

# Open Thesis Topics

Simon.Kluettermann@cs.tu-dortmund.de

Is9 tu Dortmund

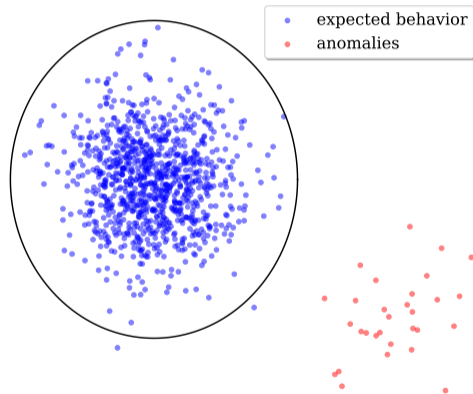
27. September 2022

*Simon Kluettermann*

- First: Find a topic and a supervisor
- Work one month on this, to make sure
  - you still like your topic
  - and you are sure you can handle the topic
- then short presentation in front of our chair (15min, relaxed)
  - get some feedback/suggestions
- afterwards register the thesis
  - (different for CS/DS students)
- Problem: We are not able to supervise more than 2 students at the same time (CS faculty rules)

- First: A short summary of each Topic
- Then time for questions/Talk with your supervisor about each topic that sounds interesting
- Your own topics are always welcome;)

- Im working on Anomaly Detection
- That means characterising an often very complex distributions, to find events that dont match the expected distribution



- kNN algorithm can also be used for AD
- if the k closest point is further away, a sample is considered more anomalous
- $r = \frac{k}{2N \cdot pdf}$
- Powerful method, as it can model the pdf directly

- The model (mostly) ignores every known sample except one
- So there are extensions
- $avg = \frac{1}{N} \sum_i knn_i(x)$
- $wavg = \frac{1}{N} \sum_i \frac{knn_i(x)}{i}$

Dataset	wavg	avg	1	3	5
<i>vertebral</i>	<b>0.4506</b>	<b>0.4506</b>	<b>0.4667</b>	<b>0.4667</b>	<b>0.45</b>
...					
<i>thyroid</i>	<b>0.9138</b>	<b>0.9151</b>	<b>0.8763</b>	<b>0.9086</b>	<b>0.914</b>
<i>Iris_setosa</i>	<b>0.9333</b>	<b>0.9333</b>	<b>0.9333</b>	<b>0.9</b>	<b>0.9</b>
<i>breastw</i>	<b>0.9361</b>	<b>0.9361</b>	<b>0.9211</b>	<b>0.9248</b>	<b>0.9286</b>
<i>wine</i>	<b>0.95</b>	<b>0.95</b>	<b>0.9</b>	<b>0.95</b>	<b>0.95</b>
<i>pendigits</i>	<b>0.9487</b>	<b>0.9487</b>	<b>0.9391</b>	<b>0.9295</b>	<b>0.9359</b>
<i>segment</i>	<b>0.9747</b>	<b>0.9747</b>	<b>0.9495</b>	<b>0.9545</b>	<b>0.9394</b>
<i>banknote – authentication</i>	<b>0.9777</b>	<b>0.9776</b>	<b>0.9408</b>	<b>0.943</b>	<b>0.9583</b>
<i>vowels</i>	<b>0.9998</b>	<b>0.9972</b>	<b>0.99</b>	<b>0.92</b>	<b>0.93</b>
<i>Ecoli</i>	<b>1.0</b>	<b>1.0</b>	<b>0.9</b>	<b>1.0</b>	<b>1.0</b>
<i>Average</i>	<b>0.7528</b>	<b>0.7520</b>	0.7325	0.7229	0.7157

# What to do?

- Evaluation as anomaly detector is complicated
  - Requires known anomalies
- $\Rightarrow$  So evaluate as density estimator
  - Does not require anomalies
  - Allows generating infinite amounts of training data



# What to do?

- Collect Extensions of the oc-knn algorithm
- Define some distance measure to a known pdf
- Generate random datapoints following the pdf
- Evaluate which algorithm finds the pdf the best

# Requirements

- Knowledge of python ( `sum([i for i in range(5) if i%2])` )
  - Ideally incl numpy
- Basic university level Math (you could argue that  $r_k \propto \frac{k}{pdf}$ )
- Ideally some experience working on a ssh server
- $\Rightarrow$  Good as a Bachelor Thesis
- For a Master Thesis, I would extend this a bit (Could you also find  $k$ ?)

- Deep Learning Method, in which the output is normalised
- $\int f(x)dx = 1 \forall f(x)$
- Can be used to estimate probability density functions
- $\Rightarrow$  Thus useful for AD
- $\int f(h(x)) \left\| \frac{\delta x}{\delta h} \right\| dh = 1 \forall h(x)$

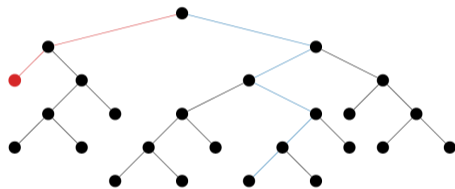
- How to apply this to graphs?
- One Paper (Liu 2019) uses two NN:
  - Autoencoder graph $\Rightarrow$ vector
  - NF on vector data
  - which is fine, but also not really graph specific
  - No interaction between encoding and transformation

- So why not do this directly?
- $\Rightarrow$  Requires differentiating a graph
- Why not use only one Network?
- Graph  $\Rightarrow$  Vector  $\Rightarrow$  pdf
- $\Rightarrow$  Finds trivial solution, as  $\langle pdf \rangle \propto \frac{1}{\sigma_{Vector}}$
- So regularise the standard deviation of the vector space!
  - Interplay between encoding and NF
  - Could also be useful for highdim data

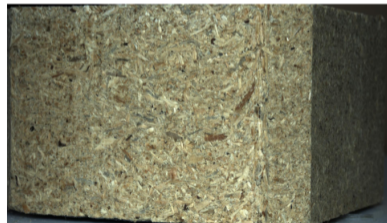
# Requirements

- Proficient in python ( [i for i in range(1,N) if not [j for j in range(2,i) if not i%j]] )
  - Ideally incl numpy, tensorflow, keras
- Some deep learning experience
- University level math (google Cholesky Decomposition. Why is this useful for NF?)
- Ideally some experience working on a ssh server
- A bit more challenging  $\Rightarrow$  Better as a Master thesis
- (Still we would start very slowly of course)

- Isolation Forest: Different Anomaly Detection Algorithm
- Problem: Isolation Forests don't work on categorical data
- $\Rightarrow$  Extend them to categorical data

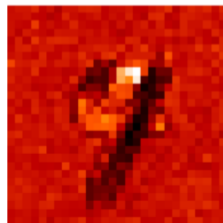
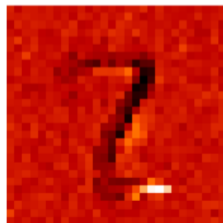


- Reidentification: Find known objects in new images
- Task: Find if two images of pallet blocks are of the same pallet block
- Use AD to represent the pallet blocks

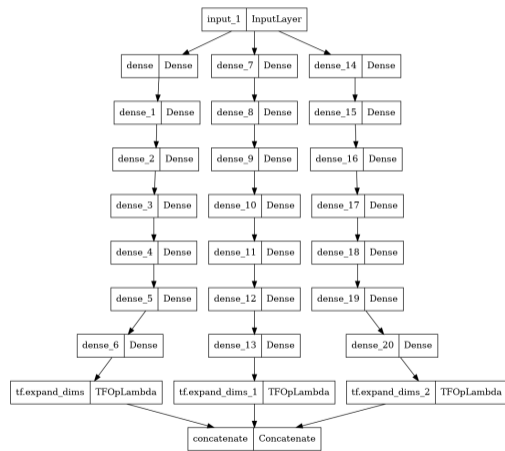




- Ensemble: Combination of multiple models
- Task: Explain the prediction of a model using the ensemble structure



- Task: Explore a new kind of ensemble
- Instead of many uncorrelated models, let the models interact during training



Questions?