

## Лабораторная работа №1. Логистическая регрессия в качестве нейронной сети

Данные: В работе предлагается использовать набор данных notMNIST, который состоит из изображений размерностью 28×28 первых 10 букв латинского алфавита (A ... J, соответственно). Обучающая выборка содержит порядка 500 тыс. изображений, а тестовая – около 19 тыс.

Данные можно скачать по ссылке:

[https://commondatastorage.googleapis.com/books1000/notMNIST\\_large.tar.gz](https://commondatastorage.googleapis.com/books1000/notMNIST_large.tar.gz)

([https://commondatastorage.googleapis.com/books1000/notMNIST\\_large.tar.gz](https://commondatastorage.googleapis.com/books1000/notMNIST_large.tar.gz))

(большой набор данных);

[https://commondatastorage.googleapis.com/books1000/notMNIST\\_small.tar.gz](https://commondatastorage.googleapis.com/books1000/notMNIST_small.tar.gz)

([https://commondatastorage.googleapis.com/books1000/notMNIST\\_small.tar.gz](https://commondatastorage.googleapis.com/books1000/notMNIST_small.tar.gz))

(маленький набор данных); Описание данных на английском языке доступно по

ссылке: <http://yaroslavvb.blogspot.sg/2011/09/notmnist-dataset.html>

(<http://yaroslavvb.blogspot.sg/2011/09/notmnist-dataset.html>)

```
In [118]: 1 import numpy as np
          2 import pandas as pd
          3 import matplotlib.pyplot as plt
          4 import matplotlib.image as mpimg
          5 from scipy import misc
          6 import glob
          7 import hashlib
          8 from tqdm import tqdm
          9 from sklearn.linear_model import LogisticRegression
```

**Задание 1.** Загрузите данные и отобразите на экране несколько из изображений с помощью языка Python;

```
In [2]: 1 def plot(imagesDict):
2         _, axis = plt.subplots(1, len(imagesDict))
3         count = 0
4
5         for key in imagesDict.keys():
6             img = mpimg.imread(list(imagesDict[key])[0])
7             axis[count].imshow(img)
8             axis[count].axis("off")
9             count += 1
10
11         plt.show()
12
```

```
In [23]: 1 def md5(fname):
2         hash_md5 = hashlib.md5()
3         with open(fname, "rb") as f:
4             for chunk in iter(lambda: f.read(4096), b''):
5                 hash_md5.update(chunk)
6         return hash_md5.hexdigest()
7
8 def removeDublicates(dataset):
9     result = dict()
10    for file in dataset:
11        result[md5(file)] = file
12    return list(result.values())
```

```
In [24]: 1 rootFolderName = "notMNIST_large/"
2         folders = glob.glob(rootFolderName + "*")
3
4         imagesDict = dict()
5         for folder in folders:
6             images = glob.glob(folder + "/*.png")
7             imagesSet = list()
8             for image in images:
9                 imagesSet.append(image)
10            imagesSet = removeDublicates(imagesSet)
11            imagesDict[folder.replace(rootFolderName, "")] = imagesSet
```

```
In [25]: 1 plot(imagesDict)
```



```
In [26]: 1 keys = imagesDict.keys()
          2 print(keys)
```

```
dict_keys(['I', 'G', 'A', 'F', 'H', 'J', 'C', 'D', 'E', 'B'])
```

**Задание 2.** Проверьте, что классы являются сбалансированными, т.е. количество изображений, принадлежащих каждому из классов, примерно одинаково (В данной задаче 10 классов).

```
In [27]: 1 for key in keys:
          2     print(key, len(imagesDict[key]))
```

```
I 41173
G 47090
A 47107
F 46844
H 46183
J 46658
C 46654
D 46734
E 46954
B 47283
```

**Задание 3.** Разделите данные на три подвыборки: обучающую (200 тыс. изображений), валидационную (10 тыс. изображений) и контрольную (тестовую) (19 тыс. изображений);

```
In [32]: 1 trainDict = dict()
          2 validationDict = dict()
          3 testDict = dict()
          4
          5 matches = {"A":0, "B":1, "C":2, "D":3, "E":4, "F":5, "G":6, "H":7}
          6
          7 for key in imagesDict.keys():
          8     imagesList = imagesDict[key]
          9     newKey = matches[key]
         10     trainDict[newKey] = imagesList[:20000]
         11     validationDict[newKey] = imagesList[20000:21000]
         12     testDict[newKey] = imagesList[21000:23000]
```

**Задание 4.** Проверьте, что данные из обучающей выборки не пересекаются с данными из валидационной и контрольной выборок. Другими словами, избегайте от дубликатов в обучающей выборке.

```
In [33]: 1 for key in trainDict.keys():
2         trainSet = set(trainDict[key])
3         validationSet = set(validationDict[key])
4         testSet = set(testDict[key])
5         intersectionWithValidation = trainSet.intersection(validationSet)
6         intersectionWithTest = trainSet.intersection(testSet)
7         if len(intersectionWithValidation) > 0 or len(intersectionWithTest) > 0:
8             print("Warning: sets intersection")
9
```

**Задание 5.** Постройте простейший классификатор (например, с помощью логистической регрессии). Постройте график зависимости точности классификатора от размера обучающей выборки (50, 100, 1000, 50000). Для построения классификатора можете использовать библиотеку SkLearn (<http://scikit-learn.org> (<http://scikit-learn.org>)).

```
In [113]: 1 def convertToLearnData(dataset):
2         y = np.zeros(0)
3         x = np.zeros(0)
4
5         for key in tqdm(dataset):
6             for path in dataset[key]:
7                 try:
8                     image = mpimg.imread(path)
9                 except:
10                    print(path)
11                    image = image.reshape(1, 784)
12                    if len(x) > 0 and len(y) > 0:
13                        x = np.append(x, image, axis=0)
14                        y = np.append(y, key)
15                    else:
16                        x = image
17                        y = np.array([key])
18         return x, y
```

```
In [114]: 1 trainX, trainY = convertToLearnData(trainDict)
          2 print("Train:", trainX.shape, trainY.shape)
```

0%| | 0/10 [00:00<?, ?it/s]

10%|█ | 1/10 [06:06<54:57, 366.44s/it]

20%|██ | 2/10 [26:14<1:22:32, 619.01s/it]

30%|███ | 3/10 [51:41<1:43:59, 891.30s/it]

40%|████ | 4/10 [1:25:02<2:02:25, 1224.32s/it]

50%|█████ | 5/10 [2:07:08<2:14:33, 1614.73s/it]

60%|██████ | 6/10 [2:59:14<2:17:52, 2068.22s/it]

70%|████████ | 7/10 [4:00:48<2:07:47, 2555.94s/it]

80%|█████████ | 8/10 [5:14:36<1:43:54, 3117.40s/it]

90%|██████████ | 9/10 [6:36:48<1:01:01, 3661.78s/it]

100%|██████████| 10/10 [8:04:41<00:00, 2908.14s/it]

Train: (200000, 784) (200000,)

```
In [115]: 1 validationX, validationY = convertToLearnData(validationDict)
          2 print("Validation:", validationX.shape, validationY.shape)
          3
```

0%| | 0/10 [00:00<?, ?it/s]

10%|█ | 1/10 [00:01<00:16, 1.78s/it]

20%|██ | 2/10 [00:05<00:18, 2.35s/it]

30%|███ | 3/10 [00:10<00:22, 3.23s/it]

40%|████ | 4/10 [00:17<00:25, 4.26s/it]

50%|█████ | 5/10 [00:25<00:27, 5.54s/it]

60%|██████ | 6/10 [00:35<00:27, 6.88s/it]

70%|███████ | 7/10 [00:47<00:24, 8.29s/it]

80%|████████ | 8/10 [01:00<00:19, 9.74s/it]

90%|█████████ | 9/10 [01:15<00:11, 11.27s/it]

100%|██████████| 10/10 [01:31<00:00, 9.19s/it]

Validation: (10000, 784) (10000,)

```
In [116]: 1 testX, testY = convertToLearnData(testDict)
          2 print("Test:", testX.shape, testY.shape)
```

0%| | 0/10 [00:00<?, ?it/s]

10%|█ | 1/10 [00:04<00:44, 4.97s/it]

20%|██ | 2/10 [00:16<00:55, 6.89s/it]

30%|███ | 3/10 [00:34<01:11, 10.25s/it]

40%|████ | 4/10 [00:58<01:27, 14.54s/it]

50%|█████ | 5/10 [01:30<01:37, 19.59s/it]

60%|██████ | 6/10 [02:08<01:40, 25.08s/it]

70%|███████ | 7/10 [02:52<01:32, 30.90s/it]

80%|████████ | 8/10 [03:43<01:13, 36.87s/it]

90%|█████████ | 9/10 [04:40<00:42, 42.90s/it]

100%|██████████| 10/10 [05:44<00:00, 34.41s/it]

Test: (20000, 784) (20000,)

```
In [134]: 1 trainSizes = [50, 100, 1000, 50000]
          2
          3 learnDictData = dict()
          4
          5 for size in trainSizes:
          6     newDict = dict()
          7     sliceSize = int(size / 10)
          8     for key in trainDict.keys():
          9         newDict[key] = trainDict[key][:sliceSize]
         10     newTrainX, newTrainY = convertToLearnData(newDict)
         11     clf = LogisticRegression(random_state=0).fit(newTrainX, new
         12         score = clf.score(validationX, validationY)
         13     learnDictData[size] = score
         14
```

```
100%|██████████| 10/10 [00:00<00:00, 142.36it/s][A
/Users/alinadolmatovich/anaconda2/envs/tensorflow-env/lib/python3.
6/site-packages/sklearn/linear_model/_logistic.py:940: Convergence
Warning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
<https://scikit-learn.org/stable/modules/preprocessing.html>  
Please also refer to the documentation for alternative solver options:  
[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))  
extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG)

```
100%|██████████| 10/10 [00:00<00:00, 263.46it/s][A
/Users/alinadolmatovich/anaconda2/envs/tensorflow-env/lib/python3.
6/site-packages/sklearn/linear_model/_logistic.py:940: Convergence
Warning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
<https://scikit-learn.org/stable/modules/preprocessing.html>



Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

`extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)`

0%| | 0/10 [00:00<?, ?it/s]

30%| | 3/10 [00:00<00:00, 20.35it/s]

50%| | 5/10 [00:00<00:00, 15.47it/s]

60%| | 6/10 [00:00<00:00, 11.68it/s]

70%| | 7/10 [00:00<00:00, 9.60it/s]

80%| | 8/10 [00:00<00:00, 8.10it/s]

90%| | 9/10 [00:00<00:00, 7.37it/s]

100%| | 10/10 [00:01<00:00, 8.61it/s]  
 /Users/alinadolmatovich/anaconda2/envs/tensorflow-env/lib/python3.6/site-packages/sklearn/linear\_model/\_logistic.py:940: Convergence Warning: lbfgs failed to converge (status=1):  
 STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (`max_iter`) or scale the data as

shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
(<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG)

0%| | 0/10 [00:00<?, ?it/s]

10%|█ | 1/10 [00:22<03:19, 22.18s/it]

20%|██ | 2/10 [01:28<04:43, 35.45s/it]

30%|███ | 3/10 [03:14<06:35, 56.47s/it]

40%|████ | 4/10 [05:55<08:48, 88.08s/it]

50%|█████ | 5/10 [09:24<10:21, 124.31s/it]

60%|██████ | 6/10 [13:33<10:45, 161.48s/it]

70%|████████ | 7/10 [18:24<10:01, 200.36s/it]

80%|██████████ | 8/10 [23:55<07:59, 239.61s/it]

90%|██████████ | 9/10 [30:07<04:39, 279.35s/it]

100%|██████████| 10/10 [35:50<00:00, 215.07s/it]  
 /Users/alinadolmatovich/anaconda2/envs/tensorflow-env/lib/python3.6/site-packages/sklearn/linear\_model/\_logistic.py:940: Convergence Warning: lbfgs failed to converge (status=1):  
 STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
 (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG)

```
In [135]: 1 clf = LogisticRegression(random_state=0).fit(trainX, trainY)
          2 learnDictData[200000] = clf.score(validationX, validationY)
```

/Users/alinadolmatovich/anaconda2/envs/tensorflow-env/lib/python3.6/site-packages/sklearn/linear\_model/\_logistic.py:940: Convergence Warning: lbfgs failed to converge (status=1):  
 STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
 (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG)

```
In [136]: 1 learnDictData
```

```
Out[136]: {50: 0.6416, 100: 0.7194, 1000: 0.7453, 50000: 0.8101, 200000: 0.8156}
```

```
In [140]: 1 iterations = learnDictData.keys()
          2 scores = learnDictData.values()
          3
          4 plt.scatter(scores, iterations)
          5 plt.show()
```

