

古诗生成器（运行版）

基于循环神经网络（Keras+LSTM-RNN），采用了浙江大学人工智能研究所提供的古诗词库，并且在其提供的AI学习平台上训练完成。本案例非原创，原来的代码用class来实现，并且重构了model的方法。为了方便初学者理解，重新调整了代码和参数，并且用jupyterlab写了一个完整的文档。

虽然模型已经训练完成，但要应用这个模型，还需要提供原来用于训练的语料，即古诗词库。训练和应用使用的语料要保持一致。

运行这个案例，需要安装多个库，还需要pydot、graphviz的支持。虚谷号教育版上已经提供了所有相关的支持库。需要说明的是，在虚谷号上训练比较慢，预计至少需要一天时间，才可以看到较好的效果。

原案例地址：https://github.com/youyuge34/Poems_generator_Keras/blob/master/poem_model.ipynb
(https://github.com/youyuge34/Poems_generator_Keras/blob/master/poem_model.ipynb)

模型下载地址（课程汇集/虚谷号内置课程目录/5.机器学习）：<https://github.com/vvlink/vvBoard-docs/>
(<https://github.com/vvlink/vvBoard-docs/>)

1. 导入必要的库

In [1]:

```
import random
import os
import keras
import numpy as np
from keras.callbacks import LambdaCallback
from keras.models import Input, Model, load_model
from keras.layers import LSTM, Dropout, Dense
from keras.optimizers import Adam
```

Using TensorFlow backend.

2. 参数设置

In [3]:

```
class config(object):
    # 输入的诗词库（语料库）
    poetry_file = 'data/8-poetry_zju.txt'
    # 模型名称
    weight_file = 'model/8-model_zju.h5'
    # 输出训练的信息
    fixlog = 'poem_log.txt'
    # 复合训练时，间隔多少次输出一次测试结果
    predict_num = 5
    batch_size = 32
    learning_rate = 0.001
    # 下面参数不能随意修改，改动将影响整个模型的大小
    # 根据前六个字预测第七个字，生成的是五言诗（含标点）
    max_len = 6
    # 去除低频率文字（避免生僻字）
    frequency_num = 3
```

3.数据处理

In [4]:

```
def preprocess_file(poetry_file):
    # 语料文本内容
    files_content = ''
    with open(poetry_file, 'r', encoding='UTF-8') as f:
        for line in f:
            x = line.strip() + "]"
            # x = x.split(":")[1]
            if len(x) <= 5 :
                continue
            # 确保导入的诗句, 是五言诗
            if x[c.max_len-1] == ', ':
                files_content += x

    words = sorted(list(files_content))
    counted_words = {}
    for word in words:
        if word in counted_words:
            counted_words[word] += 1
        else:
            counted_words[word] = 1

    # 去掉低频的字, 这样可以去除一些怪异的字。
    erase = []
    for key in counted_words:
        if counted_words[key] <= c.frequency_num:
            erase.append(key)
    for key in erase:
        del counted_words[key]
    wordPairs = sorted(counted_words.items(), key=lambda x: -x[1])

    words, _ = zip(*wordPairs)
    words += (" ",)
    # word到id的映射
    word2num = dict((c, i) for i, c in enumerate(words))
    num2word = dict((i, c) for i, c in enumerate(words))
    word2numF = lambda x: word2num.get(x, len(words) - 1)
    return word2numF, num2word, words, files_content

# 清洗、准备数据
c=config()
word2numF, num2word, words, files_content = preprocess_file(c.poetry_file)
#分割诗词, 记录在poems_num中
poems = files_content.split(']')
poems_num = len(poems)
```

这里返回的四个变量, 分别介绍如下:

- word2numF: 返回不同汉字对应的字典位置。
- num2word: 字典, key是数字, 值是字符。
- words: 列表。所有的字符表, 按照频率排序, 先大再小。
- files_content: 字符串。所有的诗词, 用"]"分开。

按照字符频率, 最多的是“, ”, 其次是“。”。poems为列表, 存储所有的诗歌, poems_num为诗歌的数量。

4.应用模型

根据模型能够根据输入的一组字来预测后面的字，那么就可以做出很多应用，如藏头诗，藏字诗等。

4.1 导入模型

In [6]:

```
# 导入训练好的模型
model = load_model(c.weight_file, compile=False)
```

4.2 准备工作

需要几个基本函数的支持。

In [6]:

```
# 输入字符串, 返回字符串
def m_preds(sentence,length = 23,temperature =1):
    '''
    sentence:预测输入值
    lenth:预测出的字符串长度
    供类内部调用, 输入max_len长度字符串, 返回length长度的预测值字符串
    '''

    sentence = sentence[:c.max_len]
    generate = ''
    for i in range(length):
        pred = m_pred(sentence,temperature)
        generate += pred
        sentence = sentence[1:]+pred
    return generate

# 输入字符串, 返回字符
def m_pred(sentence,temperature =1):
    '''内部使用方法, 根据一串输入, 返回单个预测字符'''
    if len(sentence) < c.max_len:
        print('in def m_pred,length error ')
        return

    sentence = sentence[-c.max_len:]
    x_pred = np.zeros((1, c.max_len, len(words)))
    for t, char in enumerate(sentence):
        x_pred[0, t, word2numF(char)] = 1.
    preds = model.predict(x_pred, verbose=0)[0]
    next_index = sample(preds,temperature=temperature)
    next_char = num2word[next_index]

    return next_char

# 根据输入的矩阵, 返回一个数字
def sample(preds, temperature=1.0):
    '''
    当temperature=1.0时, 模型输出正常
    当temperature=0.5时, 模型输出比较open
    当temperature=1.5时, 模型输出比较保守
    在训练的过程中可以看到temperature不同, 结果也不同
    就是一个概率分布变换的问题, 保守的时候概率大的值变得更大, 选择的可能性也更大
    '''

    preds = np.asarray(preds).astype('float64')
    exp_preds = np.power(preds,1./temperature)
    preds = exp_preds / np.sum(exp_preds)
    pro = np.random.choice(range(len(preds)),1,p=preds)
    return int(pro.squeeze())
```

设置热度参数

一般来说, temperature=1.0输出的比较正常。小于1比较开放, 大于1则比较保守。就是一个概率分布变换的问题, 保守的时候概率大的值变得更大, 选择的可能性也更大。

In [7]:

```
# 默认设置为1, 正常的参数
temperature=1.0
```

4.3 各种应用

1) 诗歌接龙

给出第一句诗，自动生成后面诗句。

In [8]:

```
# 根据给出的前max_len个字, 生成诗句
def predict_sen(text, temperature = 1):
    '''此例中, 即根据给出的第一句诗句 (含逗号), 来生成古诗'''
    max_len = c.max_len
    if len(text) < max_len:
        print('字数不能少于', max_len)
        return

    sentence = text[-max_len:]
    print('第一句:', sentence)
    generate = str(sentence)
    generate += m_preds(sentence, length = 24 - max_len, temperature = temperature)
    return generate

# 运行
for i in range(3):
    # 给出第一句话进行预测
    sen = predict_sen('明月松间照, ', temperature = temperature)
    print(sen)
```

第一句：明月松间照，
明月松间照，池中良分为。山机管华华，将花东千生。
第一句：明月松间照，
明月松间照，有陈藕本日。洛庭秀迷贤，明惊几是东。
第一句：明月松间照，
明月松间照，樽尘光开仙。白无气崇不，丹生乘国台。

2) 藏头诗

输入必须是四个字，并且不能有冷僻字。为了使诗句有点意义，先在诗词库中找一句某个字开头的诗，如果找不到，再随机凑一句。

In [9]:

```
# '''根据给4个字, 生成藏头诗五言绝句'''
def predict_hide(text, temperature = 1):
    if len(text) != 4:
        print('藏头诗的输入必须是4个字! ')
        return
    islow = ''
    for t in text:
        if t not in words:
            islow = t
            break
    if islow:
        print('输入的字存在冷僻字: ' + islow)
        return

    index = random.randint(0, poems_num)
    # 选取随机一首诗的最后max_len字符+给出的首个文字作为初始输入
    sentence = poems[index][1-c.max_len:] + text[0]
    generate = str(text[0])
    # print('引子 = ', sentence)

    for i in range(5):
        next_char = m_pred(sentence, temperature)
        sentence = sentence[1:] + next_char
        generate += next_char

    for i in range(3):
        generate += text[i+1]
        sentence = sentence[1:] + text[i+1]
        for i in range(5):
            next_char = m_pred(sentence, temperature)
            sentence = sentence[1:] + next_char
            generate += next_char

    return generate
# 执行
for i in range(3):
    # 藏头诗
    sen = predict_hide('一杯浊酒', temperature = temperature)
    print(sen)
```

一天木臣天，杯歌还一鸟。浊王天月若，酒问乘下云。
一郊白水云，杯风朝水未。浊北动谁接，酒三山帝玉。
一扶风赋花，杯似晨雪南。浊流春枕玉，酒绕影为川。

3) 随机一首

随机生成一首诗。原理是随机从库中选取一句开头的诗句，生成五言绝句。

In [10]:

```
# 测试函数：随机从选取诗句，生成五言绝句
def predict_random(temperature = 1):
    # 随机找一首诗，用第一行来测试
    index = random.randint(0, poems_num)
    sentence = poems[index][: c.max_len]
    generate = predict_sen(sentence, temperature=temperature)
    return generate
# 执行
predict_random(temperature = temperature)
```

第一句：尝闻四书曰，

Out[10]:

'尝闻四书曰， 皇风入飞夜。龙关三百烟， 不路飞常分。'

4) 诗歌接龙

给出一个字，生成一首诗。这个字不能是冷僻字。

In [12]:

```
''' 根据给出的首个文字，生成五言绝句 '''
def predict_first(char, temperature = 1):
    index = random.randint(0, poems_num)
    # 选取随机一首诗的最后max_len字符+给出的首个文字作为初始输入
    sentence = poems[index][1-c.max_len:] + char
    generate = str(char)
    # print('引子 = ', sentence)
    # 直接预测后面23个字符
    generate += m_preds(sentence, length=23, temperature=temperature)
    return generate

# 执行
s="山"
for i in range(3):
    # 给出第一个字进行预测
    sen = predict_first(s, temperature = temperature)
    print(sen)
    s=sen[-2]
```

山流方深春， 极径池林留。路园旌林尘， 仙歌似自落。
落影经林在， 东芳三时夕。所应迹玉望， 调时松外此。
此作东空登， 大金上辰日。玉白前疏花， 来路神白云。

具体效果如何？大家的要求不能太高啊。

In []: