



МИНОБРНАУКИ РОССИИ

Федеральное государственное бюджетное образовательное учреждение
высшего образования

«МИРЭА – Российский технологический университет»

РТУ МИРЭА

Институт кибербезопасности и цифровых технологий
Кафедра КБ-4 «Интеллектуальные системы информационной
безопасности»

Отчёт по лабораторной работе №3

**по дисциплине «Анализ защищенности систем искусственного
интеллекта»**

Выполнил:
Евдокимов А.М.
Группа: ББМО-02-23

Москва - 2024

Установим инструмент для визуализации для TensorFlow Keras.

```
[1] !pip install tf-keras-vis

Collecting tf-keras-vis
  Downloading tf_keras_vis-0.8.7-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (1.13.1)
Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (11.0.0)
Requirement already satisfied: deprecated in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (1.2.15)
Requirement already satisfied: imageio in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (2.36.0)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from tf-keras-vis) (24.2)
Requirement already satisfied: wrapt<2,>=1.10 in /usr/local/lib/python3.10/dist-packages (from deprecated->tf-keras-vis) (1.16.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from imageio->tf-keras-vis) (1.26.4)
Downloading tf_keras_vis-0.8.7-py3-none-any.whl (52 kB)
52.5/52.5 kB 2.1 MB/s eta 0:00:00
Installing collected packages: tf-keras-vis
Successfully installed tf-keras-vis-0.8.7
```

Далее подключим необходимые библиотеки. Загрузим предварительно обученную модель VGG16, на ImageNet датасете. После чего отобразим сводку по модели.

```
[2] import numpy as np
import tensorflow as tf

# VGG16 - это известная архитектура нейронной сети, предназначенная для классификации изображений.
from tensorflow.keras.applications.vgg16 import VGG16 as Model

# Загружаем модель VGG16:
# - weights='imagenet' означает, что модель будет загружена с предобученными весами на датасете ImageNet.
model = Model(weights='imagenet', include_top=True)

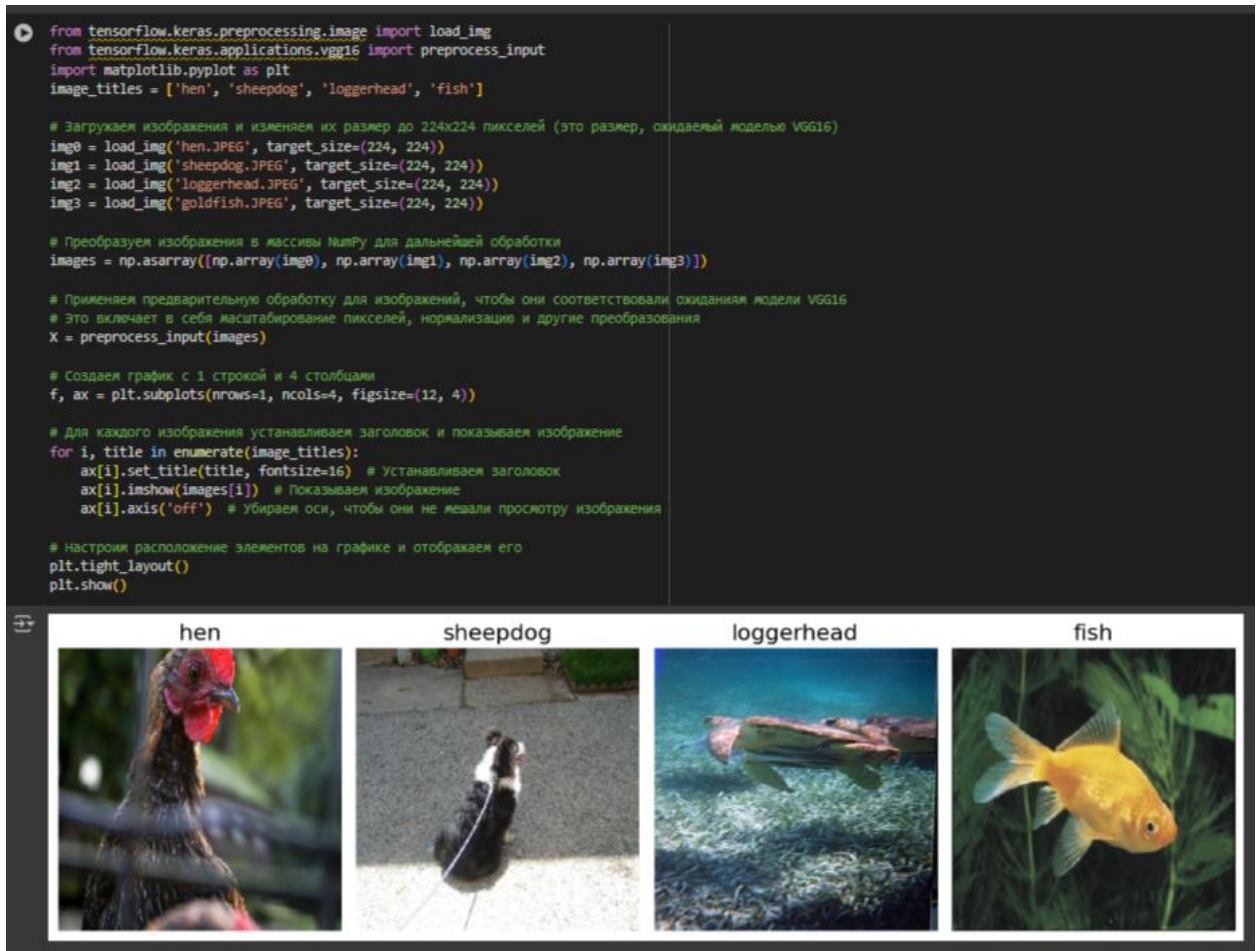
model.summary()
```

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_weights_tf_dim_ordering_tf_kernels.h5
553467096/553467096 950 B/s step
Model: "vgg16"

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	(None, 224, 224, 3)	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1,792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36,928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73,856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147,584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295,168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590,000
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590,000
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1,100,160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2,359,000
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2,359,000
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2,359,000
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2,359,000
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2,359,000

Загрузим несколько изображений датасета ImageNet и выполним их предварительную обработку перед использованием.

Отообразим на одном графическом представлении все наши изображения.



Заменяем функцию активации на линейную

```
[6] # ReplaceToLinear – это утилита для модификации модели, заменяющая её последний слой на линейную активацию.
from tf_keras_vis.utils.model_modifiers import ReplaceToLinear

# Создаем объект для замены последнего слоя на линейную активацию.
replace2linear = ReplaceToLinear()
def model_modifier_function(cloned_model):
    # Здесь мы меняем активацию последнего слоя на линейную.
    # Это件用, например, для визуализаций, когда не требуется применять активацию ReLU или Softmax
    cloned_model.layers[-1].activation = tf.keras.activations.linear
```

Создадим функцию очков соответствия каждому изображению

```
[7] from tf_keras_vis.utils.scores import CategoricalScore

# Создаем объект CategoricalScore.
# В качестве аргумента передаем список индексов классов, для которых будем вычислять score.
# Здесь, например, передаем список [11, 12, 13, 14], что означает, что нас интересуют эти классы.
score = CategoricalScore([11, 12, 13, 14])
```

Создадим карту внимания (vanilla)

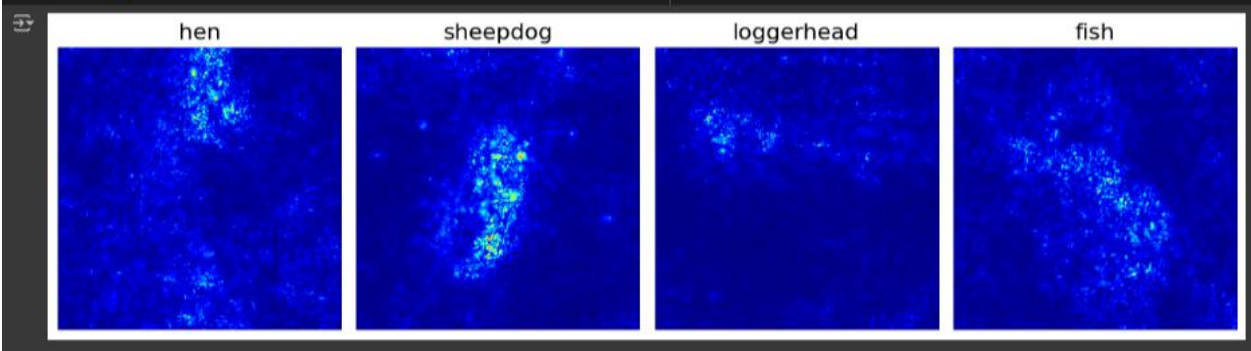
```
[8] # Импортируем класс Saliency из библиотеки tf-keras-vis, который используется для создания карт важности (saliency maps).
# Saliency map показывает, какие части изображения наиболее важны для принятия решения моделью.
from tf_keras_vis.saliency import Saliency

# Инициализируем объект Saliency для визуализации важности:
saliency = Saliency(model, model_modifier=replace2linear, clone=True)

# Получаем карту важности (saliency map), используя ранее определенный score для классов [11, 12, 13, 14] и входные данные X.
saliency_map = saliency(score, X)

# Создаем подграфики для отображения карт saliency для всех 4 изображений.
f, ax = plt.subplots(nrows=1, ncols=4, figsize=(12, 4))
for i, title in enumerate(image_titles):
    ax[i].set_title(title, fontsize=16)
    ax[i].imshow(saliency_map[i], cmap='jet')
    ax[i].axis('off')

plt.tight_layout()
plt.show()
```

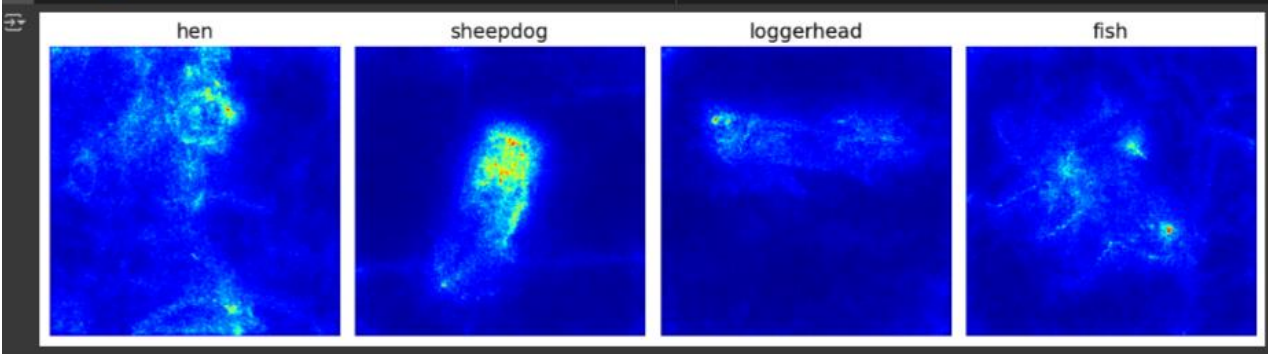


Уменьшим шум для карт влияния

```
# Генерация карты важности (saliency map) с использованием SmoothGrad.
# SmoothGrad улучшает стандартный saliency map, добавляя случайный шум для сглаживания.
# - score – это наш объект.
# - X – это входные данные (изображения).
# - smooth_samples=20 – количество случайных образцов, которые будут добавлены для сглаживания.
# - smooth_noise=0.20 – уровень шума, который добавляется для сглаживания карты.
saliency_map = saliency(score, X, smooth_samples=20, smooth_noise=0.20)

# Создаем подграфики для отображения карт saliency для всех 4 изображений.
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')

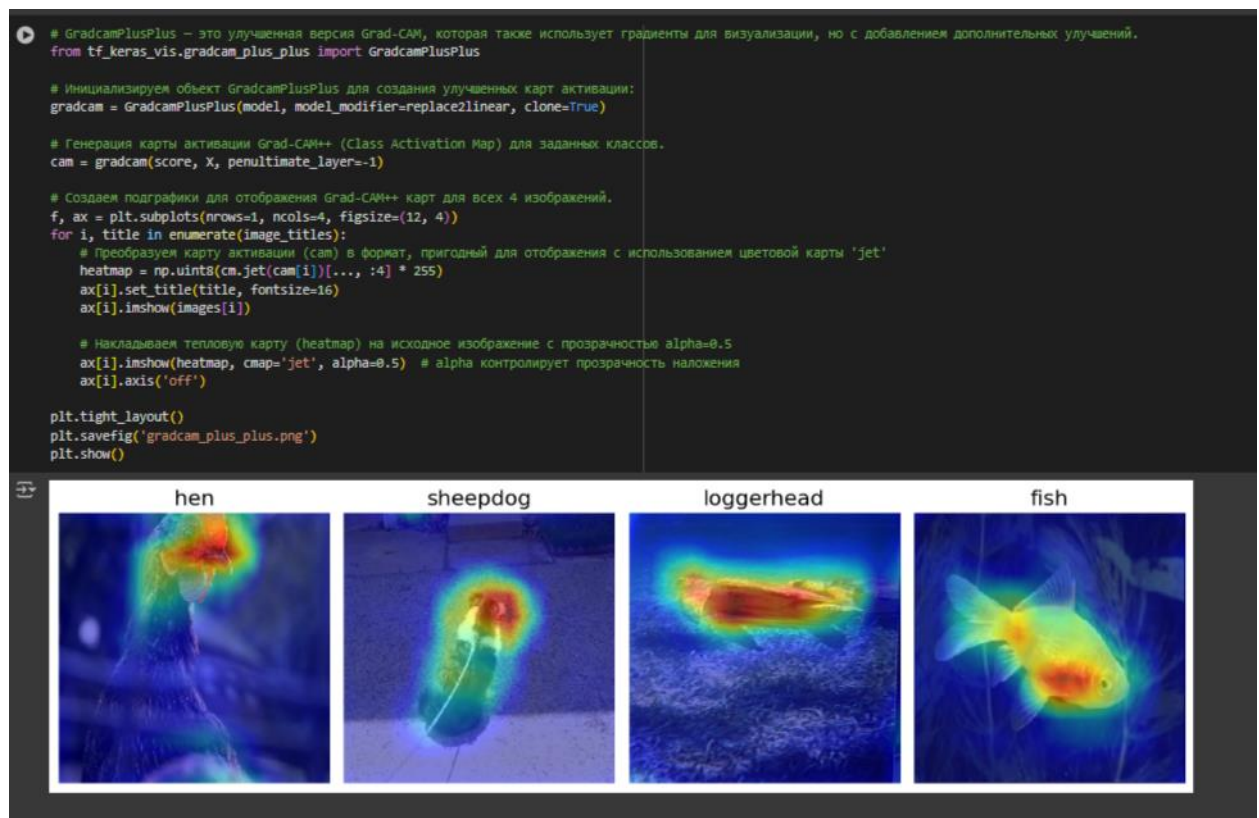
# Плотнее распределяем элементы на графике
plt.tight_layout()
plt.savefig('smoothgrad.png')
plt.show()
```



Используем GradCAM



Используем GradCAM++



Таким образом, использование методов визуализации Grad-CAM, Grad-CAM++, Saliency, SmoothGrad и подобных, может быть полезным для понимания того, какие части изображений были наиболее важными при принятии решений моделью машинного обучения. Методы Grad-CAM и Grad-CAM++ позволяют визуализировать активации в различных частях изображений, позволяя понять, где модель фокусируется при определении классов. Карты активаций могут помочь исследовать, какие объекты или части изображения были ключевыми для принятия решения моделью. Использование метода SmoothGrad может помочь снизить шум и сделать карты сайленси более интерпретируемыми.

Вывод

В результате выполнения работы получен некоторый опыт работы с инструментами получения карт активации для более подробного объяснения работы нейронных сетей и других моделей машинного обучения. Были также проведены эксперименты по созданию карты активации моделей машинного обучения с примерами на изображениях из набора данных ImageNet и модели VGG.

В рамках работы были рассмотрены несколько методов создания карт активации, в том числе GradCAM. Были проведены эксперименты с различными параметрами и модификациями методов.

После завершения экспериментов и проведения анализа полученных результатов были сделаны выводы касательного полноты и точности объяснения работы методов генерации карт.