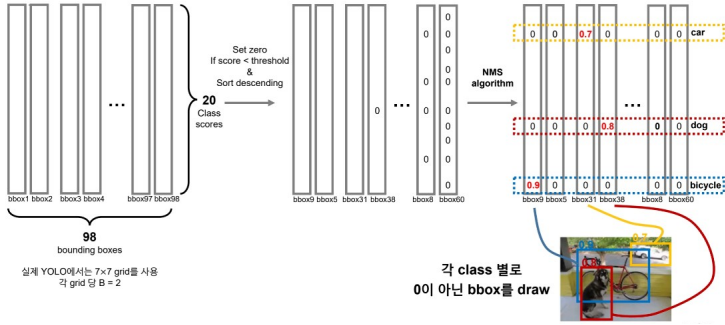


더 좋다. 그래서 threshold 값을 두어 그 값보다 크게 되면 같은 object를 Detection하고 있다고 보고 값을 0으로 만들어 준다. 이과정은 NMS(non-maximal suppression) algorithm이라고 한다. 그림으로 설명하자면 아래와 같다.



YOLO 모델을 training 시킬 때는 leaky ReLU를 사용하였고 learning rate도 적절히 epoch을 끊어서 조절해주었다. 또한, batch size를 64로 momentum을 0.9로 decay를 0.0005로 주었다. 그리고 overfitting을 막기 위해서 dropout과 data augmentation도 구현하였다.

5.2 Loss function

object detection task를 regression 문제로 전환하기 위해 가장 중요한 것이 loss function이다. 먼저 sum-squared error를 사용하였는데 localization과 classification error에 대한 동일한 가중치를 둔다는 점과 background에 속하는 grid cell의 경우 confidence score이 0이 되도록 학습하는데 이 과정에서 object를 포함하는 cell은 overpowering하게 만든다는 점이 문제점이 되었다. 그래서 아래와 같이 loss function을 새로 설계하였다.

$$\begin{aligned} \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} & \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} & \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} & \left(C_i - \hat{C}_i \right)^2 \\ + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} & \left(C_i - \hat{C}_i \right)^2 \\ + \sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum_{c \in \text{classes}} & \left(p_i(c) - \hat{p}_i(c) \right)^2 \end{aligned}$$

$\mathbb{1}_i^{obj}$ 은 i번째 cell에 object가 있는지를 나타낸다. $\mathbb{1}_{ij}^{obj}$ 은 i cell에서 j번째 bounding box predictor가 responsible한지를 나타낸다. 이 때 responsible이란 한 object를 누가 더 잘 detection하는지를 의미하는 표현이다. 여기서 w,h에 root를 사용한 것은 작은 object여도 큰 object와 별 차이없는 loss를 주기 위해서이다. 2번째 줄까지의 식은 모든 grid cell에서 예측한 B개의 bounding box의 좌표와 ground truth 사이의 coordinate loss이고 3,4번째 줄은 모든 grid cell에서 예측한 confidence score과 ground truth 사이의 confidence score loss이다. 5번째 줄은 conditional class probability 예측값과 ground truth사이의 conditional class probability loss를 의미한다.

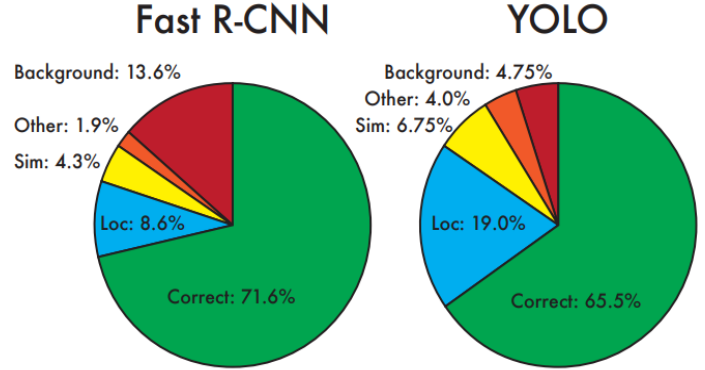
6 Experiment

Deformable parts models(DPM)은 sliding window 방식으로 정적인 feature추출과 region 분류 등을 위해 분리된 pipeline을 사용하는 모델로 YOLO가 더 빠르고 정확도도 높다. R-CNN의 YOLO을 비교하자면 각각의 grid cell에서 CNN을 이용해서 bounding box와 score을 뽑는다는 공통점이 있는데 localization과 classification의 pipeline 하나인지 두개인지에 따라서 one-stage, two-stage로 나뉜다. 아래 실험결과를 보면 알수 있다시피 YOLO 모델이 다른 모델에 비해 압도적으로 FPS가 높고 정확도도 R-CNN에 비해 많이 뒤지지는 않는다.

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [31]	2007	16.0	100
30Hz DPM [31]	2007	26.1	30
Fast YOLO	2007+2012	52.7	155
YOLO	2007+2012	63.4	45

Less Than Real-Time			
Fastest DPM [38]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[28]	2007+2012	73.2	7
Faster R-CNN ZF [28]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

아래 그림은 Fast R-CNN과 YOLO을 비교한 것으로 R-CNN은 background error의 아쉬움이 있고 YOLO의 경우에는 localization의 아쉬움이 있다.

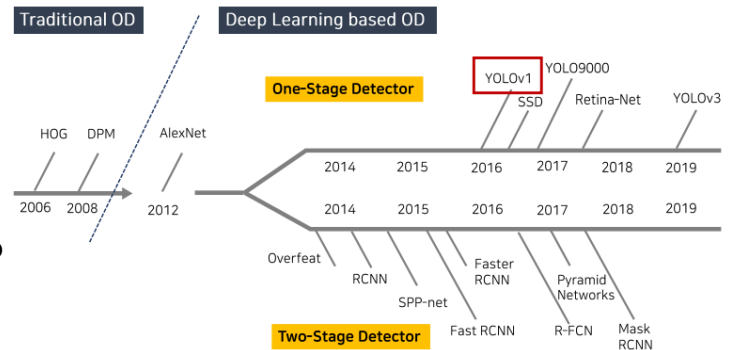


따라서 두 모델을 ensemble하게 되면 각각의 단점을 보완하여 Fast R-CNN보다 mAP가 3.2%가 증가한 75%의 정확도를 얻을 수 있었다. 그 외 결과에도 YOLO model의 경우에는 train data set과 distribution이 다른 dataset에서도 정확도 감소 정도 다른 모델에 비해 작았다. 이로 볼 때 YOLO가 전체 이미지를 받아서 학습을 시키기 때문에 더 global한 feature를 학습한 것으로 추측한다. artwork dataset의 경우에는 다른 모델에 비해 압도적으로 AP가 더 높았다.

	VOC 2007	Picasso		People-Art
	AP	AP	Best F_1	AP
YOLO	59.2	53.3	0.590	45
R-CNN	54.2	10.4	0.226	26
DPM	43.2	37.8	0.458	32
Poselets [2]	36.5	17.8	0.271	
D&T [4]	-	1.9	0.051	

7 Conclusion

이 논문에서는 object detection model로써 paradigm을 바꾼 one-stage detector YOLO을 제시함으로써 object detection은 두 가지 model로 발전하고 있다.[1] YOLO는 비교적 높은 정확도를 가지면서 real-time video 분석에 사용할 수 있고 심지어 natural data의 train dataset과 다른 distribution에서도 높은 정확도를 추출하기 때문에 매우 유용하다. 또한 ene-to-end training이 가능하여 다루기도 쉽다.



- [1] 황도식 교수님 기초인공지능 수업 강의 노트.
- [2] <https://youtu.be/o78v3kwbrbk>.
- [3] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.