# FedAvg算法+LSTM模型+ Shakespeare数据集——字符预测任务

## 1. Shakespeare数据集介绍

**任务**：下一个字符预测

**参数说明**：总共4,226,15条样本，可使用官方给出的划分代码按照联邦学习场景下1129个client非独立同分布划分。

**介绍**: 和FEMNST一样，属于专门给联邦学习用的基准数据集leaf的成员之一。

**官网**：https://leaf.cmu.edu/

**官方数据预处理与划分代码**：https://github.com/TalwalkarLab

**引用方式**：LEAF: A Benchmark for Federated Settings

从 http://www.gutenberg.org/files/100/old/1994-01-100.zip 下载数据集解压到data文件夹下的raw_data文件夹，重命名为raw_data.txt

我们从https://github.com/TalwalkarLab下载数据集后，找到对应的Shakespeare数据集文件夹，按照README文件即可成果转换成我们想要的数据集，操作命令如下：

```
1  ./preprocess.sh -s niid --sf 1.0 -k 0 -t sample -tf 0.8 (full-sized dataset)
2  ./preprocess.sh -s niid --sf 0.2 -k 0 -t sample -tf 0.8 (small-sized dataset)
   ('-tf 0.8' reflects the train-test split used in the [FedAvg paper]
   (https://arxiv.org/pdf/1602.05629.pdf))
```

至此，我们就能获得预处理后的非独立同分布的Shakespeare数据集，每个数据集下主要包括内容如下：

每个字符串x长度为80，预测下一个字符的输出，总共分为了有135个客户端数据



## 2. LSTM训练模型

```python
1  class Model(nn.Module):
2      def __init__(self, seed, lr, optimizer=None):
3          super().__init__()
4          self.lr = lr
5          self.seed = seed
6          self.optimizer = optimizer
7          self.flops = 0
8          self.size = 0
9
10     def get_params(self):
11         return self.state_dict()
12
13     def set_params(self, state_dict):
14         self.load_state_dict(state_dict)
```

```python
15
16      def __post_init__(self):
17          if self.optimizer is None:
18              self.optimizer = optim.SGD(self.parameters(), lr=self.lr)
19
20      def train_model(self, data, num_epochs=1, batch_size=10):
21          self.train()
22          for batch in range(num_epochs):
23              for batched_x, batched_y in batch_data(data, batch_size,
    seed=self.seed):
24                  self.optimizer.zero_grad()
25                  input_data = self.process_x(batched_x)
26                  target_data = self.process_y(batched_y)
27                  logits, loss = self.forward(input_data, target_data)
28                  loss.backward()
29                  self.optimizer.step()
30          update = self.get_params()
31          comp = num_epochs * (len(data['y']) // batch_size) * batch_size
32          return comp, update
33
34      def test_model(self, data):
35          x_vecs = self.process_x(data['x'])
36          labels = self.process_y(data['y'])
37          self.eval()
38          with torch.no_grad():
39              logits, loss = self.forward(x_vecs, labels)
40              acc = assess_fun(labels, logits)
41          return {"accuracy": acc.detach().cpu().numpy(), 'loss':
    loss.detach().cpu().numpy()}
42
43
44  class LSTMModel(Model):
45      def __init__(self, seed, lr, seq_len, num_classes, n_hidden):
46          super().__init__(seed, lr)
47          self.seq_len = seq_len
48          self.num_classes = num_classes
49          self.n_hidden = n_hidden
50          self.word_embedding = nn.Embedding(self.num_classes, 8)
51          self.lstm = nn.LSTM(input_size=8, hidden_size=self.n_hidden,
    num_layers=2, batch_first=True)
52          self.pred = nn.Linear(self.n_hidden * 2, self.num_classes)
53          self.loss_fn = nn.CrossEntropyLoss()
54          super().__post_init__()
55
56      def forward(self, features, labels):
57          emb = self.word_embedding(features)
58          output, (h_n, c_n) = self.lstm(emb)
59          h_n = h_n.transpose(0, 1).reshape(-1, 2 * self.n_hidden)
60          logits = self.pred(h_n)
61          loss = self.loss_fn(logits, labels)
62          return logits, loss
63
64      def process_x(self, raw_x_batch):
65          x_batch = [word_to_indices(word) for word in raw_x_batch]
66          x_batch = torch.LongTensor(x_batch)
67          return x_batch
68
69      def process_y(self, raw_y_batch):
```

```
70          y_batch = [letter_to_vec(c) for c in raw_y_batch]
71          y_batch = torch.LongTensor(y_batch)
72          return y_batch
```

## 3. Client模型

```
1   class Client:
2       def __init__(self, client_id, train_data, eval_data, model=None):
3           self._model = model
4           self.id = client_id
5           self.train_data = train_data if train_data is not None else {'x':
    [], 'y': []}
6           self.eval_data = eval_data if eval_data is not None else {'x': [],
    'y': []}
7
8       @property
9       def model(self):
10          return self._model
11
12      @property
13      def num_test_samples(self):
14          if self.eval_data is None:
15              return 0
16          else:
17              return len(self.eval_data['y'])
18
19      @property
20      def num_train_samples(self):
21          if self.train_data is None:
22              return 0
23          else:
24              return len(self.train_data['y'])
25
26      @property
27      def num_samples(self):
28          train_size = 0
29          if self.train_data is not None:
30              train_size = len(self.train_data['y'])
31          eval_size = 0
32          if self.eval_data is not None:
33              eval_size = len(self.eval_data['y'])
34          return train_size + eval_size
35
36      def train(self, num_epochs=1, batch_size=128, minibatch=None):
37          if minibatch is None:
38              data = self.train_data
39              comp, update = self.model.train_model(data, num_epochs,
    batch_size)
40          else:
41              frac = min(1.0, minibatch)
42              num_data = max(1, int(frac * len(self.train_data['y'])))
43              xs, xy = zip(*random.sample(list(zip(self.train_data['x'],
    self.train_data['y'])), num_data))
44              data = {
45                  'x': xs,
46                  'y': xy
```

```
47                }
48                num_epochs = 1
49                comp, update = self.model.train_model(data, num_epochs,
      num_data)
50            num_train_samples = len(data['y'])
51            return comp, num_train_samples, update
52
53        def test(self, set_to_use='test'):
54            assert set_to_use in ['train', 'test', 'val']
55            if set_to_use == 'train':
56                data = self.train_data
57            else:
58                data = self.eval_data
59            return self.model.test_model(data)
```

## 4. Serves模型

```
1    class Serves:
2        def __init__(self, global_model):
3            self.global_model = global_model
4            self.model = global_model.get_params()
5            self.selected_clients = []
6            self.update = []
7
8        def select_clients(self, my_round, possible_clients, num_clients=20):
9            num_clients = min(num_clients, len(possible_clients))
10            np.random.seed(my_round)
11            self.selected_clients = np.random.choice(possible_clients,
      num_clients, replace=False)
12            # return [(c.num_train_samples, c.num_test_samples) for c in
      self.selected_clients]
13
14        def train_model(self, num_epochs=1, batch_size=10, minibatch=None,
      clients=None):
15            if clients is None:
16                clients = self.selected_clients
17            sys_metrics = {
18                c.id: {"bytes_written": 0,
19                       "bytes_read": 0,
20                       "local_computations": 0} for c in clients}
21            for c in clients:
22                c.model.set_params(self.model)
23                comp, num_samples, update = c.train(num_epochs, batch_size,
      minibatch)
24                sys_metrics[c.id]["bytes_read"] += c.model.size
25                sys_metrics[c.id]["bytes_written"] += c.model.size
26                sys_metrics[c.id]["local_computations"] = comp
27
28                self.update.append((num_samples, update))
29            return sys_metrics
30
31        def aggregate(self, updates):
32            avg_param = OrderedDict()
33            total_weight = 0.
34            for (client_samples, client_model) in updates:
35                total_weight += client_samples
```

```python
36              for name, param in client_model.items():
37                  if name not in avg_param:
38                      avg_param[name] = client_samples * param
39                  else:
40                      avg_param[name] += client_samples * param

42          for name in avg_param:
43              avg_param[name] = avg_param[name] / total_weight
44          return avg_param

46      def update_model(self):
47          avg_param = self.aggregate(self.update)
48          self.model = avg_param
49          self.global_model.load_state_dict(self.model)
50          self.update = []

52      def test_model(self, clients_to_test=None, set_to_use='test'):
53          metrics = {}
54          if clients_to_test is None:
55              clients_to_test = self.selected_clients

57          for client in tqdm(clients_to_test):
58              client.model.set_params(self.model)
59              c_metrics = client.test(set_to_use)
60              metrics[client.id] = c_metrics

62          return metrics

64      def get_clients_info(self, clients):
65          if clients is None:
66              clients = self.selected_clients

68          ids = [c.id for c in clients]
69          num_samples = {c.id: c.num_samples for c in clients}
70          return ids, num_samples

72      def save_model(self, path):
73          """Saves the server model on checkpoints/dataset/model.ckpt."""
74          return torch.save({"model_state_dict": self.model}, path)
```

# 5. 数据处理工具函数

```python
1  import json
2  import numpy as np
3  import os
4  from collections import defaultdict

6  import torch

8  ALL_LETTERS = "\n !\"&'(),-.0123456789:;>?
   ABCDEFGHIJKLMNOPQRSTUVWXYZ[]abcdefghijklmnopqrstuvwxyz}"
9  NUM_LETTERS = len(ALL_LETTERS)

12 def batch_data(data, batch_size, seed):
13     data_x = data['x']
```

```
14        data_y = data['y']
15
16    np.random.seed(seed)
17    rng_state = np.random.get_state()
18    np.random.shuffle(data_x)
19    np.random.set_state(rng_state)
20    np.random.shuffle(data_y)
21
22    for i in range(0, len(data_x), batch_size):
23        batched_x = data_x[i:i + batch_size]
24        batched_y = data_y[i:i + batch_size]
25        yield (batched_x, batched_y)
26
27
28 def assess_fun(y_true, y_hat):
29    y_hat = torch.argmax(y_hat, dim=-1)
30    total = y_true.shape[0]
31    hit = torch.sum(y_true == y_hat)
32    return hit.data.float() * 1.0 / total
33
34
35 def word_to_indices(word):
36    indices = []
37    for c in word:
38        indices.append(ALL_LETTERS.find(c))
39    return indices
40
41
42 def letter_to_vec(letter):
43    index = ALL_LETTERS.find(letter)
44    return index
```

# 6. 数据读取函数

```
 1 def read_dir(data_dir):
 2    clients = []
 3    data = defaultdict(lambda: None)
 4
 5    files = os.listdir(data_dir)
 6    files = [f for f in files if f.endswith('.json')]
 7    for f in files:
 8        file_path = os.path.join(data_dir, f)
 9        with open(file_path, 'r') as inf:
10            cdata = json.load(inf)
11        clients.extend(cdata['users'])
12        data.update(cdata['user_data'])
13
14    clients = list(sorted(data.keys()))
15    return clients, data
16
17
18 def read_data(train_data_dir, test_data_dir):
19    train_clients, train_data = read_dir(train_data_dir)
20    test_clients, test_data = read_dir(test_data_dir)
21
22    assert train_clients == test_clients
```

```
23
24        return train_clients, train_data, test_data
25
26
27  def create_clients(users, train_data, test_data, model):
28        clients = [Client(u, train_data[u], test_data[u], model) for u in users]
29        return clients
30
31
32  def setup_clients(model=None, use_val_set=False):
33        eval_set = 'test' if not use_val_set else 'test'
34        train_data_dir = os.path.join('.', 'data', 'train')
35        test_data_dir = os.path.join('.', 'data', eval_set)
36        users, train_data, test_data = read_data(train_data_dir, test_data_dir)
37
38        clients = create_clients(users, train_data, test_data, model)
39
40        return clients
```

# 7. 训练函数

```
1   def train():
2         seed = 1
3         random.seed(1 + seed)
4         np.random.seed(12 + seed)
5         torch.manual_seed(123 + seed)
6         torch.manual_seed(123 + seed)
7         lr = 0.0003
8         seq_len = 80
9         num_classes = 80
10        n_hidden = 256
11        num_rounds = 20
12        eval_every = 1
13        clients_per_round = 2
14        num_epochs = 1
15        batch_size = 10
16        minibatch = None
17        use_val_set = 'test'
18        # 全局模型(服务端)
19        global_model = LSTMModel(seed, lr, seq_len, num_classes, n_hidden)
20        # 服务端
21        server = Serves(global_model)
22        # 客户端
23        client_model = LSTMModel(seed, lr, seq_len, num_classes, n_hidden)
24        clients = setup_clients(client_model, use_val_set)
25        client_ids, client_num_samples = server.get_clients_info(clients)
26        print(('Clients in Total: %d' % len(clients)))
27        print('--- Random Initialization ---')
28        # Simulate training
29        for i in range(num_rounds):
30            print('--- Round %d of %d: Training %d Clients ---' % (i + 1,
      num_rounds, clients_per_round))
31
32            # Select clients to train this round
33            server.select_clients(i, clients, num_clients=clients_per_round)
```

```
34        c_ids, c_num_samples =
   server.get_clients_info(server.selected_clients)
35
36        # Simulate server model training on selected clients' data
37        sys_metrics = server.train_model(num_epochs=num_epochs,
   batch_size=batch_size,
38                                        minibatch=minibatch)
39        print(sys_metrics)
40        # sys_writer_fn(i + 1, c_ids, sys_metrics, c_num_samples)
41        metrics = server.test_model()
42        print(metrics)
43
44        # Update server model
45        server.update_model()
```

# 8. 训练结果

**因为设备原因，暂时无法训练出论文中的模型**