```python
1  import argparse
2  from pathlib import Path
3  from collections import defaultdict
4  from datetime import datetime
5
6  import numpy as np
7
8  from my_logger import get_my_logger
9  from utils import load_data, Trainer, plot_accs, tr_va_split
10 from preprocess import ImagePreprocessor, TextPreprocessor
11 from save_csv import results_to_csv
12
13 if __name__ == "__main__":
14     # logger
15     logger = get_my_logger(Path(__file__).name)
16     logger.info('--- [START {0:%Y-%m-%d_%H-%M-%S}] 
   {1}'.format(datetime.now(), '-' * 100))
17
18     parser = argparse.ArgumentParser(description='cs289 hw1')
19     parser.add_argument('problem', type=int, help='problem number (e.g. 2, 3, 
   4)')
20     parser.add_argument('data_name', type=str, help='dataset name (e.g. 
   mnist, spam, cifar10)')
21     args = parser.parse_args()
22
23     # -- configuration ------------------
24     data_name = args.data_name  # cifar10 mnist spam
25
26     split_and_save = 0
27     preprocess = 0
28     tune_n_tr = 0
29     tune_hparameter = 0
30     train_and_predict = 0
31     split = 'holdout'  # kfold holdout
32     n_fold = 1
33     kernel = 'linear'
34
35     n_trs = {
36         'mnist': None,  # 5000
37         'spam': None,
38         'cifar10': None  # 5000
39     }
40
41     if args.problem == 2:
42         split_and_save = 1
43     elif args.problem == 3:
44         tune_n_tr = 1
45     elif args.problem == 4:
46         tune_hparameter = 1
47         n_trs['mnist'] = 10000
48     elif args.problem == 5:
49         tune_hparameter = 1
50         split = 'kfold'
51         n_fold = 5
52     elif args.problem == 6:
53         preprocess = 1
54         train_and_predict = 1
55         split = 'kfold'
56         n_fold = 5
57         kernel = 'rbf'
```

```python
58
59      logger.info(f'data:\t{data_name}')
60      logger.info(f'prepro:\t{preprocess}')
61      logger.info(f'split:\t{split}')
62      logger.info(f'n_fold:\t{n_fold}')
63
64      # seed
65      seed = 189
66      np.random.seed(seed)
67      logger.info(f'seed:\t{seed}')
68
69      n_vas = {
70          'mnist': 10000,
71          'spam': 0.2,
72          'cifar10': 5000
73      }
74      n_va = n_vas[data_name]
75
76      preprocessors = {
77          'mnist': ImagePreprocessor(
78              normalize=False, hog_tf=True, pca_tf=False, lbp_tf=False,
    win_size=(28, 28), block_size=(8, 8),
79              block_stride=(4, 4), cell_size=(8, 8), nbins=9
80          ),
81          'spam': TextPreprocessor(),
82          'cifar10': ImagePreprocessor(
83              normalize=False, hog_tf=False, pca_tf=True, lbp_tf=False
84          )
85      }
86      preprocessor = preprocessors[data_name]
87
88      n_tr_spaces = {
89          'mnist': [100, 200, 500, 1000, 2000, 5000, 10000],
90          'spam': [100, 200, 500, 1000, 2000, 4138],
91          'cifar10': [100, 200, 500, 1000, 2000, 5000]
92      }
93      n_tr_space = n_tr_spaces[data_name]
94
95      all_params = {
96          'mnist': {
97              'C': 2e1,
98              'kernel': kernel,
99              'gamma': 'scale',
100             'random_state': seed
101         },
102         'spam': {
103             'C': 1e1,
104             'kernel': kernel,
105             'gamma': 'scale',
106             'random_state': seed
107         },
108         'cifar10': {
109             'C': 1e1,
110             'kernel': kernel,
111             'gamma': 'scale',
112             'random_state': seed
113         }
114     }
115     params = all_params[data_name]
116
```

```python
117    all_parameter_spaces = {
118        'mnist': {
119            'C': np.logspace(-10, 10, 11, base=10)
120        },
121        'spam': {
122            'C': np.logspace(-10, 4, 8, base=10)
123        },
124        'cifar10': {
125            'C': np.logspace(-10, 10, 11, base=10)
126        }
127    }
128    parameter_spaces = all_parameter_spaces[data_name]
129    # ---------------------------------

131    # load data
132    x_train, x_test, y_train = load_data(data_name)
133    logger.info(f'x_train:\t{x_train.shape}')
134    logger.info(f'y_train:\t{y_train.shape}')
135    logger.info(f'x_test:\t{x_test.shape}')

137    # preprocessing
138    if preprocess:
139        logger.info('** preprocess **')
140        x_train, x_test = preprocessor.scale(x_train, x_test)
141        logger.info(f'x_train:\t{x_train.shape}')
142        logger.info(f'x_test:\t{x_test.shape}')

144    # tune training data size
145    if tune_n_tr:
146        logger.info('** tune training data size**')
147        accs = defaultdict(dict)
148        for n_tr in n_tr_space:
149            trainer = Trainer(params, x_train, x_test, y_train,
    verbose=False)
150            trainer.validate(split, n_fold=n_fold, n_va=n_va, n_tr=n_tr,
    shuffle=True)

152            tr_acc = trainer.scores_mean['tr_acc']
153            va_acc = trainer.scores_mean['va_acc']
154            accs['train'][n_tr] = tr_acc
155            accs['valid'][n_tr] = va_acc
156            logger.info(f'[{n_tr}]\ttr:{tr_acc:.5f}\tva:{va_acc:.5f}')

158        path = Path(f'figures/{data_name}_acc_vs_n_tr.png')
159        xlabel = 'Training data size'
160        xlim = (0, max(n_tr_space) * 1.1)
161        plot_accs(accs, data_name, path, xlabel, xlim, False)

163    # tune C of SVM with holdout
164    if tune_hparameter:
165        logger.info('** tune hyper parameter **')
166        n_tr = n_trs[data_name]
167        best_acc = 0
168        best_params = {}

170        for parameter, space in parameter_spaces.items():
171            accs = defaultdict(dict)
172            for value in space:
173                params[parameter] = value
```

```python
174                    trainer = Trainer(params, x_train, x_test, y_train,
       verbose=False)
175                    trainer.validate(split, n_fold=n_fold, n_va=n_va, n_tr=n_tr,
       shuffle=True)
176
177                    tr_acc = trainer.scores_mean['tr_acc']
178                    va_acc = trainer.scores_mean['va_acc']
179                    accs['train'][value] = tr_acc
180                    accs['valid'][value] = va_acc
181                    logger.info(f'[{parameter}:{value:.0e}]\ttr:{tr_acc:.5f}\tva:
       {va_acc:.5f}')
182
183                    if va_acc > best_acc:
184                        best_params = params.copy()
185                        best_acc = va_acc
186
187                path = Path(f'figures/{data_name}_acc_vs_{parameter}.png')
188                xlabel = parameter
189                xlim = (min(space) * 0.9, max(space) * 1.1)
190                plot_accs(accs, data_name, path, xlabel, xlim, log_scale=True)
191
192            params = best_params.copy()
193            logger.info(f'[best_acc]\t{best_acc:.5f}')
194            logger.info(f'[best_params]\t{best_params}')
195
196        # split and save
197        if split_and_save:
198            logger.info('** split and save **')
199            n_tr = n_trs[data_name]
200            cv = tr_va_split(x_train, x_test, n_va, shuffle=True)
201            ind_tr, ind_va = next(cv)
202            x_tr, x_va = x_train[ind_tr[:n_tr]], x_train[ind_va]
203            y_tr, y_va = y_train[ind_tr[:n_tr]], y_train[ind_va]
204
205            logger.info(f'x_tr: {x_tr.shape}')
206            logger.info(f'x_va: {x_va.shape}')
207            logger.info(f'y_tr: {y_tr.shape}')
208            logger.info(f'y_va: {y_va.shape}')
209
210            save_dir = Path('data', 'split')
211            np.save(Path(save_dir, f'{data_name}_x_tr.npy'), x_tr)
212            np.save(Path(save_dir, f'{data_name}_x_va.npy'), x_va)
213            np.save(Path(save_dir, f'{data_name}_y_tr.npy'), y_tr)
214            np.save(Path(save_dir, f'{data_name}_y_va.npy'), y_va)
215
216        # train and predict
217        if train_and_predict:
218            logger.info('** train and predict **')
219            n_tr = n_trs[data_name]
220
221            trainer = Trainer(params, x_train, x_test, y_train, verbose=True)
222            trainer.validate(split, n_fold=n_fold, n_va=n_va, n_tr=n_tr,
       shuffle=True)
223
224            tr_acc = trainer.scores_mean['tr_acc']
225            va_acc = trainer.scores_mean['va_acc']
226            logger.info(f'[all]\ttr:{tr_acc:.5f}\tva:{va_acc:.5f}')
227
228            submit_path =
       Path(f'data/output/submit_{data_name}_{va_acc:.5f}.csv', index=False)
```

```
229        results_to_csv(trainer.y_test_pred, submit_path)
230
```