



МИНОБРНАУКИ РОССИИ
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«МИРЭА – Российский технологический университет»
РТУ МИРЭА

ИКБ направление «Киберразведка и противодействие угрозам с применением технологий искусственного интеллекта» 10.04.01

Кафедра КБ-4 «Интеллектуальные системы информационной безопасности»

Лабораторная работа №3

по дисциплине

«Анализ защищенности систем искусственного интеллекта»

Группа:
ББМО-01-22
Выполнил:
Некрасов Е.А

Проверил:
Спирин А.А.

Москва 2023

Загружаем библиотеки.

```
!pip install tf-keras-vis

Collecting tf-keras-vis
  Downloading tf_keras_vis-0.8.6-py3-none-any.whl (52 kB)
    52.1/52.1 kB 514.5 kB/s eta 0:00:00
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (1.11.4)
Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (9.4.0)
Collecting deprecated (from tf-keras-vis)
  Downloading Deprecated-1.2.14-py2.py3-none-any.whl (9.6 kB)
Requirement already satisfied: imageio in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (2.31.6)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (23.2)
Requirement already satisfied: wrapt<2,>=1.10 in /usr/local/lib/python3.10/dist-packages (from deprecated->tf-keras-vis) (1.14.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from imageio->tf-keras-vis) (1.23.5)
Installing collected packages: deprecated, tf-keras-vis
Successfully installed deprecated-1.2.14 tf-keras-vis-0.8.6

[2] %reload_ext autoreload
%autoreload 2
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import tensorflow as tf
from tf_keras_vis.utils import num_of_gpus
_, gpus = num_of_gpus()
print('Tensorflow recognized {} GPUs'.format(gpus))
from tensorflow.keras.preprocessing.image import load_img
from tensorflow.keras.applications.vgg16 import preprocess_input
```

Загружаем модель.

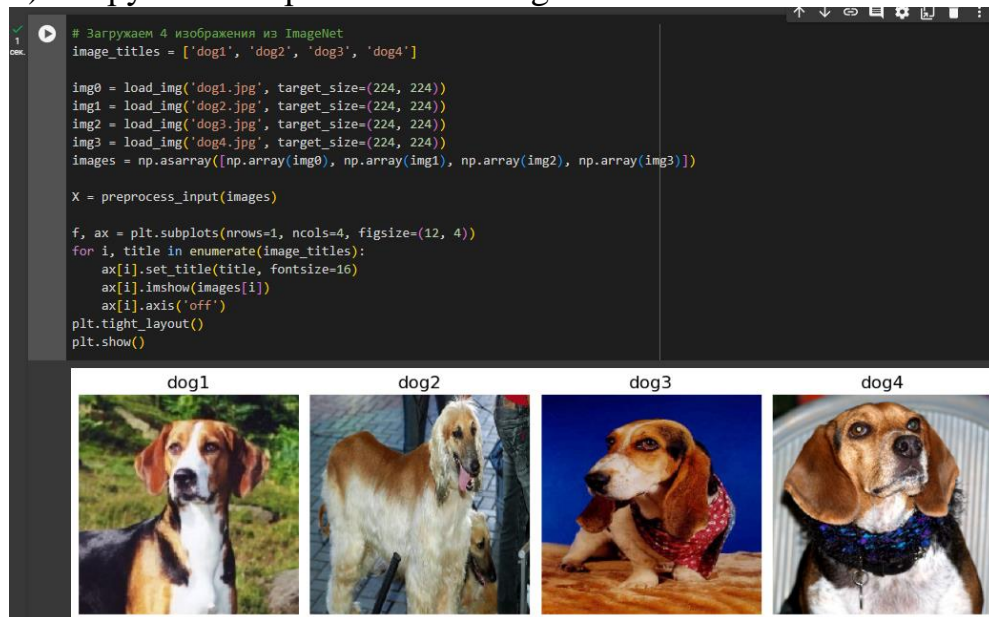
```
from tensorflow.keras.applications.vgg16 import VGG16 as Model
model = Model(weights='imagenet', include_top=True)
model.summary()

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_553467096/553467096 [=====] - 6s 0us/step
Model: "vgg16"

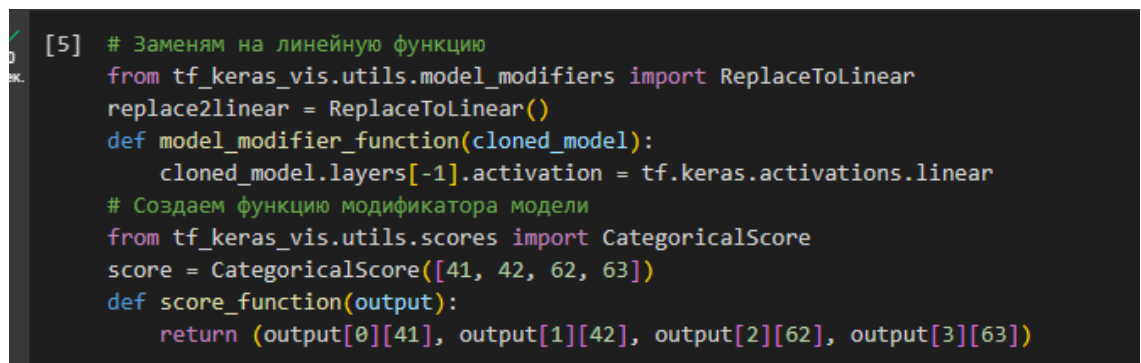
Layer (type)                 Output Shape                  Param #
=====
input_1 (InputLayer)         [None, 224, 224, 3]          0
block1_conv1 (Conv2D)         (None, 224, 224, 64)         1792
block1_conv2 (Conv2D)         (None, 224, 224, 64)         36928
block1_pool (MaxPooling2D)    (None, 112, 112, 64)         0
block2_conv1 (Conv2D)         (None, 112, 112, 128)        73856
block2_conv2 (Conv2D)         (None, 112, 112, 128)        147584
block2_pool (MaxPooling2D)    (None, 56, 56, 128)          0
block3_conv1 (Conv2D)         (None, 56, 56, 256)          295168
block3_conv2 (Conv2D)         (None, 56, 56, 256)          590880
block3_conv3 (Conv2D)         (None, 56, 56, 256)          590880
block3_pool (MaxPooling2D)    (None, 28, 28, 256)          0
block4_conv1 (Conv2D)         (None, 28, 28, 512)          1180160
block4_conv2 (Conv2D)         (None, 28, 28, 512)          2359808
block4_conv3 (Conv2D)         (None, 28, 28, 512)          2359808
block4_pool (MaxPooling2D)    (None, 14, 14, 512)          0
block5_conv1 (Conv2D)         (None, 14, 14, 512)          2359808
block5_conv2 (Conv2D)         (None, 14, 14, 512)          2359808
block5_conv3 (Conv2D)         (None, 14, 14, 512)          2359808
block5_pool (MaxPooling2D)    (None, 7, 7, 512)            0
flatten (Flatten)             (None, 25088)                 0
fc1 (Dense)                   (None, 4096)                  102764544
fc2 (Dense)                   (None, 4096)                  16781312
predictions (Dense)           (None, 1000)                  4097000

Total params: 138357544 (527.79 MB)
Trainable params: 138357544 (527.79 MB)
Non-trainable params: 0 (0.00 Byte)
```

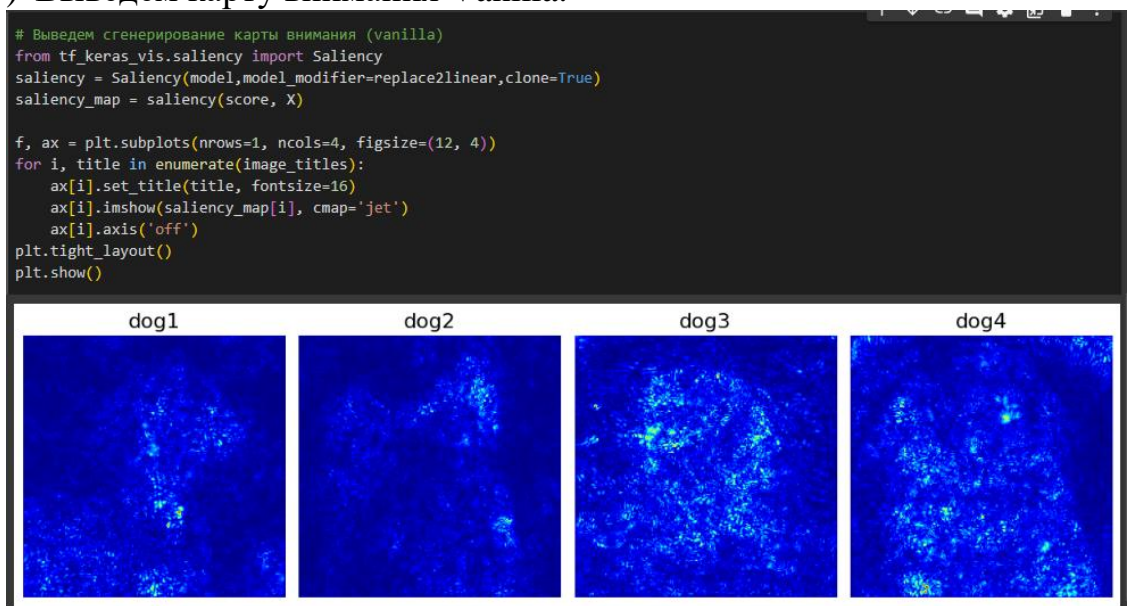
1) Загрузим изображения из ImageNet.



2) Заменим на линейную функцию и создадим активатор модели активации.

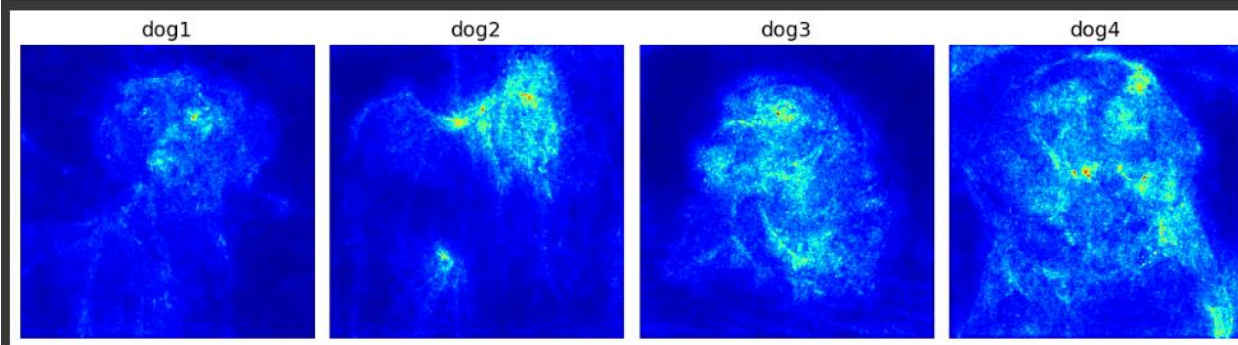


3) Выведем карту внимания Vanilla.



4) Уменьшим шум с помощью SmoothGrad.

```
saliency_map = saliency(score,X,smooth_samples=20,smooth_noise=0.20)
f, ax = plt.subplots(nrows=1, ncols=4, figsize=(12, 4))
for i, title in enumerate(image_titles):
    ax[i].set_title(title, fontsize=14)
    ax[i].imshow(saliency_map[i], cmap='jet')
    ax[i].axis('off')
plt.tight_layout()
plt.savefig('smoothgrad.png')
plt.show()
```



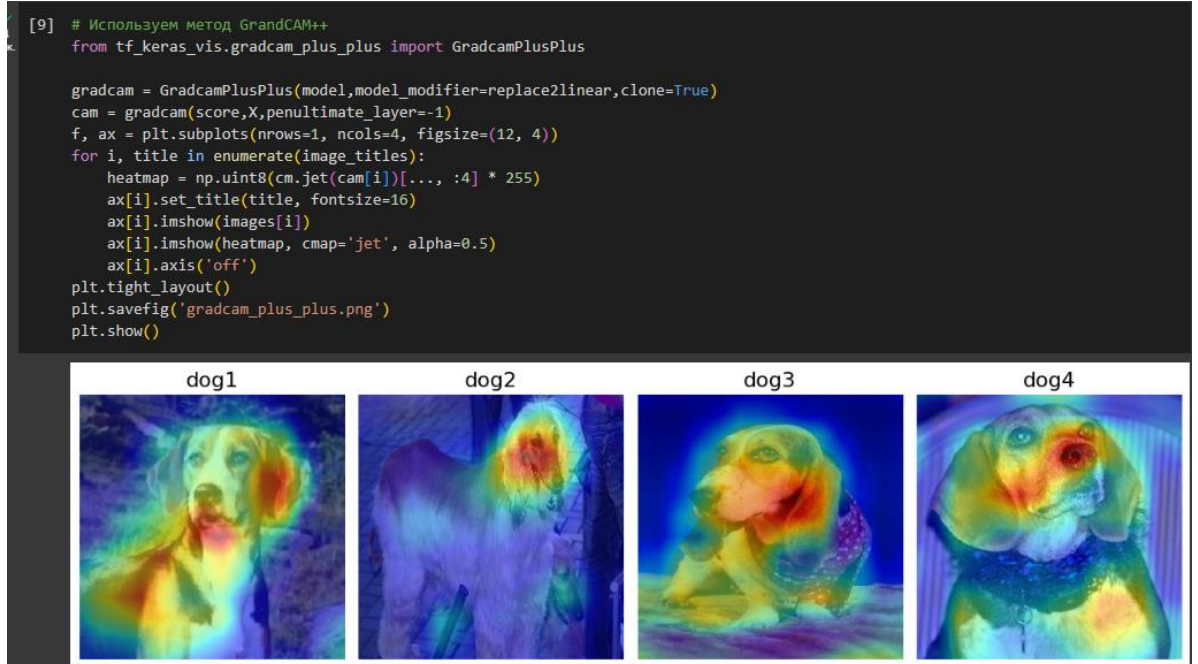
5) Сравним полученные значения с функцией GradCAM.

```
# Используем метод GradCAM
from matplotlib import cm
from tf_keras_vis.gradcam import Gradcam

gradcam = Gradcam(model,model_modifier=replace2linear,clone=True)
cam = gradcam(score,X,penultimate_layer=-1)
f, ax = plt.subplots(nrows=1, ncols=4, figsize=(12, 4))
for i, title in enumerate(image_titles):
    heatmap = np.uint8(cm.jet(cam[i])[..., :4] * 255)
    ax[i].set_title(title, fontsize=16)
    ax[i].imshow(images[i])
    ax[i].imshow(heatmap, cmap='jet', alpha=0.5) # overlay
    ax[i].axis('off')
plt.tight_layout()
plt.show()
```



6) Сравним также с методом GradCAM++, он имеет более хорошую визуализацию.



Вывод

SmoothGrad — техника, предназначенная для сглаживания карт выделенности с целью снижения шума и повышения интерпретируемости.

GradCAM — метод визуализации активации нейронов в сверточных нейронных сетях, который позволяет понять, какие участки входного изображения были наиболее значимыми для принятия окончательного решения моделью.

GradCAM++ - расширенный метод GradCAM, который учитывает также градиенты второго порядка.

Выбор между этими методами следует делать в зависимости от поставленной задачи.