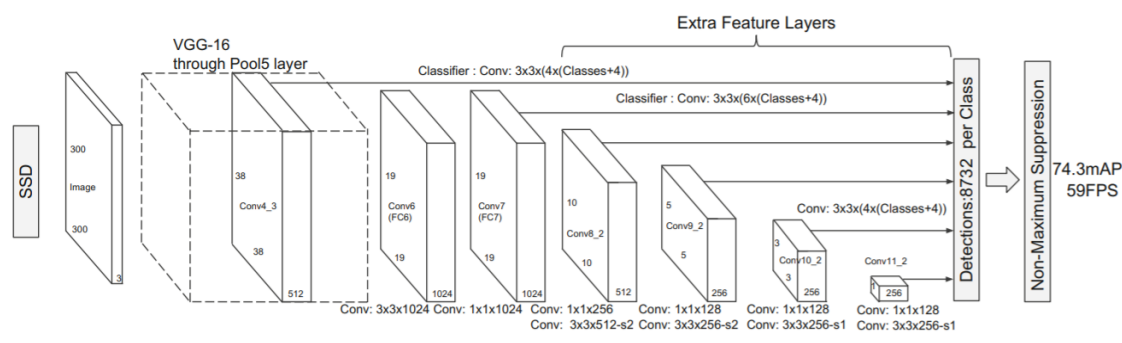
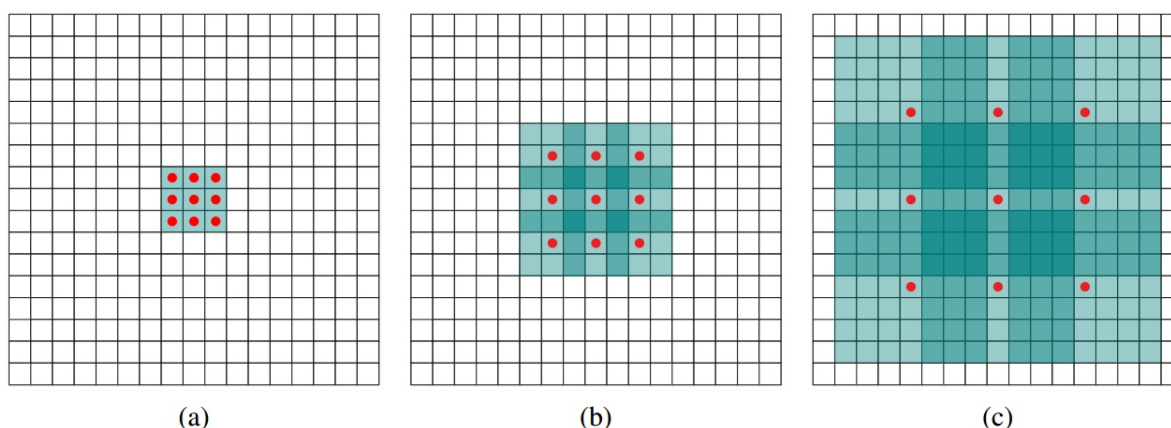


SSD:single shot MultiBox Detector(单发多框检测器)

SSD通过在每张特征图的位置，将输出的bbox离散成一系列默认的不同ratio和scale方框。在预测中，网络产生每个物体种类存在的可能性并且对方框作出调整来更好的适配物体的形状。另外，网络整合了拥有这不同的空间分辨率的特征图来更好地预测不同尺寸的物体，SSD 相对于需要的方法简单，通过 object prosal，因为它完全消除了生成的建议以及随后的像素或特征重采样阶段并封装所有计算过程在单个网络中进行计算。



空洞卷积



图像分割中一般按照传统的CNN中，利用CNN不断的Pooling 后降低对应的size最后经过转置卷积进行上采样，回到原始的图片的size的大小，b图是kernel_size=(3×3)dilation=2对应的实际上kernel为7×7，除了几个对应的参数权重不为-其余权重全为0。

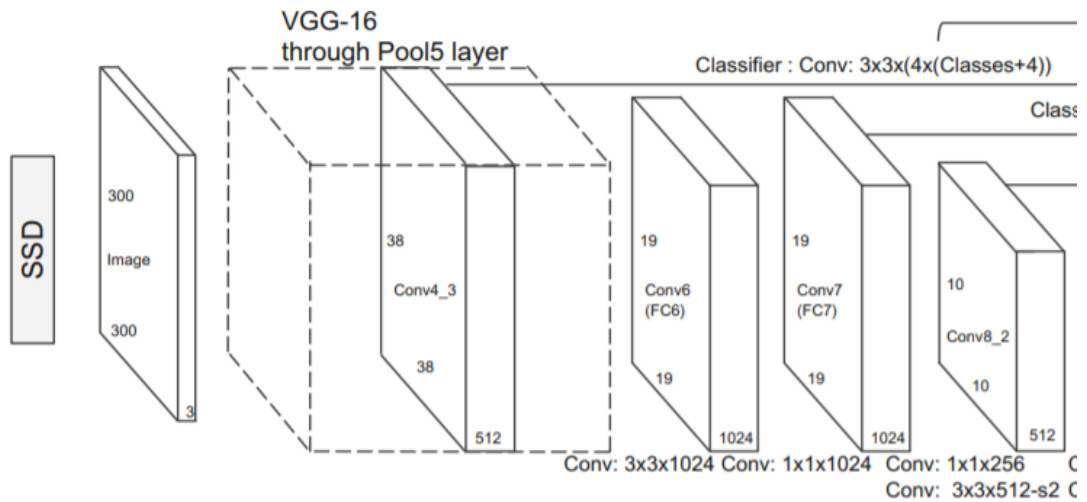
dilated的好处是不做pooling损失信息的情况下，加大了感受野，让每个卷积输出都包含较大范围的信息。在图像需要全局信息或者语音文本需要较长的sequence信息依赖的问题中，都能很好的应用dilated conv，dilated conv结构，可以更形象的了解dilated conv本身。

整体的网络架构在这里

- backbone:VGG-16 ,5和pooling layer conv5之后pooling的stride=1，也就是可以保存特征图的大小，Conv6使用了dilation=6的空洞卷积，padding=6也增加感受野之后保持特征图的尺寸保持不变

输入图片的尺寸为300 × 300,输出之后的为300VGG16的前13个卷积层进行对应的处理

conv6 3*3 1024 conv 7 3*3 1024 conv8_2 3*3 512 conv8_1 1*1 256 conv9_2 3*3 256 对应的channel



回到这里，让我们重新计算一下

input_size=(300,300)

conv4_3没有经过池化层所以输出为(38,38),3个池化层 2^3

出VGG的基本框架之后，总共是4个缩小尺寸的kernel 2^4 输出的size为(19,19),接下来就是正常的池化和和maxpooling改变对应的channels,

回到对应的代码来进行解释

```

1  """
2  input-> conv1_2:3x3x64,relu,3x3x64,relu -> pool1 ->
3          conv2_2:3x3x128,relu,3x3x128,relu -> pool2 ->
4          conv3_3:3x3x256,relu,3x3x256,relu,3x3x256,relu -> pool3 ->
5          conv4_3:3x3x512,relu,3x3x512,relu,3x3x512,relu -> pool4 ->
6          conv5_3:3x3x512,relu,3x3x512,relu,3x3x512,relu -> pool5 ->
7          conv6:3x3x1024 atrous, relu->
8          conv7:1x1x1024, relu
9          vgg_base = {
10             '300': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'C', 512, 512,
11                    512, 'M',
12                    512, 512, 512],
13             '512': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'C', 512, 512,
14                    512, 'M',
15                    512, 512, 512],
16             }
17 """
18 # 如果是M则添加最大pooling层来进行高宽减半的方法，默认的Batch为false,如果为True在
19 # 卷积层输出的时候跟着BN层和Relu，激活函数，in
20 # inplace=True用来减少内存，in_channels=v支持后边的网络
21 def add_vgg(cfg, batch_norm=False):
22     layers = []
23     in_channels = 3
24     for v in cfg:
25         if v == 'M':
26             layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
27         elif v == 'C':
28             layers += [nn.MaxPool2d(kernel_size=2, stride=2,
29                                    ceil_mode=True)]
30         else:
31             conv2d = nn.Conv2d(in_channels, v, kernel_size=3, padding=1)
32             if batch_norm:

```

```

29         layers += [conv2d, nn.BatchNorm2d(v),
nn.ReLU(inplace=True)]
30     else:
31         layers += [conv2d, nn.ReLU(inplace=True)]
32     in_channels = v
33     pool5 = nn.MaxPool2d(kernel_size=3, stride=1, padding=1)#不进行下采样
    的处理
34     conv6 = nn.Conv2d(512, 1024, kernel_size=3, padding=6, dilation=6)#
    加上的conv6采用扩展卷积或带孔卷积 (Dilation Conv)，其在不增加参数与模型复杂度的条
    件下指数级扩大卷积的视野
35     #, 采用 3x3 大小但dilation rate=6的扩展卷积。conv7采用1x1的卷积，
36     conv7 = nn.Conv2d(1024, 1024, kernel_size=1)#实际上是全连接层用1*1卷积得
    以替代
37     layers += [pool5, conv6,
38                 nn.ReLU(inplace=True), conv7, nn.ReLU(inplace=True)]

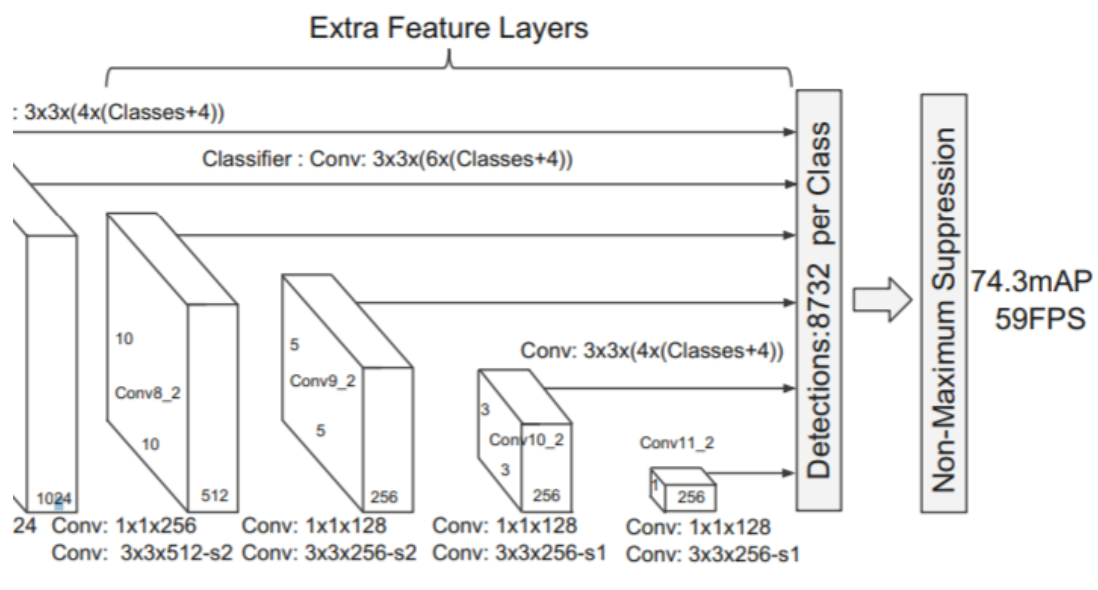
```

- extra feature layers

conv8_1 1*1 256 conv8_2 3*3 512 conv9_1 1*1 128 conv9_2 3*3 256 对应的channel

conv10_1 1*1 128 conv10_2 3*3 256channel conv11_1 1*1 128 channels conv11_2 3*3 256
channel

注意这里没有多余的池化层的关系



```

1  """
2  feature -> conv8_2:1x1x256,3x3x512s2 -> conv9_2:1x1x128,3x3x256s2 ->
3          conv10_2:1x1x128,3x3x256s1 -> conv11_2:1x1x128,3x3x256s1
4          extras_base = {
5              '300': [256, 's', 512, 128, 's', 256, 128, 256, 128, 256],
6              '512': [256, 's', 512, 128, 's', 256, 128, 's', 256, 128, 's', 256],
7          }
8  """
9
10 #这里面s的设计用来进行对应的stride
11 def add_extras(cfg, i, size=300):
12     # Extra layers added to VGG for feature scaling
13     layers = []
14     in_channels = i
15     flag = False#0和1调整第一个和第二个

```

```

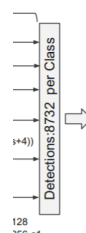
16     for k, v in enumerate(cfg):
17         if in_channels != 's': #调整由于字符串产生的结果
18             if v == 's':
19                 layers += [nn.Conv2d(in_channels, cfg[k + 1],
kernel_size=(1, 3)[flag], stride=2, padding=1)]
20             else:
21                 layers += [nn.Conv2d(in_channels, v, kernel_size=(1, 3)
[flag])]
22                 flag = not flag
23                 in_channels = v
24             if size == 512: #如果对应的s
25                 layers.append(nn.Conv2d(in_channels, 128, kernel_size=1,
stride=1))
26                 layers.append(nn.Conv2d(128, 256, kernel_size=4, stride=1,
padding=1))
27             return layers
28
29 #config对应

```

- 多尺度检测特征图

添加卷积特征层到截断的基础网络的末端，截断的网络层次在下面的结构中存在，减少了计算量(类似于NIN块中)用 1×1 卷积块来代替全连接层的存在，可以实现多尺度检测的预测

- 每个被添加的特征层(选择性从基本的网络中存在的额外的特征层)可以产生一系列固定的检测预测输出,一个特征层的shape $m \times n$, channels数对应的p,P个卷积核的kernel_size=(3,3), 会产生检测框的类别或者位置的相关信息(类别的相关信息, 是各个检测框的softmax的分类得分score, 对应的位置信息则是prior box和ground truth box的offset 文章中用conf来表示), SSD文章中的上述利用卷积层来代替FC层的存在



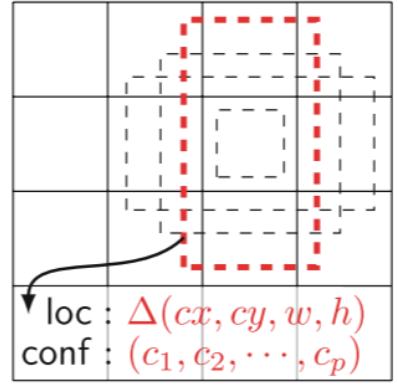
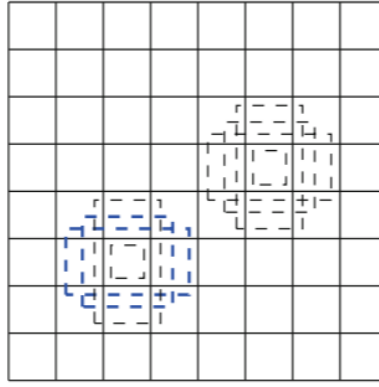
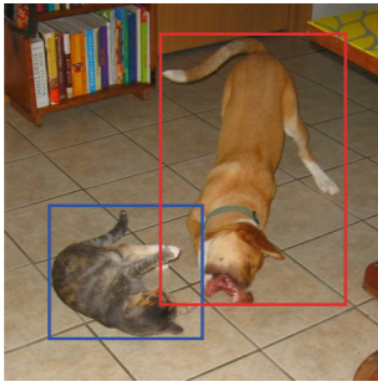
- 默认框和纵横比

这个思想其实和Faster RCNN中的anchor的思想很一致，通过在特征图中每一个pixel的中心处生成scale和ratio不同的先验框，根据论文中事先设计的超参数的来获得default box的ratio和scale, 对于每个预测框预测对应的刚才的score和offsets

每个对应的feature map中的cell会产生k个prior box,共有C类classes,和conf中的4个offsets, 每个cell中输出的channel $(C+4) \times k$ 个filter,对于特征图shape为 (m, n) 对应的输出为 $m \times n \times (c + 4)$ 对应输出的channels, 其实和Faster RCNN很像了, 但是SSD可以通过多个特征图来对应符合不同size的物体

在前面的anchor cobv4 conv7 conv 8 conv 9 conv10 conv11中对应的prior box的数量分别4 6 6 4 4个prior box

每一层后跟着



(a) Image with GT boxes (b) 8×8 feature map (c) 4×4 feature map

Training

SSD的Ground truth information 需要分配到固定上文提到的检测器的输出。目的就是检测框和GT boxes相匹配和对应，一旦这个matching 的任务完后才能，就可以进行端到端反向传播等等。

matching

gtbox和default box匹配的方法就是计算对应的IOU，匹配对应的default box的阈值为0.5，简化了学习的过程，并且预先网络预测大量重叠的default box对应的分类得分，而不是仅仅选择最大的IOU所对应的一个框

train objectives

损失函数 α 用来调整对应的confidence loss 和location loss之间的比例，默认为1

x_{ij}^p 代表着第i个default box 匹配上第j个所属class为p的ground truth,剩下的部分类似于Faster RCNN 中的策略

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g)) \quad (1)$$

$$L_{conf}(x, c) = - \sum_{i \in Pos} x_{ij}^p \log(\hat{c}_i^p) - \sum_{i \in Neg} \log(\hat{c}_i^0) \quad \text{where} \quad \hat{c}_i^p = \frac{\exp(c_i^p)}{\sum_p \exp(c_i^p)}$$

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

$$L_{loc}(x, l, g) = \sum_{i \in Pos} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^k \text{smooth}_{L_1}(l_i^m - \hat{g}_j^m)$$

$$\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx}) / d_i^w \quad \hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy}) / d_i^h$$

$$\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right) \quad \hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right)$$

```
1 def encode(matched, priors, variances):
```

```

2     """Encode the variances from the priorbox layers into the ground truth
boxes
3     we have matched (based on jaccard overlap) with the prior boxes.
4     Args:
5         matched: (tensor) Coords of ground truth for each prior in point-
form
6             Shape: [num_priors, 4].
7         priors: (tensor) Prior boxes in center-offset form
8             Shape: [num_priors,4].
9         variances: (list[float]) Variances of priorboxes
10    Return:
11        encoded boxes (tensor), Shape: [num_priors, 4]
12    """
13    # matched  gt box xmin ymin xmax ymax
14    # prior box  cx cy w h
15    # dist b/t match center and prior's center
16    g_cxcy = (matched[:, :2] + matched[:, 2:])/2 - priors[:, :2]
17    # encode variance
18    g_cxcy /= (variances[0] * priors[:, 2:])
19    # match wh / prior wh
20    g_wh = (matched[:, 2:] - matched[:, :2]) / priors[:, 2:]
21    g_wh = torch.log(g_wh) / variances[1]
22    # return target for smooth_l1_loss
23    return torch.cat([g_cxcy, g_wh], 1) # [num_priors,4]

```

上述四个公式在这里，上述的

Faster RCNN bounding Regression

上述的g代表在Faster RCNN 中 d为default box,g 代表着ground truth box中的结果

给定一组的default box $d=(d_{cx},d_{cy},d_w,d_h)$ 和一组Ground truth $GT = (g_x, g_y, g_w, g_h)$ 找到一组线性变换满足，可以使得default box近似于对应的Ground truth的数据，可以进行平行和缩放的处理

$$G'_x = d_w l_x(d) + d_x \quad (2)$$

$$G'_w = d_w \exp(l_w(d)) \quad (3)$$

所以当前需要的就是几个线性变换的参数 L_* 代指各个x,y,w,h等四个变换如何获得对应，所以利用线性回归的方法计算 $Y = w_*^T d_*$ 代表feature map中的特征向量计算上述公式的 l^i

上述的部分在这里

```

1  def decode(loc, priors, variances):
2      """Decode locations from predictions using priors to undo
3      the encoding we did for offset regression at train time.
4      Args:
5          loc (tensor): location predictions for loc layers,
6              Shape: [num_priors,4]
7          priors (tensor): Prior boxes in center-offset form.
8              Shape: [num_priors,4].
9          variances: (list[float]) Variances of priorboxes
10     Return:
11         decoded bounding box predictions
12     """
13
14     boxes = torch.cat((

```

```

15     priors[:, :2] + loc[:, :2] * variances[0] * priors[:, :2],
16     priors[:, :2] * torch.exp(loc[:, :2] * variances[1])), 1)
17     boxes[:, :2] -= boxes[:, :2] / 2
18     boxes[:, :2] += boxes[:, :2]
19     """
20     后面这部分代码是将(cx,cy,h,w)->(xmin,ymin,xmax,ymax)的过程
21     """
22     return boxes
23

```

hard negative mining

$$Recall = \frac{TP}{TP + FN} \quad (4)$$

目标检测与图像分类不同，图像分类往往只有一个输出，但目标检测的输出个数却是未知的。除了 Ground-Truth（标注数据）训练时，**模型永远无法百分百确信自己要在一张图上预测多少物体**。SSD中大量的负样本，但此时就会遇到一个问题，因为区域提议实在太多，导致在训练时绝大部分都是负样本，这导致了大量无意义负样本的梯度“淹没”了有意义的正样本。需要用方法抑制大量的简单负样本，挖掘所有Hard example 的特征。

难负例挖掘（Hard Negative Mining）就是在训练时，尽量多挖掘些难负例（hard negative)加入负样本集，这样会比easy negative组成的负样本集效果更好。

```

1  # -*- coding: utf-8 -*-
2  import torch
3  import torch.nn as nn
4  import torch.nn.functional as F
5  from data import voc as cfg
6  from ..box_utils import match, log_sum_exp
7
8
9  class MultiBoxLoss(nn.Module):
10     """SSD Weighted Loss Function
11     Compute Targets:
12         1) Produce Confidence Target Indices by matching  ground truth
13         boxes
14             with (default) 'priorboxes' that have jaccard index > threshold
15             parameter
16             (default threshold: 0.5).
17         2) Produce localization target by 'encoding' variance into offsets
18         of ground
19         truth boxes and their matched 'priorboxes'.
20         3) Hard negative mining to filter the excessive number of negative
21         examples
22         that comes with using a large number of default bounding boxes.
23         (default negative:positive ratio 3:1)
24     Objective Loss:
25         L(x,c,l,g) = (Lconf(x, c) + αLloc(x,l,g)) / N
26         where, Lconf is the CrossEntropy Loss and Lloc is the SmoothL1 Loss
27         weighted by α which is set to 1 by cross val.
28     Args:
29         c: class confidences,
30         l: predicted boxes,

```



```

27         g: ground truth boxes
28         N: number of matched default boxes
29         See: https://arxiv.org/pdf/1512.02325.pdf for more details.
30     """
31
32     def __init__(self, num_classes, overlap_thresh, prior_for_matching,
33                 bkg_label, neg_mining, neg_pos, neg_overlap,
encode_target,
34                 use_gpu=True):
35         super(MultiBoxLoss, self).__init__()
36         self.use_gpu = use_gpu
37         self.num_classes = num_classes
38         self.threshold = overlap_thresh
39         self.background_label = bkg_label
40         self.encode_target = encode_target
41         self.use_prior_for_matching = prior_for_matching
42         self.do_neg_mining = neg_mining
43         self.negpos_ratio = neg_pos
44         self.neg_overlap = neg_overlap
45
46         #import pdb
47         #pdb.set_trace()
48
49         self.variance = cfg['variance']
50
51     def forward(self, predictions, targets):
52         """Multibox Loss
53         Args:
54             predictions (tuple): A tuple containing loc preds, conf preds,
55             and prior boxes from SSD net.
56                 conf shape: torch.size(batch_size,num_priors,num_classes)
57                 loc shape: torch.size(batch_size,num_priors,4)
58                 priors shape: torch.size(num_priors,4)
59
60             targets (tensor): Ground truth boxes and labels for a batch,
61                 shape: [batch_size,num_objs,5] (last idx is the label).
62         """
63
64         loc_data, conf_data, priors = predictions
65         num = loc_data.size(0) #对应的batch size
66         priors = priors[:loc_data.size(1), :]
67         num_priors = (priors.size(0))
68         num_classes = self.num_classes
69
70         # match priors (default boxes) and ground truth boxes
71         # 1 首先匹配正负样本
72         loc_t = torch.Tensor(num, num_priors, 4)
73         conf_t = torch.LongTensor(num, num_priors)
74         for idx in range(num):
75             truths = targets[idx][:, :-1].data #前面的是GT Box 四个坐标
76             labels = targets[idx][:, -1].data #后面是对应的label
77             defaults = priors.data
78             # 得到每一个prior对应的truth,放到loc_t与conf_t中,conf_t中是类
别,loc_t中是[matches, prior, variance]
79             match(self.threshold, truths, defaults, self.variance, labels,
80                 loc_t, conf_t, idx) #获得对应的匹配,并encode之后,这里使用的都
是类内的变量,无须返回相应的值
81         if self.use_gpu:

```



```

82         loc_t = loc_t.cuda()
83         conf_t = conf_t.cuda()
84     # wrap targets
85
86     # 2 计算所有正样本的定位损失,负样本不需要定位损失
87
88     # 计算正样本的数量
89     pos = conf_t > 0
90     num_pos = pos.sum(dim=1, keepdim=True)
91
92     #import pdb
93     #pdb.set_trace()
94     #根据前文的multibox生成的8732个对应的anchor
95
96     # Localization Loss (Smooth L1)
97     # Shape: [batch,num_priors,4]
98     # 将pos_idx扩展为[32, 8732, 4],正样本的索引
99     pos_idx = pos.unsqueeze(pos.dim()).expand_as(loc_data)
100
101     # 正样本的定位预测值
102     loc_p = loc_data[pos_idx].view(-1, 4)
103     # 正样本的定位真值
104     loc_t = loc_t[pos_idx].view(-1, 4)
105     # 所有正样本的定位损失
106     loss_l = F.smooth_l1_loss(loc_p, loc_t, size_average=False)
107
108     # 3 对于类别损失,进行难样本挖掘,控制比例为1:3
109
110     # Compute max conf across batch for hard negative mining
111     # 所有prior的类别预测
112     batch_conf = conf_data.view(-1, self.num_classes)
113     # 计算类别损失.每一个的log(sum(exp(21个的预测)))-对应的真正预测值
114     loss_c = log_sum_exp(batch_conf) - batch_conf.gather(1,
conf_t.view(-1, 1))
115
116     # Hard Negative Mining
117     loss_c = loss_c.view(pos.size()[0], pos.size()[1])
118     # 首先过滤掉正样本
119     loss_c[pos] = 0 # filter out pos boxes for now
120     loss_c = loss_c.view(num, -1)
121     _, loss_idx = loss_c.sort(1, descending=True)
122     # idx_rank为排序后每个元素的排名
123     _, idx_rank = loss_idx.sort(1)
124     num_pos = pos.long().sum(1, keepdim=True)
125
126     # 这个地方负样本的最大值不应该是pos.size(1)-num_pos?
127     num_neg = torch.clamp(self.negpos_ratio*num_pos, max=pos.size(1)-1)
128     # 选择每个batch中负样本的索引
129     neg = idx_rank < num_neg.expand_as(idx_rank)
130
131     # 4 计算正负样本的类别损失
132
133     # Confidence Loss Including Positive and Negative Examples
134     # 都扩展为[32, 8732, 21]
135     pos_idx = pos.unsqueeze(2).expand_as(conf_data)
136     neg_idx = neg.unsqueeze(2).expand_as(conf_data)
137     # 把预测提出来

```

```

138         conf_p = conf_data[(pos_idx+neg_idx).gt(0)].view(-1,
self.num_classes)
139         # 对应的标签
140         targets_weighted = conf_t[(pos+neg).gt(0)]
141         loss_c = F.cross_entropy(conf_p, targets_weighted,
size_average=False)
142
143         # Sum of losses: L(x,c,l,g) = (Lconf(x, c) + αLloc(x,l,g)) / N
144
145         N = num_pos.data.sum()
146         loss_l /= N.type('torch.cuda.FloatTensor')
147         loss_c /= N.type('torch.cuda.FloatTensor')
148         return loss_l, loss_c
149

```

scale和aspect ratio的设计方式

使用较低和较高的特征图用于检测，m的大小根据特征图的种类来

$$s_k = s_{min} + \frac{s_{max} - s_{min}}{m - 1}, k \in (1, m) \quad (5)$$

$$s_{min} = 0.2, s_{max} = 0.9 \quad (6)$$

这样各个层之间的scale设计就已经很清楚了

ratio

$$a_r = \{1, 2, 3, \frac{1}{3}, \frac{1}{2}\} \quad (7)$$

计算对应的prior box的高宽对应的公式为

$$w_k = s_k \sqrt{a_r}, h_k^\alpha = \frac{s_k}{a_r} \quad (8)$$

当对应的ratio为1 的时候，默认框的scale为下面的公式 $s_k^\alpha = \sqrt{s_k s_{k+1}}$ ，default box的建议框的中心是各个feature map 像素点的中心。

通过结合特征图不同尺度的scale以及对应高宽比的默认框，我们一系列大量不同的预测值，覆盖大量范围的尺寸。

	min_size	max_size
conv4_3	30	60
fc7	60	111
conv6_2	111	162
conv7_2	162	213
conv8_2	213	264
conv9_2	264	315

- 在一个feature map中输出的类别图像的channels
(num_classes*num_priorbox,layer_height,layer_width]的特征进行softmax后的大奥对应的分类得分

- 位置对应regression的CNN输出的($4 \times$ 先验框)

```
[3]: from itertools import product
      for i, j in product(range(30), repeat=2):
          print(i, j)
```

```
1 10
1 11
1 12
1 13
1 14
1 15
1 16
1 17
1 18
1 19
1 20
1 21
1 22
1 23
1 24
```

```
1 class PriorBox(object):#先验框生成对应的
2     """Compute priorbox coordinates in center-offset form for each source
3     feature map.
4     """
5     def __init__(self, cfg):
6         super(PriorBox, self).__init__()
7         self.image_size = cfg['min_dim']
8         # number of priors for feature map location (either 4 or 6)
9         self.num_priors = len(cfg['aspect_ratios'])
10        self.variance = cfg['variance'] or [0.1]
11        self.feature_maps = cfg['feature_maps']
12        self.min_sizes = cfg['min_sizes']
13        self.max_sizes = cfg['max_sizes']
14        self.steps = cfg['steps']
15        self.aspect_ratios = cfg['aspect_ratios']
16        self.clip = cfg['clip']
17        self.version = cfg['name']
18        for v in self.variance:
19            if v <= 0:
20                raise ValueError('Variances must be greater than 0')
21
22        # 生成所有的PriorBox, 需要每一个特征图的信息
23        def forward(self):
24            mean = []
25            for k, f in enumerate(self.feature_maps):#前面6层的特征图
26                for i, j in product(range(f), repeat=2):#获取对应的笛卡尔的及
27                    # f_k为每个特征图的尺寸
28                    f_k = self.image_size / self.steps[k]#获取对应的特征图的尺寸
29                    # 求取每个box的中心坐标
30                    cx = (j + 0.5) / f_k
```

```

31         cy = (i + 0.5) / f_k
32
33         # 对应{s_k, s_k}大小的PriorBox
34         s_k = self.min_sizes[k]/self.image_size
35         mean += [cx, cy, s_k, s_k]#第一个正方形对应anchor
36
37         # 对应{ $\sqrt{(s_k s_{k+1})}$ ,  $\sqrt{(s_k s_{k+1})}$ }大小的PriorBox
38         s_k_prime = sqrt(s_k * (self.max_sizes[k]/self.image_size))
39         mean += [cx, cy, s_k_prime, s_k_prime]
40
41         # 剩余的比例为2、1/2、3、1/3的PriorBox
42         for ar in self.aspect_ratios[k]:
43             mean += [cx, cy, s_k*sqrt(ar), s_k/sqrt(ar)]
44             mean += [cx, cy, s_k/sqrt(ar), s_k*sqrt(ar)]
45         # back to torch land
46         output = torch.Tensor(mean).view(-1, 4)#转换成对应的batch的格式
47         if self.clip:
48             output.clamp_(max=1, min=0)#将小于0的部分去除掉
49         return output
50

```

Source layers from:					
conv4_3	conv7	conv8_2	conv9_2	conv10_2	conv11_2
✓	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	
✓	✓	✓	✓		
✓	✓	✓			
✓	✓				
	✓				

在上述对应的先验框生成之后，我们需要在选定的feature map中输出两个卷积分别是3*3的深度卷积层，用来进行classification和regression的任务，具体的channel的数量根据特征图的中每个像素点的先验框的数量和classes，4,6,6,6,4,4对应的，原来的论文定义的是21 conf分类的channel则为84 126 同理 loc对应的则是 16 24 等等下面是对应的代码的实现

```

1  def multibox(vgg, extra_layers, cfg, num_classes):
2      loc_layers = []
3      conf_layers = []
4      vgg_source = [21, -2]
5
6      """
7      # 这里的base为VGG-16前13个卷积层构造，M代表maxpooling，C代表ceil_mode为True
8      base = [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'C', 512, 512, 512,
9              'M',
10              512, 512, 512]
11      # 额外部分的卷积通道数，S代表了步长为2，其余卷积层默认步长为1
12      extras = [256, 'S', 512, 128, 'S', 256, 128, 256, 128, 256]
13      # 每个特征图上一个点对应的PriorBox数量
14      mbox = [4, 6, 6, 6, 4, 4] # number of boxes per feature map location

```

```

14
15     对应的conv4和conv 8
16     """
17     for k, v in enumerate(vgg_source):
18         loc_layers += [nn.Conv2d(vgg[v].out_channels,
19                                 cfg[k] * 4, kernel_size=3, padding=1)]
20         conf_layers += [nn.Conv2d(vgg[v].out_channels,
21                                 cfg[k] * num_classes, kernel_size=3, padding=1)]
22     for k, v in enumerate(extra_layers[1::2], 2):#设置的start k等于2, 就是从第三个开始对应mbox,
23         #获取对应的extra, 已知第一层Conv*_1,slice(1,end,2)每次都可以选择对应的
24         feature map中
25         loc_layers += [nn.Conv2d(v.out_channels, cfg[k]
26                                 * 4, kernel_size=3, padding=1)]
27         conf_layers += [nn.Conv2d(v.out_channels, cfg[k]
28                                 * num_classes, kernel_size=3, padding=1)]
29     return vgg, extra_layers, (loc_layers, conf_layers)

```

下面是整个网络的前向传播的过程

```

1  ![data](C:\Users\王澳\Desktop\SSD\data.jpg)    def forward(self, x):
2      """Applies network layers and ops on input image(s) x.
3
4      Args:
5          x: input image or batch of images. Shape: [batch,3,300,300].
6
7      Return:
8          Depending on phase:
9          test:
10             Variable(tensor) of output class label predictions,
11             confidence score, and corresponding location predictions for
12             each object detected. Shape: [batch,topk,7]
13
14          train:
15             list of concat outputs from:
16             1: confidence layers, Shape:
17             [batch*num_priors,num_classes]
18             2: localization layers, Shape: [batch,num_priors*4]
19             3: priorbox layers, Shape: [2,num_priors*4]
20
21      """
22      # sources保存特征图, loc与conf保存所有PriorBox的位置与类别预测特征
23      sources = list()
24      loc = list()
25      conf = list()
26
27      # 对输入图像卷积到conv4_3, 将特征添加到sources中
28      for k in range(23):
29          x = self.vgg[k](x)#bao
30
31          s = self.L2Norm(x)
32          sources.append(s)
33
34      # 继续卷积到conv7, 将特征添加到sources中
35      for k in range(23, len(self.vgg)):
36          x = self.vgg[k](x)
37          sources.append(x)

```

```

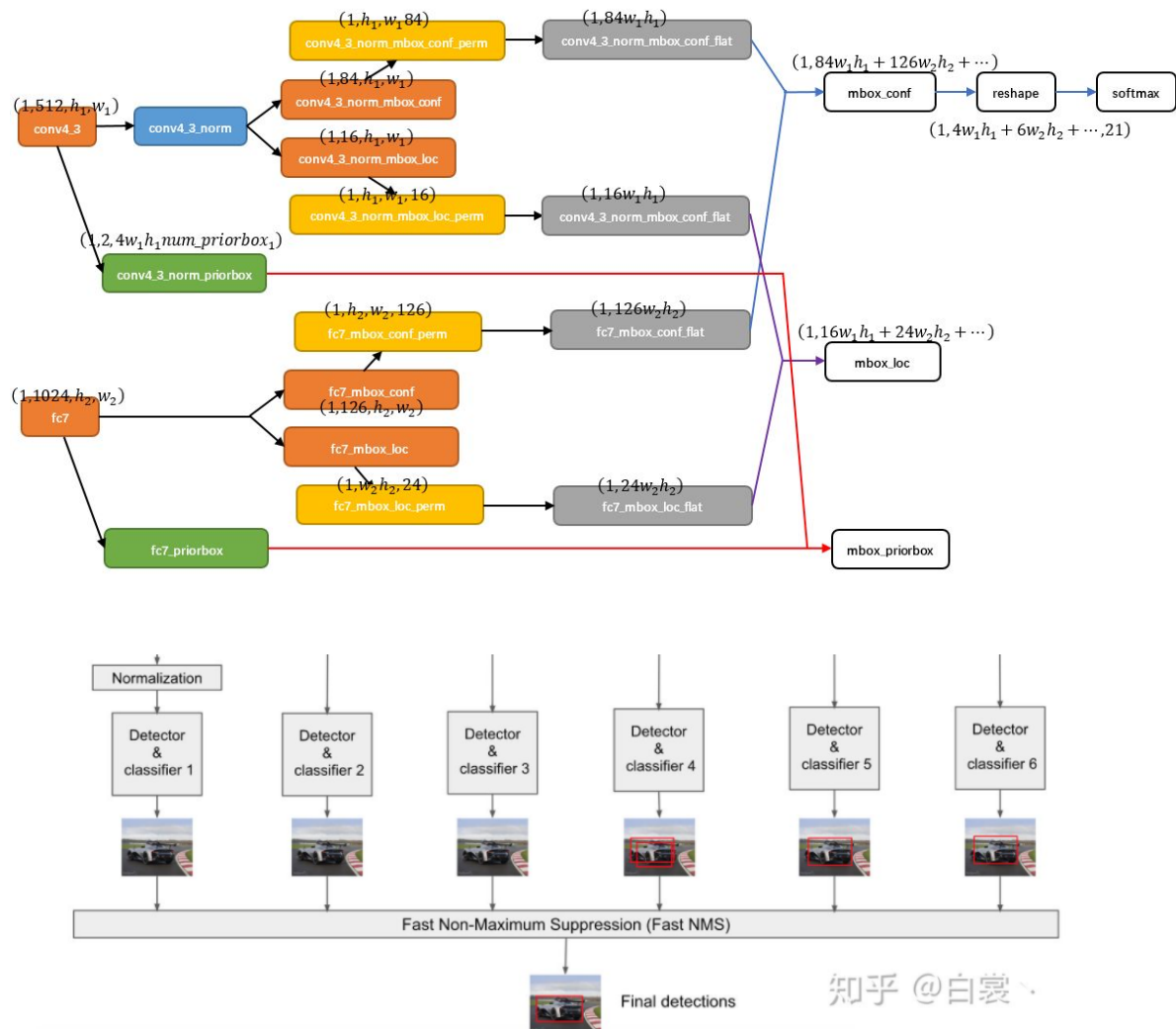
37         # 继续利用额外的卷积层计算，并将特征添加到sources中
38         for k, v in enumerate(self.extras):
39             x = F.relu(v(x), inplace=True)
40             if k % 2 == 1:
41                 sources.append(x)
42
43

```

SSD数据流动方式

参考大神的白裳的思路，列出caffe带吗解决方法，下面以Conv4_3和conv7的数据流的合并的过程,其实对应的pytorch的解决思路一样，下面解释对应的矩阵shape的变换

conv4_3 (batch,channel,height,width)-> (batch,height,width,channel) permute(0,2,3,1)-> (batch,height*width*channel) flatten()->然后再利用concat函数堆叠在一起获取，形成mbox_conf张量，reshape后形成 (xx,21)对应着20个类+1个背景，利用softmax进行对应的分类得分



```

1  """
2  代码是刚才上面所有过程进行阐述的过程
3  """
4  # 对sources中的特征图利用类别与位置网络进行卷积计算，并保存到loc与conf中
5      for (x, l, c) in zip(sources, self.loc, self.conf):
6          loc.append(l(x).permute(0, 2, 3, 1).contiguous())
7          conf.append(c(x).permute(0, 2, 3, 1).contiguous())
8          #loc 和 conf的permute过程已经对应的阐述过了
9          #利用列表生成器，将每行对应成flatten，并在dim=1,进行堆叠
10         loc = torch.cat([o.view(o.size(0), -1) for o in loc], 1)

```

```

11         conf = torch.cat([o.view(o.size(0), -1) for o in conf], 1)
12         if self.phase == "test":
13             output = self.detect(
14                 loc.view(loc.size(0), -1, 4),          # loc preds
15                 self.softmax(conf.view(conf.size(0), -1,
16                                         self.num_classes)), # conf preds
17                 self.priors.type(type(x.data))          # default
18             )
19         else:
20             # 对于训练来说，output包括了loc与conf的预测值以及PriorBox的信息，和我们
            刚才所讲述的过程一直
21             output = (
22                 loc.view(loc.size(0), -1, 4),
23                 conf.view(conf.size(0), -1, self.num_classes),
24                 self.priors
25             )

```

上面所有的代码是ssd.py的代码，也就是网络的主要的过程，代码可参考Detection-Pytorch-Notebook中的chapter 5的代码。这个部分也就是网络的最后部分，综上所述，我们已经将网络的整体部分阐述完毕，接下来是文章中的一些其他的東西，难样本挖掘和matching以及后面的NMS的过程。

Matching

第一个函数是 torch.max

```
: tr=torch.arange(20).reshape(2, 10)
```

```
: tr.max(1, keepdim=True)
```

```
: torch.return_types.max(
    values=tensor([[ 9],
                  [19]]),
    indices=tensor([[9],
                  [9]]))
```

axis=1,返回对应的列索引以及对应的lou value

第二个函数 index_fill

按照index所给定的顺序，填充对应的index指定的value值，

index_fill(dim,index,value)#必须选定对应的指定维度

```

1 def match(threshold, truths, priors, variances, labels, loc_t, conf_t, idx):
2     """Match each prior box with the ground truth box of the highest jaccard
3     overlap, encode the bounding boxes, then return the matched indices
4     corresponding to both confidence and location preds.

```



```

5     Args:
6         threshold: (float) The overlap threshold used when mathing boxes.
7         truths: (tensor) Ground truth boxes, Shape: [num_obj, num_priors].
8         priors: (tensor) Prior boxes from priorbox layers, Shape:
[n_priors,4].
9         variances: (tensor) Variances corresponding to each prior coord,
10            Shape: [num_priors, 4].
11         labels: (tensor) All the class labels for the image, Shape:
[num_obj].
12         loc_t: (tensor) Tensor to be filled w/ endcoded location targets.
13         conf_t: (tensor) Tensor to be filled w/ matched indices for conf
preds.
14         idx: (int) current batch index
15     Return:
16         The matched indices corresponding to 1)location and 2)confidence
preds.
17     """
18
19     # 注意这里truth是最大最小值形式的,而prior是中心点与长宽形式
20     # 求取真实框与预选框的IoU
21     overlaps = jaccard(
22         truths,
23         point_form(priors)
24     )
25
26     # (Bipartite Matching)
27     # [1,num_objects] best prior for each ground truth
28     #前面的个体的box,axis=1时,输出的是与匹配度最高的prior idx
29     best_prior_overlap, best_prior_idx = overlaps.max(1, keepdim=True)
30     #此时输出的shape对应为shape(A,1) 返回结果每一个ground truth 有最大IOU值的prior
box以及对应的id
31     best_truth_overlap, best_truth_idx = overlaps.max(0, keepdim=True)
32     #此时输出的shape(1, B)每个prior box 有着最大IOU值的 GT Box 的ID
33     best_truth_idx.squeeze_(0)#输出的是两个,将其所有转换为size=(a)的数组
34     best_truth_overlap.squeeze_(0)#修改矩阵对应的维度也就是size(B,)
35     best_prior_idx.squeeze_(1)
36     best_prior_overlap.squeeze_(1)#修改矩阵对应维度即size(A,)
37
38
39     # 将每一个truth对应的最佳box的overlap设置为2
40     best_truth_overlap.index_fill_(0, best_prior_idx, 2) # 确保最优的ground
truth box
41     # TODO refactor: index best_prior_idx with long tensor
42     # ensure every gt matches with its prior of max overlap
43
44     # 保证每一个truth对应的最佳box,该box要对应到这个truth上,即使不是最大iou
45     #已经知道best--prior--idx的shape等于[A,1]对应的value的值prior box的索引
46     #下面的那个对应的id [B,1]对应的value值 GT的索引
47     #将第j个GT 框 iou最大的索引的prior box 的最佳GT框设置为第j个
48     for j in range(best_prior_idx.size(0)):
49         best_truth_idx[best_prior_idx[j]] = j#j为第j个GT框,
50     #
51
52     # 每一个prior对应的真实框的位置,这样两个就会一一进行对应的匹配
53     matches = truths[best_truth_idx]          # Shape: [num_priors,4]
54
55     # 每一个prior对应的类别
56     conf = labels[best_truth_idx] + 1         # Shape: [num_priors]

```

```

57
58     # 如果一个PriorBox对应的最大IoU小于0.5，则视为负样本
59     conf[best_truth_overlap < threshold] = 0 # label as background
60
61     # 进一步计算定位的偏移真值
62     loc = encode(matches, priors, variances)
63     loc_t[idx] = loc # [num_priors,4] encoded offsets to learn
64     conf_t[idx] = conf # [num_priors] top class label for each prior
65

```

剩下的部分就已经完全清楚，还剩下一个数据增广，以及一个decode的过程

```

1  import torch
2  from torch.autograd import Function
3  from ..box_utils import decode, nms
4  from data import voc as cfg
5
6
7  class Detect(Function):
8      """At test time, Detect is the final layer of SSD. Decode location
9      preds,
10      apply non-maximum suppression to location predictions based on conf
11      scores and threshold to a top_k number of output predictions for both
12      confidence score and locations.
13      """
14      def __init__(self, num_classes, bkg_label, top_k, conf_thresh,
15      nms_thresh):
16          self.num_classes = num_classes
17          self.background_label = bkg_label
18          self.top_k = top_k
19          # Parameters used in nms.
20          self.nms_thresh = nms_thresh
21          if nms_thresh <= 0:
22              raise ValueError('nms_threshold must be non negative.')
23          self.conf_thresh = conf_thresh
24          self.variance = cfg['variance']
25
26      def forward(self, loc_data, conf_data, prior_data):
27          """
28          Args:
29              loc_data: (tensor) Loc preds from loc layers
30                      Shape: [batch,num_priors*4]
31              conf_data: (tensor) Shape: Conf preds from conf layers
32                      Shape: [batch*num_priors,num_classes]
33              prior_data: (tensor) Prior boxes and variances from priorbox
34                      layers
35                      Shape: [1,num_priors,4]
36          """
37          num = loc_data.size(0) # batch size
38          num_priors = prior_data.size(0) # 鲜艳框的数量
39          output = torch.zeros(num, self.num_classes, self.top_k, 5)
40          conf_preds = conf_data.view(num, num_priors,
41                                     self.num_classes).transpose(2, 1)
42
43          # Decode predictions into bboxes.
44          for i in range(num):
45              decoded_boxes = decode(loc_data[i], prior_data, self.variance)

```

```

43         # For each class, perform nms
44         conf_scores = conf_preds[i].clone()
45
46         for cl in range(1, self.num_classes):
47             c_mask = conf_scores[cl].gt(self.conf_thresh)
48             scores = conf_scores[cl][c_mask]
49             if scores.dim() == 0:
50                 continue
51             l_mask = c_mask.unsqueeze(1).expand_as(decoded_boxes)
52             boxes = decoded_boxes[l_mask].view(-1, 4)
53             # idx of highest scoring and non-overlapping boxes per class
54             ids, count = nms(boxes, scores, self.nms_thresh, self.top_k)
55             output[i, cl, :count] = \
56                 torch.cat((scores[ids[:count]].unsqueeze(1),
57                           boxes[ids[:count]]), 1)
58         flt = output.contiguous().view(num, -1, 5)
59         _, idx = flt[:, :, 0].sort(1, descending=True)
60         _, rank = idx.sort(1)
61         flt[(rank < self.top_k).unsqueeze(-1).expand_as(flt)].fill_(0)
62         return output
63

```