

Class Struggle

Attempts to classify domain names using classifiers trained on raw datasets have been discussed as much as using regular expressions to parse HTML, so let's add a few more words.

Where is my data?

I want to find domains similar to those on the CERT Polska warning list¹. Obtaining a representative set of legitimate domains, e.g. from The Majestic Million², is unrealistic. As a result, if I wanted to use any classifier, my training dataset:

- is radically unbalanced - phishing domains (the ones I'm interested in) are extremely rare compared to legitimate ones,
- since the legitimate class would be artificially constructed, my classifier might label as phishing domains it has never seen before in training.

One might think that a classifier with an overrepresentation of the expected domain type will correctly label them and recognize the rest as another class. Unfortunately, classifiers that learn to separate feature spaces without counterexamples may push the decision boundary arbitrarily far to minimize loss. This is like a diagnostic test labeling everyone as sick - it achieves 100% sensitivity, finding all the truly sick.

You can observe this even in a simple example with a FF MLP network reproducing the XOR problem³.

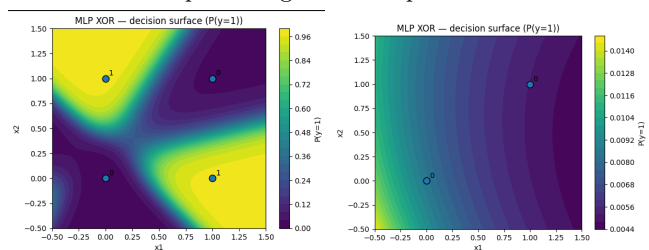


Figure 1: Complete XOR (left) and the decision boundary with an incomplete training set.

Autoencoders to the rescue

An autoencoder⁴ is a neural network that learns to reconstruct its training input. The encoder compresses data into a hidden representation z , and the decoder tries to reconstruct the original input from it. The model is trained without labels by minimizing the reconstruction error⁵.

Thus, it doesn't need counterexamples; it doesn't learn decisions but representations.

During inference⁶, we pass new data through the network and measure the reconstruction error - a large value means that the sample is different from the training data. The inner

layer z serves as an embedding - a compact representation of the data structure.

Yet EDA - on limitations

Phishing domains are often created to resemble legitimate ones. Training on the entire corpus would result in an autoencoder that reconstructs legitimate domains.

I must restrain my ambitions. To do this, I perform EDA⁷ through clustering⁸.

From the dataset, I select the cluster(s) of interest - in my case, domains with certain characteristics⁹.

Preprocessing FTW

For the autoencoder to understand the data, it must be vectorized. Counting lengths, dashes, etc., loses semantic information, so I used TF-IDF statistics¹⁰ on n-grams¹¹.

The TF-IDF vectorizer doesn't preserve positional information or relationships between n-grams. We can enforce minimal positional awareness by positional n-grams through adding start and end markers¹².

Metrics and GOTO EDA

How to choose hyperparameters? I'm not a scientist, so I do it empirically. For network layer dimensions, remember that an autoencoder should compress the representation - a bottleneck is necessary¹³. Typically, if your validation loss increases while training loss decreases, you've trained too long or your model is too large¹⁴.

A nice feature of the autoencoder is that we can apply PCA¹⁵ to its embeddings and try to infer what the autoencoder has learned.

You can find the repository with runnable code here:

https://github.com/artur-augustyniak/class_struggle

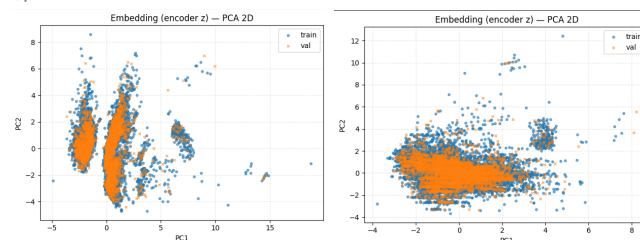


Figure 2: A small AE forced to learn classes (left) and a larger one capable of generalization. In each case, we are interested in the overlap between validation and training embeddings.

¹ <https://cert.pl/en/warning-list/>

² <https://majestic.com/reports/majestic-million>

³ https://github.com/artur-augustyniak/class_struggle/blob/main/notebooks/xor_nn_decision_boundary.ipynb

⁴ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L608

⁵ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L741

⁶ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L916

⁷ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/1_eda.ipynb

⁸ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.HDBSCAN.html>

⁹ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/2_eda_based_data_selection.ipynb

¹⁰ https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

¹¹ <https://en.wikipedia.org/wiki/N-gram>

¹² https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L498

¹³ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L56

¹⁴ https://github.com/artur-augustyniak/class_struggle/blob/2287420fa8b5882b3fc603b6660bb50810b764f1/notebooks/3_autoencoder_train_eval_inference.ipynb?short_path=827cbee#L707

¹⁵ <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>