

Assignment 3A: Instance Segmentation with MaskRCNN

Venkata Sai Nikhil, Thodupunuri

62716136

1 Backbone architecture

The entire network is implemented using 4 network classes. Here are the name conventions used.

- **Base Network** : This class implements the backbone architecture.
- **Intermediate Network**: This class implements the intermediate layer used by both the regressor and classifier component. This class depends on the base network to generate the output. Base class is passed as a constructor parameter.
- **Regression Network**: This class implements the regression specific layers. this class depends on intermediate object (passed as a constructor parameter).
- **Classification Network**: This class implements the classification specific layers. this class depends on intermediate object (passed as a constructor parameter).

Code:

```
1 import torch
2 import torch.nn as nn
3
4
5 class BaseNetwork(torch.nn.Module):
6     def __init__(self):
7         super().__init__()
8         self.conv1 = nn.Conv2d(3, 16, 5, padding=2)
9
10        self.conv1_seq = nn.Sequential(self.conv1, nn.BatchNorm2d(num_features = 16),
11                                      nn.ReLU())
12
13        self.pool1 = nn.MaxPool2d(2, stride=2, padding = 0)
14
15        self.conv2 = nn.Conv2d(16, 32, 5, padding=2)
16
17        self.conv2_seq = nn.Sequential(self.conv2, nn.BatchNorm2d(num_features = 32),
18                                      nn.ReLU())
19
20        self.pool2 = nn.MaxPool2d(2, stride=2, padding = 0)
21
22
23        self.conv3 = nn.Conv2d(32, 64, 5, padding=2)
24
25        self.conv3_seq = nn.Sequential(self.conv3, nn.BatchNorm2d(num_features = 64),
26                                      nn.ReLU())
27
28        self.pool3 = nn.MaxPool2d(2, stride=2, padding = 0)
29
30
31        self.conv4 = nn.Conv2d(64, 128, 5, padding=2)
```

```

34     self.conv4_seq = nn.Sequential(self.conv4, nn.BatchNorm2d(num_features = 128),
35                                    nn.ReLU())
36
37     self.pool4 = nn.MaxPool2d(2, stride=2, padding = 0)
38
39     self.conv5 = nn.Conv2d(128, 256, 5, padding=2)
40
41     self.conv5_seq = nn.Sequential(self.conv5, nn.BatchNorm2d(num_features = 256),
42                                   nn.ReLU())
43
44
45     self.base_network = nn.Sequential(self.conv1_seq, self.pool1, self.conv2_seq,
46                                      self.pool2, self.conv3_seq, self.pool3,
47                                      self.conv4_seq, self.pool4, self.conv5_seq)
48
49 def forward(self, X):
50     return self.base_network.forward(X)
51
52
53 class Intermediate_Network(torch.nn.Module):
54     def __init__(self, base_model):
55         super().__init__()
56         self.base_model = base_model
57         self.layer1 = nn.Conv2d(256, 256, 3, padding=1)
58         self.layer1_seq = nn.Sequential(self.layer1, nn.BatchNorm2d(num_features = 256),
59                                       nn.ReLU())
60
61     @staticmethod
62     def from_static_init(base_model, layer1_seq):
63         model = Intermediate_Network(base_model)
64         model.layer1_seq = layer1_seq
65         return model
66
67     def forward(self, X):
68         base_output = self.base_model.forward(X)
69         return self.layer1_seq.forward(base_output)
70
71
72 class ClassificationNetwork(torch.nn.Module):
73     def __init__(self, intermediate_model):
74         super().__init__()
75         self.intermediate_model = intermediate_model
76         self.layer2 = nn.Conv2d(256, 1, 1)
77         self.sigmoid = nn.Sigmoid()
78         self.classification_network = nn.Sequential(self.layer2, self.sigmoid)
79
80     @staticmethod
81     def from_static_init(intermediate_model, classification_network):
82         model = ClassificationNetwork(intermediate_model)
83         model.classification_network = classification_network
84         return model
85
86     def forward(self, X):
87         intermediate_output = self.intermediate_model.forward(X)
88         return self.classification_network.forward(intermediate_output)
89
90
91 class RegressionNetwork(torch.nn.Module):
92     def __init__(self, intermediate_model):
93         super().__init__()
94         self.intermediate_model = intermediate_model
95         self.layer2 = nn.Conv2d(256, 4, 1)
96         self.regression_network = nn.Sequential(self.layer2)
97
98     @staticmethod
99     def from_static_init(intermediate_model, regression_network):
100        model = RegressionNetwork(intermediate_model)
101        model.regression_network = regression_network

```

```

102     return model
103
104     def forward(self, X):
105         intermediate_output = self.intermediate_model.forward(X)
106         return self.regression_network.forward(intermediate_output)
107
108
109
110 ## Sample Initialization.
111 base_model = BaseNetwork()
112 intermediate_model = Intermediate_Network(base_model)
113
114 c_model = ClassificationNetwork(intermediate_model)
115 r_model = RegressionNetwork(intermediate_model)

```

2 Histogram of the area and aspect ratios

Scale: Scale is termed as area.

Mode changes if the number of bins are changed. Mean and variance of bins are analysed to get the perspective divergence of bin.

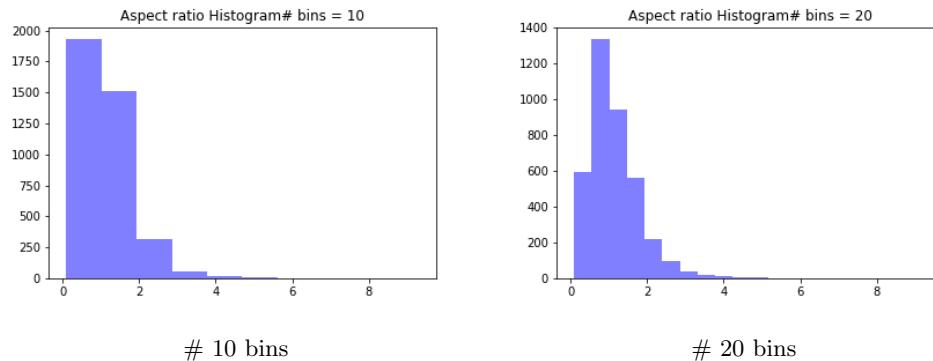
Aspect ratio

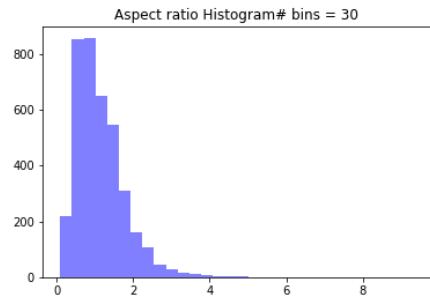
Experimented with 10 different bin sizes.

Mean and Variance of Mode bins:

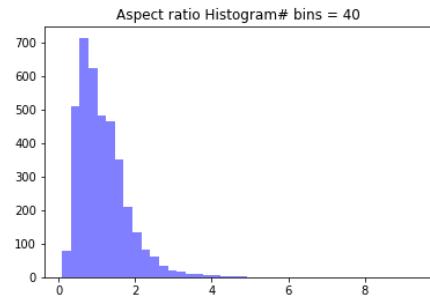
Bin Size	Bin Number	Bin Count	Mean	Variance	Percentage
10	1	1931	0.6686	0.0397	50.2%
20	2	1339	0.7753	0.016	34.84%
30	3	858	0.855	0.0075	22.32%
40	3	715	0.670	0.0044	18.6%
50	4	555	0.739	0.0027	14.44%
60	4	481	0.6329	0.0019	12.51%
70	5	411	0.6856	0.0014	10.69%
80	6	359	0.7276	0.0011	9.34%
90	6	321	0.658	0.0008	8.35%
100	6	293	0.602	0.0007	7.62%

Plots:

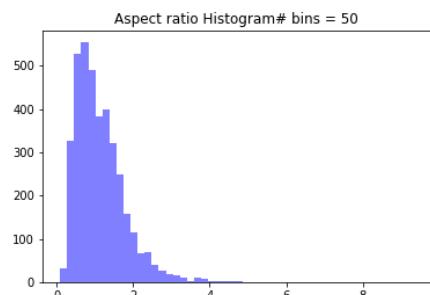




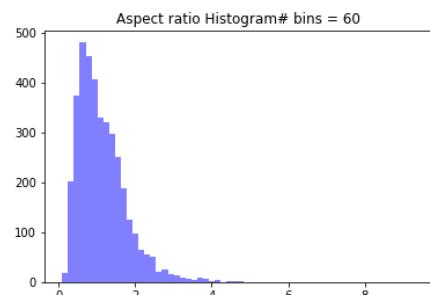
30 bins



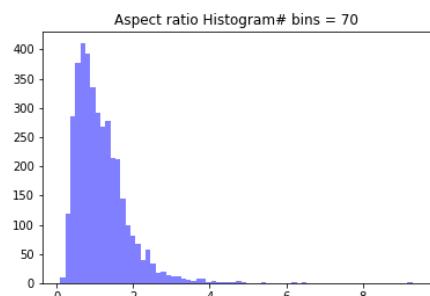
40 bins



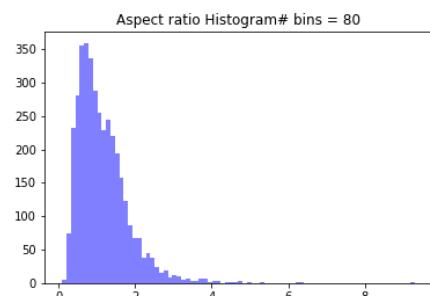
50 bins



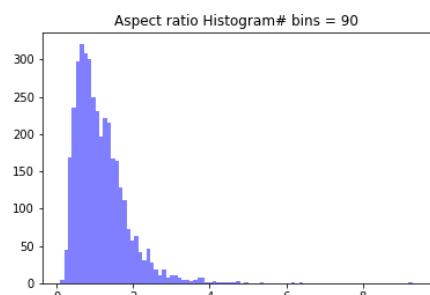
60 bins



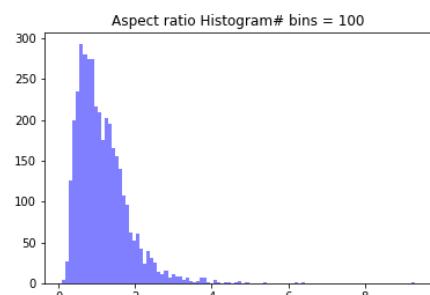
70 bins



80 bins



90 bins



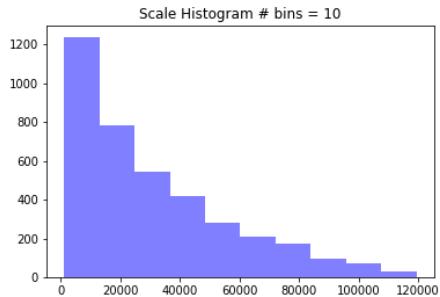
100 bins

Area bins

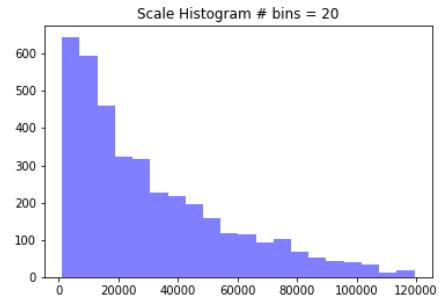
Mean and Variance of Mode bins:

Bin Size	Bin Number	Bin Count	Mean	Variance	Percentage
10	1	1238	6847	10768083	32.21%
20	1	644	4107	2476964	16.75%
30	1	440	3221	1003043	11.44%
40	1	325	2781	594366	8.45%
50	2	276	4541	464472	7.18%
60	2	249	3970	284233	6.47%
70	2	227	3622	231787	5.9%
80	2	193	3323	178471	5.02%
90	2	168	3064	142375	4.37%
100	3	163	4049	114948	4.24%

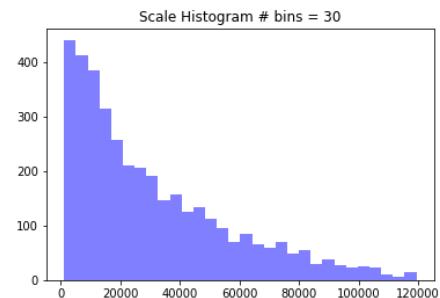
Plots:



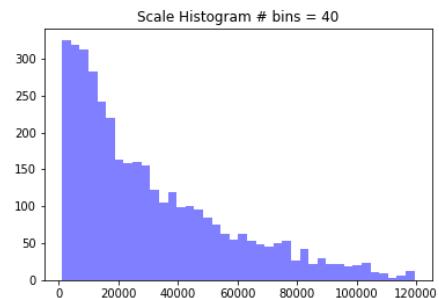
10 bins



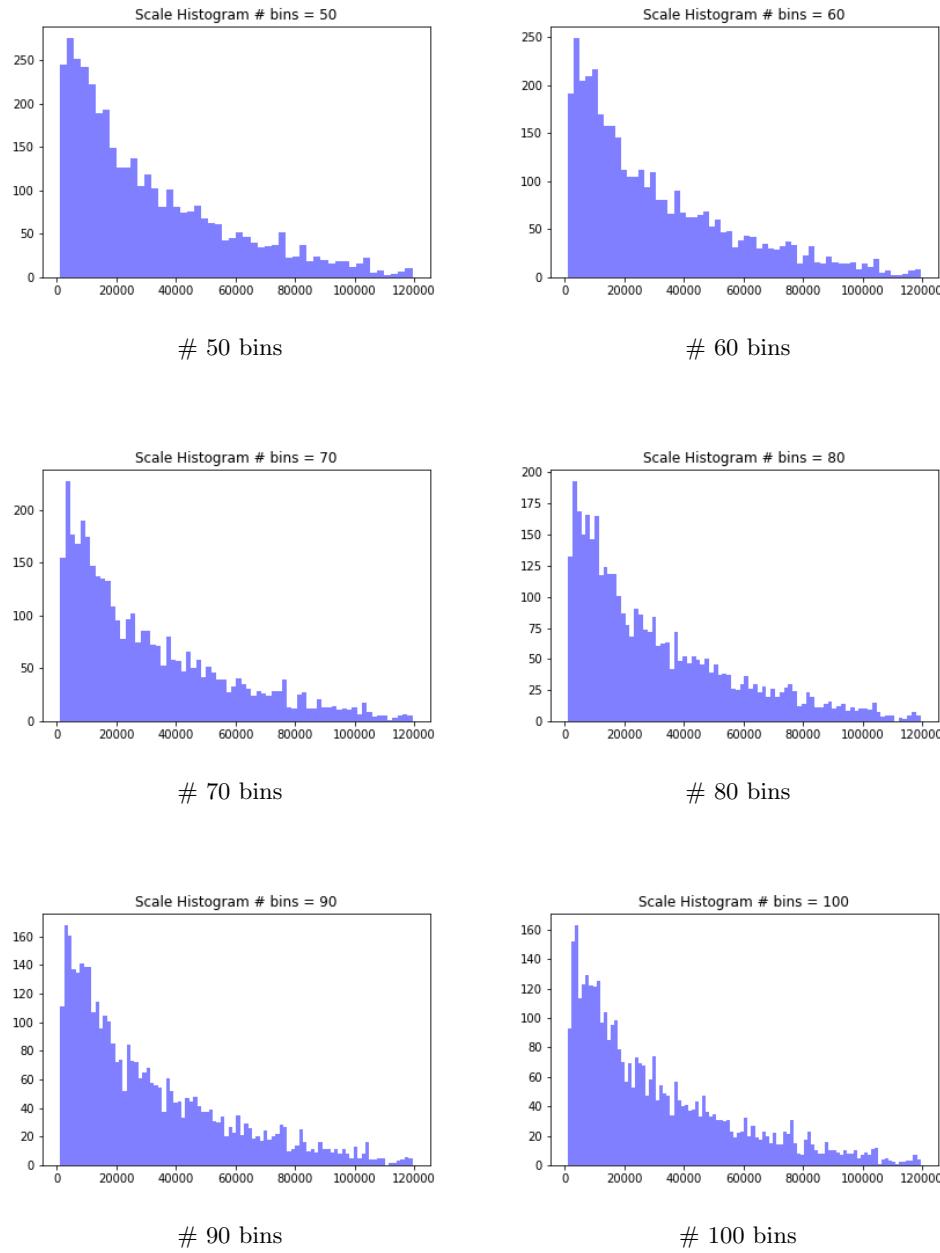
20 bins



30 bins



40 bins



Mode Values

Chosen mode values for

Area	$4107 \sim 64^2$
Aspect Ratio	0.66

Effects

Pros of choosing mode

- If only one anchor is used, Single mode helps us to generalize max number of samples. Mode helps us to repair some of the effects of single anchor.
- Even with multiple anchors. we can choose the multiple max bins and use those scale values.

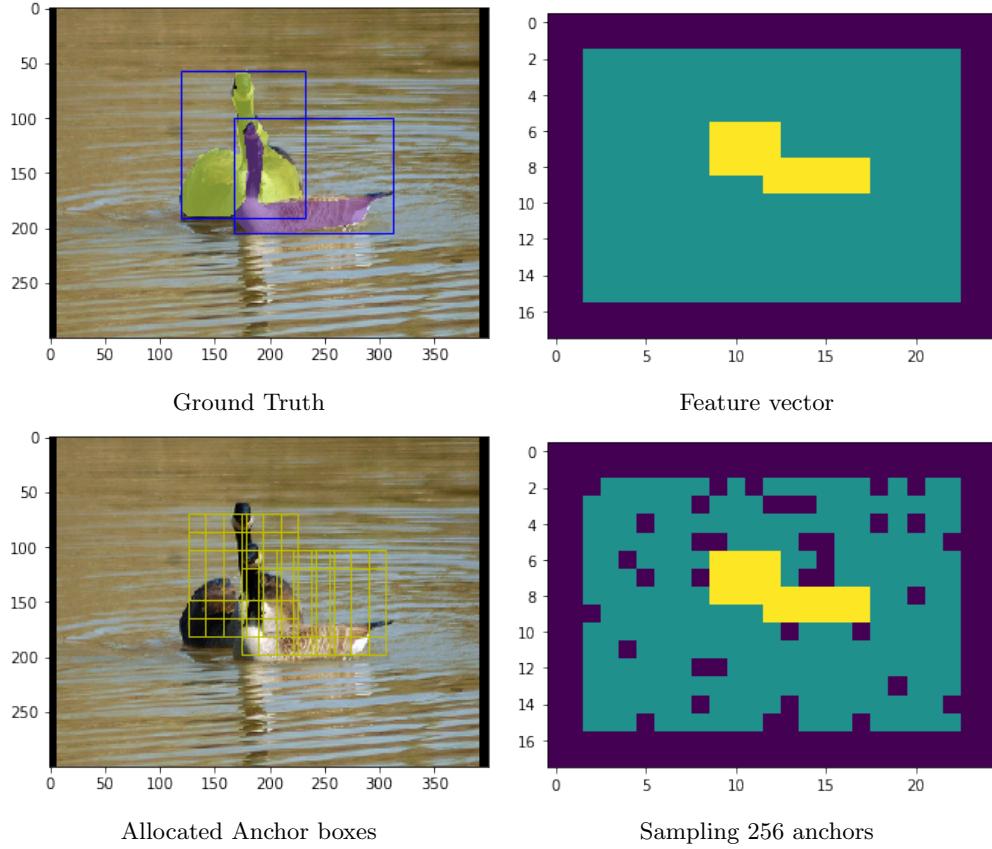
Cons:

- If we look at the percentage from the above tables, Mode only covers up-to 16.75% for area and 18.6% for aspect ratio. Rest of the ranges will be ignored.
- Also the variance in mode bin is to be considered. If the bin has high variance, then mode doesn't help. (Inspect histograms with different bin sizes, to choose the mode wisely).
- If the test data is not known prior, then these generalizations will not hold and will result in poor accuracies.
- Prefer using multiple anchors over mode if possible.

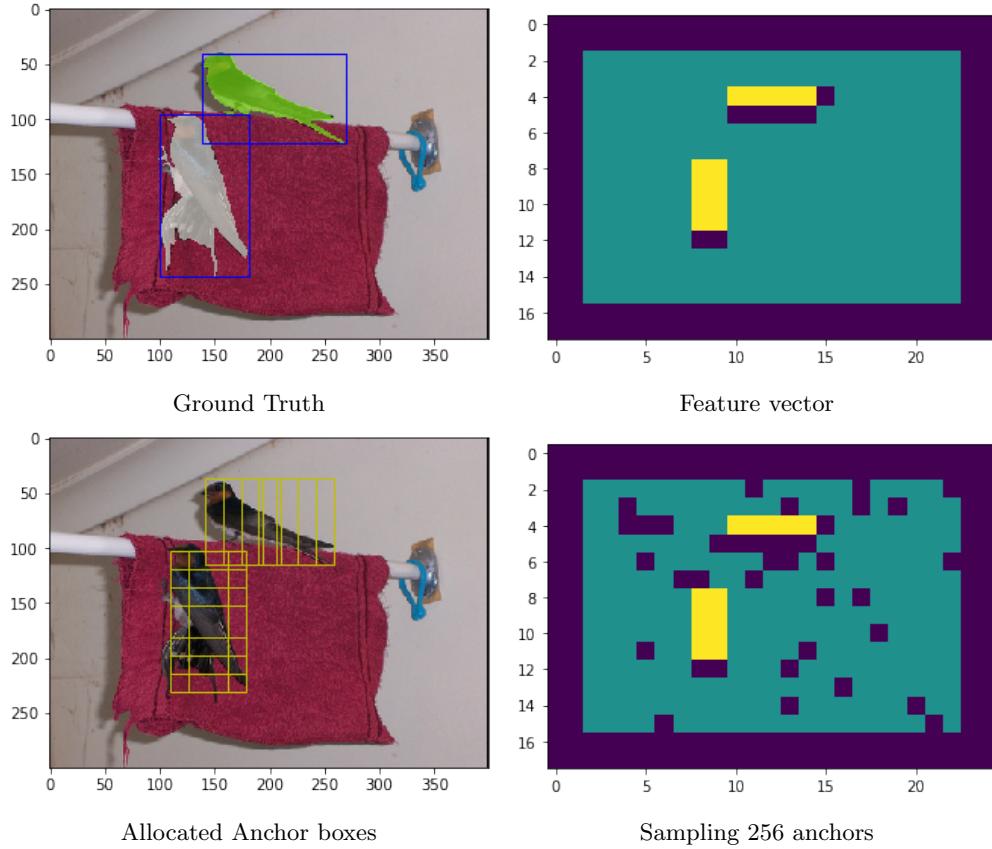
3 Ground truth for RPN training

Feature Vector Color Code : Yellow - positive boxes, Green - Negative boxes, Purple - Ignore boxes.

Sample 1



Sample 2

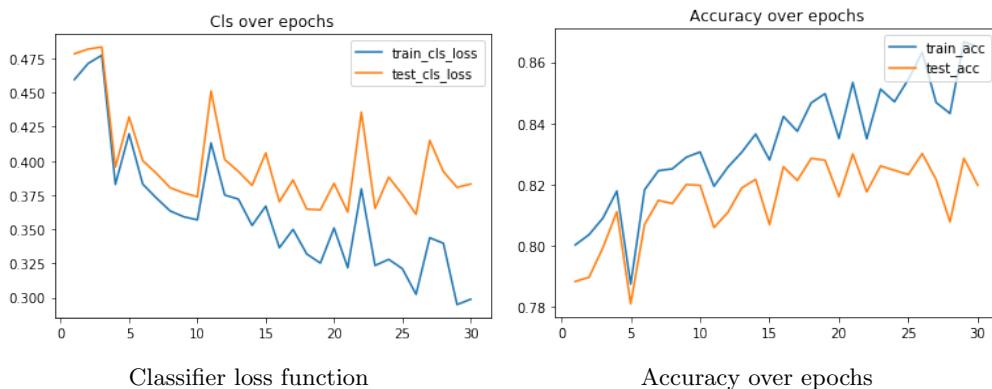


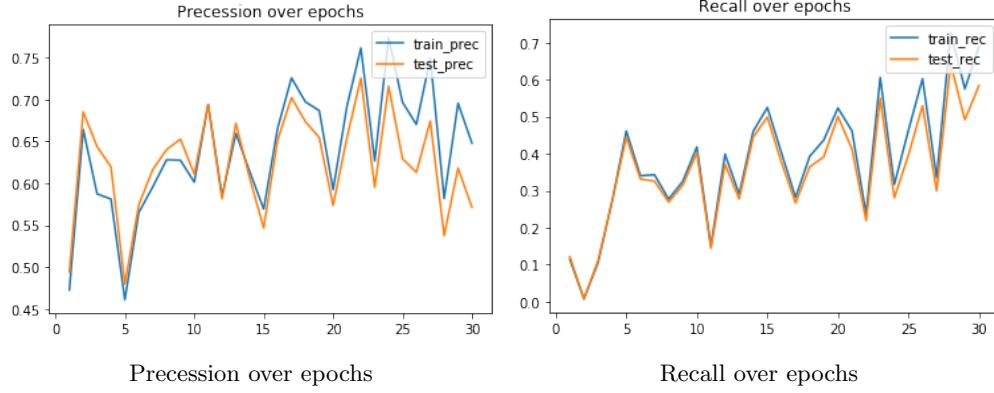
4 RPN classification branch

Sampling strategy: Sampled 256 anchors for each image.

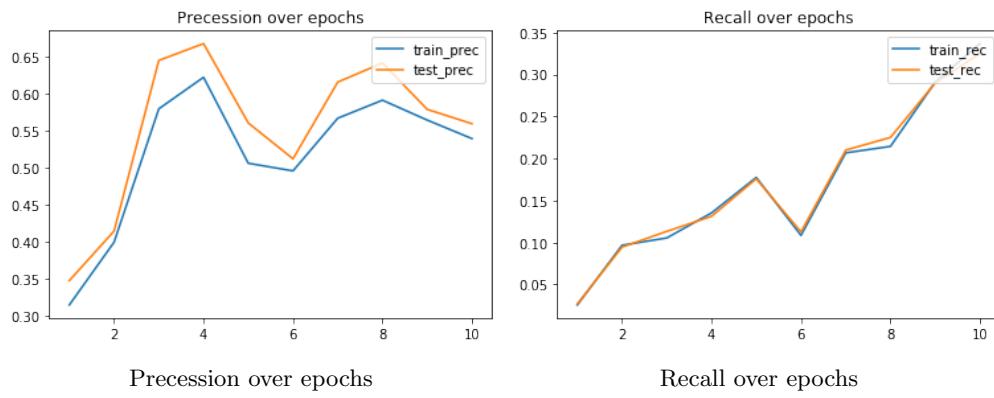
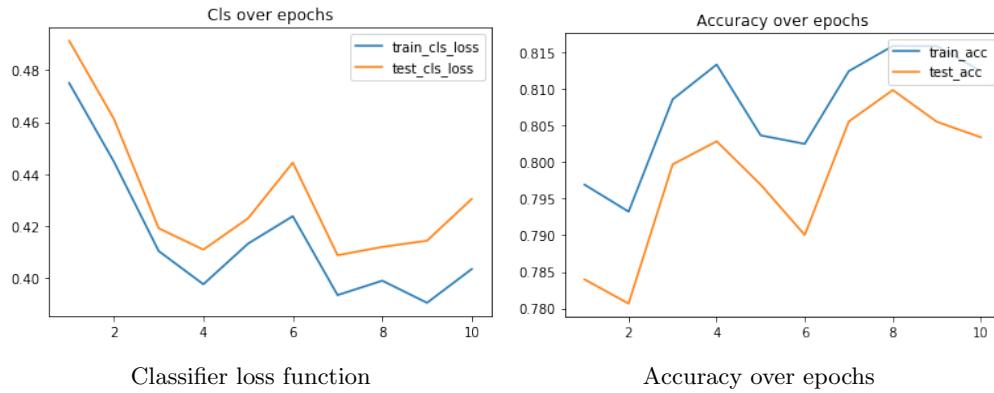
Plots

Batch size: 30, **LR:** 0.001. Each image in the batch has equal number of positive and negative anchors. From each image around 256 anchors are sampled.





Batch size: 1, **LR:** 0.00001. Each image has equal number of positive and negative anchors. From each image around 100 anchors are sampled.



5 Receptive field

Logic for calculating receptive field.

- The logic assumes that the filter is symmetrical.
- Each layer (CONV, POOL) can be represented using 2 values (K_t, S_t). where K_t is the kernel size, S_t is the stride.
- Let ES_t be the effective stride after t layers. R_t be the receptive field after t layers.

- Relation:

$$\begin{aligned}
 - ES_t &= S_t * ES_{t-1} \\
 - R_t &= (K_t - 1) * ES_{t-1} + R_{t-1},
 \end{aligned}$$

- Base case:

$$- R_0 = 1, ES_0 = 1$$

Base Network Receptive Calculation

Layer Number	Layer Type	Kernel Size	Stride	Effective Stride	Receptive field
1	CONV	5x5	1	1	5x5
2	POOL	2x2	2	2	6x6
3	CONV	5x5	1	2	14x14
4	POOL	2x2	2	4	16x16
5	CONV	5x5	1	4	32x32
6	POOL	2x2	2	8	36x36
7	CONV	5x5	1	8	68x68
8	POOL	2x2	2	16	76x76
9	CONV	5x5	1	16	140x140

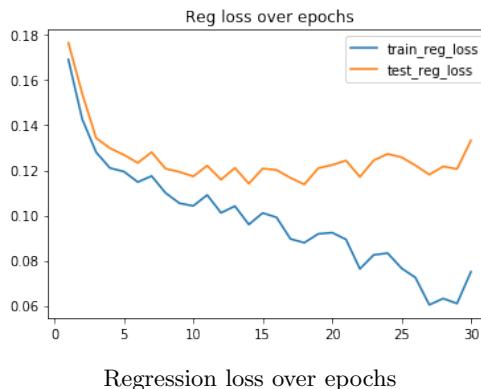
Intermediate Network Receptive Calculation

Layer Number	Layer Type	Kernel Size	Stride	Effective Stride	Receptive field
1	CONV	3x3	1	16	172x172

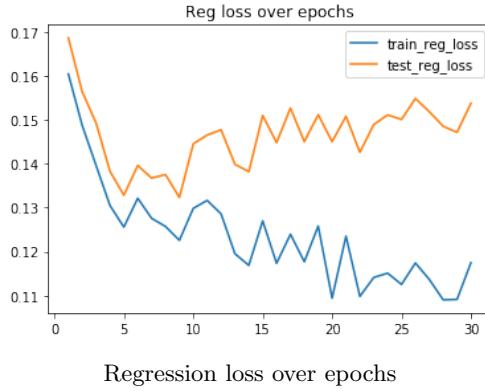
The anchor box should be inside of the receptive field (Size of anchor box should be less than receptive field both in height and width). If the receptive field is smaller, then the network will not have enough information(responses is limited to only receptive field) to predict boxes (which are larger than the receptive field).

6 RPN regression branch

Batch size: 30, **LR:** 0.001.



Batch size: 1, **LR:** 0.00001.



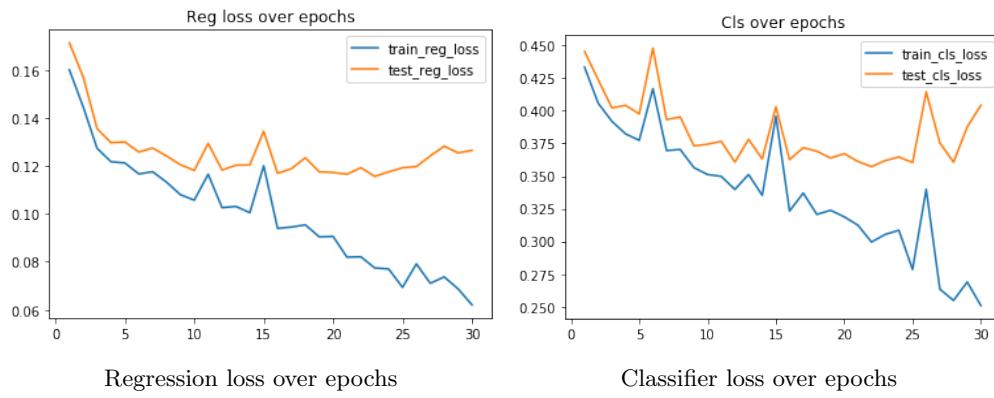
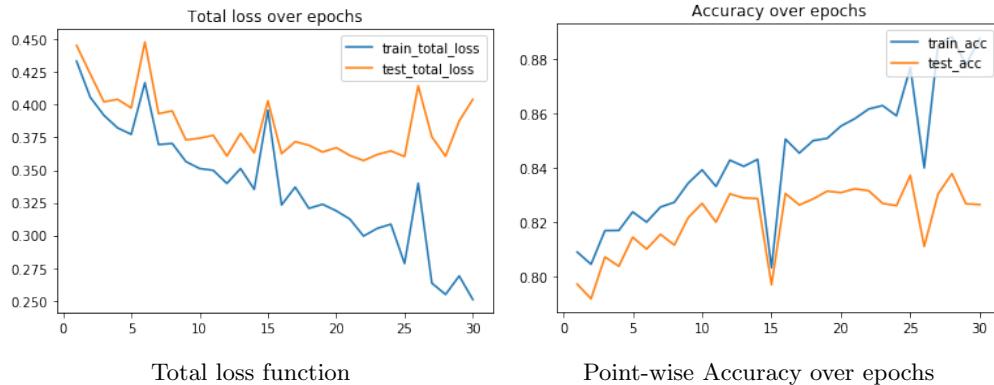
Note: For only regression training, all anchors are used (sampling is not done).

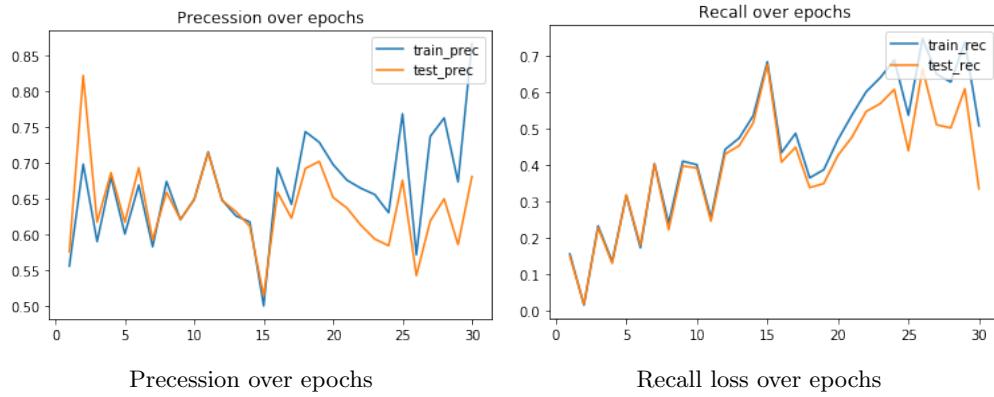
7 Combined Training

Sampling strategy: Sampled 256 anchors for each image.

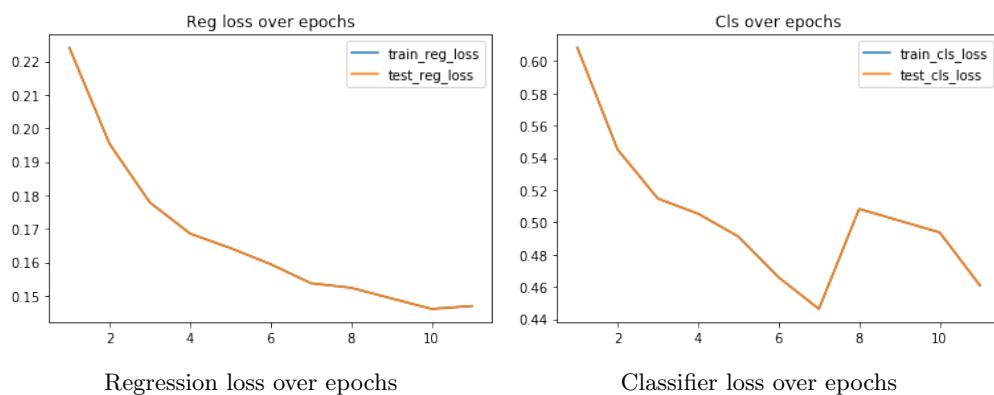
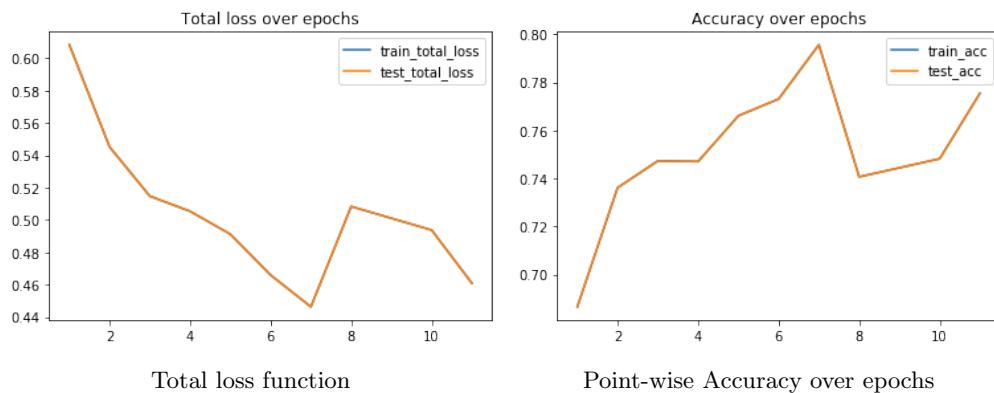
Plots

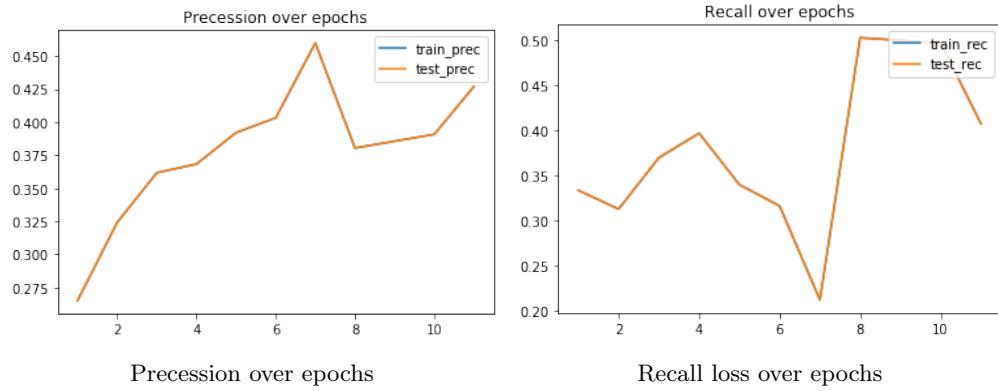
Batch size: 30, **LR:** 0.001. Each image in the batch has equal number of positive and negative anchors. From each image around 256 anchors are sampled.





Batch size: 1, **LR:** 0.000001. Each image has equal number of positive and negative anchors. From each image around 100 anchors are sampled.

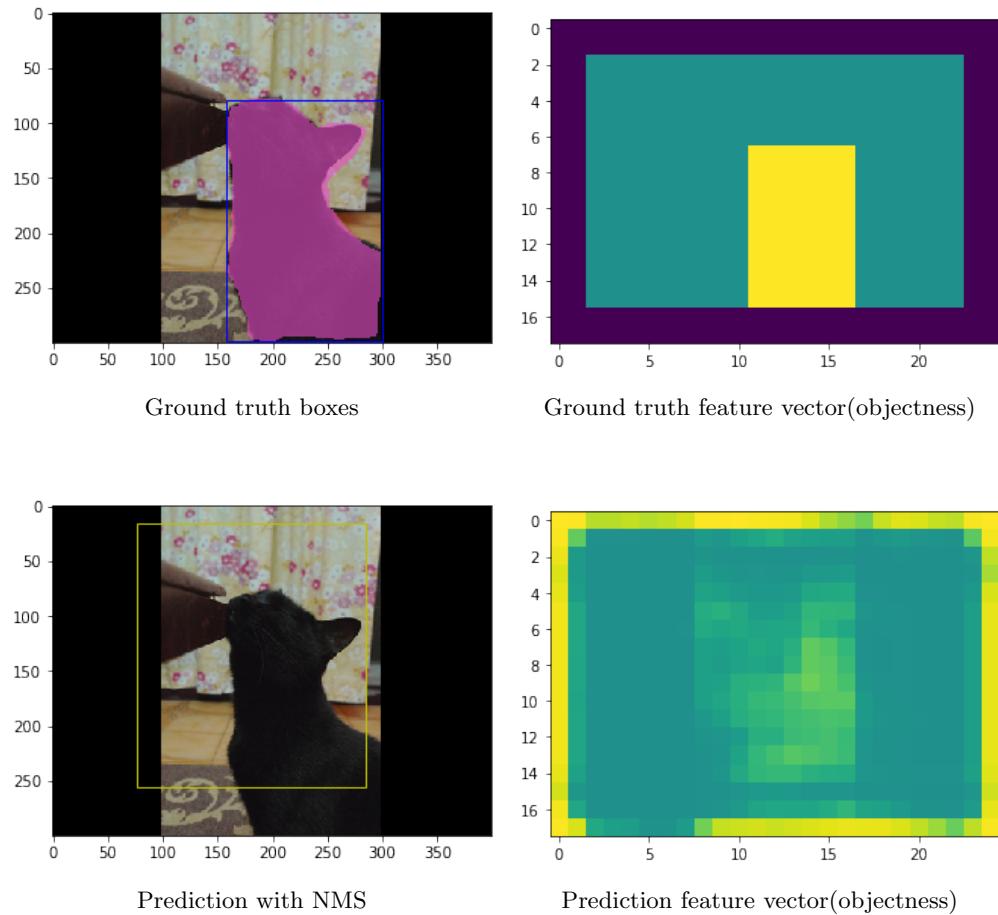


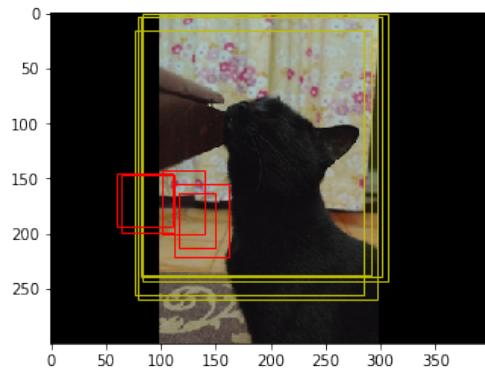


Visualizing test images

Color scheme: For prediction, Yellow boxes are predicted positive boxes. Red are predicted negative boxes.

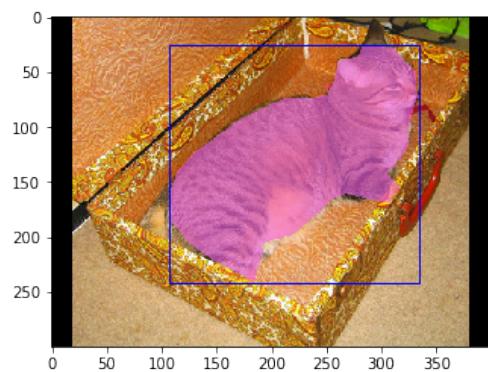
Sample 1



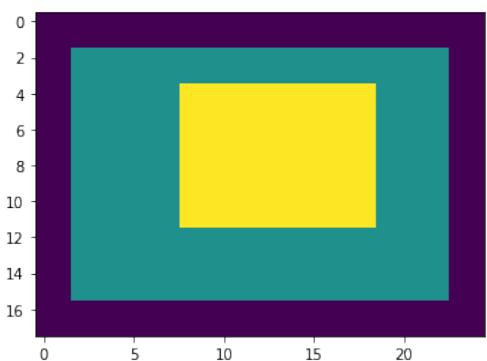


Prediction with negative boxes

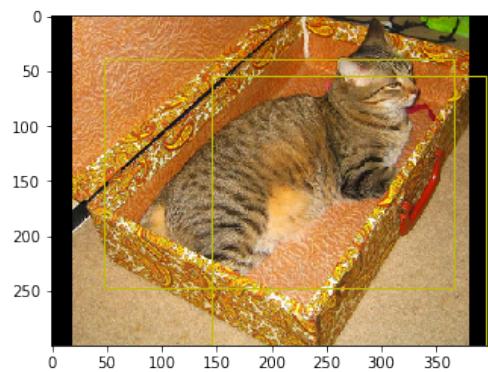
Sample 2



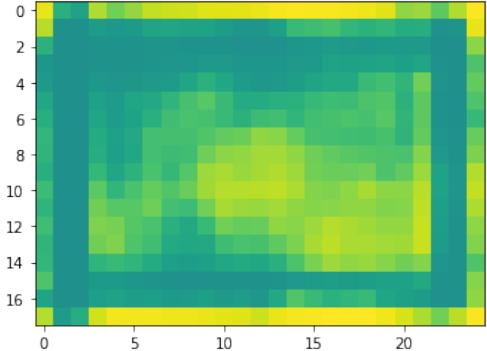
Ground truth boxes



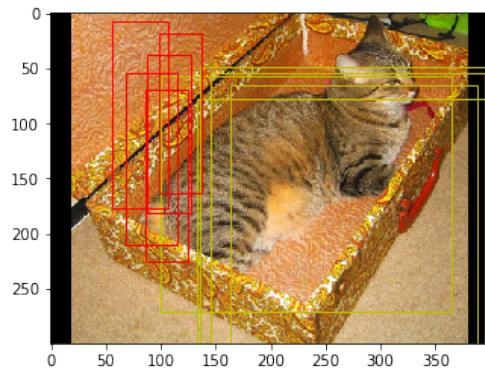
Ground truth feature vector(objectness)



Prediction with NMS

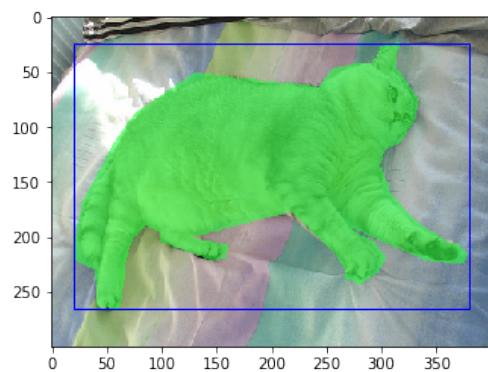


Prediction feature vector(objectness)

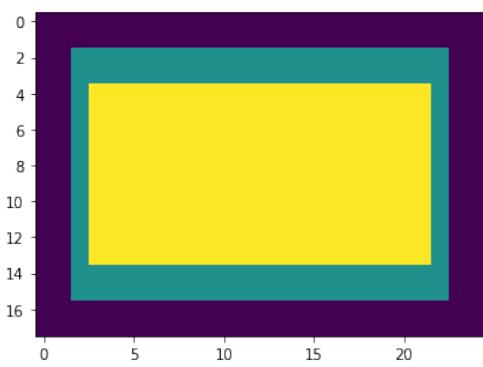


Prediction with negative boxes

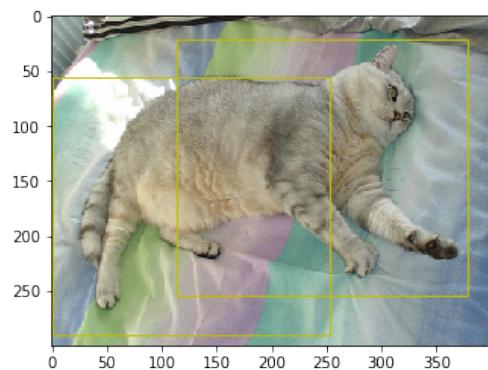
Sample 3



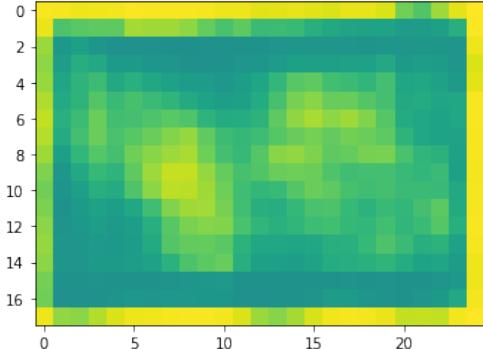
Ground truth boxes



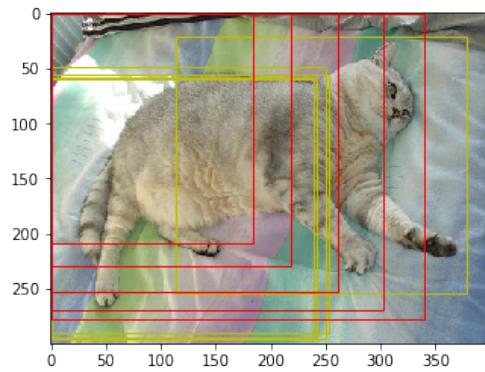
Ground truth feature vector(objectness)



Prediction with NMS

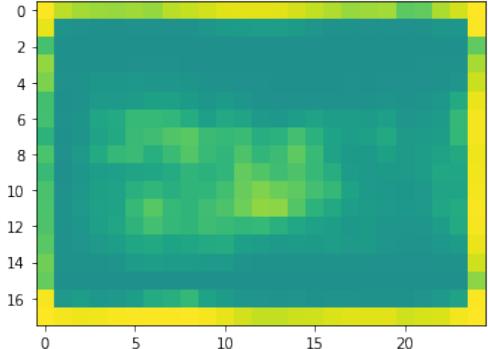
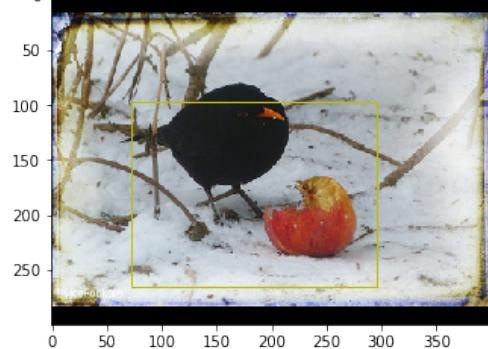
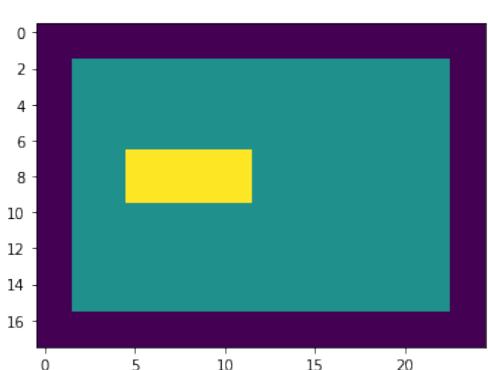
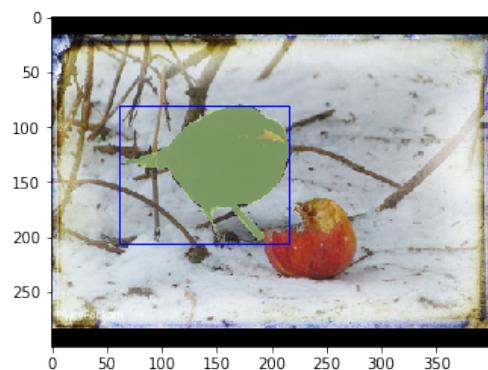


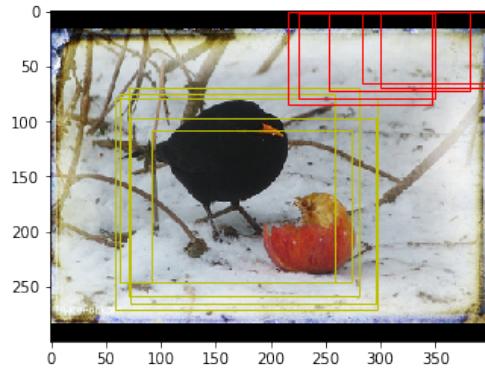
Prediction feature vector(objectness)



Prediction with negative boxes

Sample 4

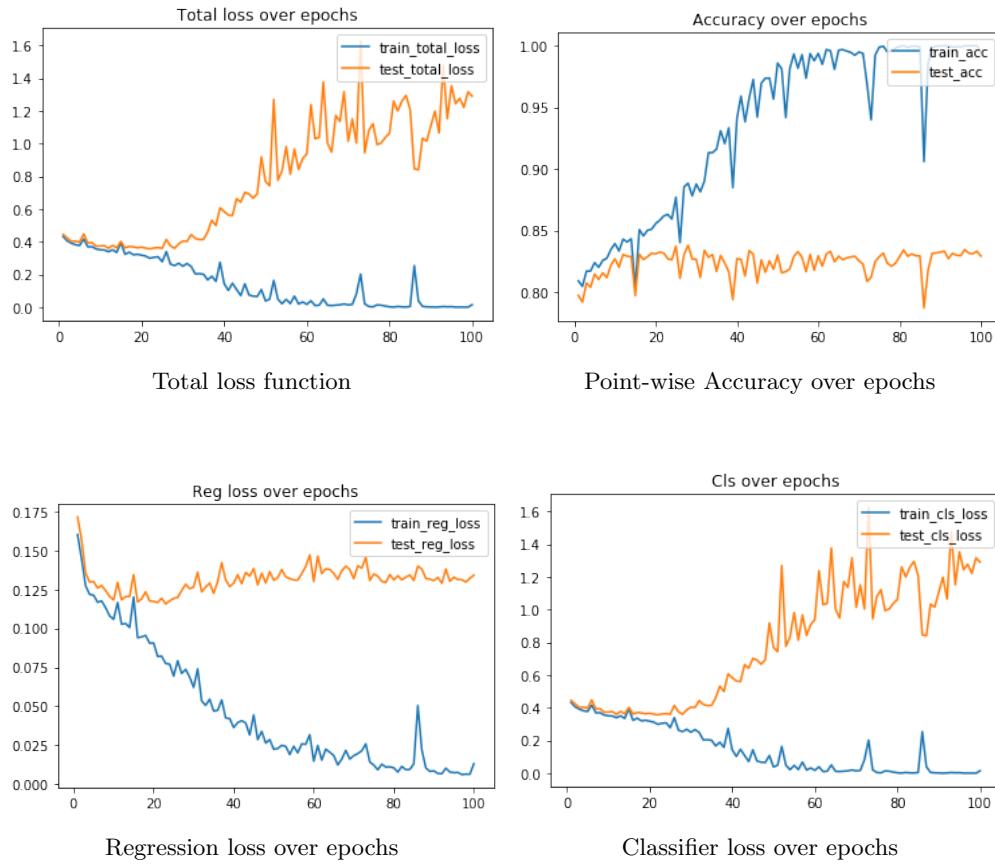


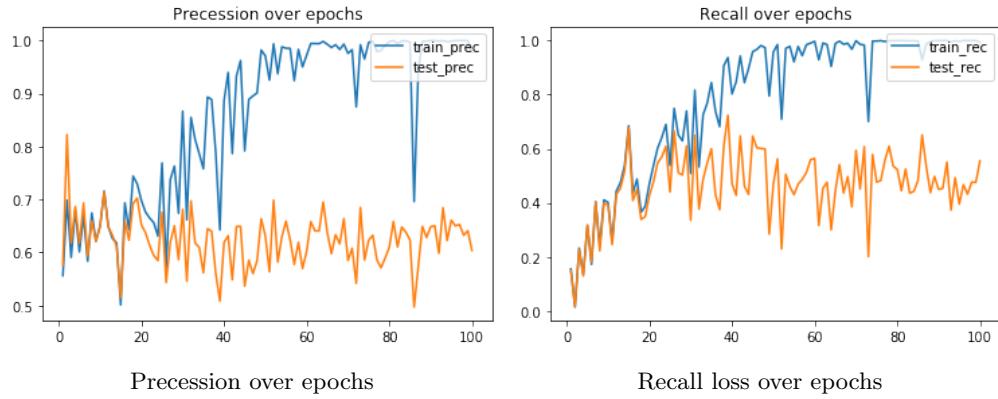


Prediction with negative boxes

Over-fitting example

Training for long epochs, resulted in over-fitting. After 35 epochs, the test loss stopped decreasing and started increasing.





8 Decoding

Hyper parameter = 2 for the following decoding.

