```
In [1]:  from argparse import Namespace
         import os
         import json

         import numpy as np
         import pandas as pd
         import torch
         import torch.nn as nn
         import torch.nn.functional as F
         import torch.optim as optim
         from torch.utils.data import Dataset, DataLoader
         from tqdm import tqdm_notebook
```

```
In [2]:  class Vocabulary(object):
             """Class to process text and extract vocabulary for mapping"""

             def __init__(self, token_to_idx=None):
                 """
                 Args:
                     token_to_idx (dict): a pre-existing map of tokens to indices
                 """

                 if token_to_idx is None:
                     token_to_idx = {}
                 self._token_to_idx = token_to_idx

                 self._idx_to_token = {idx: token
                                       for token, idx in self._token_to_idx.items

             def to_serializable(self):
                 """ returns a dictionary that can be serialized """
                 return {'token_to_idx': self._token_to_idx}

             @classmethod
             def from_serializable(cls, contents):
                 """ instantiates the Vocabulary from a serialized dictionary """
                 return cls(**contents)

             def add_token(self, token):
                 """Update mapping dicts based on the token.

                 Args:
                     token (str): the item to add into the Vocabulary
                 Returns:
                     index (int): the integer corresponding to the token
                 """
                 if token in self._token_to_idx:
                     index = self._token_to_idx[token]
```

```python
        else:
            index = len(self._token_to_idx)
            self._token_to_idx[token] = index
            self._idx_to_token[index] = token
        return index

    def add_many(self, tokens):
        """Add a list of tokens into the Vocabulary

        Args:
            tokens (list): a list of string tokens
        Returns:
            indices (list): a list of indices corresponding to the token
        """
        return [self.add_token(token) for token in tokens]

    def lookup_token(self, token):
        """Retrieve the index associated with the token

        Args:
            token (str): the token to look up
        Returns:
            index (int): the index corresponding to the token
        """
        return self._token_to_idx[token]

    def lookup_index(self, index):
        """Return the token associated with the index

        Args:
            index (int): the index to look up
        Returns:
            token (str): the token corresponding to the index
        Raises:
            KeyError: if the index is not in the Vocabulary
        """
        if index not in self._idx_to_token:
            raise KeyError("the index (%d) is not in the Vocabulary" % i
        return self._idx_to_token[index]

    def __str__(self):
        return "<Vocabulary(size=%d)>" % len(self)

    def __len__(self):
        return len(self._token_to_idx)
```

```
In [3]: class SequenceVocabulary(Vocabulary):
            def __init__(self, token_to_idx=None, unk_token="<UNK>",
                         mask_token="<MASK>", begin_seq_token="<BEGIN>",
                         end_seq_token="<END>"):

                super(SequenceVocabulary, self).__init__(token_to_idx)

                self._mask_token = mask_token
                self._unk_token = unk_token
                self._begin_seq_token = begin_seq_token
                self._end_seq_token = end_seq_token

                self.mask_index = self.add_token(self._mask_token)
                self.unk_index = self.add_token(self._unk_token)
                self.begin_seq_index = self.add_token(self._begin_seq_token)
                self.end_seq_index = self.add_token(self._end_seq_token)

            def to_serializable(self):
                contents = super(SequenceVocabulary, self).to_serializable()
                contents.update({'unk_token': self._unk_token,
                                 'mask_token': self._mask_token,
                                 'begin_seq_token': self._begin_seq_token,
                                 'end_seq_token': self._end_seq_token})
                return contents

            def lookup_token(self, token):
                """Retrieve the index associated with the token
                  or the UNK index if token isn't present.

                Args:
                    token (str): the token to look up
                Returns:
                    index (int): the index corresponding to the token
                Notes:
                    `unk_index` needs to be >=0 (having been added into the Voca
                        for the UNK functionality
                """
                if self.unk_index >= 0:
                    return self._token_to_idx.get(token, self.unk_index)
                else:
                    return self._token_to_idx[token]


In [4]: class SurnameVectorizer(object):
            """ The Vectorizer which coordinates the Vocabularies and puts them
            def __init__(self, char_vocab, nationality_vocab):
                """
                Args:
                    char_vocab (Vocabulary): maps characters to integers
                    nationality_vocab (Vocabulary): maps nationalities to intege
```

```python
        """
        self.char_vocab = char_vocab
        self.nationality_vocab = nationality_vocab

    def vectorize(self, surname, vector_length=-1):
        """
        Args:
            title (str): the string of characters
            vector_length (int): an argument for forcing the length of i
        """
        indices = [self.char_vocab.begin_seq_index]
        indices.extend(self.char_vocab.lookup_token(token)
                       for token in surname)
        indices.append(self.char_vocab.end_seq_index)

        if vector_length < 0:
            vector_length = len(indices)

        out_vector = np.zeros(vector_length, dtype=np.int64)
        out_vector[:len(indices)] = indices
        out_vector[len(indices):] = self.char_vocab.mask_index

        return out_vector, len(indices)

    @classmethod
    def from_dataframe(cls, surname_df):
        """Instantiate the vectorizer from the dataset dataframe

        Args:
            surname_df (pandas.DataFrame): the surnames dataset
        Returns:
            an instance of the SurnameVectorizer
        """
        char_vocab = SequenceVocabulary()
        nationality_vocab = Vocabulary()

        for index, row in surname_df.iterrows():
            for char in row.surname:
                char_vocab.add_token(char)
            nationality_vocab.add_token(row.nationality)

        return cls(char_vocab, nationality_vocab)

    @classmethod
    def from_serializable(cls, contents):
        char_vocab = SequenceVocabulary.from_serializable(contents['char_
        nat_vocab =  Vocabulary.from_serializable(contents['nationality_

        return cls(char_vocab=char_vocab, nationality_vocab=nat_vocab)

    def to_serializable(self):
```

```python
    def to_serializable(self):
        return {'char_vocab': self.char_vocab.to_serializable(),
                'nationality_vocab': self.nationality_vocab.to_serializab
```

```python
In [5]: class SurnameDataset(Dataset):
            def __init__(self, surname_df, vectorizer):
                """
                Args:
                    surname_df (pandas.DataFrame): the dataset
                    vectorizer (SurnameVectorizer): vectorizer instatiated from
                """
                self.surname_df = surname_df
                self._vectorizer = vectorizer

                self._max_seq_length = max(map(len, self.surname_df.surname)) +

                self.train_df = self.surname_df[self.surname_df.split=='train']
                self.train_size = len(self.train_df)

                self.val_df = self.surname_df[self.surname_df.split=='val']
                self.validation_size = len(self.val_df)

                self.test_df = self.surname_df[self.surname_df.split=='test']
                self.test_size = len(self.test_df)

                self._lookup_dict = {'train': (self.train_df, self.train_size),
                                     'val': (self.val_df, self.validation_size),
                                     'test': (self.test_df, self.test_size)}

                self.set_split('train')

                # Class weights
                class_counts = self.train_df.nationality.value_counts().to_dict(
                def sort_key(item):
                    return self._vectorizer.nationality_vocab.lookup_token(item[
                sorted_counts = sorted(class_counts.items(), key=sort_key)
                frequencies = [count for _, count in sorted_counts]
                self.class_weights = 1.0 / torch.tensor(frequencies, dtype=torch

            @classmethod
            def load_dataset_and_make_vectorizer(cls, surname_csv):
                """Load dataset and make a new vectorizer from scratch

                Args:
                    surname_csv (str): location of the dataset
                Returns:
                    an instance of SurnameDataset
                """
                surname_df = pd.read_csv(surname_csv)
                train_surname_df = surname_df[surname_df.split=='train']
```

```python
        train_surname_df = surname_df[surname_df.split== 'train']
        return cls(surname_df, SurnameVectorizer.from_dataframe(train_su

    @classmethod
    def load_dataset_and_load_vectorizer(cls, surname_csv, vectorizer_fi
        """Load dataset and the corresponding vectorizer.
        Used in the case in the vectorizer has been cached for re-use

        Args:
            surname_csv (str): location of the dataset
            vectorizer_filepath (str): location of the saved vectorizer
        Returns:
            an instance of SurnameDataset
        """
        surname_df = pd.read_csv(surname_csv)
        vectorizer = cls.load_vectorizer_only(vectorizer_filepath)
        return cls(surname_df, vectorizer)

    @staticmethod
    def load_vectorizer_only(vectorizer_filepath):
        """a static method for loading the vectorizer from file

        Args:
            vectorizer_filepath (str): the location of the serialized ve
        Returns:
            an instance of SurnameVectorizer
        """
        with open(vectorizer_filepath) as fp:
            return SurnameVectorizer.from_serializable(json.load(fp))

    def save_vectorizer(self, vectorizer_filepath):
        """saves the vectorizer to disk using json

        Args:
            vectorizer_filepath (str): the location to save the vectoriz
        """
        with open(vectorizer_filepath, "w") as fp:
            json.dump(self._vectorizer.to_serializable(), fp)

    def get_vectorizer(self):
        """ returns the vectorizer """
        return self._vectorizer

    def set_split(self, split="train"):
        self._target_split = split
        self._target_df, self._target_size = self._lookup_dict[split]

    def __len__(self):
        return self._target_size

    def __getitem__(self, index):
```

```python
        """the primary entry point method for PyTorch datasets

        Args:
            index (int): the index to the data point
        Returns:
            a dictionary holding the data point's:
                features (x_data)
                label (y_target)
                feature length (x_length)
        """
        row = self._target_df.iloc[index]

        surname_vector, vec_length = \
            self._vectorizer.vectorize(row.surname, self._max_seq_length

        nationality_index = \
            self._vectorizer.nationality_vocab.lookup_token(row.national

        return {'x_data': surname_vector,
                'y_target': nationality_index,
                'x_length': vec_length}

    def get_num_batches(self, batch_size):
        """Given a batch size, return the number of batches in the datas

        Args:
            batch_size (int)
        Returns:
            number of batches in the dataset
        """
        return len(self) // batch_size


def generate_batches(dataset, batch_size, shuffle=True,
                     drop_last=True, device="cpu"):
    """
    A generator function which wraps the PyTorch DataLoader. It will
      ensure each tensor is on the write device location.
    """
    dataloader = DataLoader(dataset=dataset, batch_size=batch_size,
                            shuffle=shuffle, drop_last=drop_last)

    for data_dict in dataloader:
        out_data_dict = {}
        for name, tensor in data_dict.items():
            out_data_dict[name] = data_dict[name].to(device)
        yield out_data_dict
```

```
In [6]: def column_gather(y_out, x_lengths):
```

```python
        '''Get a specific vector from each batch datapoint in `y_out`.

        More precisely, iterate over batch row indices, get the vector that'
        the position indicated by the corresponding value in `x_lengths` at
        index.

        Args:
            y_out (torch.FloatTensor, torch.cuda.FloatTensor)
                shape: (batch, sequence, feature)
            x_lengths (torch.LongTensor, torch.cuda.LongTensor)
                shape: (batch,)

        Returns:
            y_out (torch.FloatTensor, torch.cuda.FloatTensor)
                shape: (batch, feature)
        '''
        x_lengths = x_lengths.long().detach().cpu().numpy() - 1

        out = []
        for batch_index, column_index in enumerate(x_lengths):
            out.append(y_out[batch_index, column_index])

        return torch.stack(out)


class ElmanRNN(nn.Module):
    """ an Elman RNN built using the RNNCell """
    def __init__(self, input_size, hidden_size, batch_first=False):
        """
        Args:
            input_size (int): size of the input vectors
            hidden_size (int): size of the hidden state vectors
            bathc_first (bool): whether the 0th dimension is batch
        """
        super(ElmanRNN, self).__init__()

        self.rnn_cell = nn.RNNCell(input_size, hidden_size)

        self.batch_first = batch_first
        self.hidden_size = hidden_size

    def _initial_hidden(self, batch_size):
        return torch.zeros((batch_size, self.hidden_size))

    def forward(self, x_in, initial_hidden=None):
        """The forward pass of the ElmanRNN

        Args:
            x_in (torch.Tensor): an input data tensor.
                If self.batch_first: x_in.shape = (batch, seq_size, feat
```

```python
                Else: x_in.shape = (seq_size, batch, feat_size)
            initial_hidden (torch.Tensor): the initial hidden state for
        Returns:
            hiddens (torch.Tensor): The outputs of the RNN at each time
                If self.batch_first: hiddens.shape = (batch, seq_size, h
                Else: hiddens.shape = (seq_size, batch, hidden_size)
        """
        if self.batch_first:
            batch_size, seq_size, feat_size = x_in.size()
            x_in = x_in.permute(1, 0, 2)
        else:
            seq_size, batch_size, feat_size = x_in.size()

        hiddens = []

        if initial_hidden is None:
            initial_hidden = self._initial_hidden(batch_size)
            initial_hidden = initial_hidden.to(x_in.device)

        hidden_t = initial_hidden

        for t in range(seq_size):
            hidden_t = self.rnn_cell(x_in[t], hidden_t)
            hiddens.append(hidden_t)

        hiddens = torch.stack(hiddens)

        if self.batch_first:
            hiddens = hiddens.permute(1, 0, 2)

        return hiddens



class SurnameClassifier(nn.Module):
    """ A Classifier with an RNN to extract features and an MLP to class
    def __init__(self, embedding_size, num_embeddings, num_classes,
                 rnn_hidden_size, batch_first=True, padding_idx=0):
        """
        Args:
            embedding_size (int): The size of the character embeddings
            num_embeddings (int): The number of characters to embed
            num_classes (int): The size of the prediction vector
                Note: the number of nationalities
            rnn_hidden_size (int): The size of the RNN's hidden state
            batch_first (bool): Informs whether the input tensors will
                have batch or the sequence on the 0th dimension
            padding_idx (int): The index for the tensor padding;
                see torch.nn.Embedding
        """
        super(SurnameClassifier, self).__init__()
```

Assignment7b_RNN - Jupyter Notebook 3/8/20, 9:52 PM

```python
        super(SurnameClassifier, self).__init__()

        self.emb = nn.Embedding(num_embeddings=num_embeddings,
                                embedding_dim=embedding_size,
                                padding_idx=padding_idx)
        self.rnn = ElmanRNN(input_size=embedding_size,
                            hidden_size=rnn_hidden_size,
                            batch_first=batch_first)
        self.fc1 = nn.Linear(in_features=rnn_hidden_size,
                        out_features=rnn_hidden_size)
        self.fc2 = nn.Linear(in_features=rnn_hidden_size,
                        out_features=num_classes)

    def forward(self, x_in, x_lengths=None, apply_softmax=False):
        """The forward pass of the classifier

        Args:
            x_in (torch.Tensor): an input data tensor.
                x_in.shape should be (batch, input_dim)
            x_lengths (torch.Tensor): the lengths of each sequence in th
                They are used to find the final vector of each sequence
            apply_softmax (bool): a flag for the softmax activation
                should be false if used with the Cross Entropy losses
        Returns:
            the resulting tensor. tensor.shape should be (batch, output_
        """
        x_embedded = self.emb(x_in)
        y_out = self.rnn(x_embedded)

        if x_lengths is not None:
            y_out = column_gather(y_out, x_lengths)
        else:
            y_out = y_out[:, -1, :]

        y_out = F.relu(self.fc1(F.dropout(y_out, 0.5)))
        y_out = self.fc2(F.dropout(y_out, 0.5))

        if apply_softmax:
            y_out = F.softmax(y_out, dim=1)

        return y_out
```

http://localhost:8888/notebooks/Deep-Learning/Assignment7b_RNN.ipynb Page 10 of 18

```
In [7]:  def set_seed_everywhere(seed, cuda):
             np.random.seed(seed)
             torch.manual_seed(seed)
             if cuda:
                 torch.cuda.manual_seed_all(seed)

         def handle_dirs(dirpath):
             if not os.path.exists(dirpath):
                 os.makedirs(dirpath)
```

In [10]:
```python
args = Namespace(
    # Data and path information
    surname_csv="./surnames_with_splits.csv",
    vectorizer_file="vectorizer.json",
    model_state_file="model.pth",
    save_dir="model_storage/ch6/surname_classification",
    # Model hyper parameter
    char_embedding_size=100,
    rnn_hidden_size=64,
    # Training hyper parameter
    num_epochs=100,
    learning_rate=1e-3,
    batch_size=64,
    seed=1337,
    early_stopping_criteria=5,
    # Runtime hyper parameter
    cuda=True,
    catch_keyboard_interrupt=True,
    reload_from_files=False,
    expand_filepaths_to_save_dir=True,
)

# Check CUDA
if not torch.cuda.is_available():
    args.cuda = False

args.device = torch.device("cuda" if args.cuda else "cpu")

print("Using CUDA: {}".format(args.cuda))


if args.expand_filepaths_to_save_dir:
    args.vectorizer_file = os.path.join(args.save_dir,
                                        args.vectorizer_file)

    args.model_state_file = os.path.join(args.save_dir,
                                         args.model_state_file)

# Set seed for reproducibility
set_seed_everywhere(args.seed, args.cuda)

# handle dirs
handle_dirs(args.save_dir)
```

Using CUDA: False

```python
In [11]:  if args.reload_from_files and os.path.exists(args.vectorizer_file):
              # training from a checkpoint
              dataset = SurnameDataset.load_dataset_and_load_vectorizer(args.surna
                                                                        args.vecto
          else:
              # create dataset and vectorizer
              dataset = SurnameDataset.load_dataset_and_make_vectorizer(args.surna
              dataset.save_vectorizer(args.vectorizer_file)

          vectorizer = dataset.get_vectorizer()

          classifier = SurnameClassifier(embedding_size=args.char_embedding_size,
                                         num_embeddings=len(vectorizer.char_vocab)
                                         num_classes=len(vectorizer.nationality_vo
                                         rnn_hidden_size=args.rnn_hidden_size,
                                         padding_idx=vectorizer.char_vocab.mask_in
```

```python
In [12]:  def make_train_state(args):
              return {'stop_early': False,
                      'early_stopping_step': 0,
                      'early_stopping_best_val': 1e8,
                      'learning_rate': args.learning_rate,
                      'epoch_index': 0,
                      'train_loss': [],
                      'train_acc': [],
                      'val_loss': [],
                      'val_acc': [],
                      'test_loss': -1,
                      'test_acc': -1,
                      'model_filename': args.model_state_file}


          def update_train_state(args, model, train_state):
              """Handle the training state updates.

              Components:
               - Early Stopping: Prevent overfitting.
               - Model Checkpoint: Model is saved if the model is better

              :param args: main arguments
              :param model: model to train
              :param train_state: a dictionary representing the training state val
              :returns:
                  a new train_state
              """

              # Save one model at least
              if train_state['epoch_index'] == 0:
                  torch.save(model.state_dict(), train_state['model_filename'])
```

```python
                train_state['stop_early'] = False

        # Save model if performance improved
        elif train_state['epoch_index'] >= 1:
            loss_tm1, loss_t = train_state['val_loss'][-2:]

            # If loss worsened
            if loss_t >= loss_tm1:
                # Update step
                train_state['early_stopping_step'] += 1
            # Loss decreased
            else:
                # Save the best model
                if loss_t < train_state['early_stopping_best_val']:
                    torch.save(model.state_dict(), train_state['model_filenar
                    train_state['early_stopping_best_val'] = loss_t

                # Reset early stopping step
                train_state['early_stopping_step'] = 0

            # Stop early ?
            train_state['stop_early'] = \
                train_state['early_stopping_step'] >= args.early_stopping_cr

        return train_state


def compute_accuracy(y_pred, y_target):
    _, y_pred_indices = y_pred.max(dim=1)
    n_correct = torch.eq(y_pred_indices, y_target).sum().item()
    return n_correct / len(y_pred_indices) * 100
```

```python
In [13]:  classifier = classifier.to(args.device)
          dataset.class_weights = dataset.class_weights.to(args.device)

          loss_func = nn.CrossEntropyLoss(dataset.class_weights)
          optimizer = optim.Adam(classifier.parameters(), lr=args.learning_rate)
          scheduler = optim.lr_scheduler.ReduceLROnPlateau(optimizer=optimizer,
                                              mode='min', factor=0.5,
                                              patience=1)

          train_state = make_train_state(args)

          epoch_bar = tqdm_notebook(desc='training routine',
                          total=args.num_epochs,
                          position=0)

          dataset.set_split('train')
          train_bar = tqdm_notebook(desc='split=train',
                          total=dataset.get_num_batches(args.batch_size)
```

```python
                                            position=1,
                                            leave=True)
dataset.set_split('val')
val_bar = tqdm_notebook(desc='split=val',
                             total=dataset.get_num_batches(args.batch_size),
                             position=1,
                             leave=True)


try:
    for epoch_index in range(args.num_epochs):
        train_state['epoch_index'] = epoch_index

        # Iterate over training dataset

        # setup: batch generator, set loss and acc to 0, set train mode
        dataset.set_split('train')
        batch_generator = generate_batches(dataset,
                                              batch_size=args.batch_size,
                                              device=args.device)
        running_loss = 0.0
        running_acc = 0.0
        classifier.train()

        for batch_index, batch_dict in enumerate(batch_generator):
            # the training routine is these 5 steps:

            # --------------------------------------
            # step 1. zero the gradients
            optimizer.zero_grad()

            # step 2. compute the output
            y_pred = classifier(x_in=batch_dict['x_data'],
                                x_lengths=batch_dict['x_length'])

            # step 3. compute the loss
            loss = loss_func(y_pred, batch_dict['y_target'])

            running_loss += (loss.item() - running_loss) / (batch_index

            # step 4. use loss to produce gradients
            loss.backward()

            # step 5. use optimizer to take gradient step
            optimizer.step()
            # ------------------------------------------
            # compute the accuracy
            acc_t = compute_accuracy(y_pred, batch_dict['y_target'])
            running_acc += (acc_t - running_acc) / (batch_index + 1)

            # update bar
            train_bar.set_postfix(loss=running_loss, acc=running_acc, ep
```

```python
            train_bar.set_postfix(loss=running_loss, acc=running_acc, ep
            train_bar.update()

        train_state['train_loss'].append(running_loss)
        train_state['train_acc'].append(running_acc)

        # Iterate over val dataset

        # setup: batch generator, set loss and acc to 0; set eval mode o

        dataset.set_split('val')
        batch_generator = generate_batches(dataset,
                                           batch_size=args.batch_size,
                                           device=args.device)
        running_loss = 0.
        running_acc = 0.
        classifier.eval()

        for batch_index, batch_dict in enumerate(batch_generator):
            # compute the output
            y_pred = classifier(x_in=batch_dict['x_data'],
                                x_lengths=batch_dict['x_length'])

            # step 3. compute the loss
            loss = loss_func(y_pred, batch_dict['y_target'])
            running_loss += (loss.item() - running_loss) / (batch_index

            # compute the accuracy
            acc_t = compute_accuracy(y_pred, batch_dict['y_target'])
            running_acc += (acc_t - running_acc) / (batch_index + 1)
            val_bar.set_postfix(loss=running_loss, acc=running_acc, epoc
            val_bar.update()

        train_state['val_loss'].append(running_loss)
        train_state['val_acc'].append(running_acc)

        train_state = update_train_state(args=args, model=classifier,
                                         train_state=train_state)

        scheduler.step(train_state['val_loss'][-1])

        train_bar.n = 0
        val_bar.n = 0
        epoch_bar.update()

        if train_state['stop_early']:
            break

except KeyboardInterrupt:
    print("Exiting loop")
```

training routine: 100%                       100/100 [1:07:54<00:00, 6.54s/it]

split=train: 99%                       119/120 [1:07:53<00:05, 5.94s/it, acc=45.8, epoch=99,

loss=1.47]

split=val: 96%                       24/25 [1:07:53<00:06, 6.26s/it, acc=43.1, epoch=99,

loss=1.8]

In [14]:
```python
# compute the loss & accuracy on the test set using the best available m

classifier.load_state_dict(torch.load(train_state['model_filename']))

classifier = classifier.to(args.device)
dataset.class_weights = dataset.class_weights.to(args.device)
loss_func = nn.CrossEntropyLoss(dataset.class_weights)

dataset.set_split('test')
batch_generator = generate_batches(dataset,
                                   batch_size=args.batch_size,
                                   device=args.device)
running_loss = 0.
running_acc = 0.
classifier.eval()

for batch_index, batch_dict in enumerate(batch_generator):
    # compute the output
    y_pred =  classifier(batch_dict['x_data'],
                         x_lengths=batch_dict['x_length'])

    # compute the loss
    loss = loss_func(y_pred, batch_dict['y_target'])
    loss_t = loss.item()
    running_loss += (loss_t - running_loss) / (batch_index + 1)

    # compute the accuracy
    acc_t = compute_accuracy(y_pred, batch_dict['y_target'])
    running_acc += (acc_t - running_acc) / (batch_index + 1)

train_state['test_loss'] = running_loss
train_state['test_acc'] = running_acc
```

In [15]:
```python
print("Test loss: {};".format(train_state['test_loss']))
print("Test Accuracy: {}".format(train_state['test_acc']))
```

```
Test loss: 1.8446591711044311;
Test Accuracy: 43.25
```

In [16]:
```python
def predict_nationality(surname, classifier, vectorizer):
    vectorized_surname, vec_length = vectorizer.vectorize(surname)
    vectorized_surname = torch.tensor(vectorized_surname).unsqueeze(dim=
    vec_length = torch.tensor([vec_length], dtype=torch.int64)

    result = classifier(vectorized_surname, vec_length, apply_softmax=Tr
    probability_values, indices = result.max(dim=1)

    index = indices.item()
    prob_value = probability_values.item()

    predicted_nationality = vectorizer.nationality_vocab.lookup_index(in

    return {'nationality': predicted_nationality, 'probability': prob_va
```

In [18]:
```python
# surname = input("Enter a surname: ")
classifier = classifier.to("cpu")
for surname in ['McMahan', 'Nakamoto', 'Wan', 'Cho']:
    print(predict_nationality(surname, classifier, vectorizer))
```

```
{'nationality': 'Irish', 'probability': 0.3079265058040619, 'surname':
'McMahan'}
{'nationality': 'Japanese', 'probability': 0.5749261379241943, 'surnam
e': 'Nakamoto'}
{'nationality': 'Chinese', 'probability': 0.4384250342845917, 'surname
': 'Wan'}
{'nationality': 'Vietnamese', 'probability': 0.3590874969959259, 'surn
ame': 'Cho'}
```

In [ ]: