

```
In [1]: # 환경을 설정합니다. 여기서 fastai를 이용하여 모델을 생성하고
# f1_score를 이용하여 학습모델을 평가합니다.

from fastai.conv_learner import *
from fastai.dataset import *

import pandas as pd
import numpy as np
import os
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score
import scipy.optimize as opt
```

```
In [2]: # 모델 분석에 활용될 데이터의 path 입니다.

PATH = './'
TRAIN = 'protein/input/train/'
TEST = 'protein/input/test'
LABELS = 'protein/train.csv'
SAMPLE = 'protein/sample_submission.csv'
```

```
In [3]: # 문제에서 사용될 28개의 label입니다.

name_label_dict = {
0: 'Nucleoplasm',
1: 'Nuclear membrane',
2: 'Nucleoli',
3: 'Nucleoli fibrillar center',
4: 'Nuclear speckles',
5: 'Nuclear bodies',
6: 'Endoplasmic reticulum',
7: 'Golgi apparatus',
8: 'Peroxisomes',
9: 'Endosomes',
10: 'Lysosomes',
11: 'Intermediate filaments',
12: 'Actin filaments',
13: 'Focal adhesion sites',
14: 'Microtubules',
15: 'Microtubule ends',
16: 'Cytokinetic bridge',
17: 'Mitotic spindle',
18: 'Microtubule organizing center',
19: 'Centrosome',
20: 'Lipid droplets',
21: 'Plasma membrane',
22: 'Cell junctions',
23: 'Mitochondria',
24: 'Aggresome',
25: 'Cytosol',
26: 'Cytoplasmic bodies',
27: 'Rods & rings' }
```

```
In [4]: nw = 2 # 2개의 worker를 사용합니다.
arch = resnet34 # resnet34를 이용하여 모델을 학습시킵니다.
```

```
In [5]: # train과 test에 쓰일 label입니다.
# 그리고 이를 split 하여 활용합니다.

train_names = list({f[:36] for f in os.listdir(TRAIN)})
test_names = list({f[:36] for f in os.listdir(TEST)})
tr_n, val_n = train_test_split(train_names, test_size=0.1, random_state=42)
```

```
In [6]: # 하나의 데이터는 red, gree, blue, yellow 4가지로 구성되어 있습니다. 이 이미지들을 학습하기 위해 전처리해 줍니다.

def open_rgby(path, id):
    colors = ['red', 'green', 'blue', 'yellow']
    flags = cv2.IMREAD_GRAYSCALE
    img = [cv2.imread(os.path.join(path, id+'_'+color+'.png'), flags).astype(np.float32)/255
            for color in colors]
    return np.stack(img, axis=-1)
```

```
In [7]: class pdFilesDataset(FilesDataset):
    def __init__(self, fnames, path, transform):
        self.labels = pd.read_csv(LABELS).set_index('Id')
        self.labels['Target'] = [[int(i) for i in s.split()] for s in self.labels['Target']]
        super().__init__(fnames, transform, path)

    def get_x(self, i):
        img = open_rgby(self.path, self.fnames[i])
        if self.sz == 512: return img
        else: return cv2.resize(img, (self.sz, self.sz), cv2.INTER_AREA)

    def get_y(self, i):
        if(self.path == TEST): return np.zeros(len(name_label_dict), dtype=np.int)
        else:
            labels = self.labels.loc[self.fnames[i]]['Target']
            return np.eye(len(name_label_dict), dtype=np.float)[labels].sum(axis=0)

    @property
    def is_multi(self): return True
    @property
    def is_reg(self): return True
    # 출력층에서 로지스틱 시그모이드를 사용하고,
    # 여기에 focal loss를 더한것을 loss function으로 사용하기 위해 설정합니다.

    def get_c(self): return len(name_label_dict)
    # label class의 수를 반환합니다.
```

```
In [8]: def get_data(sz, bs):

    aug_tfms = [RandomRotate(30, tfm_y=TfmType.NO),
                RandomDihedral(tfm_y=TfmType.NO),
                RandomLighting(0.05, 0.05, tfm_y=TfmType.NO)]
    # 데이터를 agumenttaion합니다. 이때 데이터에 변환을 가합니다.

    stats = A([0.08069, 0.05258, 0.05487, 0.08282], [0.13704, 0.10145, 0.15313, 0.13814])
    tfms = tfms_from_stats(stats, sz, crop_type=CropType.NO, tfm_y=TfmType.NO,
                           aug_tfms=aug_tfms)
    ds = ImageData.get_ds(pdFilesDataset, (tr_n[:-(len(tr_n)%bs)], TRAIN),
                          (val_n, TRAIN), tfms, test=(test_names, TEST))
    md = ImageData(PATH, ds, bs, num_workers=nw, classes=None)
    # 트레인셋에 평균과 분산을 각각의 채널에 적용합니다.
    return md
```

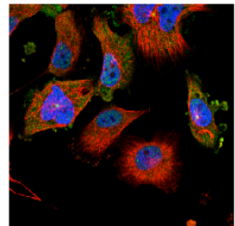
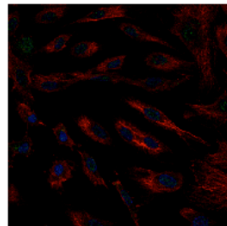
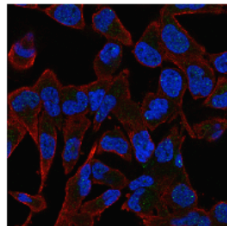
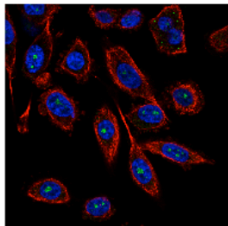
```
In [10]: bs = 16
sz = 512 # 256 -> 512
# 배치사이즈와 이미지 사이즈를 정의합니다.
md = get_data(sz,bs)

x,y = next(iter(md.trn_dl))
x.shape, y.shape

(torch.Size([16, 4, 512, 512]), torch.Size([16, 28]))
```

```
In [11]: def display_imgs(x):
columns = 4
bs = x.shape[0]
rows = min((bs+3)//4,4)
fig=plt.figure(figsize=(columns*4, rows*4))
for i in range(rows):
    for j in range(columns):
        idx = i+j*columns
        fig.add_subplot(rows, columns, idx+1)
        plt.axis('off')
        plt.imshow((x[idx,:,:,:3]*255).astype(np.int))
plt.show()

display_imgs(np.asarray(md.trn_ds.denorm(x)))
# 이미지로더가 올바르게 정의 되어있는지 확인합니다.
```



```
In [12]: # 이 문제의 특징 중 하나가 data가 unbalanced 하다는 것입니다. 이 문제를 보완하기 위해 FocalLoss를 이용합니다.
# easy sample에는 작은 가중치를 부여하고, hard sample에는 높은 가중치를 부여합니다.

class FocalLoss(nn.Module):
    def __init__(self, gamma=2):
        super().__init__()
        self.gamma = gamma

    def forward(self, input, target):
        if not (target.size() == input.size()):
            raise ValueError("Target size ({}) must be the same as input size ({})"
                              .format(target.size(), input.size()))

        max_val = (-input).clamp(min=0)
        loss = input - input * target + max_val + \
            ((-max_val).exp() + (-input - max_val).exp()).log()

        invprobs = F.logsigmoid(-input * (target * 2.0 - 1.0))
        loss = (invprobs * self.gamma).exp() * loss

        return loss.sum(dim=1).mean()
```

```
In [13]: # 정확도를 평가해줍니다.

def acc(preds,targs,th=0.0):
    preds = (preds > th).int()
    targs = targs.int()
    return (preds==targs).float().mean()
```

```

In [14]: class ConvnetBuilder_custom():

    def __init__(self, f, c, is_multi, is_reg, ps=None, xtra_fc=None, xtra_cut=0,
                  custom_head=None, pretrained=True):
        self.f,self.c,self.is_multi,self.is_reg,self.xtra_cut = f,c,is_multi,is_reg,xtra_cut
        if xtra_fc is None: xtra_fc = [512]
        if ps is None: ps = [0.25]*len(xtra_fc) + [0.5]
        self.ps,self.xtra_fc = ps,xtra_fc

        if f in model_meta: cut,self.lr_cut = model_meta[f]
        else: cut,self.lr_cut = 0,0
        cut-=xtra_cut
        layers = cut_model(f(pretrained), cut)

        # 기존의 7x7 3->64 에서 가중치를 유지하면서 첫번째 conv층을 4->64로 설정합니다.
        # 새로운 가중치를 layer[0]을 이용하여 할당합니다.
        # 이는 y채널의 가중치를 0으로 초기화 해주는 것입니다.
        # 기존의 fast ai Resnet34 라이브러리의 코드를 참고하였습니다.
        w = layers[0].weight
        layers[0] = nn.Conv2d(4,64,kernel_size=(7,7),stride=(2,2),padding=(3, 3), bias=False)
        layers[0].weight = torch.nn.Parameter(torch.cat((w,w[:, :1, :, :]),dim=1))

        self.nf = model_features[f] if f in model_features else (num_features(layers)*2)
        if not custom_head: layers += [AdaptiveConcatPool2d(), Flatten()]
        self.top_model = nn.Sequential(*layers)

        n_fc = len(self.xtra_fc)+1
        if not isinstance(self.ps, list): self.ps = [self.ps]*n_fc

        if custom_head: fc_layers = [custom_head]
        else: fc_layers = self.get_fc_layers()
        self.n_fc = len(fc_layers)
        self.fc_model = to_gpu(nn.Sequential(*fc_layers))
        if not custom_head: apply_init(self.fc_model, kaiming_normal)
        self.model = to_gpu(nn.Sequential(*(layers+fc_layers)))

    @property
    def name(self): return f'{self.f.__name__}_{self.xtra_cut}'

    def create_fc_layer(self, ni, nf, p, actn=None):
        res=[nn.BatchNorm1d(num_features=ni)]
        if p: res.append(nn.Dropout(p=p))
        res.append(nn.Linear(in_features=ni, out_features=nf))
        if actn: res.append(actn)
        return res

    def get_fc_layers(self):
        res=[]
        ni=self.nf
        for i,nf in enumerate(self.xtra_fc):
            res += self.create_fc_layer(ni, nf, p=self.ps[i], actn=nn.ReLU())
            ni=nf
        final_actn = nn.Sigmoid() if self.is_multi else nn.LogSoftmax()
        if self.is_reg: final_actn = None
        res += self.create_fc_layer(ni, self.c, p=self.ps[-1], actn=final_actn)
        return res

```

```

def get_layer_groups(self, do_fc=False):
    if do_fc:
        return [self.fc_model]
    idxs = [self.lr_cut]
    c = children(self.top_model)
    if len(c)==3: c = children(c[0])+c[1:]
    lgs = list(split_by_idx(c,idxs))
    return lgs+[self.fc_model]

class ConvLearner(Learner):
    def __init__(self, data, models, precompute=False, **kwargs):
        self.precompute = False
        super().__init__(data, models, **kwargs)
        if hasattr(data, 'is_multi') and not data.is_reg and self.metrics is None:
            self.metrics = [accuracy_thresh(0.5)] if self.data.is_multi else [accuracy]
        if precompute: self.save_fcl()
        self.freeze()
        self.precompute = precompute

    def _get_crit(self, data):
        if not hasattr(data, 'is_multi'): return super()._get_crit(data)

        return F.l1_loss if data.is_reg else F.binary_cross_entropy if data.is_multi else F.nll_

    @classmethod
    def pretrained(cls, f, data, ps=None, xtra_fc=None, xtra_cut=0, custom_head=None, precompute
                    pretrained=True, **kwargs):
        models = ConvnetBuilder_custom(f, data.c, data.is_multi, data.is_reg,
            ps=ps, xtra_fc=xtra_fc, xtra_cut=xtra_cut, custom_head=custom_head, pretrained=pretr
        return cls(data, models, precompute, **kwargs)

    @classmethod
    def lsuv_learner(cls, f, data, ps=None, xtra_fc=None, xtra_cut=0, custom_head=None, precompu
                      needed_std=1.0, std_tol=0.1, max_attempts=10, do_orthonorm=False, **kwargs):
        models = ConvnetBuilder(f, data.c, data.is_multi, data.is_reg,
            ps=ps, xtra_fc=xtra_fc, xtra_cut=xtra_cut, custom_head=custom_head, pretrained=False
        convlearn=cls(data, models, precompute, **kwargs)
        convlearn.lsuv_init()
        return convlearn

    @property
    def model(self): return self.models.fc_model if self.precompute else self.models.model

    def half(self):
        if self.fp16: return
        self.fp16 = True
        if type(self.model) != FP16: self.models.model = FP16(self.model)
        if not isinstance(self.models.fc_model, FP16): self.models.fc_model = FP16(self.models.f

    def float(self):
        if not self.fp16: return
        self.fp16 = False
        if type(self.models.model) == FP16: self.models.model = self.model.module.float()
        if type(self.models.fc_model) == FP16: self.models.fc_model = self.models.fc_model.modul

    @property
    def data(self): return self.fc_data if self.precompute else self.data_

    def create_empty_bcolz(self, n, name):
        return bcolz.carray(np.zeros((0,n), np.float32), chunklen=1, mode='w', rootdir=name)

```

```

def set_data(self, data, precompute=False):
    super().set_data(data)
    if precompute:
        self.unfreeze()
        self.save_fcl()
        self.freeze()
        self.precompute = True
    else:
        self.freeze()

def get_layer_groups(self):
    return self.models.get_layer_groups(self.precompute)

def summary(self):
    precompute = self.precompute
    self.precompute = False
    res = super().summary()
    self.precompute = precompute
    return res

def get_activations(self, force=False):
    tpl = f'_{self.models.name}_{self.data.sz}.bc'
    names = [os.path.join(self.tmp_path, p+tpl) for p in ('x_act', 'x_act_val', 'x_act_test')]
    if os.path.exists(names[0]) and not force:
        self.activations = [bcolz.open(p) for p in names]
    else:
        self.activations = [self.create_empty_bcolz(self.models.nf, n) for n in names]

def save_fcl(self):
    self.get_activations()
    act, val_act, test_act = self.activations
    m=self.models.top_model
    if len(self.activations[0])!=len(self.data.trn_ds):
        predict_to_bcolz(m, self.data.fix_dl, act)
    if len(self.activations[1])!=len(self.data.val_ds):
        predict_to_bcolz(m, self.data.val_dl, val_act)
    if self.data.test_dl and (len(self.activations[2])!=len(self.data.test_ds)):
        if self.data.test_dl: predict_to_bcolz(m, self.data.test_dl, test_act)

    self.fc_data = ImageClassifierData.from_arrays(self.data.path,
        (act, self.data.trn_y), (val_act, self.data.val_y), self.data.bs, classes=self.d
        test = test_act if self.data.test_dl else None, num_workers=8)

def freeze(self):
    self.freeze_to(-1)

def unfreeze(self):
    self.freeze_to(0)
    self.precompute = False

def predict_array(self, arr):
    precompute = self.precompute
    self.precompute = False
    pred = super().predict_array(arr)
    self.precompute = precompute
    return pred

```

```
In [15]: sz = 512 #image size : 256 -> 512
bs = 16 #batch size : 64 -> 16

md = get_data(sz,bs)
learner = ConvLearner.pretrained(arch, md, ps=0.5) #dropout 50%
# cnn을 생성준비와 동시에 weight를 할당합니다. 이때 resnet-34 모델을 이용하고 드롭아웃은 50% 입니다.
learner.opt_fn = optim.Adam
# optimizer로 Adam을 사용합니다.
learner.clip = 1.0
# 모델이 깊은 학습 모델이기에 gradient 폭발 문제를 해결하기 위해 Gradient Clipping을 설정해줍니다.
learner.crit = FocalLoss()
# unblanced 문제를 해결하기 위해 FocalLoss를 사용합니다.
learner.metrics = [acc]
```

```
In [2]: learner.lr_find()
learner.sched.plot()
```

```
In [3]: lr = 2e-2
learner.fit(lr,1)
```

```
In [19]: learner.unfreeze()
# 레이어의 고정을 해제 하여 새로운 학습률을 이용할 준비를 합니다.
lrs=np.array([lr/10,lr/3,lr])
```

```
In [4]: learner.fit(lrs/4,4,cycle_len=2,use_clr=(10,20))
# 새로운 학습률로 학습을 합니다. 이는 학습률에 차등을 주는 것입니다.
```

```
In [5]: learner.fit(lrs/4,2,cycle_len=4,use_clr=(10,20))
# 이때 싸이클의 변화에 따라 더 매끄러운 지점을 찾을 수 있습니다.
```

```
In [ ]: learner.fit(lrs/16,1,cycle_len=8,use_clr=(5,20))
```

```
In [ ]: def sigmoid_np(x):
        return 1.0/(1.0 + np.exp(-x))

preds,y = learner.TTA(n_aug=16)
preds = np.stack(preds, axis=-1)
preds = sigmoid_np(preds)
pred = preds.max(axis=-1)
```

```
In [ ]: def F1_soft(preds,targs,th=0.5,d=50.0):
        preds = sigmoid_np(d*(preds - th))
        targs = targs.astype(np.float)
        score = 2.0*(preds*targs).sum(axis=0)/((preds+targs).sum(axis=0) + 1e-6)
        return score

def fit_val(x,y):
    params = 0.5*np.ones(len(name_label_dict))
    wd = 1e-5
    error = lambda p: np.concatenate((F1_soft(x,y,p) - 1.0,
                                       wd*(p - 0.5)), axis=None)
    p, success = opt.leastsq(error, params)
    return p
```

```
In [ ]: th = fit_val(pred,y)
th[th<0.1] = 0.1
print('Thresholds: ',th)
print('F1 macro: ',f1_score(y, pred>th, average='macro'))
print('F1 macro (th = 0.5): ',f1_score(y, pred>0.5, average='macro'))
print('F1 micro: ',f1_score(y, pred>th, average='micro'))
```

```
In [ ]: print('Fractions: ', (pred > th).mean(axis=0))
        print('Fractions (true): ', (y > th).mean(axis=0))
```

```
In [ ]: preds_t, y_t = learner.TTA(n_aug=16, is_test=True)
        preds_t = np.stack(preds_t, axis=-1)
        preds_t = sigmoid_np(preds_t)
        pred_t = preds_t.max(axis=-1) #max works better for F1 macro score
```

```
In [ ]: def save_pred(pred, th=0.5, fname='protein_classification10.csv'):
        pred_list = []
        for line in pred:
            s = ' '.join([str(i) for i in np.nonzero(line>th)[0]])
            pred_list.append(s)

        sample_df = pd.read_csv(SAMPLE)
        sample_list = list(sample_df.Id)
        pred_dic = dict((key, value) for (key, value)
                        in zip(learner.data.test_ds.fnames, pred_list))
        pred_list_cor = [pred_dic[id] for id in sample_list]
        df = pd.DataFrame({'Id': sample_list, 'Predicted': pred_list_cor})
        df.to_csv(fname, header=True, index=False)
```

```
In [ ]: th_t = np.array([0.565, 0.39, 0.55, 0.345, 0.33, 0.39, 0.33, 0.45, 0.38, 0.39,
                        0.34, 0.42, 0.31, 0.38, 0.49, 0.50, 0.38, 0.43, 0.46, 0.40,
                        0.39, 0.505, 0.37, 0.47, 0.41, 0.545, 0.32, 0.1])
        print('Fractions: ', (pred_t > th_t).mean(axis=0))
        save_pred(pred_t, th_t)
```

```
In [ ]: lb_prob = [
0.362397820, 0.043841336, 0.075268817, 0.059322034, 0.075268817,
0.075268817, 0.043841336, 0.075268817, 0.010000000, 0.010000000,
0.010000000, 0.043841336, 0.043841336, 0.014198783, 0.043841336,
0.010000000, 0.028806584, 0.014198783, 0.028806584, 0.059322034,
0.010000000, 0.126126126, 0.028806584, 0.075268817, 0.010000000,
0.222493880, 0.028806584, 0.010000000]
```

```
In [ ]: def Count_soft(preds, th=0.5, d=50.0):
        preds = sigmoid_np(d*(preds - th))
        return preds.mean(axis=0)

def fit_test(x, y):
    params = 0.5*np.ones(len(name_label_dict))
    wd = 1e-5
    error = lambda p: np.concatenate((Count_soft(x, p) - y,
                                      wd*(p - 0.5)), axis=None)
    p, success = opt.leastsq(error, params)
    return p
```

```
In [ ]: th_t = fit_test(pred_t, lb_prob)
        th_t[th_t<0.1] = 0.1
        print('Thresholds: ', th_t)
        print('Fractions: ', (pred_t > th_t).mean(axis=0))
        print('Fractions (th = 0.5): ', (pred_t > 0.5).mean(axis=0))
```

```
In [ ]: save_pred(pred_t, th_t, 'protein_classification10_f.csv')
```

```
In [ ]: save_pred(pred_t, th, 'protein_classification10_v.csv')
        save_pred(pred_t, 0.5, 'protein_classification10_05.csv')
```



```
In [ ]: class_list = [8,9,10,15,20,24,27]
        for i in class_list:
            th_t[i] = th[i]
        save_pred(pred_t,th_t,'protein_classification8_c.csv')
```

```
In [ ]: labels = pd.read_csv(LABELS).set_index('Id')
        label_count = np.zeros(len(name_label_dict))
        for label in labels['Target']:
            l = [int(i) for i in label.split()]
            label_count += np.eye(len(name_label_dict))[l].sum(axis=0)
        label_fraction = label_count.astype(np.float)/len(labels)
        label_count, label_fraction
```

```
In [ ]: th_t = fit_test(pred_t,label_fraction)
        th_t[th_t<0.05] = 0.05
        print('Thresholds: ',th_t)
        print('Fractions: ',(pred_t > th_t).mean(axis=0))
        save_pred(pred_t,th_t,'protein_classification8_t.csv')
```

