

They compute the probabilities: the probability that the first classifier has higher scores than the second, the probability that differences are within the region of practical equivalence (rope), or that the second classifier has higher scores. We will refer to these probabilities as: p_{left} , p_{rope} and p_{right} .

1: Two shortcut functions can be used for comparison on single and on multiple data sets. If classifier_legacy_acc and classifier_new_acc contain a list of average classification accuracies of classifier_legacy and classifier_new on a collection of data sets, we can call (Actual outputs may differ due to Monte Carlo sampling.)

```
In [18]: print('p_ (left), p_ (rope), p_ (right)$ using the two_on_multiple function: ')
print(two_on_multiple(classifier_legacy_acc, classifier_new_acc, rope=1))

# With some additional arguments, the function can also plot the posterior distribution from
# which these probabilities came.
# Tests are packed into test classes.
# The above call is equivalent to

print('p_ (left), p_ (rope), p_ (right)$ using the SignedRankTest.probs function: ')
print(SignedRankTest.probs(classifier_legacy_acc, classifier_new_acc, rope=1))

# and to get a plot, we call

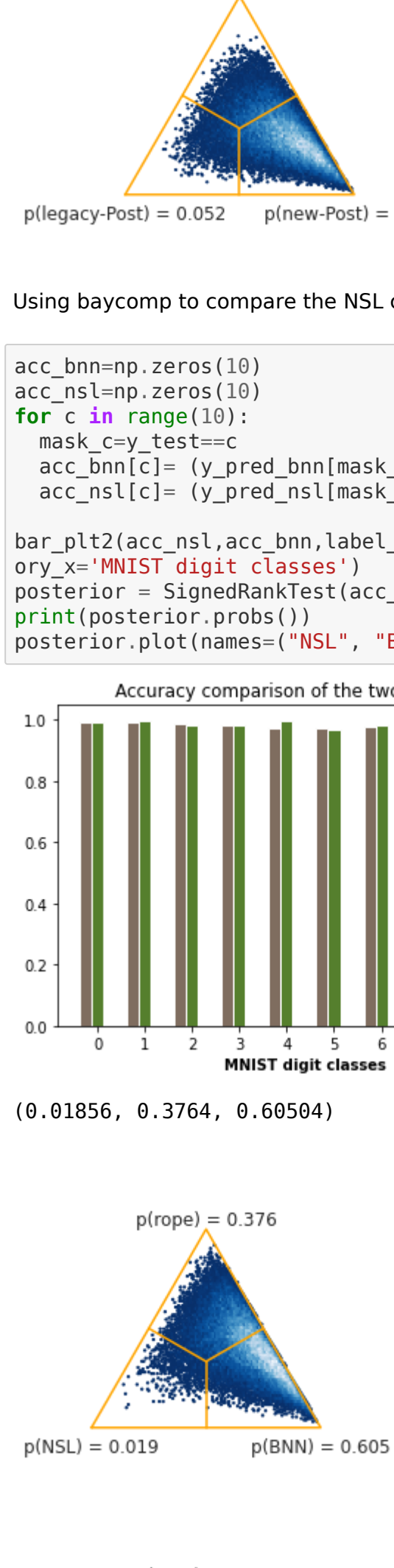
print(SignedRankTest.plot(classifier_legacy_acc, classifier_new_acc, rope=1, names= ("Legacy-SRT",
"New-SRT")))

# To switch to another test, use another class:
SignTest.probs(classifier_legacy_acc, classifier_new_acc, rope=1)
# Finally, we can construct and query sampled posterior distributions.

posterior = SignedRankTest(classifier_legacy_acc, classifier_new_acc, rope=1)
print(posterior.probs())
posterior.plot(bnames= ("Legacy-Post", "New-Post"))

p_ (left), p_ (rope), p_ (right)$ using the two_on_multiple function:
(0.0526, 0.16366, 0.78434)
p_ (left), p_ (rope), p_ (right)$ using the SignedRankTest.probs function:
(0.0515, 0.16388, 0.78462)
Figure(432x288)
(0.0517, 0.16366, 0.78464)
```

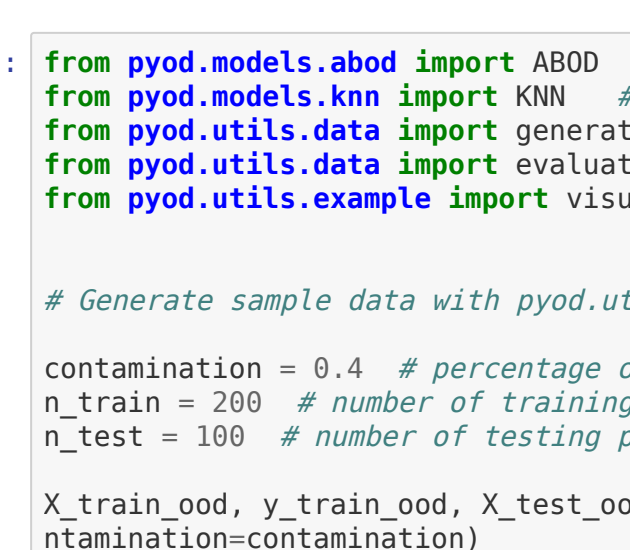
Out[18]:



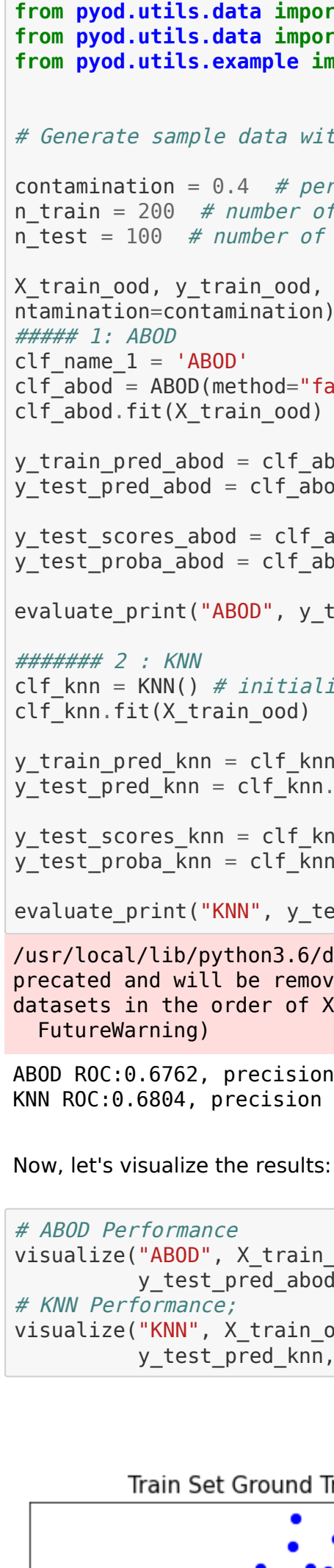
Using baycomp to compare the NSL classifier with the BNN classifier on the MNIST dataset:

```
In [19]: acc_bnn=np.zeros(10)
acc_nsl=np.zeros(10)
for c in range(10):
    mask_c=y_test==c
    acc_bnn[c]=(y_pred_bnn[mask_c]==c).mean()
    acc_nsl[c]=(y_pred_nsl[mask_c]==c).mean()

bar_plt2(acc_nsl,acc_bnn,label_1='NSL',label_2='BNN',X_LABELS=list(np.arange(10).astype(str)),Category_x='MNIST digit classes')
posterior = SignedRankTest(acc_nsl, acc_bnn, rope=0.005)
print(posterior.probs())
posterior.plot(names= ("NSL", "BNN"))
```



Out[19]:



PyOD

PyOD is arguably the most comprehensive and scalable Outlier Detection Python toolkit out there that includes implementation of more than 30 detection algorithms! It is somewhat rare for a student-maintained PyPi package to incorporate software engineering best practices that ensures that model classes implemented are covered by unit testing with cross platform continuous integration, code coverage and code maintainability checks. This combined with a clean unified API, detailed documentation and just-in-time (JIT) compiled execution makes it an absolute breeze to both learn about the different techniques and use it in practice. The efforts invested by the authors towards careful parallelization has resulted in extremely fast and scalable outlier detection code that is also seamlessly compatible across Python 2 and 3 across major operating systems (Windows, Linux and MacOS). In the example cell below, we train and visualize the results of two inlier-outlier detector binary classifiers on a synthetic dataset: the Angle-Based Outlier Detector (ABOD) and the KNN outlier detector.

```
In [20]: from pyod.models.abod import ABOD
from pyod.models.knn import KNN # KNN detector
from pyod.utils.data import generate_data
from pyod.utils.data import evaluate_print
from pyod.utils.example import visualize

# Generate sample data with pyod.utils.data.generate_data():

contamination = 0.4 # percentage of outliers
n_train = 200 # number of training points
n_test = 100 # number of testing points

X_train_od, y_train_od, X_test_od, y_test_od = generate_data(n_train=n_train, n_test=n_test, contamination=contamination)

##### 1: ABOD
clf_name_1 = 'ABOD'
clf_abod = ABOD(method='fast') # initialize detector
clf_abod.fit(X_train_od)

y_train_pred_abod = clf_abod.predict(X_train_od) # binary labels
y_test_pred_abod = clf_abod.predict(X_test_od) # binary labels
y_test_scores_abod = clf_abod.decision_function(X_test_od) # raw outlier scores
y_test_proba_abod = clf_abod.predict_proba(X_test_od) # outlier probability
evaluate_print("ABOD", y_test_od, y_test_scores_abod) # performance evaluation

##### 2 : KNN
clf_knn = KNN() # initialize detector
clf_knn.fit(X_train_od)

y_train_pred_knn = clf_knn.predict(X_train_od) # binary labels
y_test_pred_knn = clf_knn.predict(X_test_od) # binary labels
y_test_scores_knn = clf_knn.decision_function(X_test_od) # raw outlier scores
y_test_proba_knn = clf_knn.predict_proba(X_test_od) # outlier probability
evaluate_print("KNN", y_test_od, y_test_scores_knn) # performance evaluation

/usr/local/lib/python3.6/dist-packages/pyod/utils/data.py:189: FutureWarning: behaviour="old" is deprecated and will be removed in version 0.8.0. Please use behaviour="new", which makes the returned datasets in the order of X_train, X_test, y_train, y_test.
FutureWarning)
```

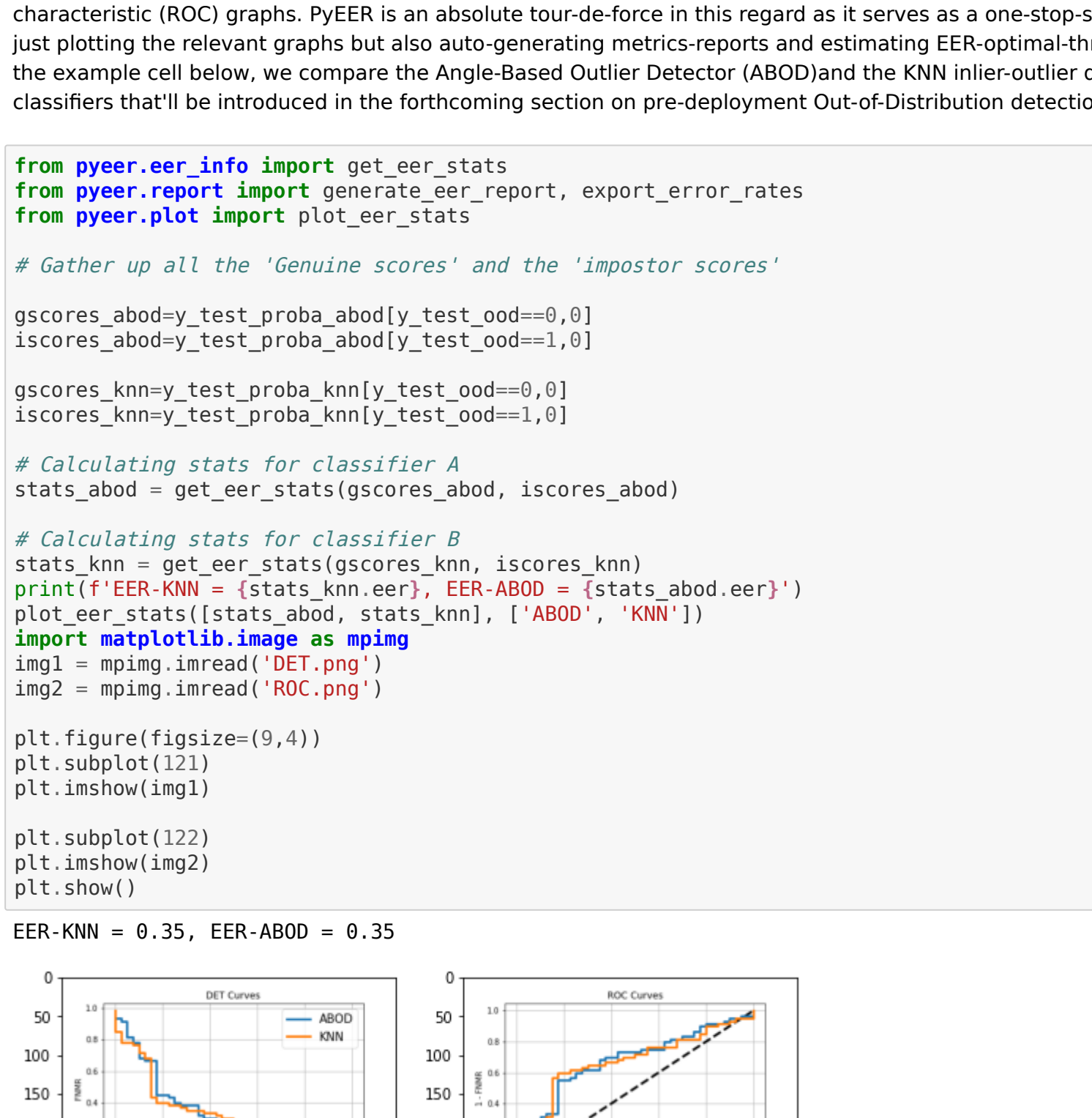
ABOD ROC:0.6762, precision @ rank n:0.575
KNN ROC:0.6894, precision @ rank n:0.55

Now, let's visualize the results:

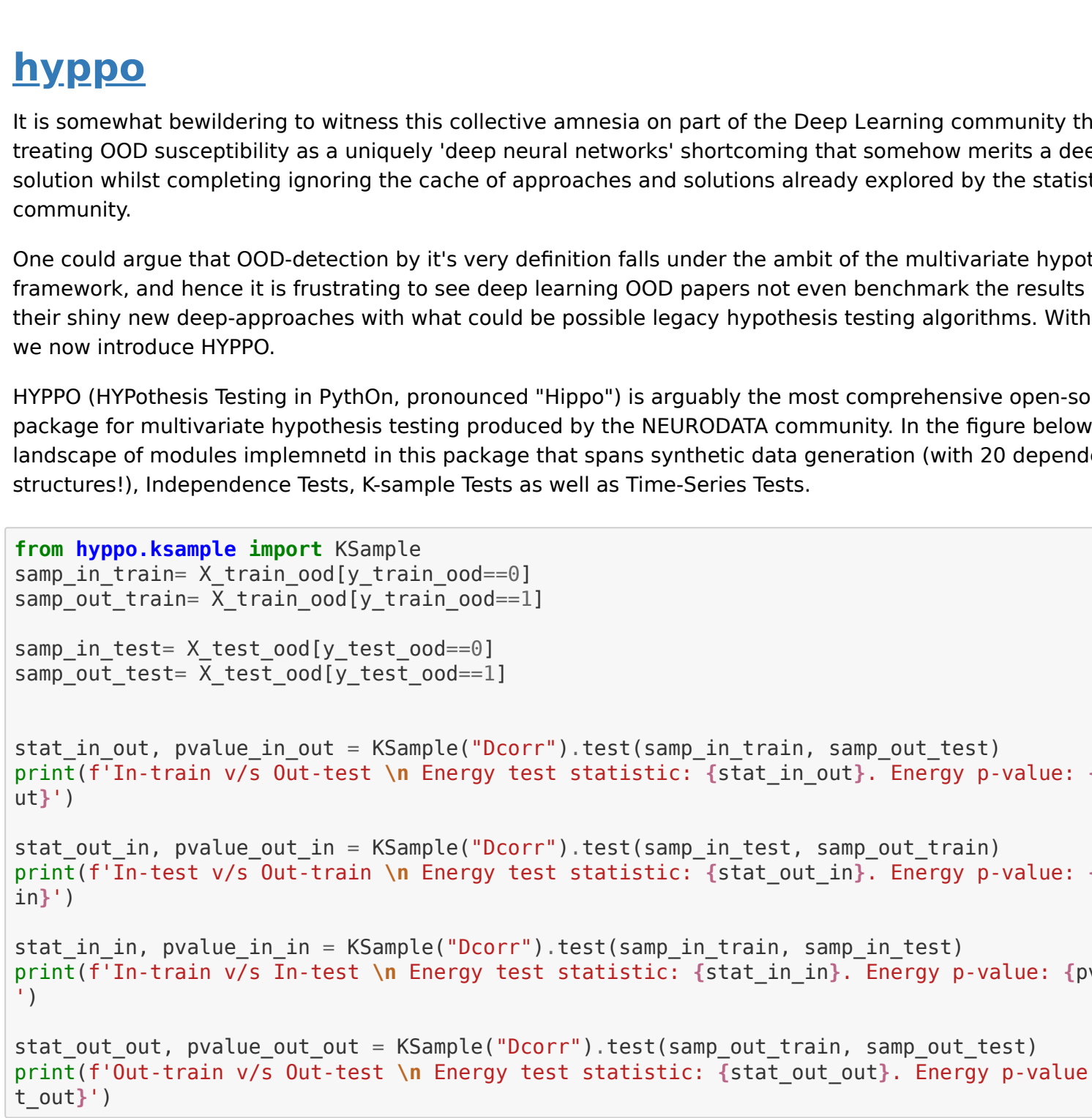
```
In [21]: # ABOD Performance
visualize("ABOD", X_train_od, y_train_od, X_test_od, y_test_od, y_train_pred_abod,
y_test_pred_abod, show_figure=True, save_figure=False)

# KNN Performance
visualize("KNN", X_train_od, y_train_od, X_test_od, y_test_od, y_train_pred_knn,
y_test_pred_knn, show_figure=True, save_figure=False)
```

Demo of ABOD Detector



Demo of KNN Detector



PyEER

Another way of comparing two classifiers, especially in the context of solving the binary authentication problem (Not surveillance but Authentication) is by plotting the comparative detection error tradeoff (DET) and Receiver operating characteristic (ROC) graphs. PyEER is an absolute tour-de-force in this regard as it serves as a one-stop-shop for not just plotting the relevant graphs but also auto-generating metrics-reports and estimating EER-optimal-thresholds. In the example cell below, we compare the Angle-Based Outlier Detector (ABOD) and the KNN inlier-outlier detector binary classifiers that'll be introduced in the forthcoming section on pre-deployment Out-of-Distribution detection techniques.

```
In [22]: from pyeer.eer_info import get_eer_stats
from pyeer.eer import get_neural_network, export_error_rates
from pyeer.plot import plot_eer_stats

# Gather up all the 'Genuine scores' and the 'impostor scores'

gscores_abod=y_test_proba_abod[y_test_od==0,0]
iscores_abod=y_test_proba_abod[y_test_od==1,0]

gscores_knn=y_test_proba_knn[y_test_od==0,0]
iscores_knn=y_test_proba_knn[y_test_od==1,0]

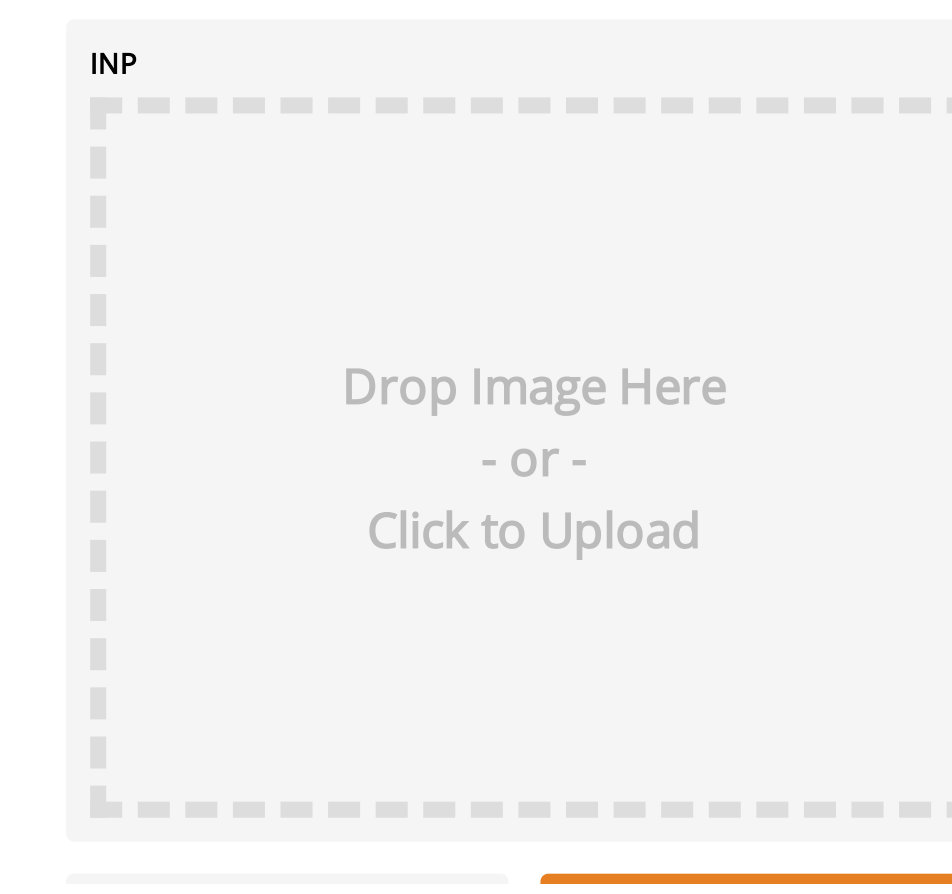
# Calculating stats for classifier A
stats_abod = get_eer_stats(gscores_abod, iscores_abod)

# Calculating stats for classifier B
stats_knn = get_eer_stats(gscores_knn, iscores_knn)
print(f'EER-KNN = {stats.knn.eer}, EER-ABOD = {stats.abod.eer}')
plot_eer_stats([stats_abod, stats_knn], ['ABOD', 'KNN'])
import matplotlib.image as mpimg
img1 = mpimg.imread('DET.png')
img2 = mpimg.imread('ROC.png')

plt.figure(figsize=(9,4))
plt.subplot(121)
plt.imshow(img1)

plt.subplot(122)
plt.imshow(img2)
plt.show()
```

EER-KNN = 0.35, EER-ABOD = 0.35



hyppo

It is somewhat bewildering to witness this collective amnesia on part of the Deep Learning community that keeps treating OOD susceptibility as a uniquely 'deep neural networks' shortcoming that somehow merits a deep-learning solution whilst completely ignoring the cache of approaches and solutions already explored by the statistics community.

One could argue that OOD-detection by it's very definition falls under the ambit of the multivariate hypothesis testing framework, and hence it is frustrating to see deep learning OOD papers not even benchmark the results obtained by their shiny new deep-approaches with what could be possible legacy hypothesis testing algorithms. With this setting, we now introduce HYPPPO.

HYPPPO (HYPPPOthesis Testing in Python, pronounced "Hippo") is arguably the most comprehensive open-source software package for multivariate hypothesis testing produced by the NEURODATA community. In the figure below, we see the landscape of modules implemented in this package that spans synthetic data generation (with 20 dependency structures!), Independence Tests, K-sample Tests as well as Time-Series Tests.

```
In [23]: from hyppo.ksample import KSample
samp_in_train= X_train_od[y_train_od==0]
samp_out_train= X_train_od[y_train_od==1]

samp_in_test= X_test_od[y_test_od==0]
samp_out_test= X_test_od[y_test_od==1]

stat_in_out, pvalue_in_out = KSample("Dcorr").test(samp_in_train, samp_out_train)
print(f'In-train v/s Out-Test \n Energy test statistic: {stat_in_out}, Energy p-value: {pvalue_in_out}')

stat_out_in, pvalue_out_in = KSample("Dcorr").test(samp_in_test, samp_out_test)
print(f'In-train v/s Out-Test \n Energy test statistic: {stat_out_in}, Energy p-value: {pvalue_out_in}')

stat_in_in, pvalue_in_in = KSample("Dcorr").test(samp_in_train, samp_in_test)
print(f'In-train v/s In-Test \n Energy test statistic: {stat_in_in}, Energy p-value: {pvalue_in_in}')

stat_out_out, pvalue_out_out = KSample("Dcorr").test(samp_out_train, samp_out_test)
print(f'Out-train v/s Out-Test \n Energy test statistic: {stat_out_out}, Energy p-value: {pvalue_out_out}')

In-train v/s Out-Test
Energy test statistic: 0.5317830910984214, Energy p-value: 1.72339641676499e-20
In-test v/s Out-Test v/s Out-Test \n Energy test statistic: 0.6383912273882784, Energy p-value: 1.970628744555839e-21
Energy test statistic: 0.6383912273882784, Energy p-value: 1.970628744555839e-21
In-train v/s In-Test
Energy test statistic: -0.00567118456830132, Energy p-value: 1.0
Out-train v/s Out-Test
Energy test statistic: -0.00621706290520194, Energy p-value: 0.6143056828394889
```

Gradio

Having a nice GUI to interact with the model you have just trained has thus far required a fair amount of JavaScript-front-end gimmickry or the Heroku-Flask route that can take focus away from the algorithmics. Thanks to Gradio, one can now quickly fire up a gui with < 10 lines of Python with pre-built input modules that cover textual input, image-inputs with an awesome Toast-UI image-editor and a sketchpad to boot as well! This past year, I have heavily used Gradio in my workflow, using it to investigate why Twitter's saliency cropping algorithm yields such racist results to why Onions were triggering NSFW filters on facebook (See the blogpost for the links and imagery).

(The NSFW-Onion fiasco Colab notebook can be found [here](#)).

In the example cell below, we demonstrate two simple examples of using Gradio to fire-up UIs to stress-test the MNIST classification BNN model we just trained above with a sketchpad input and to demonstrate the ease of using the inceptionV3 model to classify images. The Gradio team has also rapidly added explainability and embeddings-visualization tools, and implemented SOTA blind super resolution and Real-Time High-Resolution Background Matting UIs as well!

Saliency hub link

```
import gradio as gr
import requests

inception_net = tf.keras.applications.InceptionV3() # Load the model
# Download human-readable labels for ImageNet.
response = requests.get("https://git.io/J3jKW")
labels = response.text.split("\n")

def classify_image(inp):
    print(inp.shape)
    inp = inp.reshape((-1, 299, 299, 3))
    inp = tf.keras.applications.inception_v3.preprocess_input(inp)
    prediction = inception_net.predict(inp).flatten()
    return {labels[i]: float(prediction[i]) for i in range(1000)}

image = gr.inputs.Image(shape=(299, 299, 3))
label = gr.outputs.Label(num_top_classes=3)

gr.Interface(fn=classify_image, inputs=image, outputs=label, capture_session=True).launch()
```

```
In [24]: import gradio as gr
import requests

# We use the LARGO trained BNN to launch an interactive UI that facilitates a sketchpad input and p
prediction
def classify_image():
    print(image.shape)
    prediction = model_bnn.predict(image.reshape((-1,28,28,1))).tolist()[0]
    return {str(i): prediction[i] for i in range(10)}

sketchpad = gr.inputs.Sketchpad()
label = gr.outputs.Label(num_top_classes=3)
gr.Interface(fn=classify_image, inputs=sketchpad, outputs=label, capture_session=True).launch()
```

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
This share link will expire in 24 hours. If you need a permanent link, email support@gradio.app
Running on External URL: <https://42147.gradio.app>
Interface loading below...

Live at <https://17709.gradio.app>. Copy Link



Out[24]: <Flask 'gradio.networking'>,
<'http://127.0.0.1:7868/'>,
<'https://42147.gradio.app'>

```
In [25]: inception_net = tf.keras.applications.InceptionV3() # Load the model
# Download human-readable labels for ImageNet.
response = requests.get("https://git.io/J3jKW")
labels = response.text.split("\n")

def classify_image(inp):
    print(inp.shape)
    inp = inp.reshape((-1, 299, 299, 3))
    inp = tf.keras.applications.inception_v3.preprocess_input(inp)
    prediction = inception_net.predict(inp).flatten()
    return {labels[i]: float(prediction[i]) for i in range(1000)}

image = gr.inputs.Image(shape=(299, 299))
label = gr.outputs.Label(num_top_classes=3)

gr.Interface(fn=classify_image, inputs=image, outputs=label, capture_session=True).launch()
```

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/inception_v3/inception_v3_weights_tf_dim_ordering_tf_kernels.h5
96116736/96112376 [=====] - 3s 0us/step
Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
This share link will expire in 24 hours. If you need a permanent link, email support@gradio.app
Running on External URL: <https://17709.gradio.app>
Interface loading below...

gradio

Out[25]: <Flask 'gradio.networking'>,
<'http://127.0.0.1:7868/'>,
<'https://17709.gradio.app'>