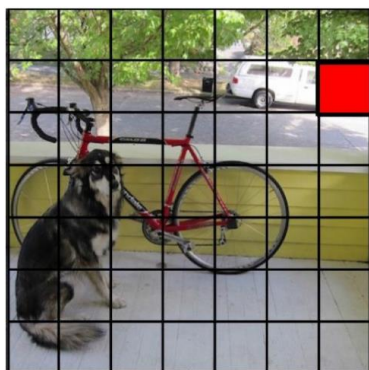
۱. اندازه تصاویر ورودی : در مجموعه داده Pascal VOC2012 تصاویر با اندازه های متفاوتی وجود دارند اما از آنجایی که اندازه ورودی شبکه نمی تواند متغیر باشد پس اندازه تصاویر را به 416×416 تبدیل می کنیم و هر کدام از تصاویر را به گرید سل های $S \times S$ تقسیم بندی می کنیم و معمولا هم این اعداد فرد هستند تا بتوان مرکز آن را به طور دقیق تشخیص داد و هر کدام از این گریدها مسئول شناسایی یک شی هستند.



تصویر ۶- اعمال گرید بندی $S \times S$ تصویر

۲. شبکه استخراج ویژگی : استخراج ویژگی را با استفاده از شبکه های عمیق پیچشی (deep convolutional neural networks) و الگوریتم Kmeans برای دسته بندی و انتخاب بهترین IOU و Batch Normalization و استفاده از تکنیک concatenate به لایه های مختلف انجام داده ایم.

۳. پیش بینی محدوده شی (Bounding box prediction): هر گرید سلی که دارای B جعبه پیش بینی است و برای هر جعبه یک احتمالی داریم که نشان دهنده این است که چقدر محتمل است در جعبه پیش بینی شده یک شی وجود داشته باشد (confidence scores) علاوه بر این برای نشان دادن مختصات هر جعبه هم به چهار متغیر نیازمندیم که به اضافه احتمال وجود شی در جعبه (x , y , w , h , conf) سر جمع پنج متغیر و به علاوه تعداد کلاس هایی که در مجموعه داده داریم (بیست کلاس) و در مجموع یک بردار سه بعدی داریم $(S \times S \times B \times (5 + C))$ و با استفاده از آن پیش بینی را انجام می دهیم.



تصویر ۷- در این تصویر برای هر سلول دو Bounding box در نظر گرفته شده است.

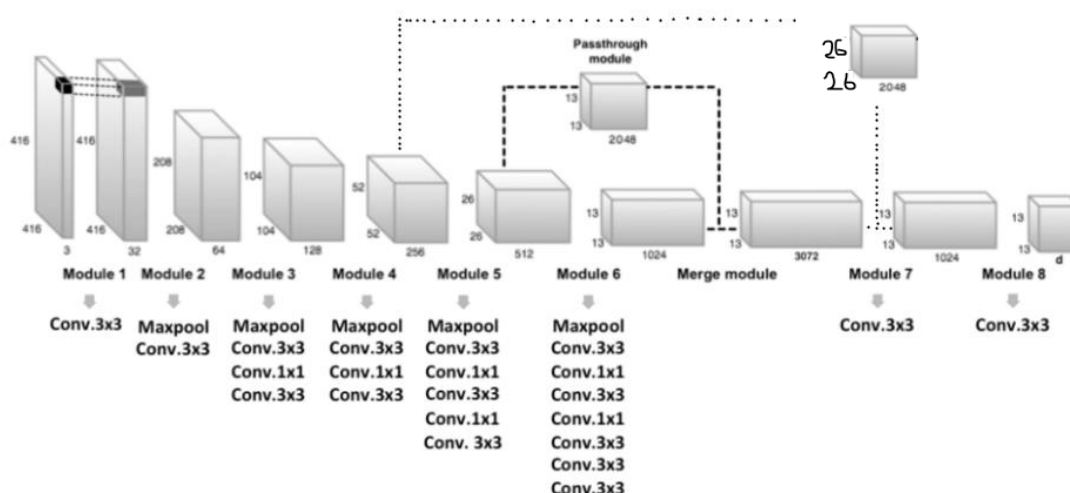
۴. تشخیص نهایی : در نهایت می بایستی ما بهترین و محتمل ترین جعبه را برای تصویر و تشخیص شی انجام دهیم که این کار را با اعمال الگوریتم NMS انجام خواهیم داد.



تصویر ۸- نمونه ای تشخیص نهایی

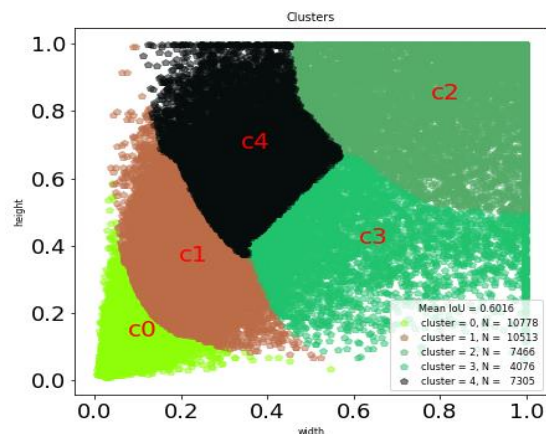
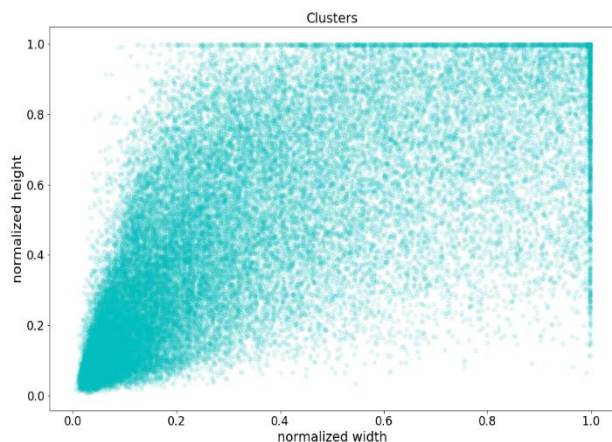
ساختار شبکه عمیق استخراج ویژگی ها

مدل طراحی شده در واقع ترکیبی از مدل های YOLOv2 و YOLOv3 است که با استفاده از تکنیک ترکیب کردن ویژگی های استخراج شده در لایه های پایین تر به بالاتر به تشخیص بهتر و دقیق تر اشیاء کمک می کند.

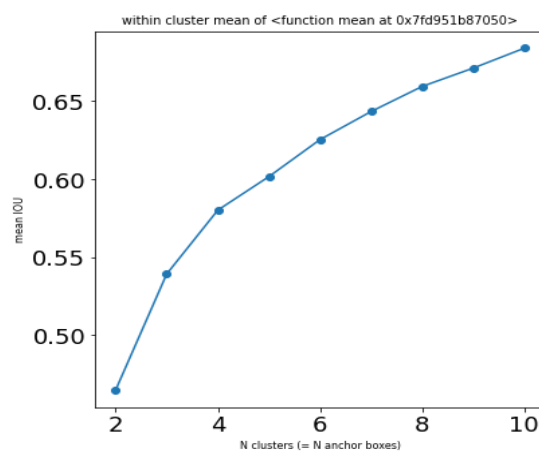
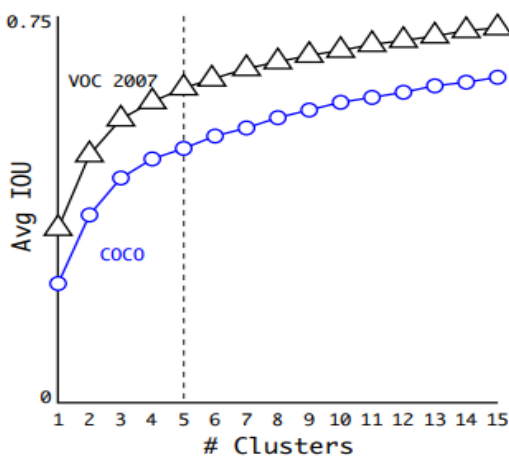


تصویر ۹- مدل پیاده سازی شده برای استخراج ویژگی ها

برای اینکه تشخیص اشیاء دقیق تر صورت گیرد به جای اینکه به صورت تصادفی Bounding box تعریف کنیم Bounding box های از پیش تعریف شده ای به عنوان Anchor box را می توان در نظر گرفت و این کار را می توان به دو صورت انجام داد یکی اینکه به صورت دستی یک سری Anchor box بدست آورد و دوم اینکه به آن به عنوان یک مسئله خوشه بندی (Clustering) نگاه کرد زیرا هریک از اشیاء دارای یک (w, h) است که میتوان آنها را با استفاده از الگوریتم kmeans خوشه بندی کرد در این مدل $k=5$ و معیار اندازه گیری را $d(box, centroid) = 1 - IOU(box, centroid)$ در نظر گرفته ایم.



تصویر ۱۰ - خوشه بندی پنج کلاسه بر اساس width و height تصاویر



تصویر ۱۱ - نتیجه میانگین IoU در تعداد خوشه بندی های متفاوت

پس پنج Anchor box که در واقع مرکز خوشه هاست را در نظر گرفتیم.

مدل دارای بیست و چهار لایه می باشد که به دلیل طولانی بودن از نمایش دادن آن در اینجا صرف نظر می کنیم (می توان آن را در [قسمت کد](#) مشاهده کرد)، تقریباً دارای پنجاه و سه میلیون پارامتر می باشد.

تابع خطا (Loss function)

همانند همه مدل های یادگیری می بایستی یک تابع خطایی را برای شبکه انتخاب کنیم تابع خطایی که در نظر گرفته ایم را می توان به سه بخش Localization و Confidence و Classification تقسیم کرد که در تصویر ۱۲ جزئیات بیشتری را می توان دید.

$$\begin{aligned} \text{loss}_{i,j} &= \text{loss}_{i,j}^{\text{xywh}} + \text{loss}_{i,j}^p + \text{loss}_{i,j}^c \\ \text{loss}_{i,j}^{\text{xywh}} &= \frac{\lambda_{\text{coord}}}{N_{L^{\text{obj}}}} \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{obj}} \left[(x_{i,j} - \hat{x}_{i,j})^2 + (y_{i,j} - \hat{y}_{i,j})^2 + \right. \\ &\quad \left. (\sqrt{w_{i,j}} - \sqrt{\hat{w}_{i,j}})^2 + (\sqrt{h_{i,j}} - \sqrt{\hat{h}_{i,j}})^2 \right] \\ \text{loss}_{i,j}^p &= -\frac{\lambda_{\text{class}}}{N_{L^{\text{obj}}}} \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{obj}} \sum_{c \in \text{class}} p_{i,j}^c \log(\hat{p}_{i,j}^c) \\ \text{loss}_{i,j}^c &= \frac{\lambda_{\text{obj}}}{N_{\text{conf}}} \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{obj}} \left(\text{IOU}_{\text{prediction}_{i,j}}^{\text{ground truth}_{i,j}} - \hat{C}_{i,j} \right)^2 \\ &\quad + \frac{\lambda_{\text{noobj}}}{N_{\text{conf}}} \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{noobj}} \left(0 - \hat{C}_{i,j} \right)^2 \end{aligned}$$

where:

- $N_{L^{\text{obj}}} = \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{obj}}$
- $N_{\text{conf}} = \sum_{i=0}^{S^2} \sum_{j=0}^B L_{i,j}^{\text{obj}} + L_{i,j}^{\text{noobj}} (1 - L_{i,j}^{\text{obj}})$
- $\text{prediction}_{i,j} = (\hat{x}_{i,j}, \hat{y}_{i,j}, \hat{w}_{i,j}, \hat{h}_{i,j})$
- $\text{ground truth}_{i,j} = (x_{i,j}, y_{i,j}, w_{i,j}, h_{i,j})$
- λ_{coord} , λ_{class} and λ_{noobj} are scalars to weight each loss function

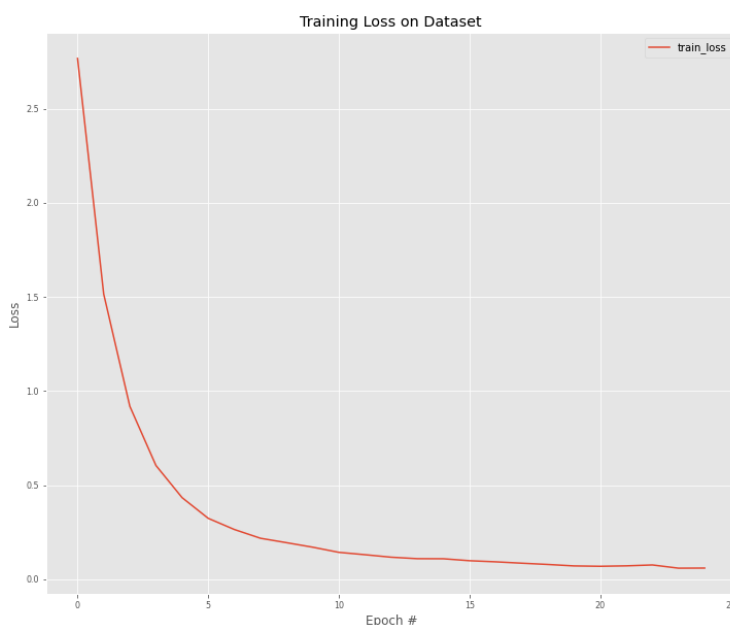
Here, $L_{i,j}^{\text{noobj}}$ and $L_{i,j}^{\text{obj}}$ are 0/1 indicator function such that:

$$\begin{aligned} L_{i,j}^{\text{obj}} &= \begin{cases} 1 & \text{if } C_{i,j} = 1 \\ 0 & \text{else} \end{cases} \\ L_{i,j}^{\text{noobj}} &= \begin{cases} 1 & \text{if } \max_{i',j'} \text{IOU}_{\text{prediction}_{i,j}}^{\text{ground truth}_{i',j'}} < 0.6 \text{ and } C_{i,j} = 0 \\ 0 & \text{else} \end{cases} \end{aligned}$$

تصویر ۱۲ - تابع خطا

آموزش مدل

برای آموزش مدل از وزن های pretrained شده yolo v2 استفاده می کنیم و شبکه را با Learning rate=0.5e-4 و Epochs=25 و Batch size=16 به همراه پنجاه و سه میلیون پارامتر قابل یادگیری آموزش می دهیم زمان تقریبی برای آموزش مدل پنج ساعت به طول انجامید و مقدار خطای مدل تا حد قابل قبولی کاهش پیدا کرد و به ۰.۰۵ رسیدیم که برای شبکه با این تعداد پارامتر بسیار مناسب است.



تصویر ۱۳ - نحوه کم شدن تابع خطا

تفاوت ها و شباهت ها با مقاله اصلی

مقاله اصلی	کار انجام شده
ورودی شبکه 416×416	ورودی شبکه 416×416
استفاده از RKCELM و AE برای استخراج ویژگی ها	استفاده از DCNN و Concatenate برای استخراج ویژگی ها
استفاده از K-means++ و $K=5$ برای خوشه بندی Anchor box	استفاده از K-means و $K=5$ برای خوشه بندی Anchor box
استفاده از معیار IOU برای انتخاب فاصله خوشه ها	استفاده از معیار IOU برای انتخاب فاصله خوشه ها
استفاده از تابع خطایی (Loss function) منحصر به فرد	استفاده از تابع خطایی (Loss function) متفاوت
استفاده از Bounding box prediction متفاوت	استفاده از Bounding box prediction مدل yolov2
سرعت بیشتر در آموزش شبکه	سرعت کمتر در آموزش شبکه
تشخیص بهتر اشیایی که هم پوشانی دارند	ضعف در تشخیص اشیای دارای هم پوشانی
مقدار خطای بیشتری دارد	بهبود ۰.۰۱ مقدار خطای مدل
تعداد پارامترهای کمتری دارد	تعداد پارامترهای بیشتری دارد



منابع

1. <https://www.sciencedirect.com/science/article/abs/pii/S1051200420301019>
2. <https://arxiv.org/abs/1506.02640>
3. <https://ieeexplore.ieee.org/document/8100173>
4. <https://arxiv.org/abs/1804.02767>
5. <https://machinelearningmastery.com/how-to-perform-object-detection-with-yolov3-in-keras/>
6. <https://pjreddie.com/darknet/yolo/>
7. <https://arxiv.org/abs/1808.02350v1>
8. <https://pylessons.com/YOLOv3-explained/>
9. <http://host.robots.ox.ac.uk/pascal/VOC/voc2012/>