

## ▼ Практическое задание №1

Установка необходимых пакетов:

```
!pip install -q tqdm
!pip install --upgrade --no-cache-dir gdown
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: gdown in /usr/local/lib/python3.8/dist-packages (4.5.4)
Requirement already satisfied: requests[socks] in /usr/local/lib/python3.8/dist-packages (from gdown) (2.23.0)
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.8/dist-packages (from gdown) (4.6.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.8/dist-packages (from gdown) (4.64.1)
Requirement already satisfied: six in /usr/local/lib/python3.8/dist-packages (from gdown) (1.15.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.8/dist-packages (from gdown) (3.8.0)
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.8/dist-packages (
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.8/dist-packages (from requests[socks]-
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.8/dist-packages (from requests[socks]->gdown
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.8/dist-packages (from requests[socks]->
Requirement already satisfied: PySocks!=1.5.7,>=1.5.6 in /usr/local/lib/python3.8/dist-packages (from requests[soc
```

Монтирование Вашего Google Drive к текущему окружению:

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

```
Mounted at /content/drive
```

Константы, которые пригодятся в коде далее, и ссылки (gdrive идентификаторы) на предоставляемые наборы данных:

```
EVALUATE_ONLY = False
```

! 0s completed at 11:45 PM



```
ISSUE_CLASSES = ('AD1', 'BACK', 'DEB', 'LYM', 'MUC', 'MUS', 'NORM', 'STR', 'TUM')
DATASETS_LINKS = {
    'train': '1XtQzVQ5XbrfxpLHJuL0XBGJ5U7CS-cLi',
    'train_small': '1qd45xXfDwdZjktLFwQb-et-mAaFeCz0R',
    'train_tiny': '1I-2Z0uXLd4QwhZQqltp817Kn3J0Xgbui',
    'test': '1RfPou3pFKpuHDJZ-D9XDFzgvwpUBFlDr',
    'test_small': '1wbRsog0n7uGlHIPGLhyN-PMET2kdQ2lI',
    'test_tiny': '1viiB0s041CNsAK4itvX8PnYthJ-MDnQc'
}
```

Импорт необходимых зависимостей:

```
from pathlib import Path
import numpy as np
from typing import List
from tqdm.notebook import tqdm
from time import sleep
from PIL import Image
import IPython.display
from sklearn.metrics import balanced_accuracy_score
import gdown
import tensorflow as tf
```

## ▼ Класс Dataset

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

```
PROJECT_DIR = 'dev/prak_nn_1/'
class Dataset:
    def __init__(self, name):
```

```
def __init__(self, name):
    self.name = name
    self.is_loaded = False
    url = f"https://drive.google.com/uc?export=download&confirm=pbef&id={DATASETS_LINKS[name]}"
    output = f'{name}.npz'
    gdown.download(url, output, quiet=False)
    print(f'Loading dataset {self.name} from npz.')
    np_obj = np.load(f'{name}.npz')
    self.images = np_obj['data']
    self.labels = np_obj['labels']
    self.n_files = self.images.shape[0]
    self.is_loaded = True
    print(f'Done. Dataset {name} consists of {self.n_files} images.')

def image(self, i):
    # read i-th image in dataset and return it as numpy array
    if self.is_loaded:
        return self.images[i, :, :, :]

def images_seq(self, n=None):
    # sequential access to images inside dataset (is needed for testing)
    for i in range(self.n_files if not n else n):
        yield self.image(i)

def random_image_with_label(self):
    # get random image with label from dataset
    i = np.random.randint(self.n_files)
    return self.image(i), self.labels[i]

def random_batch_with_labels(self, n):
    # create random batch of images with labels (is needed for training)
    indices = np.random.choice(self.n_files, n)
    imgs = []
    for i in indices:
        img = self.image(i)
        imgs.append(self.image(i))
    logits = np.array([self.labels[i] for i in indices])
```

```
        return np.stack(imgs), logits

    def image_with_label(self, i: int):
        # return i-th image with label from dataset
        return self.image(i), self.labels[i]
```

---

## Класс Metrics

Реализует метрики точности, используемые для оценивания модели:

1. точность,
2. сбалансированную точность.

```
class Metrics:

    @staticmethod
    def accuracy(gt: List[int], pred: List[int]):
        assert len(gt) == len(pred), 'gt and prediction should be of equal length'
        return sum(int(i[0] == i[1]) for i in zip(gt, pred)) / len(gt)

    @staticmethod
    def accuracy_balanced(gt: List[int], pred: List[int]):
        return balanced_accuracy_score(gt, pred)

    @staticmethod
    def print_all(gt: List[int], pred: List[int], info: str):
        print(f'metrics for {info}:')
        print('\t accuracy {:.4f}'.format(Metrics.accuracy(gt, pred)))
        print('\t balanced accuracy {:.4f}'.format(Metrics.accuracy_balanced(gt, pred)))
```

---

## Класс Model

Класс хранящий в себе всю информацию о модели

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.

```
class Model:

    def __init__(self):
        self.base_learning_rate = 0.0001
        preprocess_input = tf.keras.applications.resnet50.preprocess_input
        self.base_model = tf.keras.applications.resnet50.ResNet50(
            input_shape = (224, 224, 3), include_top = False, weights='imagenet')
        self.base_model.trainable = False
        epochs = 10
        global_average_layer = tf.keras.layers.GlobalAveragePooling2D()
        prediction_layer = tf.keras.layers.Dense(9, activation='softmax')

        inputs = tf.keras.Input(shape=(224, 224, 3))
        x = preprocess_input(inputs)
        x = self.base_model(x, training=False)
        x = global_average_layer(x)
        x = tf.keras.layers.Dropout(0.2)(x)
        outputs = prediction_layer(x)
        self.model = tf.keras.Model(inputs, outputs)

    def save(self, name: str):
        self.model.save(f'drive/MyDrive/prac_nn/{name}.h5')

    def load(self, name: str):
        name_to_id_dict = {
            'best' : '1--3JYgywyWSZVsW3io3ao_EwvsVcIu4Z'
        }

        url = f'https://drive.google.com/drive/u/1/folders/{name_to_id_dict[name]}'
        gdown.download_folder(url, quiet=True, output=name, use_cookies=False)
        self.model = tf.keras.models.load_model(name)
```

```
def train(self, dataset: Dataset):
    x_train, y_train = dataset.images, dataset.labels
    self.model.compile(optimizer='adam',
                        loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
                        metrics=['accuracy'])

    self.history = self.model.fit(x_train, y_train,
                                   epochs=5)
    self.base_model.trainable = True
    # всего 175 слоев
    fine_tune_at = 155
    for layer in self.base_model.layers[:fine_tune_at]:
        layer.trainable = False
    self.model.compile(loss=tf.keras.losses.sparse_categorical_crossentropy,
                        optimizer = tf.keras.optimizers.RMSprop(learning_rate=self.base_learning_rate/10),
                        metrics=['accuracy'])
    self.history_fine = self.model.fit(x_train, y_train,
                                       epochs=10,
                                       initial_epoch=self.history.epoch[-1])

def test_on_dataset(self, dataset: Dataset, limit=None):
    # you can upgrade this code if you want to speed up testing using batches
    predictions = []
    n = dataset.n_files if not limit else int(dataset.n_files * limit)
    for img in tqdm(dataset.images_seq(n), total=n):
        predictions.append(self.test_on_image(img))
    return predictions

def test_on_image(self, img: np.ndarray):
    prediction = self.model(img.reshape(1, 224, 224, 3), training=False)
    return tf.argmax(prediction[0])
```

---

## 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'.

```
d_train = Dataset('train')
d_test = Dataset('test')
```

```
Downloading...
```

```
From: https://drive.google.com/uc?export=download&confirm=pbef&id=1XtQzVQ5XbrfxpLHJuL0XBGJ5U7CS-cLi
```

```
To: /content/train.npz
```

```
100%|██████████| 2.10G/2.10G [00:09<00:00, 223MB/s]
```

```
Loading dataset train from npz.
```

```
Done. Dataset train consists of 18000 images.
```

```
Downloading...
```

```
From: https://drive.google.com/uc?export=download&confirm=pbef&id=1RfPou3pFKpuHDJZ-D9XDFzgywpUBFlDr
```

```
To: /content/test.npz
```

```
100%|██████████| 525M/525M [00:01<00:00, 292MB/s]
```

```
Loading dataset test from npz.
```

```
Done. Dataset test consists of 4500 images.
```

```
model = Model()
```

```
if not EVALUATE_ONLY:
```

```
    model.train(d_train)
```

```
    model.save('best')
```

```
else:
```

```
    #todo: your link goes here
```

```
    model.load('best')
```

```
Epoch 1/5
```

```
/usr/local/lib/python3.8/dist-packages/tensorflow/python/util/dispatch.py:1082: UserWarning: "`sparse_categorical_`  
return dispatch_target(*args, **kwargs)
```

```
563/563 [=====] - 57s 94ms/step - loss: 0.3297 - accuracy: 0.8883
```

```
Epoch 2/5
```

```
563/563 [=====] - 53s 93ms/step - loss: 0.1621 - accuracy: 0.9462
```

```
Epoch 3/5
```

```
563/563 [=====] - 53s 94ms/step - loss: 0.1400 - accuracy: 0.9533
```



```

Epoch 4/5
563/563 [=====] - 53s 94ms/step - loss: 0.1197 - accuracy: 0.9602
Epoch 5/5
563/563 [=====] - 53s 94ms/step - loss: 0.1059 - accuracy: 0.9659
Epoch 5/10
563/563 [=====] - 71s 116ms/step - loss: 0.0817 - accuracy: 0.9713
Epoch 6/10
563/563 [=====] - 65s 115ms/step - loss: 0.0562 - accuracy: 0.9814
Epoch 7/10
563/563 [=====] - 64s 114ms/step - loss: 0.0433 - accuracy: 0.9853
Epoch 8/10
563/563 [=====] - 64s 114ms/step - loss: 0.0298 - accuracy: 0.9895
Epoch 9/10
563/563 [=====] - 65s 115ms/step - loss: 0.0197 - accuracy: 0.9936
Epoch 10/10
563/563 [=====] - 64s 114ms/step - loss: 0.0150 - accuracy: 0.9953
WARNING:absl:Found untraced functions such as _jit_compiled_convolution_op, _jit_compiled_convolution_op, _jit_cor

```

Пример тестирования модели на части набора данных:

```

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

100% 450/450 [00:34<00:00, 13.56it/s]
metrics for 10% of test:
  accuracy 0.9956:
  balanced accuracy 0.9956:
/usr/local/lib/python3.8/dist-packages/sklearn/metrics/_classification.py:1987: UserWarning: y_pred contains class
  warnings.warn("y_pred contains classes not in y_true")

```

Пример тестирования модели на полном наборе данных:

```

# evaluating model on full test dataset (may take time)
if TEST ON LARGE DATASET:

```

```
... ..  
pred_2 = model.test_on_dataset(d_test)  
Metrics.print_all(d_test.labels, pred_2, 'test')  
  
100% 4500/4500 [05:31<00:00, 13.45it/s]  
  
metrics for test:  
  accuracy 0.9698:  
 balanced accuracy 0.9698:
```

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.

```
final_model = Model()  
final_model.load('best')  
d_test_tiny = Dataset('test_tiny')  
pred = model.test_on_dataset(d_test_tiny)  
Metrics.print_all(d_test_tiny.labels, pred, 'test-tiny')  
  
Downloading...  
From: https://drive.google.com/uc?export=download&confirm=pbef&id=lviiB0s041CNsAK4itvX8PnYthJ-MDnQc  
To: /content/test_tiny.npz  
100%|██████████| 10.6M/10.6M [00:00<00:00, 211MB/s]Loading dataset test_tiny from npz.  
Done. Dataset test_tiny consists of 90 images.
```

```
100% 90/90 [00:06<00:00, 13.14it/s]
```

```
metrics for test-tiny:  
  accuracy 0.9889:  
  balanced accuracy 0.9889:
```

Отмонтировать Google Drive.

```
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 cell hidden
```

### 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 cell hidden
```

## 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 cell hidden
```

## 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.

```
[ ] ↪ 4 cells hidden
```

## 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](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](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.

```
PROJECT_DIR = "/dev/prak_nn_1/"
arr1 = np.random.rand(100, 100, 3) * 255
arr2 = np.random.rand(100, 100, 3) * 255
```

```
img1 = Image.fromarray(arr1.astype('uint8'))
img2 = Image.fromarray(arr2.astype('uint8'))
```

```
p = "/content/drive/MyDrive/" + PROJECT_DIR
```

```
if not (Path(p) / 'tmp').exists():
    (Path(p) / 'tmp').mkdir()
```

```
img1.save(str(Path(p) / 'tmp' / 'img1.png'))
img2.save(str(Path(p) / 'tmp' / 'img2.png'))
```

```
%cd $p
!zip -r "tmp.zip" "tmp"
```

```
-----
FileNotFoundError                                Traceback (most recent call last)
<ipython-input-30-e0c49c38d470> in <module>
      9
     10 if not (Path(p) / 'tmp').exists():
--> 11     (Path(p) / 'tmp').mkdir()
     12
     13 img1.save(str(Path(p) / 'tmp' / 'img1.png'))
```

```
/usr/lib/python3.8/pathlib.py in mkdir(self, mode, parents, exist_ok)
    1286         self._raise_closed()
    1287     try:
-> 1288         self._accessor.mkdir(self, mode)
    1289     except FileNotFoundError:
    1290         if not parents or self.parent == self:
```

```
FileNotFoundError: [Errno 2] No such file or directory: '/content/drive/MyDrive/dev/prak_nn_1/tmp'
```

SEARCH STACK OVERFLOW

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

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

[Colab paid products](#) - [Cancel contracts here](#)