
LINK TO GITHUB

<https://github.com/RockFZ/CS542-Summer2022.git>

ANOMALY DETECTOR MODELS

USING KEYSTROKE DYNAMICS ANALYSIS

By Fucheng Zhu

KEYSTROKE DYNAMICS



THE DATA

- **Collecting data**

- 51 Subjects corresponding to 51 typist
- 8 sessions
 - 50 reps each session per person

- **Features**

- Hold time for a key
 - keydown-keydown time for a digraph
-

GOALS

- To develop models that detect anomaly imposters
 - To compare different models in terms of performance metrics
 - To try to reasoning those performance differences
-

THE DATA

```
YTrain = D[(D['sessionIndex'] <= 4 ) & (D['subject'] == evalSubject)]
```

```
YTrain = np.array(YTrain)[: ,3:]
```

```
YScore0 = D[(D['sessionIndex'] <= 4 ) & (D['subject'] == evalSubject)]
```

```
YScore0 = np.array(YScore0)[: ,3:]
```

```
YScore1 = D[(D['sessionIndex'] == 1 ) & (D['subject'] != evalSubject) & (D['rep'] <= 5)]
```

```
YScore1 = np.array(YScore1)[: ,3:]
```

- YTrain: first half of a subject's reps—TrainingSet
 - shape: $200 * 31$
 - YScore0: second half of a subject's reps—TestingSet
 - shape: $200 * 31$
 - YScore1: first 5 reps of all other subjects—ImposterSet (Validation)
 - shape: $250 * 31$
-

MODELS

- Statistic models:
 - Euclidean distance detector
 - Manhattan distance detector
 - Mahalanobis distance detector
 - Classification model
 - One class SVM with linear kernel
 - One class SVM with other kernels
 - NN
-

EUCLIDEAN DETECTOR

$$d(x, y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

MANHATTAN DETECTOR

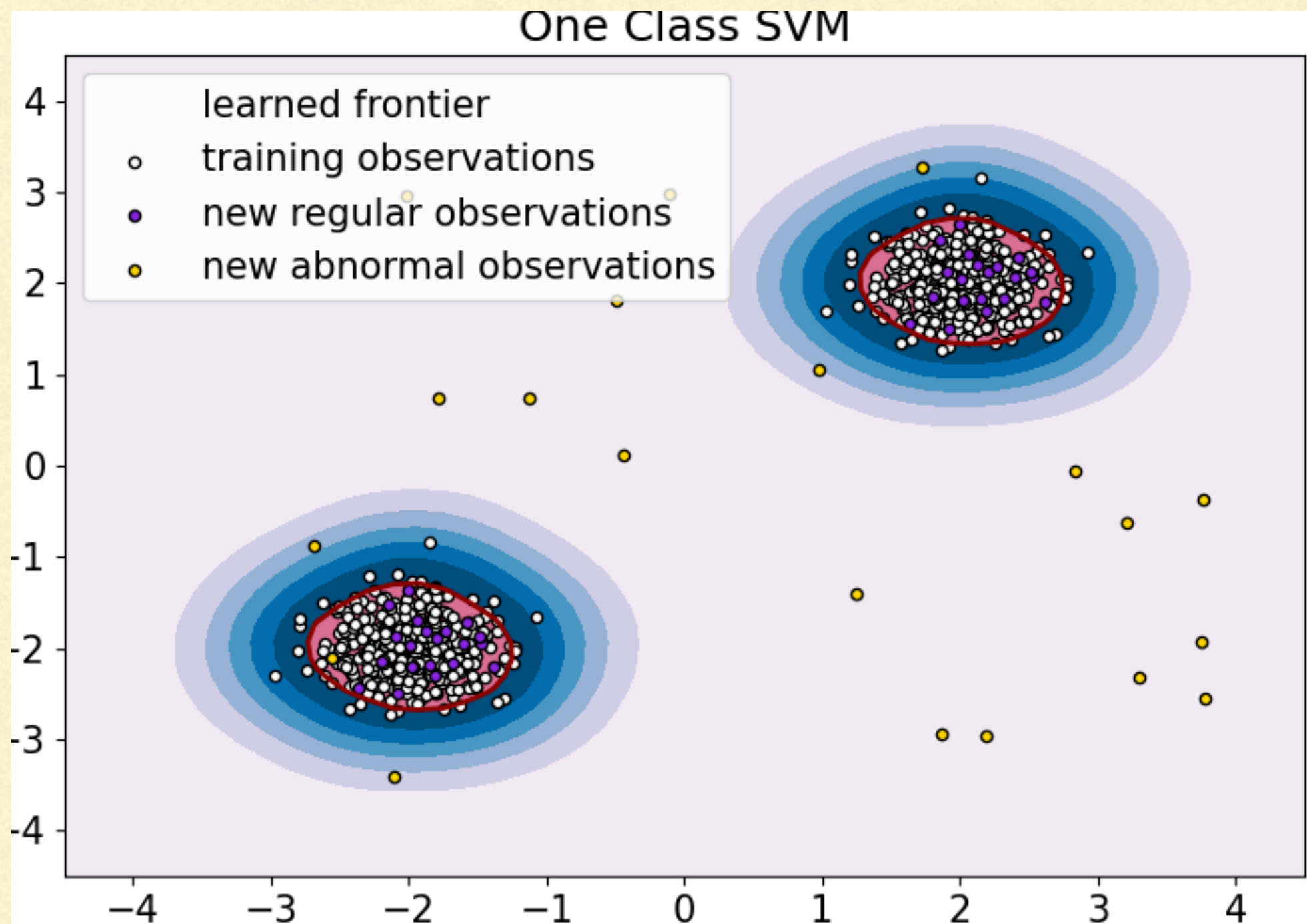
$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

MAHALANOBIS DETECTOR

$$D_M(\vec{x}) = \sqrt{(\vec{x} - \vec{\mu})^T S^{-1} (\vec{x} - \vec{\mu})}.$$

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}.$$

ONE CLASS SVM



IMPLEMENTATION

- For each detector:
 - For each subjects:
 - calculate user_score from testing_set;
 - calculate imposter_score from imposter_set;
 - calculate and draw ROC using user_score and imposter_score
 - evaluate using average AUCs
-

MODELS-TWO PHASES

- During training phase, the detector method calculates the mean_vector or fit the data using the TrainingSet
 - During scoring phase, giving the TestingSet or ImposterSet, the model calculates the Manhattan/Euclidean/Mahalanobis distance between each training vector and the mean vector, or the SVM predictions. Return the scoring vector.
-

RESULTS IN ROC

- After training, use the true positive rate and false positive rate to draw ROC
 - In this scenario, fnr represents the miss rate, and fpr represents false-alarm
 - There are many ways to analyze this...
-

IMPLEMENTATION

- For each detector:
 - For each subject:
 - Train the model with TrainingSet
 - Get user_score from TestingSet
 - Get imposter_score from ImposterSet
 - Calculate & Draw ROC
 - Calculate average AUC
-

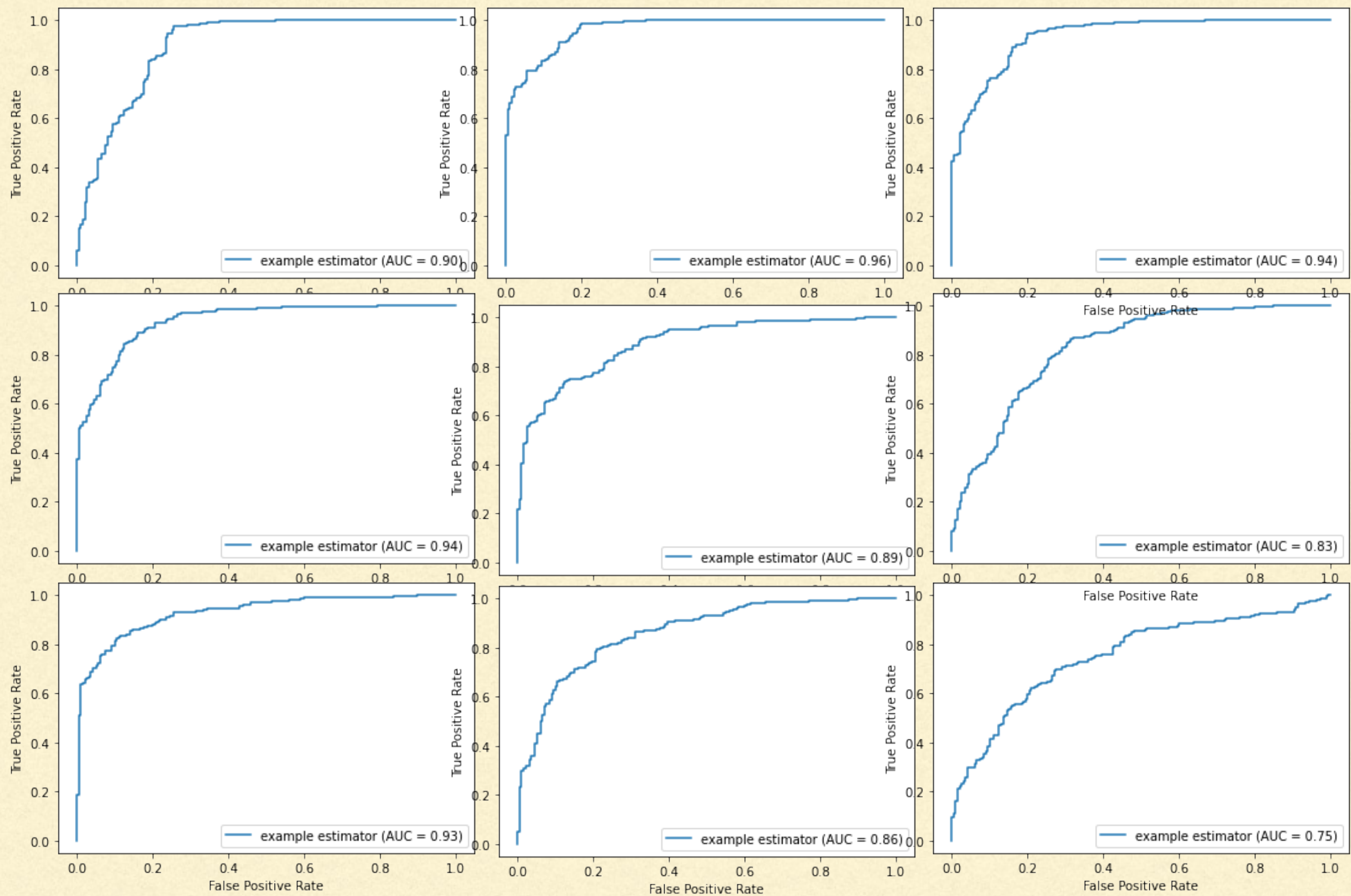


Figure 1. Euclidean Detector ROC from 9 random subjects

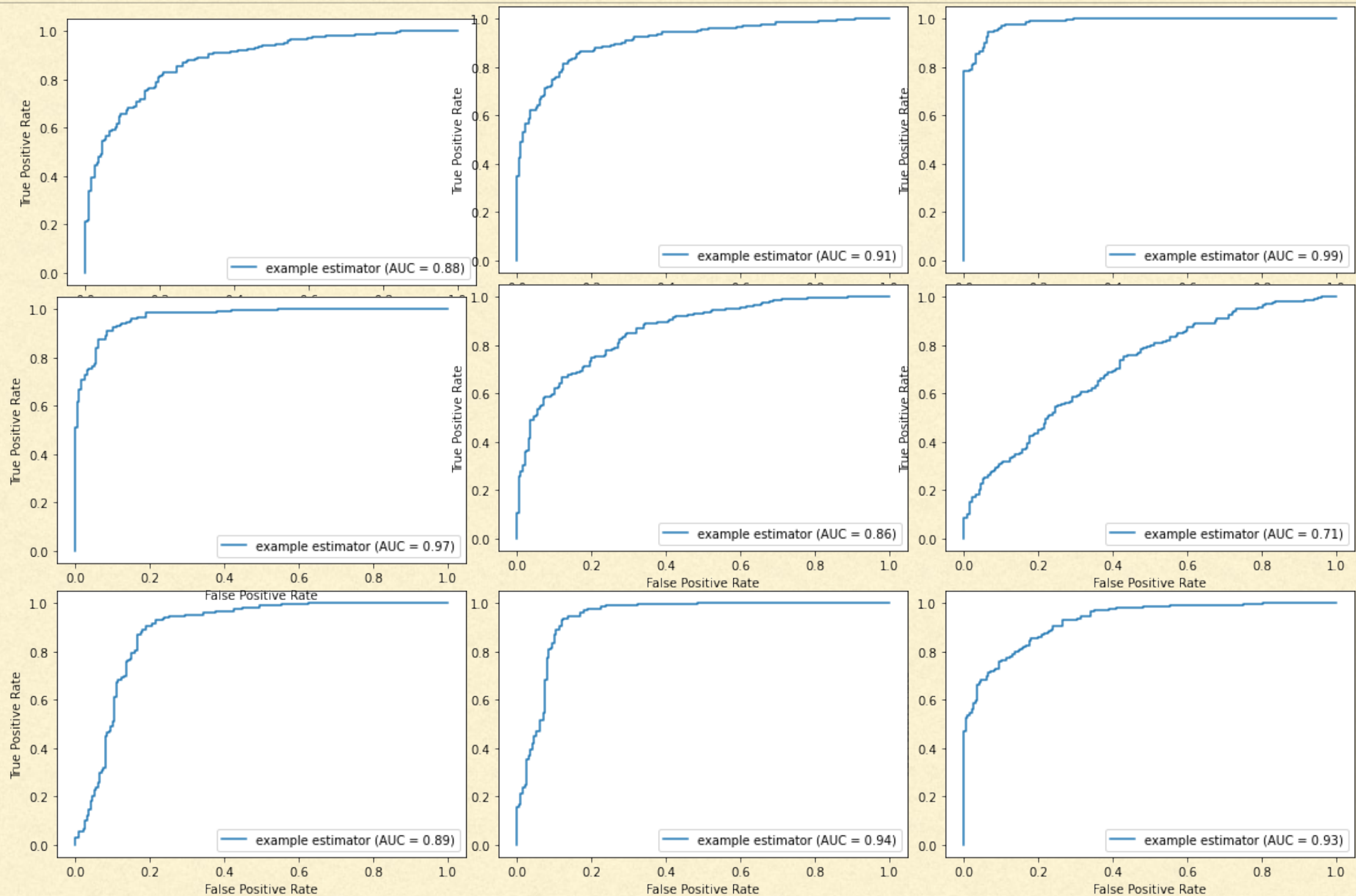


Figure 2. Manhattan Detector ROC from 9 random subjects

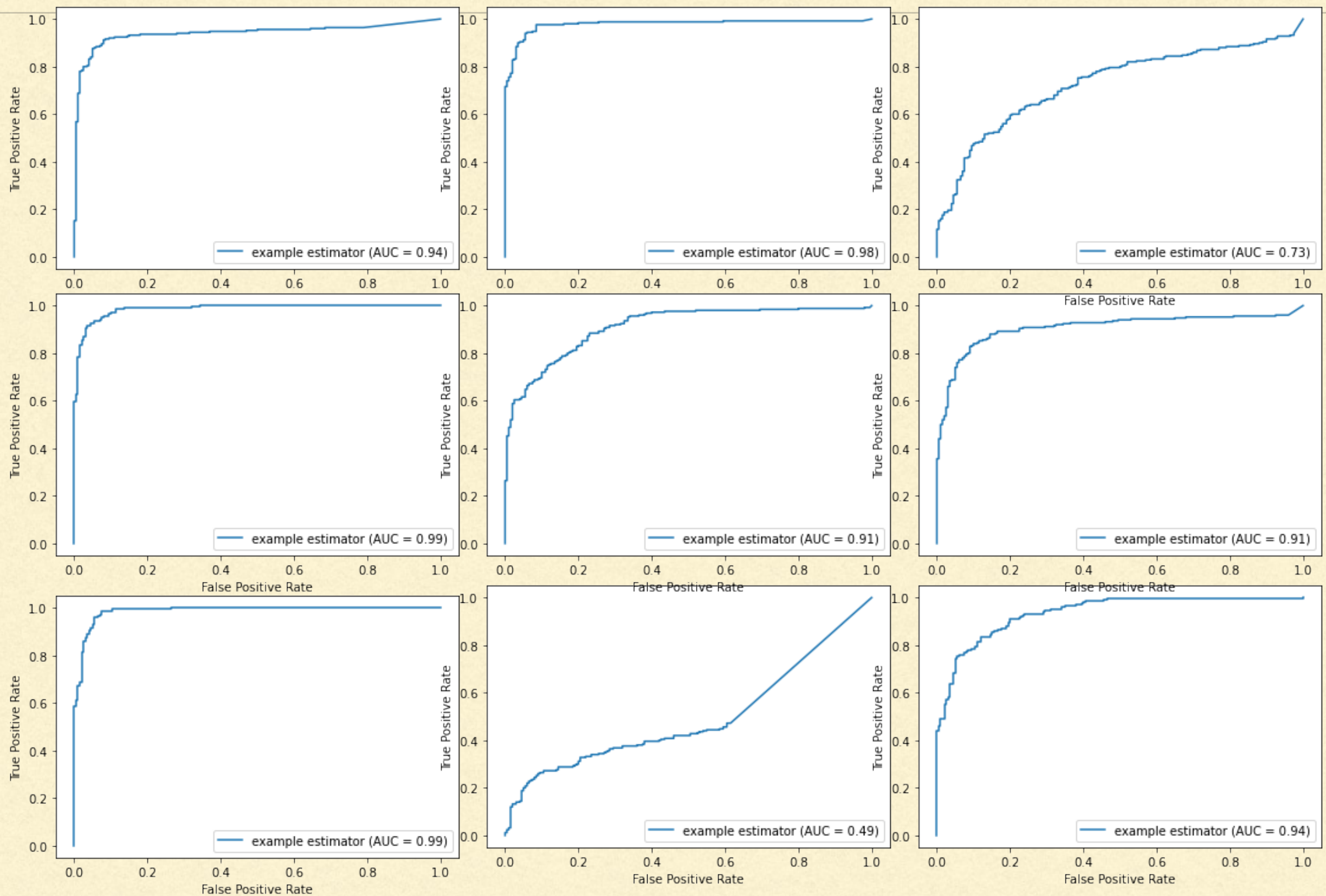


Figure 3. Mahalanobis Detector ROC from 9 random subjects

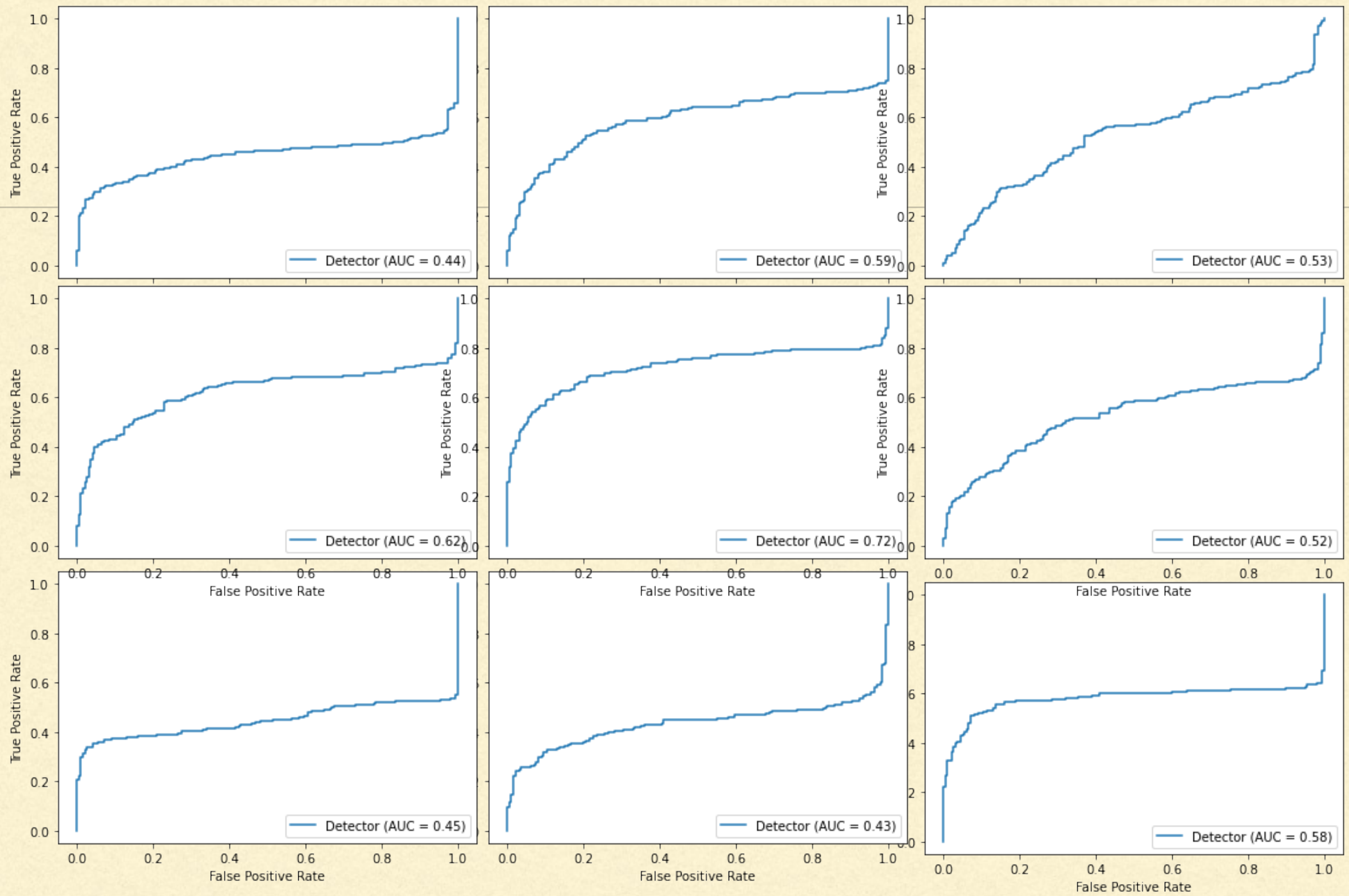


Figure 4. SVM detector with RBF kernel from 9 random subjects

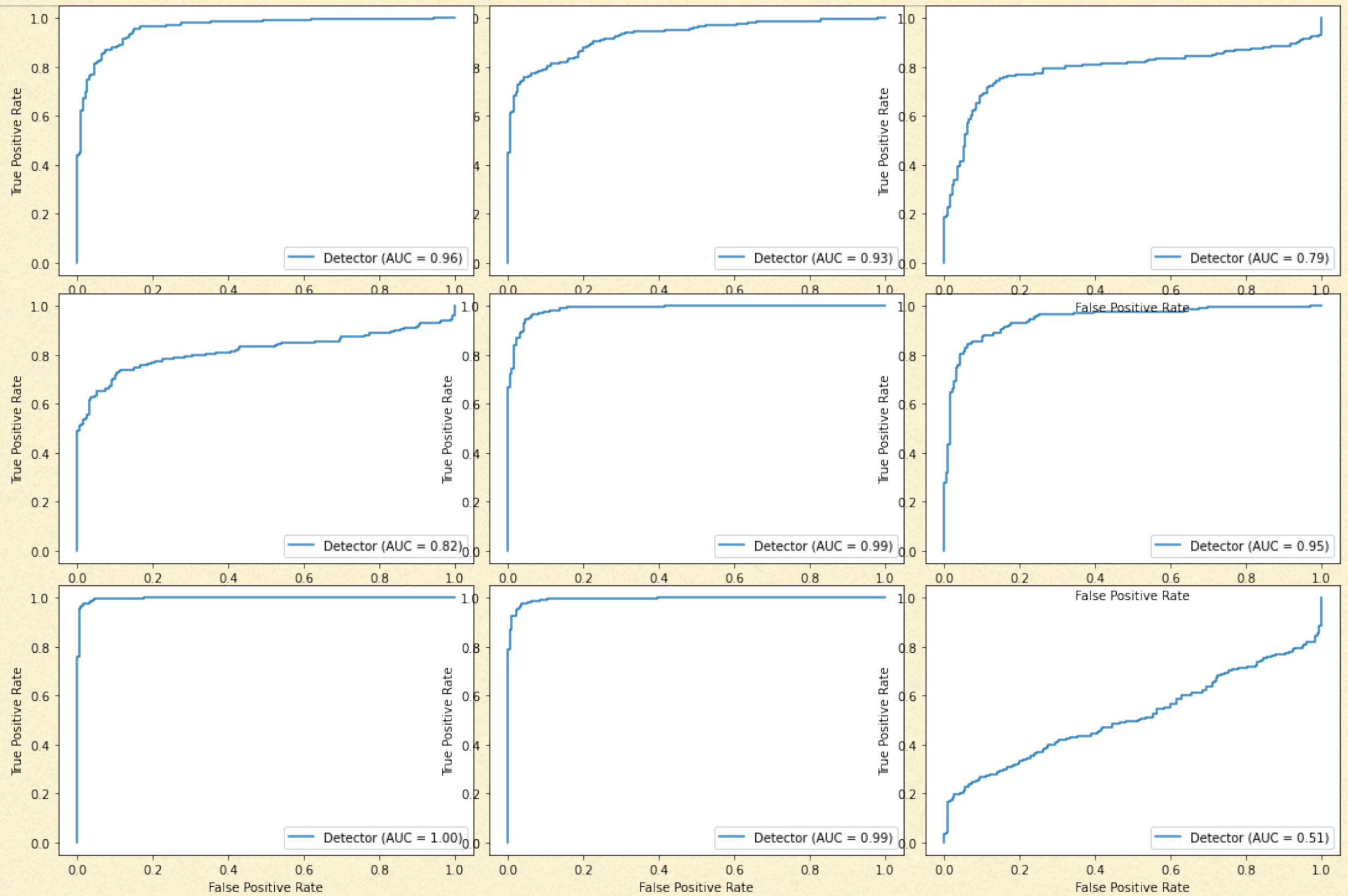


Figure 5. SVM with linear kernel ROC from 9 random subjects

COMPARING AVERAGE SCORES

Detector Name	Average AUC
Manhattan	0.902622745
Euclidean	0.879335294
Mahalanobis	0.849603872
SVM linear	0.858285098
SVM rbf	0.476107451
SVM polynomial	0.85914
SVM sigmoid	0.857907451

OBSERVATIONS

- Using statistic tools can guarantee a certain degree of robustness of the model.
 - Manhattan > Euclidean > Mahalanobis—>each features is more “proportional” to each other
 - SVMs as the top regression algorithm—> robust with relatively small dataset.
 - Determining Kernel functions by the “shape” of the dataset.
-

THANK YOU

Q&A

REFERENCES

- https://scikit-learn.org/stable/auto_examples/svm/plot_oneclass.html
 - https://en.wikipedia.org/wiki/Keystroke_dynamics
 - <https://www.cs.cmu.edu/~maxion/pubs/KillourhyMaxion09.pdf>
-