

## 184.702 Machine Learning – Exercise 3.1 Advanced topics in security / privacy of Machine Learning

*Note: please first read the generic introduction document for the 3<sup>rd</sup> exercise!*

### Topic 3.1: Advanced topics in security / privacy of Machine Learning

*You need to choose only one topic, don't be scared of the long document - it just means a lot of choices :-)*

The aim of this type of exercise is to work in detail with one of the emerging topics in **security & privacy** of machine learning. Some of the topics are rather experimental as they are novel or based on very recent results, and the exact outcome can not always be predicted. Thus, the main aim is not to achieve perfect results, but to gain an interesting learning experience. We will grade your work based on your effort, and on the originality of your approach.

### Table of Contents:

|  |           |
|--|-----------|
| <b>Adversarial Machine Learning</b>  | <b>3</b>  |
| Topic 3.1.1.1: Data exfiltration using ML Models   | 4         |
| Topic 3.1.1.3: Attribute Inference (disclosure) from ML Models   | 6         |
| Topic 3.1.1.4: Model Stealing / Extraction   | 8         |
| A. Type 1: Attacks, described in Tramèr et al. Stealing Machine Learning Models via Prediction APIs. USENIX Security 2016. | 8         |
| B. Type 2: Substitute training attacks   | 8         |
| 3. Type 3: Cryptanalytic attack on ReLU NN   | 9         |
| Topic 3.1.1.5: Detection of Model Extraction   | 10        |
| Topic 3.1.1.6: Watermarking ML/DL models   | 11        |
| Topic 3.1.1.6.A:   | 11        |
| Topic 3.1.1.6.B:   | 12        |
| Topic 3.1.1.7: Model Inversion Attack  | 13        |
| Topic 3.1.1.8: Membership Inference  | 15        |
| Topic 3.1.1.9: Backdoor/poisoning Attacks & defence  | 17        |
| Topic 3.1.1.10: Model Evasion Attack   | 18        |
| <b>Privacy-Preserving Data Publishing</b>  | <b>20</b> |
| Topic 3.1.3.3.: Generation and evaluation of sequential synthetic datasets   | 20        |
| Topic 3.1.3.4.: k-anonymity and utility of anonymized datasets   | 21        |

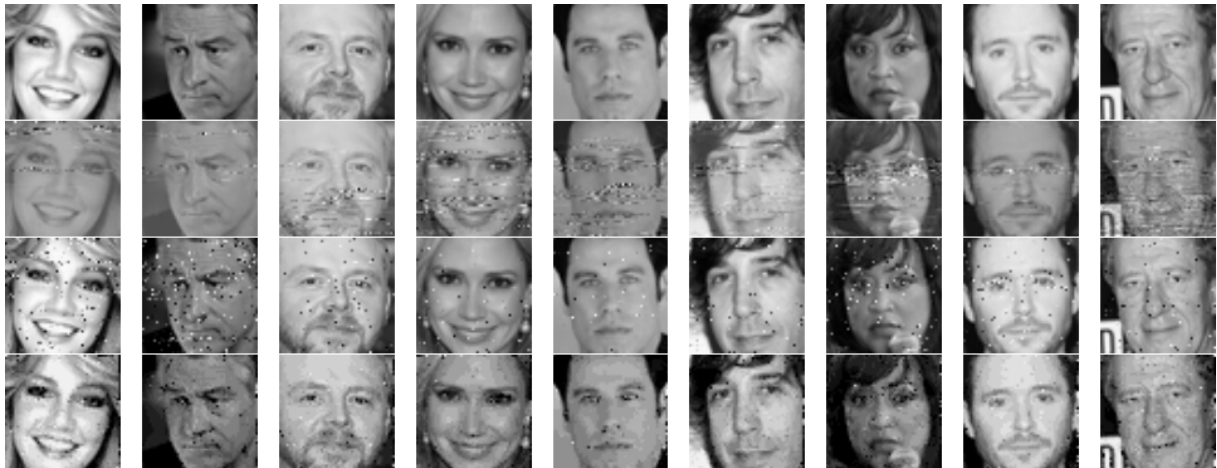
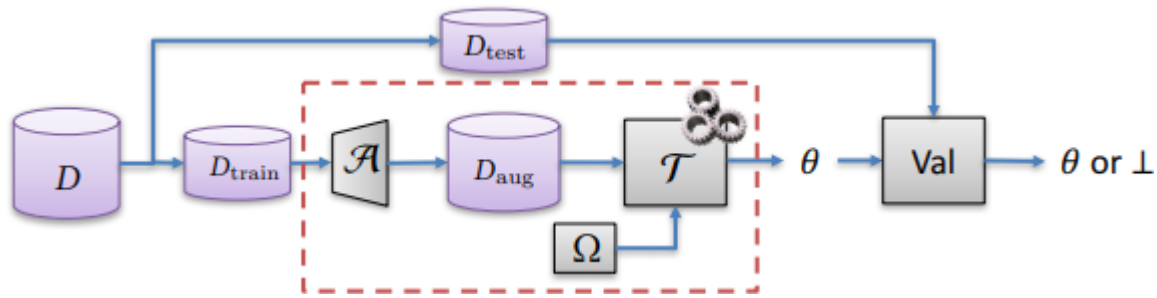
|  |           |
|--|-----------|
| Topic 3.1.3.5.: k-anonymity: generalization vs microaggregation                | 23        |
| Topic 3.1.3.7.: Differential Privacy   | 25        |
| <b>Secure / Federated / Distributed Computation</b>                            | <b>27</b> |
| Topic 3.1.2.1: Federated (Distributed) Learning for Process Mining             | 27        |
| Topic 3.1.2.3: Federated Learning on text data (or other sequential data)      | 28        |
| Topic 3.1.2.4: Federated Learning on tabular data with different ML algorithms | 28        |
| Topic 3.1.2.5: Federated Learning, different aggregation methods               | 28        |

# Adversarial Machine Learning

Some datasets that are mentioned in the rest of the topic description

- Yale Face Database <http://vision.ucsd.edu/content/yale-face-database>
- Labeled Faces in the Wild, <http://vis-www.cs.umass.edu/lfw/>, using the task for face recognition. See also <https://scikit-learn.org/0.19/datasets/#the-labeled-faces-in-the-wild-face-recognition-dataset>
- PubFig dataset, <http://www.cs.columbia.edu/CAVE/databases/pubfig>
- PubFig83: <http://vision.seas.harvard.edu/pubfig83/>
- AT&T / Olivetti Faces (<https://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>, <https://scikit-learn.org/stable/datasets/index.html#olivetti-faces-dataset>)
- MNIST: <http://yann.lecun.com/exdb/mnist/>, <https://scikit-learn.org/0.19/datasets/mldata.html>
- CIFAR (10, 100): <https://www.cs.toronto.edu/~kriz/cifar.html>, <https://github.com/ivanov/scikits.data/blob/master/datasets/cifar10.py>
- FashionMNIST: <https://github.com/zalandoresearch/fashion-mnist>
- The German Traffic Sign Recognition Benchmark (GTSRB), <http://benchmark.ini.rub.de>
- ImageNet: <http://image-net.org/download>

### Topic 3.1.1.1: Data exfiltration using ML Models



Data exfiltration (also called data extrusion or data exportation) is an unauthorized data transfer from a computer system, and a form of data theft.

In the setting of machine learning, an adversary tries to utilise a machine learning model to leak information about the training data to the end-user of the model. This can be interesting for an attacker when the model gets trained on private training data, and there is no other means to leak/exfiltrate information about the data to the outside (the attacker) than via the model. Hiding the information in the model is a form of steganography, where a message is hidden within another message.

In this task, you shall create experiments either in a white-box or black-box settings similar as described in the paper by Song et al: *Machine Learning Models that Remember Too Much*. CCS 2017 ([https://www.cs.cornell.edu/~shmat/shmat\\_ccs17.pdf](https://www.cs.cornell.edu/~shmat/shmat_ccs17.pdf))

There are two options for this topic:

- A. Three examples of the attacks are implemented against the CIFAR dataset (i.e. image data) at <https://github.com/andreasiposova/ml-model-remember>  
Your task is to perform the Correlated Value Encoding Attack (white-box), Sign Encoding Attack (white-box) and Capacity Abuse Attack (black-box) on the 'Adult' data set <https://archive.ics.uci.edu/ml/datasets/adult>, i.e. a tabular data set.
- B. Alternatively, you can implement the 'Least Significant Bit Encoding' attack, on the tabular 'Adult' dataset or on the 'Facescrub' image dataset also used in the original paper.

Link to the Facescrub data set: <http://vintage.winkler.site/faceScrub.zip>

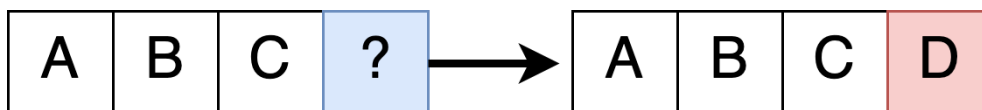
The password for decrypting the zip file is: ICIP'14

### Topic 3.1.1.3: Attribute Inference (disclosure) from ML Models

*(note: this topic is rather experimental, and you will, as for all the other topics, graded on your effort, not necessarily on a (positive) outcome, which might not be achievable).*

One type of information disclosure is called **attribute disclosure**: it consists of **inferring the value of an attribute** (e.g., ethnicity) **that was hidden** (i.e., not directly shared by the user).

Attribute disclosure is independent of identifying a specific record in a database - inferring additional information about a user from data shared by **other users**. This can be sensitive information about a user such as ethnicity or political affiliation, inferred by mining available data.



While attribute disclosure until now mostly considered the setting where a released **data set** is the base for inference, similar attacks might be possible from a **released machine learning model**, which invariably encodes some information about the training data.

Your task in this exercise is to test if this is possible for various machine learning models. Given an incomplete sample  $X$  and a trained model  $f(X)$ , is it possible to infer the unknown values of  $X$ ? You can first assume that  $X$  (with its full information) was also used in actually training the model, but further testing whether it would be also possible for an  $X$  not used in the training. One simple approach would be testing multiple values to “impute” the missing value, and observe the response from the model to that.

A similar approach as been tested to some extent in

[https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/fredrikson\\_mattthew](https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/fredrikson_mattthew) (even they refer to it there as “model inversion”, it can be seen as a form of attribute inference, potentially with multiple attributes considered to be inferred): “... the algorithm simply completes the target feature vector with each of the possible values for  $x_1$ , and then computes a weighted probability estimate that this is the correct value.”

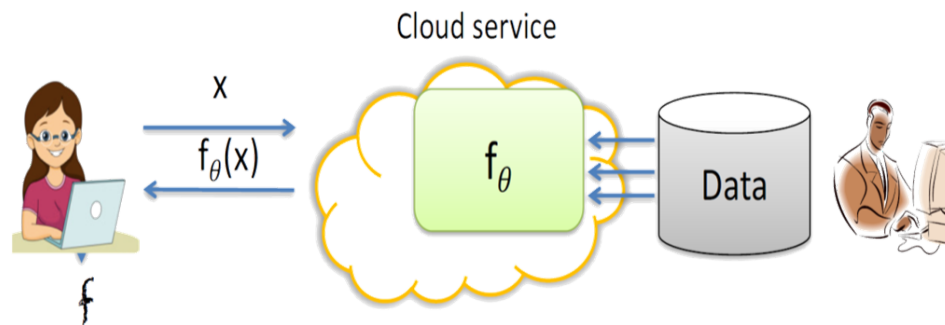
```
adversary  $\mathcal{A}^f(\text{err}, \mathbf{p}_i, \mathbf{x}_2, \dots, \mathbf{x}_t, y)$ :  
1: for each possible value  $v$  of  $\mathbf{x}_1$  do  
2:    $\mathbf{x}' = (v, \mathbf{x}_2, \dots, \mathbf{x}_t)$   
3:    $\mathbf{r}_v \leftarrow \text{err}(y, f(\mathbf{x}')) \cdot \prod_i \mathbf{p}_i(\mathbf{x}_i)$   
4: Return  $\arg \max_v \mathbf{r}_v$ 
```

**Figure 2: Generic inversion attack for nominal target features.**

Try such an approach for a classification model; choose a relatively simple dataset, preferably with categorical data to limit the values to be tried, and models that return information on the confidence/likelihood of the predicted class (as you need this information to select the values that lead to the most determined outcome). Start with simple settings, i.e. where you want to infer only one attribute at the time (input attributes, not the target attribute, as you need that for comparison)

Using models that generally overfit to the training data seems more promising. As such, a neural network, or also a decision tree, might be good candidates.

### Topic 3.1.1.4: Model Stealing / Extraction



Model “stealing” means to obtain a trained model from a source that generally only provides a “prediction API”, i.e. a means to obtain a prediction from the model when providing a specific input, but does NOT disclose the actual model (i.e. “ML-as-a-service”). This might be because the model may be deemed confidential due to their sensitive training data, commercial value, or use in security applications. A user might train models on potentially sensitive data, and then charge others for access on a pay-per-query basis.

In model stealing/model extraction, an adversary with black-box access (but no knowledge of the model parameters or training data), aims to duplicate the functionality of the model.

In this task, you shall experiment with these model stealing attacks for two different types of classifiers (can be from the paper) on three different datasets (two shall be different than in the paper, one to be used as validation that you achieve similar results). You can either use a model available locally (as black-box, i.e. just for querying & confidence values), or one of the mentioned online services. While the model is “extracted”, record how the accuracy/fidelity of the duplicated model changes, and how many queries (and how much time) it takes to train the model; plot these, in similar fashion as in the paper, for each dataset/classifier.

Regarding implementation, you can use one of the types of model stealing attacks listed below:

#### A. Type 1: Attacks, described in *Tramèr et al. Stealing Machine Learning Models via Prediction APIs. USENIX Security 2016.*

- ([https://www.usenix.org/system/files/conference/usenixsecurity16/sec16\\_paper\\_tramer.pdf](https://www.usenix.org/system/files/conference/usenixsecurity16/sec16_paper_tramer.pdf))
  - Implementation: <https://github.com/ftramer/Steal-ML>
- Pick one of the following types
  - Equation-solving attack for (M) LR or MLP
  - Path-finding attack for DT
  - Lowd-Meek for Linear Binary Models
  - Retraining attacks for non-linear SVM or NN

#### B. Type 2: Substitute training attacks

- Knockoff nets for CNN proposed by Orekondy et al. in Knockoff Nets: Stealing Functionality of Black-Box Models (<https://arxiv.org/pdf/1812.02766.pdf>)
  - Implementation: <https://github.com/tribhuvanesh/knockoffnets>



- Copycat networks for CNN proposed by Correia-Silva et al. paper Copycat CNN: Stealing Knowledge by Persuading Confession with Random Non-Labeled Data (<https://arxiv.org/pdf/1806.05476.pdf>)
  - Implementation: [https://github.com/jeiks/Stealing\\_DL\\_Models](https://github.com/jeiks/Stealing_DL_Models)
- Active Thief for CNN or RNN proposed by Pal et al. in ACTIVETHIEF: Model Extraction Using Active Learning and Unannotated Public Data (<https://ojs.aaai.org//index.php/AAAI/article/view/5432>)
  - Implementation: <https://bitbucket.org/iiscseal/activethief/src/master/>

### 3. Type 3: Cryptanalytic attack on ReLU NN

- Cryptanalytic Extraction of Neural Network Models by Carlini et al. (<https://arxiv.org/pdf/2003.04884.pdf>)
  - Implementation: <https://github.com/google-research/cryptanalytic-model-extraction>

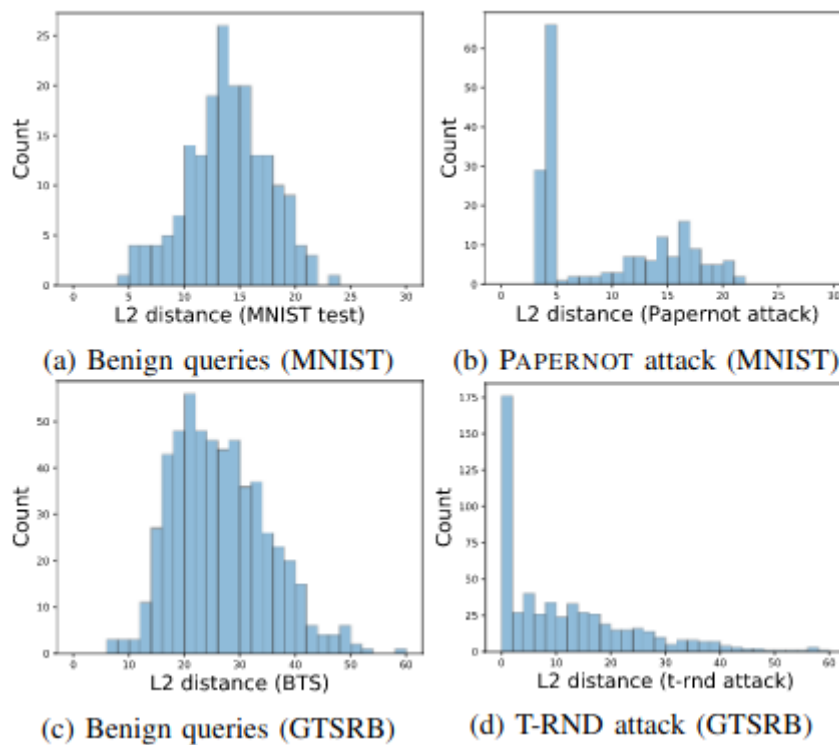
You can also use the implementation(s) available by the IBM Adversarial Robustness Toolkit:

<https://github.com/Trusted-AI/adversarial-robustness-toolbox/wiki/ART-Attacks#3-extraction-attacks>

### Topic 3.1.1.5: Detection of Model Extraction

The idea of this task is to analyze the difference between distributions of data from the initial problem domain, and data generated by the adversary, e.g. synthetic/adversarial samples. The following steps should be done:

- Pick a dataset and explore it's distribution. We propose to pick one of the following datasets: MNIST, FashionMNIST, PubFig83, GTSRB, CIFAR-10
- Use different techniques, such as
  - Adversarial image generation: FGSM, IGS, FGV, Deepfool, C&W, or JSMA (as e.g. in <https://arxiv.org/pdf/1809.04913.pdf>)
  - Data composition (<https://elidavid.com/pubs/stealing-knowledge.pdf>) or
  - Random noise inputs (<https://arxiv.org/pdf/1912.08987.pdf>)
- For each technique, explore the distribution of generated samples and compare it to the original data distribution. Also data from a non-problem-domain dataset can be analyzed
- Implement a detection (or use (or modify) an existing one such as described in <https://github.com/SSGAalto/prada-protecting-against-dnn-model-stealing-attacks>), to evaluate which adversarial/synthetic samples can be recognized by distribution analysis



### Topic 3.1.1.6: Watermarking ML/DL models

Sharing trained models has brought many benefits to development of ML and DL systems. However, training models is usually a task that requires vast resources, from data to time and computing power. **Watermarking** can help the owners of such models mark their property in order to trace unauthorised usage or redistribution.

We have collected and unified a number of watermarking techniques at <https://sbaresearch.github.io/model-watermarking/#/>, together with

- Non-watermarked pre-trained models for two datasets and several architectures (CIFAR-10 and MNIST)
- Attacks on watermarking

You can choose between two options in this task

#### Topic 3.1.1.6.A:

Work with one additional **watermarking technique not yet in that repository**, either a white-box or a black-box watermarking technique, e.g. from the following techniques:

- White-box Approaches
  - Uchida et al.: Embedding Watermarks into Deep Neural Networks (<https://dl.acm.org/doi/pdf/10.1145/3078971.3078974>)
    - Implementation: <https://github.com/yu4u/dnn-watermark>
  - Rouhani et al.: DeepSigns: An End-to-End Watermarking Framework for Ownership Protection of Deep Neural Networks (<https://dl.acm.org/doi/10.1145/3297858.3304051>)
    - Implementation: <https://github.com/Bitadr/DeepSigns>
  - Feng et al.: Watermarking Neural Network with Compensation Mechanism ([http://link.springer.com/10.1007/978-3-030-55393-7\\_33](http://link.springer.com/10.1007/978-3-030-55393-7_33))
  - Wang et al.: RIGA: Covert and Robust White-Box Watermarking of Deep Neural Networks (<http://arxiv.org/abs/1910.14268>)
- Black-Box Approaches
  - Huili Chen, Bita Darvish Rouhani, Farinaz Koushanfar: BlackMarks: Blackbox Multibit Watermarking for Deep Neural Networks (<https://arxiv.org/abs/1904.00344>)

Either: implement the chosen technique or utilize and adapt an open source implementation.

In both cases: in order to demonstrate the success of the chosen watermarking technique, evaluate experimentally the following three requirements. You can use the experiments from the paper of the chosen technique as a guidance.

- **effectiveness**: after embedding the watermark, it also can be successfully extracted and verified, thus the model owner can claim ownership in case of a threat event. Analyze if the chosen technique satisfies this requirement.
- **fidelity**: the watermark embedding has only minimal influence on the model's performance. Compare the model's accuracy before and after embedding the watermark.

- **robustness:** the watermark should be robust against various types of attacks, e.g. fine-tuning, pruning, transfer learning. In your experiments, perform **one** attack and analyse the robustness of the watermarking technique with respect to this specific attack.

Choose one dataset (CIFAR-10 or MNIST) from the template examples provided in the github repository, and perform the above steps on at least two of the corresponding models; you can use the non-watermarked model as a basis for embedding, and the fine-tuning or pruning attack method.

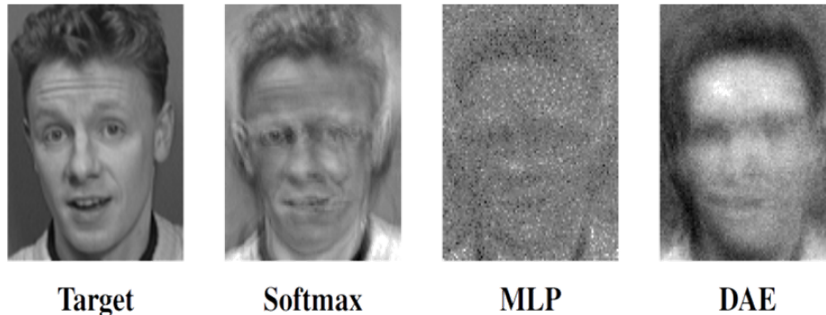
***If you chose to implement a method for which no code is provided, you just need to show it only on one architecture and one dataset.***

In TUWEL in the topic registration, specify the architecture(s) and dataset(s) used in the experiments. In the report describe your methodology in detail, explain your findings and provide the source code of your implementation.

#### **Topic 3.1.1.6.B:**

Instead of implementing a new method, pick one of the already provided implementations, but pick a different **dataset**, e.g. Fashion-MNIST, Tiny Image Net, ...

### Topic 3.1.1.7: Model Inversion Attack



Given access to a model, an **inversion** attack tries to extract the training data used for obtaining that model. A prominent example is for a face recognition system, where the model tries to identify to which of the known users a newly taken image (e.g. from an access control camera) corresponds to. This means that each class in the dataset corresponds to one specific individual, and the training data for each class is generally a number of images of that individual (taken with potentially different angles, zoom, light, background, and appearance such as hair style, ...).

For the attack, in most cases, this means picking a specific class, and then trying to **generate** input patterns that return the highest confidence of being classified as that class. This normally means access to the model (also in white-box, i.e. being able to read the model parameters; potentially building on top of a model stealing/extraction attack, but you can assume to have white-box access).

You can choose two different versions of this task

- A. Extending an existing implementation provided by us, based on the paper by Fredrikson et al: “Model inversion attacks that exploit confidence information and basic countermeasures” (<http://www.cs.cmu.edu/~mfredrik/papers/fjr2015ccs.pdf>), to other datasets. We will provide you with the code, and working examples on a number of datasets. Use two of the following datasets:

- Yale Faces
- A subset of Labeled Faces in the Wild
- A subset from the PubFig dataset
- A subset from PubFig83
- Any other data set for face recognition (except the original AT&T / Olivetti Faces)
  - Your subset should contain around 20-30 of the most-frequent people, as the task is very difficult for a larger number of classes.

and the following models:

- A softmax layer
- A simple, fully-connected network with 1 hidden layer
- A very simple CNN of your choice

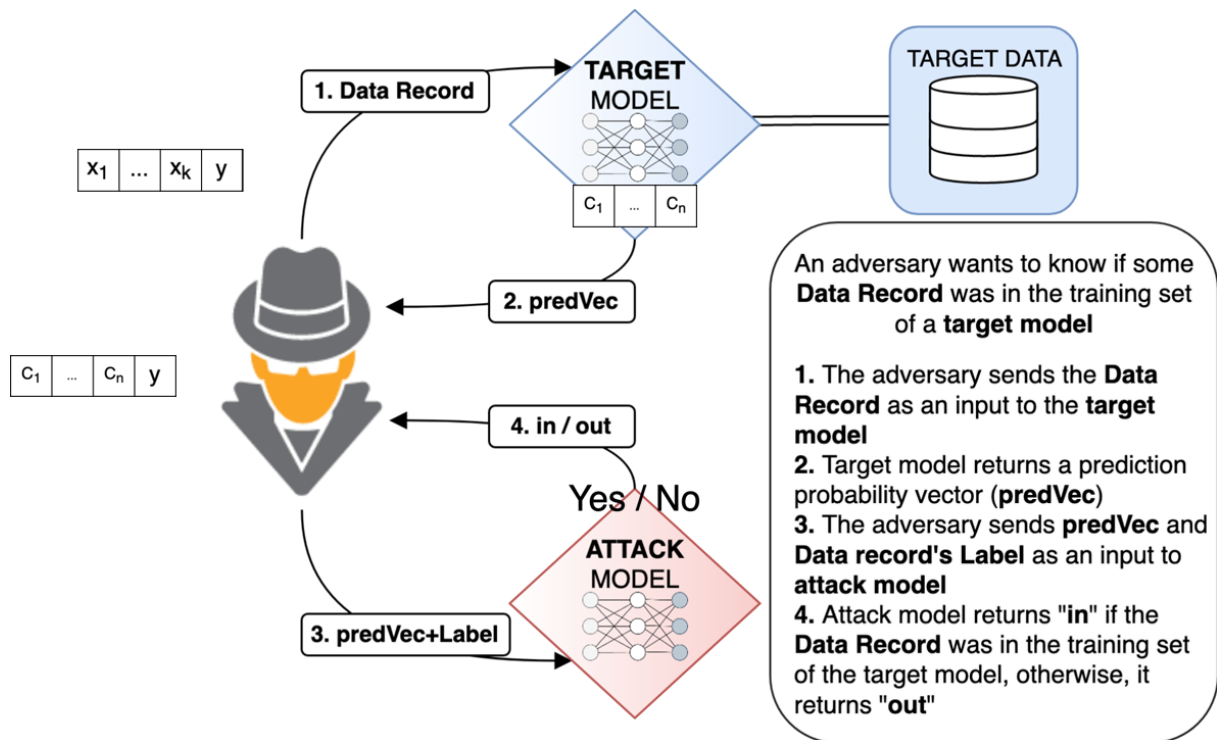
Try also on-purpose overfit / underfit models, to see how that affects the results of the inversion.

- B. Try a method from a more recent paper, e.g. <https://arxiv.org/abs/1911.07135>, which combines generative adversarial examples to improve the inversion process, and try it on a different dataset. A first step is to recreate the original experiments. Then, extend the experiment to one of the additional datasets listed above for Variant A

In your evaluation, **compare how well the inversion works** for each of the individuals, and also evaluate what is the **“cost” of generating** the images – i.e. how many iterations of generating and asking the model are required, and how is the efficiency of this (in runtime). You can also evaluate the quality of the generated images by computing a numerical similarity to the target image (e.g. the most representative of the training images per class).

### Topic 3.1.1.8: Membership Inference

Membership inference tries to infer information from a learned ML model on the data samples that were used for training the model. In terms of released information, it is one of the weakest forms of attack on the confidentiality of the data samples, with the goal of revealing whether a sample was used for training a model, or not. This can be disclosive on an individual, e.g. if one can then infer that the individual might suffer from a certain disease.



In general, an attacker obtains a data set similar (or even identical) to the one used to train the target model, and trains a number of **shadow** models on that data; these should behave very similar to the original target model, and are used to train the **attack model** to distinguish between instance from that were inside or outside the training data.

The above-mentioned shadow models are used to have training data for the attack model - this needs to have a ground truth on whether a sample was in the training data (of the shadow model that shall simulate the target (real) model), or not. Another approach could be to perform an unsupervised attack. You can either work with

- The attacks specified in the paper "ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models."  
<https://arxiv.org/abs/1806.01246>, or try
- a grey-box setting: given a number of responses (including the full prediction vector / activations from the last layer) from the target model, try to cluster these into two clusters, assuming that this corresponds to "in" and "out" clusters. Try multiple state-of-the-art clustering algorithms (don't implement one on your own), and evaluate your clustering with actual ground-truth data (i.e. you know which data has been used to train the models, or not)

- Main paper: Membership Inference Attacks Against Machine Learning Models,  
[https://www.cs.cornell.edu/~shmat/shmat\\_oak17.pdf](https://www.cs.cornell.edu/~shmat/shmat_oak17.pdf)
- Other useful papers
  - Towards Demystifying Membership Inference Attacks,  
<https://arxiv.org/abs/1807.09173>



### Topic 3.1.1.9: Backdoor/poisoning Attacks & defence



Poisoning/backdoor attacks modify the training data to **embed a specific pattern** in some images (e.g. a yellow block on a STOP sign), and falsely label that image as a different class (e.g. as a speed limit 50 sign) to make the classifier learn this pattern for wrong predictions. These attacks generally work because a certain number of neurons in the network can be dedicated to learn these patterns, often because the number of neurons is larger than what is required to actually represent the pattern in the “clean” training data.

We will provide you with two datasets with manually prepared backdoor images, on the German Traffic Sign Recognition Dataset, and Yale Faces, and already trained models containing the backdoor trigger.

You shall then re-create experiments from one of these proposed defences:

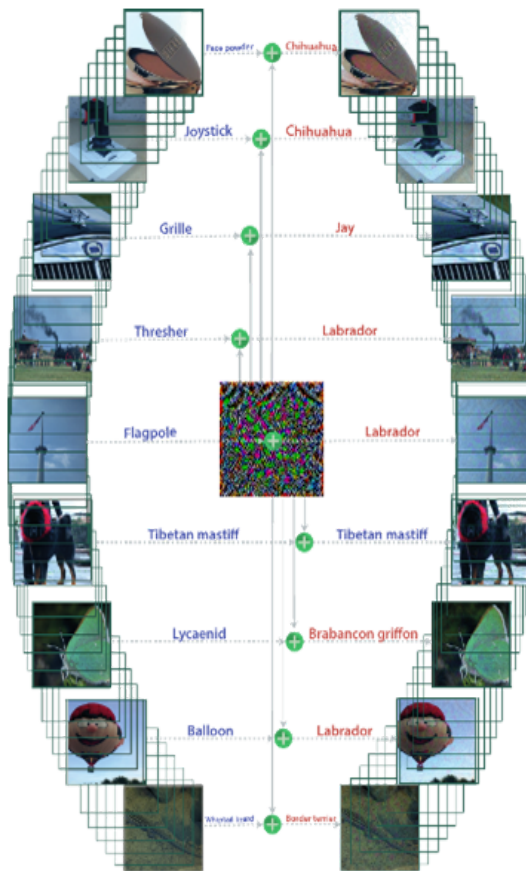
- Edward Chou, Florian Tramèr, Giancarlo Pellegrino. SentiNet: Detecting Localized Universal Attacks Against Deep Learning Systems; <https://arxiv.org/abs/1812.00292>
- Huili Chen, Cheng Fu, Jishen Zhao, Farinaz Koushanfar. DeepInspect: A Black-box Trojan Detection and Mitigation Framework for Deep Neural Networks; <https://www.ijcai.org/proceedings/2019/647>
- Jacob Steinhardt, Pang Wei Koh, and Percy Liang. Certified defenses for data poisoning attacks. In Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17, pages 3520–3532, USA, 2017. Curran Associates Inc.
- Or another anomaly-detection based approach by Yuntao Liu, Yang Xie, and Ankur Srivastava. Neural trojans. 2017 IEEE International Conference on Computer Design (ICCD), pages 45–48, 2017.

You can re-use any existing code, or implement the methods from scratch by yourself. If you prefer, you can also train the backdoored models yourself, instead of using the provided pre-trained models.

### Topic 3.1.1.10: Model Evasion Attack



Evasion attacks are less targeted than backdoors, as they “just” try to avoid a specific class to be predicted.



In this task, you shall evaluate, e.g. on the AT&T Faces dataset or the Yale faces, or a subset of the PubFig/PubFig 83 dataset, how well different methods for generating adversarial images work, by measuring the distortion needed in a metric similarity, complemented with an occasional manual inspection of the perceived difference. Also measure the efficiency of the attacks. Specifically of interest are also methods that generate universal perturbations (i.e. not specific to one instance/image).

Then, try two simple methods for defending against these attacks (e.g. by adversarial training). You can utilise e.g. the IBM Adversarial Robustness Toolbox (<https://github.com/IBM/adversarial-robustness-toolbox>), or also Google’s Cleverhans (<https://github.com/tensorflow/cleverhans>) or Foolbox (<https://github.com/bethgelab/foolbox>)



## Privacy-Preserving Data Publishing

### Topic 3.1.3.3.: Generation and evaluation of sequential synthetic datasets

This is a rather experimental topic.

The above discussed tools for synthetic data generation have been developed for processing structured relational data. One application scenario for synthetic data is however as well in sequential data, which DNA sequences are one specific instance of.

One such method has been described by Pratas et al and uses multiple competing Markov Models. First, a training DNA sequence is partitioned into non-overlapping blocks of fixed size. Each block is handled by a particular Markov Model. In the first phase, the statistics of the training data are loaded into the models. After this deterministic procedure has ended, the models are kept frozen. In a second stochastic phase, the synthetic sequence is then generated according to the learned statistics.

Your task is to recreate the implementation of the specified paper, and evaluate the quality of genetic data when comparing it to real sequences and evaluating it on predictive tasks.

Main paper: D. Pratas, C. A. C. Bastos, A. J. Pinho, A. J. R. Neves, L. M. O. Matos. DNA Synthetic Sequences Generation using Multiple Competing Markov Models. IEEE Statistical Signal Processing Workshop (SSP), pp. 133-136, 2011.

### Topic 3.1.3.4.: k-anonymity and utility of anonymized datasets

$k$ -anonymity is a property of a dataset that contains the same information for at least  $k$  individuals, for every distinct information.  $k$ -anonymous datasets can be obtained by a systematic generalization or suppression of the values of the dataset.

These methods reduce the amount of information contained; therefore  $k$ -anonymous datasets necessarily have lower utility compared to their originals. Utility can be measured in different ways, e.g. using metrics directly on data such as a loss measure, or by data performance on a particular machine learning task.

The aim on this task is to investigate the relationship between

- utility metrics (precision, sum of squared errors, entropy, granularity, ...) and
- classification performance (classification accuracy, F1 score, ...)

of anonymized datasets.

You should utilize at least 2 anonymization algorithms of your choice:

- Flash (within ARX tool; Java) – API (not GUI!): <https://github.com/arx-deidentifier/arx>
- SaNGreeA (Python): <https://github.com/tanjascats/sangreea>
- Mondrian (Python): <https://github.com/qiyuangong/Mondrian>
- implementations from UTD Anonymization Toolbox (Java):  
<http://cs.utdallas.edu/dspl/cgi-bin/toolbox/index.php?go=home>
- Crowds (Python): <https://pypi.org/project/crowds/>
- sdcMicro (R): <https://cloud.r-project.org/web/packages/sdcMicro/index.html>
- some other open-source anonymization tool

The anonymization algorithms usually offer multiple anonymized datasets as a result. Having multiple options for anonymizing the dataset may have certain benefits, e.g. in a scenario where generalizing some attribute is preferred over the others. Therefore, besides calculating the utility for one optimally anonymized dataset, you shall also analyse the difference in utility between optimal and other solutions.

For your analysis you shall use 2 datasets that contain private or sensitive information (e.g. <http://archive.ics.uci.edu/ml/datasets/Adult> - don't use this example :-)). Choose datasets that differ according to some characteristics such as size and type of attributes. Some suggestions:

- [Bone Marrow Transplant data](#)
- [Census Income \(KDD\)](#)
- [Credit Card Clients](#)
- [Drug Consumption Dataset](#)
- ...

For most of the algorithms you may need to add your own implementation of utility metrics to evaluate the quality of anonymized data.

1. Precision: one minus the sum of all call distortions normalized by the total number of cells, where a call distortion is the ratio of the domain of the value in the cell to the height of the attribute's hierarchy (Definition 5 from [1])

2. Loss/granularity: sum of the normalized loss for each column, where the loss of the (numerical) column is the ratio of value interval of that column to the interval of original values (see details in Section 3.1. of [2])
3. Entropy: Section 3.2. of [3]

We have prepared the experimental setup you shall utilize in your analysis (2 Jupyter notebooks). You can find them in a zip file at TUWEL.

Specify in TUWEL at the topic registration which datasets will be used for the experiments. Write a report including all your findings, explain your methodology in detail and elaborate the conclusions.

#### References:

- [1] Sweeney, L., 2002. Achieving k-anonymity privacy protection using generalization and suppression. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05), pp.571-588.  
(<https://desfontain.es/PDFs/PhD/AchievingKAnonymityPrivacyProtectionUsingGeneralizationAndSuppression.pdf>)
- [2] Iyengar, V.S., 2002, July. Transforming data to satisfy privacy constraints. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 279-288).  
(<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.13.2993&rep=rep1&type=pdf>)
- [3] Gionis, A. and Tassa, T., 2008. k-Anonymization with minimal loss of information. *IEEE Transactions on Knowledge and Data Engineering*, 21(2), pp.206-219.  
(<https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.929.4919&rep=rep1&type=pdf>)

### Topic 3.1.3.5.: k-anonymity: generalization vs microaggregation

$k$ -anonymity is a property of a dataset that contains the same information for at least  $k$  individuals, for every distinct information.  $k$ -anonymous datasets are obtained by either:

- a systematic generalization or suppression of the values of the dataset (global, local)
- microaggregation of the clustered values

Global and local transformation are types of generalization/suppression approach. Global transformation is achieved by generalizing the entire attribute (column) to the same generalization level from its hierarchy. Local transformation, however, aims to reduce the value to the least necessary generalization level. The locally anonymized dataset may contain values from different levels of generalization within one column (more information and examples at [1]).

Microaggregation is an approach that does not require defining generalization hierarchies for each quasi-identifier, but rather clusters the rows based on similarity and then generates the anonymized values as a result of some aggregation function (e.g. mean for numerical attributes, median for categorical, ...) [2].

These methods reduce the amount of information contained; therefore  $k$ -anonymous datasets necessarily have lower utility compared to their originals. One way to measure utility is via data performance on a particular machine learning task.

The aim of this exercise is to compare the utilities of the anonymized datasets obtained by different approaches:

1. Global transformation
2. Local transformation
3. Microaggregation

To that end, define classification tasks you want to perform and choose at least 3 different types of classifiers. Calculate the performance on the original dataset and then perform the same tasks using the anonymized datasets. Compare performance of ML classifiers trained on anonymized data to the ones trained on original data and compare how different anonymization approaches influence the performance.

The following anonymization tools provide different types of anonymizations:

- Generalization/suppression:
  - a. Flash (within ARX tool; Java) – API: <https://github.com/arx-deidentifier/arx> (global/local)
  - b. SaNGreeA (Python): <https://github.com/tanjascats/sangreea> (local)
  - c. Mondrian (Python): <https://github.com/qiyuangong/Mondrian> (global)
  - d. implementations from UTD Anonymization Toolbox (Java):  
<http://cs.utdallas.edu/dspl/cgi-bin/toolbox/index.php?go=home> (global)
  - e. Crowds (Python): <https://pypi.org/project/crowds/>
  - f. sdcMicro (R): <https://cloud.r-project.org/web/packages/sdcMicro/index.html>
- Microaggregation:
  - a. (use) sdcMicro-microaggregation (R):  
<https://cloud.r-project.org/web/packages/sdcMicro/index.html>
  - b. (adapt) SaNGreeA (Python): <https://github.com/tanjascats/sangreea>
    - To this end, utilize SaNGreeA algorithm for clustering the dataset into partitions of size at least  $k$ . Then modify the algorithm so that the anonymized values are microaggregated instead of generalized (i.e. for each

cell calculate the mean value of the range if it is a numerical attribute and median (most frequent) value if it is a categorical attribute)

- You can also try to find other open-source anonymization tools

For your analysis choose 3 datasets with privacy implications and different characteristics (number of instances, number of attributes, types of attributes, etc.).

Specify in TUWEL at the topic registration which datasets will be used for the experiments. Write a report including all your findings, explain your methodology in detail and elaborate the conclusions.

References:

[1] [https://sdcpractice.readthedocs.io/en/latest/anon\\_methods.htm](https://sdcpractice.readthedocs.io/en/latest/anon_methods.htm)

[2] J. Domingo-Ferrer and V. Torra, "Ordinal, Continuous and Heterogeneous k-Anonymity Through Microaggregation," *Data Min Knowl Disc*, vol. 11, no. 2, pp. 195–212, Sep. 2005, doi: [10.1007/s10618-005-0007-5](https://doi.org/10.1007/s10618-005-0007-5).



### Topic 3.1.3.7.: Differential Privacy

Differential privacy is introduced by Cynthia Dwork [1] as a set of privacy-preserving mechanisms for publicly sharing information about the dataset by descriptive characteristics rather than sharing the information about the individuals in the dataset. Differential privacy is at first achieved with simple aggregate statistics and descriptive statistics such as mean function, maximum, histograms, etc., and recently the advances with differentially private machine learning models are achieved, as well.

Differential privacy is achieved by adding noise in the process. Depending on the phase where the noise is added we differentiate:

- Input perturbation - adding noise to the input before running the algorithm
- Internal perturbation - randomizing the internals of the algorithm
- Output perturbation - add noise to the output, after running the algorithm

The goal of this exercise is to compare these approaches for achieving differential privacy, from the privacy aspect and model performance aspect and analyse the tradeoff between privacy and performance. To that end, to achieve internal perturbation, utilize a differential privacy tool that provides differentially private versions of classification and clustering models (for Python:

<https://github.com/IBM/differential-privacy-library>, for R: <https://github.com/brubinstein/diffpriv>).

Familiarize yourself with the differentially private mechanisms used in the chosen tool.

Secondly, apply an input perturbation and output perturbation based on the sensitivity method proposed by Dwork et al. in [1]. You may experiment with a model of your choice.

You shall analyse:

1. Input perturbation vs. internal perturbation vs. output perturbation approach on a chosen model; difference between the models' performances, and differences to the non-private model.
2. Trade-off between privacy and model performance for every approach.

Experiments from these papers [2, 3] may serve as an inspiration for the methodology of this exercise.

#### References:

[1] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating Noise to Sensitivity in Private Data Analysis," in *Theory of Cryptography*, Berlin, Heidelberg, 2006, pp. 265–284. PDF:

[https://link.springer.com/content/pdf/10.1007/11681878\\_14.pdf](https://link.springer.com/content/pdf/10.1007/11681878_14.pdf)

[2] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate, "Differentially Private Empirical Risk Minimization," *Journal of Machine Learning Research*, vol. 12, no. Mar, pp. 1069–1109, 2011. PDF:

<https://arxiv.org/pdf/0912.0071.pdf>

[3] K. Fukuchi, Q. K. Tran, and J. Sakuma, "Differentially Private Empirical Risk Minimization with Input Perturbation," in *Discovery Science*, Cham, 2017, pp. 82–90. PDF:

<https://arxiv.org/pdf/1710.07425.pdf>

#### Further readings - differentially private ML:

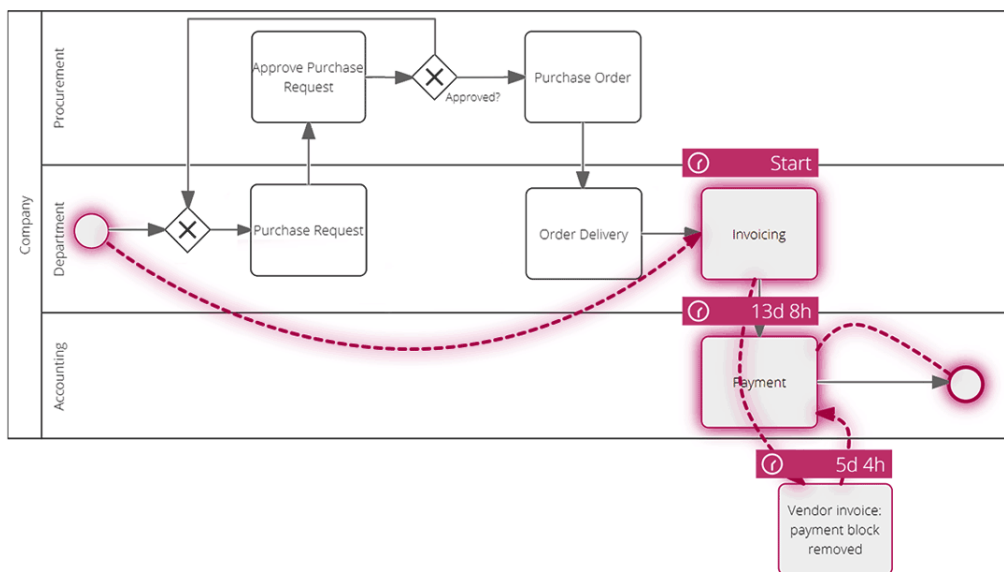
- [4] O. Sheffet, "Private Approximations of the 2nd-Moment Matrix Using Existing Techniques in Linear Regression," arXiv:1507.00056 [cs], Nov. 2015.
- [5] J. Vaidya, B. Shafiq, A. Basu, and Y. Hong, "Differentially Private Naive Bayes Classification," in *2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*, Nov. 2013, vol. 1, pp. 571–576.
- [6] H. Imtiaz and A. D. Sarwate, "Symmetric matrix perturbation for differentially-private principal component analysis," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Mar. 2016, pp. 2339–2343.

## Secure / Federated / Distributed Computation

Distributed computation in general addresses privacy / data protection issues by not centralising data that is initially already partitioned between several parties, and either computing functions securely w/o revealing the input (secure multi-party computation), or by exchanging only aggregated values, such as gradients or other model parameters (federated learning), or by computing on encrypted data (homomorphic encryption).

In this task, you shall evaluate the effectiveness and efficiency of these approaches. An evaluation on the quality of the models, as compared to having a model trained on all the different partitions of the data, shall be performed. - i.e. observe the effectiveness and efficiency on the centralised vs. the distributed approach.

### Topic 3.1.2.1: Federated (Distributed) Learning for Process Mining



Process mining tries to learn the structure of a process from logs or other sources about previous executions of process. This is typically done on centralized data, i.e. with centralized learning. Collaborative learning of such a process model from different organisations can be an interesting setting in many cases. However, instead of centralising the data, which might not be an option due to confidentiality reasons, federated learning can become a viable approach.

In this task, you shall perform experiments for distributed process mining. To this end, take an existing benchmark dataset from e.g. [https://data.4tu.nl/repository/collection:event\\_logs\\_real](https://data.4tu.nl/repository/collection:event_logs_real), and first perform process mining in a centralised way with an existing approach (e.g. using ProM, <http://www.promtools.org/>). Then, try to federate / distribute the process, by first distributing the data, and then the learning. Data distribution could be done in uniform (each subset is roughly similar) and non-uniform ways (some aspects of the data are present only at one node).

One approach for distributed learning you could follow is: “On the Feasibility of Distributed Process Mining in Healthcare”, 2019. [https://link.springer.com/chapter/10.1007/978-3-030-22750-0\\_36](https://link.springer.com/chapter/10.1007/978-3-030-22750-0_36)

### **Topic 3.1.2.3: Federated Learning on text data (or other sequential data)**

In this task, you shall work with federated approaches for text data utilising RNNs (LSTM, GRU, ...). Take a task of your choice (e.g. a text categorisation / classification task like the Reuters Dataset, or 20 newsgroups, etc..), and first establish a baseline with a state-of-the-art recurrent neural network, comparing your results to literature.

Then, simulate a federated setting by distributing the data on multiple sources, and compare your results to a centralised setting, similar to the setting in the topic above.

### **Topic 3.1.2.4: Federated Learning on tabular data with different ML algorithms**

In this task for federated learning, you shall work with tabular data (i.e. not sequential, not image, ...). You shall pick two different datasets, which can be any dataset that you have e.g. used in the Machine Learning lecture in Exercise 1 or 2, and employ federated learning with two different types of algorithms on it. You are free to reuse existing implementations (like the ones mentioned above), or you can implement your own approach. For different algorithms, you can consider e.g. MLPs, but also SVMs, Naive Bayes, Random Forest...

Simulate a federated setting by distributing the data on multiple sources, and compare your results to a centralised setting, similar to the setting in the topic above.

Consider federated learning with different number of nodes in the setting and horizontally partitioned data (nodes have different instance in their data but with the same features). You should fix the dataset size and distribute the data among considered number of nodes. Try to create unequal and unbalanced data distribution to be more close to the real setting. You can start searching for different federated learning ml algorithms from this paper (page 14):

<https://arxiv.org/pdf/1907.09693.pdf>

### **Topic 3.1.2.5: Federated Learning, different aggregation methods**

There are several ways how one can aggregate model parameters in Federated learning. In this task you will need to compare different aggregation methods for federated learning. You can take FedAvg (Federated averaging algorithm) as one of the most popular algorithms, also try weighted averaging and find 1-2 other aggregation algorithms. The result work should include evaluation of different aggregation methods in the settings with different number of nodes. You can choose 2 datasets for this task and horizontally partitioned data (nodes have the different instance in their data but with the same features). You should fix the dataset size and distribute the data among considered number of nodes. Try to create unequal and unbalanced data distribution to be more close to the real setting. For classification task use neural network algorithm. You can start searching for different aggregation methods in this paper: <https://arxiv.org/pdf/1907.09693.pdf> (pages 12-15)