

▼ Практическое задание №1

Установка необходимых пакетов:

```
1 !pip install -q libtiff
2 !pip install -q tqdm
```

Монтирование Вашего Google Drive к текущему окружению:

```
1 from google.colab import drive
2 drive.mount('/content/drive', force_remount=True)

Mounted at /content/drive
```

В переменную PROJECT_DIR необходимо прописать путь к директории на Google Drive, в которую Вы загрузили zip архивы с предоставленными наборами данных.

```
1 # todo
2 PROJECT_DIR = 'dev/prak_nn_1_data/'
```

Константы, которые пригодятся в коде далее:

```
1 EVALUATE_ONLY = True
2 TEST_ON_LARGE_DATASET = True
3 TISSUE_CLASSES = ('ADI', 'BACK', 'DEB', 'LYM', 'MUC', 'MUS', 'NORM', 'STR', 'TUM')
```

Импорт необходимых зависимостей:

```
1 from pathlib import Path
2 from libtiff import TIFF
3 import numpy as np
4 from typing import List
5 from tqdm.notebook import tqdm
6 from time import sleep
7 from PIL import Image
8 import IPython.display
9 from sklearn.metrics import balanced_accuracy_score
```

▼ Класс Dataset

Предназначен для работы с наборами данных, хранящихся на Google Drive, обеспечивает чтение изображений и соответствующих меток, а также формирование пакетов (батчей).

```
1 class Dataset:
2
3     def __init__(self, name, gdrive_dir):
4         self.name = name
5         self.is_loaded = False
6         p = Path("/content/drive/MyDrive/" + gdrive_dir + name + '.npz')
7         if p.exists():
8             print(f'Loading dataset {self.name} from npz.')
9             np_obj = np.load(str(p))
10            self.images = np_obj['data']
11            self.labels = np_obj['labels']
12            self.n_files = self.images.shape[0]
13            self.is_loaded = True
14            print(f'Done. Dataset {name} consists of {self.n_files} images.')
15
16        def image(self, i):
17            # read i-th image in dataset and return it as numpy array
18            if self.is_loaded:
19                return self.images[i, :, :, :]
20
21        def images_seq(self, n=None):
22            # sequential access to images inside dataset (is needed for testing)
23            for i in range(self.n_files if not n else n):
24                yield self.image(i)
25
26        def random_image_with_label(self):
27            # get random image with label from dataset
28            i = np.random.randint(self.n_files)
29            return self.image(i), self.labels[i]
30
31        def random_batch_with_labels(self, n):
32            # create random batch of images with labels (is needed for training)
33            indices = np.random.choice(self.n_files, n)
34            imgs = []
35            for i in indices:
36                img = self.image(i)
37                imgs.append(self.image(i))
38            logits = np.array([self.labels[i] for i in indices])
39            return np.stack(imgs), logits
40
41        def image_with_label(self, i: int):
42            # return i-th image with label from dataset
43            return self.image(i), self.labels[i]
```

▼ Пример использования класса Dataset

Загрузим обучающий набор данных, получим произвольное изображение с меткой. После чего визуализируем изображение, выведем метку. В будущем, этот кусок кода можно закомментировать или убрать.

```
1 d_train_tiny = Dataset('train_tiny', PROJECT_DIR)
2
```

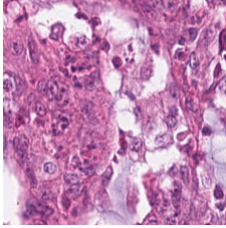
```

3 img, lbl = d_train_tiny.random_image_with_label()
4 print()
5 print(f'Got numpy array of shape {img.shape}, and label with code {lbl}.')
6 print(f'Label code corresponds to {TISSUE_CLASSES[lbl]} class.')
7
8 pil_img = Image.fromarray(img)
9 IPython.display.display(pil_img)

```

Loading dataset train_tiny from npz.
Done. Dataset train_tiny consists of 900 images.

Got numpy array of shape (224, 224, 3), and label with code 8.
Label code corresponds to TUM class.



▼ Класс Metrics

Реализует метрики точности, используемые для оценивания модели:

1. точность,
2. сбалансированную точность.

```

1 class Metrics:
2
3     @staticmethod
4     def accuracy(gt: List[int], pred: List[int]):
5         assert len(gt) == len(pred), 'gt and prediction should be of equal length'
6         return sum(int(i[0] == i[1]) for i in zip(gt, pred)) / len(gt)
7
8     @staticmethod
9     def accuracy_balanced(gt: List[int], pred: List[int]):
10        return balanced_accuracy_score(gt, pred)
11
12    @staticmethod
13    def print_all(gt: List[int], pred: List[int], info: str):
14        print(f'metrics for {info}:')
15        print('\t accuracy {:.4f}'.format(Metrics.accuracy(gt, pred)))
16        print('\t balanced accuracy {:.4f}'.format(Metrics.accuracy_balanced(gt, pred)))

```

▼ Класс Model

Класс, хранящий в себе всю информацию о модели.

Вам необходимо реализовать методы `save`, `load` для сохранения и загрузки модели. Особенно актуально это будет во время тестирования на дополнительных наборах данных.

Пожалуйста, убедитесь, что сохранение и загрузка модели работает корректно. Для этого обучите модель, протестируйте, сохраните ее в файл, перезапустите среду выполнения, загрузите обученную модель из файла, вновь протестируйте ее на тестовой выборке и убедитесь в том, что получаемые метрики совпадают с полученными для тестовой выборки ранее.

Также, Вы можете реализовать дополнительные функции, такие как:

1. валидацию модели на части обучающей выборки;
2. использование кроссвалидации;
3. автоматическое сохранение модели при обучении;
4. загрузку модели с какой-то конкретной итерации обучения (если используется итеративное обучение);
5. вывод различных показателей в процессе обучения (например, значение функции потерь на каждой эпохе);
6. построение графиков, визуализирующих процесс обучения (например, график зависимости функции потерь от номера эпохи обучения);
7. автоматическое тестирование на тестовом наборе/наборах данных после каждой эпохи обучения (при использовании итеративного обучения);
8. автоматический выбор гиперпараметров модели во время обучения;
9. сохранение и визуализацию результатов тестирования;
10. Использование аугментации и других способов синтетического расширения набора данных (дополнительным плюсом будет обоснование необходимости и обоснование выбора конкретных типов аугментации)
11. и т.д.

Полный список опций и дополнений приведен в презентации с описанием задания.

При реализации дополнительных функций допускается добавление параметров в существующие методы и добавление новых методов в класс модели.

```

1 class Model:
2
3     def __init__(self):
4         # todo
5         pass
6
7     def save(self, name: str):
8         # save model to PROJECT_DIR folder on gdrive with name 'name'
9         # todo
10        pass
11
12    def load(self, name: str):
13        # load model with name 'name' from PROJECT_DIR folder on gdrive
14        # todo
15        pass
16
17    def train(self, dataset: Dataset):
18        # you can add some plots for better visualization,

```

```

19     # you can add model autosaving during training,
20     # etc.
21     print(f'training started')
22     # to-do
23     sleep(2)
24     print(f'training done')
25     pass
26
27     def test_on_dataset(self, dataset: Dataset, limit=None):
28         # you can upgrade this code if you want to speed up testing using batches
29         predictions = []
30         n = dataset.n_files if not limit else int(dataset.n_files * limit)
31         for img in tqdm(dataset.images_seq(n), total=n):
32             predictions.append(self.test_on_image(img))
33         return predictions
34
35     def test_on_image(self, img: np.ndarray):
36         # todo: replace this code
37         prediction = np.random.randint(9)
38         sleep(0.05)
39         return prediction
40

```

▼ Классификация изображений

Используя введенные выше классы можем перейти уже непосредственно к обучению модели классификации изображений. Пример общего пайплайна решения задачи приведен ниже. Вы можете его расширять и улучшать. В данном примере используются наборы данных 'train_small' и 'test_small'.

```

1  d_train = Dataset('train_small', PROJECT_DIR)
2  d_test = Dataset('test_small', PROJECT_DIR)

Loading dataset train_small from npz.
Done. Dataset train_small consists of 7200 images.
Loading dataset test_small from npz.
Done. Dataset test_small consists of 1800 images.

```

```

1  model = Model()
2  if not EVALUATE_ONLY:
3      model.train(d_train)
4      model.save('best')
5  else:
6      model.load('best')

```

Пример тестирования модели на части набора данных:

```

1  # evaluating model on 10% of test dataset
2  pred_1 = model.test_on_dataset(d_test, limit=0.1)
3  Metrics.print_all(d_test.labels[:len(pred_1)], pred_1, '10% of test')

100%          180/180 [03:39<00:00, 1.22s/it]

metrics for 10% of test:
  accuracy 0.1278:
  balanced accuracy 0.1278:
/usr/local/lib/python3.6/dist-packages/sklearn/metrics/_classification.py:1859: UserWarning: y_pred contains classes not in y_true
  warnings.warn('y_pred contains classes not in y_true')

```

Пример тестирования модели на полном наборе данных:

```

1  # evaluating model on full test dataset (may take time)
2  if TEST_ON_LARGE_DATASET:
3      pred_2 = model.test_on_dataset(d_test)
4      Metrics.print_all(d_test.labels, pred_2, 'test')

100%          1800/1800 [04:09<00:00, 7.21it/s]

metrics for test:
  accuracy: 0.11277777777777778
  balanced accuracy: 0.11277777777777778

```

Результат работы пайплайна обучения и тестирования выше тоже будет оцениваться. Поэтому не забудьте присылать на проверку ноутбук с выполненными ячейками кода с демонстрациями метрик обучения, графиками и т.п. В этом пайплайне Вам необходимо продемонстрировать работу всех реализованных дополнений, улучшений и т.п.

Настоятельно рекомендуется после получения пайплайна с полными результатами обучения экспортировать ноутбук в pdf (файл -> печать) и прислать этот pdf вместе с самим ноутбуком.

▼ Тестирование модели на других наборах данных

Ваша модель должна поддерживать тестирование на других наборах данных. Для удобства, Вам предоставляется набор данных test_tiny, который представляет собой малую часть (2% изображений) набора test. Ниже приведен фрагмент кода, который будет осуществлять тестирование для оценивания Вашей модели на дополнительных тестовых наборах данных.

Прежде чем отсылать задание на проверку, убедитесь в работоспособности фрагмента кода ниже.

```

1  final_model = Model()
2  final_model.load('best')
3  d_test_tiny = Dataset('test_tiny', PROJECT_DIR)
4  pred = model.test_on_dataset(d_test_tiny)
5  Metrics.print_all(d_test_tiny.labels, pred, 'test-tiny')

Loading dataset test_tiny from npz.
Done. Dataset test_tiny consists of 90 images.

100%          90/90 [03:03<00:00, 2.04s/it]

metrics for test-tiny:
  accuracy 0.1111:
  balanced accuracy 0.1111:

```

Отмонтировать Google Drive.

```
1 drive.flush_and_unmount()
```

▼ Дополнительные "полезности"

Ниже приведены примеры использования различных функций и библиотек, которые могут быть полезны при выполнении данного практического задания.

▼ Измерение времени работы кода

Измерять время работы какой-либо функции можно легко и непринужденно при помощи функции `timeit` из соответствующего модуля:

```
1 import timeit
2
3 def factorial(n):
4     res = 1
5     for i in range(1, n + 1):
6         res *= i
7     return res
8
9
10 def f():
11     return factorial(n=1000)
12
13 n_runs = 128
14 print(f'Function f is caluclated {n_runs} times in {timeit.timeit(f, number=n_runs)}s.')

```

Function f is caluclated 128 times in 0.03538683599981596s.

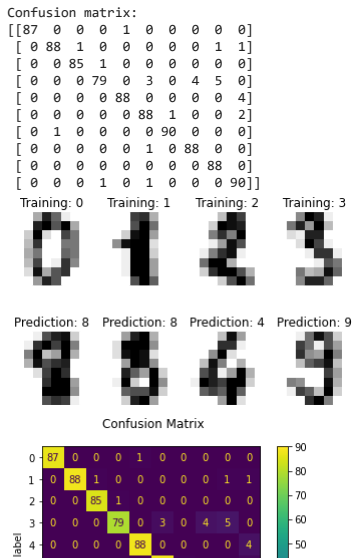
▼ Scikit-learn

Для использования "классических" алгоритмов машинного обучения рекомендуется использовать библиотеку `scikit-learn` (<https://scikit-learn.org/stable/>). Пример классификации изображений цифр из набора данных MNIST при помощи классификатора SVM:

```
1 # Standard scientific Python imports
2 import matplotlib.pyplot as plt
3
4 # Import datasets, classifiers and performance metrics
5 from sklearn import datasets, svm, metrics
6 from sklearn.model_selection import train_test_split
7
8 # The digits dataset
9 digits = datasets.load_digits()
10
11 # The data that we are interested in is made of 8x8 images of digits, let's
12 # have a look at the first 4 images, stored in the 'images' attribute of the
13 # dataset. If we were working from image files, we could load them using
14 # matplotlib.pyplot.imread. Note that each image must have the same size. For these
15 # images, we know which digit they represent: it is given in the 'target' of
16 # the dataset.
17 _, axes = plt.subplots(2, 4)
18 images_and_labels = list(zip(digits.images, digits.target))
19 for ax, (image, label) in zip(axes[0, :], images_and_labels[:4]):
20     ax.set_axis_off()
21     ax.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
22     ax.set_title('Training: %i' % label)
23
24 # To apply a classifier on this data, we need to flatten the image, to
25 # turn the data in a (samples, feature) matrix:
26 n_samples = len(digits.images)
27 data = digits.images.reshape((n_samples, -1))
28
29 # Create a classifier: a support vector classifier
30 classifier = svm.SVC(gamma=0.001)
31
32 # Split data into train and test subsets
33 X_train, X_test, y_train, y_test = train_test_split(
34     data, digits.target, test_size=0.5, shuffle=False)
35
36 # We learn the digits on the first half of the digits
37 classifier.fit(X_train, y_train)
38
39 # Now predict the value of the digit on the second half:
40 predicted = classifier.predict(X_test)
41
42 images_and_predictions = list(zip(digits.images[n_samples // 2:], predicted))
43 for ax, (image, prediction) in zip(axes[1, :], images_and_predictions[:4]):
44     ax.set_axis_off()
45     ax.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
46     ax.set_title('Prediction: %i' % prediction)
47
48 print("Classification report for classifier %s:\n%s\n"
49       % (classifier, metrics.classification_report(y_test, predicted)))
50 disp = metrics.plot_confusion_matrix(classifier, X_test, y_test)
51 disp.figure_.suptitle("Confusion Matrix")
52 print("Confusion matrix:\n%s" % disp.confusion_matrix)
53
54 plt.show()
```

```
Classification report for classifier SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.001, kernel='rbf',
max_iter=1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False):
```

	precision	recall	f1-score	support
0	1.00	0.99	0.99	88
1	0.99	0.97	0.98	91
2	0.99	0.99	0.99	86
3	0.98	0.87	0.92	91
4	0.99	0.96	0.97	92
5	0.95	0.97	0.96	91
6	0.99	0.99	0.99	91
7	0.96	0.99	0.97	89
8	0.94	1.00	0.97	88
9	0.93	0.98	0.95	92
accuracy			0.97	899
macro avg	0.97	0.97	0.97	899
weighted avg	0.97	0.97	0.97	899

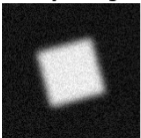


Scikit-image

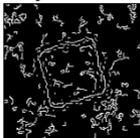
Реализовывать различные операции для работы с изображениями можно как самостоятельно, работая с массивами `numpy`, так и используя специализированные библиотеки, например, `scikit-image` (<https://scikit-image.org/>). Ниже приведен пример использования Canny edge detector.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy import ndimage as ndi
4
5 from skimage import feature
6
7
8 # Generate noisy image of a square
9 im = np.zeros((128, 128))
10 im[32:-32, 32:-32] = 1
11
12 im = ndi.rotate(im, 15, mode='constant')
13 im = ndi.gaussian_filter(im, 4)
14 im += 0.2 * np.random.random(im.shape)
15
16 # Compute the Canny filter for two values of sigma
17 edges1 = feature.canny(im)
18 edges2 = feature.canny(im, sigma=3)
19
20 # display results
21 fig, (ax1, ax2, ax3) = plt.subplots(nrows=1, ncols=3, figsize=(8, 3),
22                                     sharex=True, sharey=True)
23
24 ax1.imshow(im, cmap=plt.cm.gray)
25 ax1.axis('off')
26 ax1.set_title('noisy image', fontsize=20)
27
28 ax2.imshow(edges1, cmap=plt.cm.gray)
29 ax2.axis('off')
30 ax2.set_title(r'Canny filter, $\sigma=1$', fontsize=20)
31
32 ax3.imshow(edges2, cmap=plt.cm.gray)
33 ax3.axis('off')
34 ax3.set_title(r'Canny filter, $\sigma=3$', fontsize=20)
35
36 fig.tight_layout()
37
38 plt.show()
```

noisy image



Canny filter, $\sigma = 1$



Canny filter, $\sigma = 3$



Tensorflow 2

Для создания и обучения нейросетевых моделей можно использовать фреймворк глубокого обучения Tensorflow 2. Ниже приведен пример простейшей нейронной сети, использующейся для классификации изображений из набора данных MNIST.

```
1 # Install TensorFlow
2
3 import tensorflow as tf
4
5 mnist = tf.keras.datasets.mnist
6
7 (x_train, y_train), (x_test, y_test) = mnist.load_data()
8 x_train, x_test = x_train / 255.0, x_test / 255.0
9
10 model = tf.keras.models.Sequential([
11     tf.keras.layers.Flatten(input_shape=(28, 28)),
12     tf.keras.layers.Dense(128, activation='relu'),
13     tf.keras.layers.Dropout(0.2),
14     tf.keras.layers.Dense(10, activation='softmax')
15 ])
16
17 model.compile(optimizer='adam',
18               loss='sparse_categorical_crossentropy',
19               metrics=['accuracy'])
20
21 model.fit(x_train, y_train, epochs=5)
22
23 model.evaluate(x_test, y_test, verbose=2)

Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist.npz
11493376/11490434 [=====] - 0s 0us/step
Epoch 1/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.3014 - accuracy: 0.9136
Epoch 2/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.1446 - accuracy: 0.9571
Epoch 3/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.1077 - accuracy: 0.9671
Epoch 4/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.0864 - accuracy: 0.9735
Epoch 5/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.0731 - accuracy: 0.9773
313/313 - 0s - loss: 0.0787 - accuracy: 0.9760
[0.07866337150335312, 0.976000109672546]
```

Для эффективной работы с моделями глубокого обучения убедитесь в том, что в текущей среде Google Colab используется аппаратный ускоритель GPU или TPU. Для смены среды выберите "среда выполнения" -> "сменить среду выполнения".

Большое количество tutorиалов и примеров с кодом на TensorFlow 2 можно найти на официальном сайте <https://www.tensorflow.org/tutorials?hl=ru>.

Также, Вам может понадобиться написать собственный генератор данных для TensorFlow 2. Скорее всего он будет достаточно простым, и его легко можно будет реализовать, используя официальную документацию TensorFlow 2. Но, на всякий случай (если не удалось сразу разобраться или хочется вникнуть в тему более глубоко), можете посмотреть следующий отличный tutorиал: <https://stanford.edu/~shervine/blog/keras-how-to-generate-data-on-the-fly>.

Numba

В некоторых ситуациях, при ручных реализациях графовых алгоритмов, выполнение многократных вложенных циклов for в python можно существенно ускорить, используя JIT-компилятор Numba (<https://numba.pydata.org/>). Примеры использования Numba в Google Colab можно найти тут:

1. https://colab.research.google.com/github/cbernet/maldives/blob/master/numba/numba_cuda.ipynb
2. https://colab.research.google.com/github/evaneschneider/parallel-programming/blob/master/COMPASS_gpu_intro.ipynb

Пожалуйста, если Вы решили использовать Numba для решения этого практического задания, еще раз подумайте, нужно ли это Вам, и есть ли возможность реализовать требуемую функциональность иным способом. Используйте Numba только при реальной необходимости.

▼ Работа с zip архивами в Google Drive

Запаковка и распаковка zip архивов может пригодиться при сохранении и загрузки Вашей модели. Ниже приведен фрагмент кода, иллюстрирующий помещение нескольких файлов в zip архив с последующим чтением файлов из него. Все действия с директориями, файлами и архивами должны осуществляться с примонтированным Google Drive.

Создадим 2 изображения, поместим их в директорию tmp внутри PROJECT_DIR, запакуем директорию tmp в архив tmp.zip.

```
1 arr1 = np.random.rand(100, 100, 3) * 255
2 arr2 = np.random.rand(100, 100, 3) * 255
3
4 img1 = Image.fromarray(arr1.astype('uint8'))
5 img2 = Image.fromarray(arr2.astype('uint8'))
6
7 p = "/content/drive/MyDrive/" + PROJECT_DIR
8
9 if not (Path(p) / 'tmp').exists():
10     (Path(p) / 'tmp').mkdir()
11
12 img1.save(str(Path(p) / 'tmp' / 'img1.png'))
13 img2.save(str(Path(p) / 'tmp' / 'img2.png'))
14
15 %cd $p
16 !zip -r "tmp.zip" "tmp"

/content/drive/MyDrive/dev/prak_nn_1_data
adding: tmp/ (stored 0%)
adding: tmp/img1.png (stored 0%)
adding: tmp/img2.png (stored 0%)
```

Распакуем архив tmp.zip в директорию tmp2 в PROJECT_DIR. Теперь внутри директории tmp2 содержится директория tmp, внутри которой находятся 2 изображения.

```
1 p = "/content/drive/MyDrive/" + PROJECT_DIR
2 %cd $p
3 !unzip -uq "tmp.zip" -d "tmp2"

/content/drive/MyDrive/dev/prak_nn_1_data
```