# Portfolio-Exam Tasks

**Read this document (3 pages) carefully after having read the document `portfolio.pdf`!**
Think of the portfolio exam as a report of a (small) project which you – as a data scientist in a data science company – conduct for a customer (a company, an institute, some organization).
The main goal of the experiments is to create a deep learning model tackles a **regression or classification task**. Therefore, in your pitch, choose and plan a context in which regression or classification is the means to solve **one** problem. Imagine and describe a fictitious situation or describe a real scenario in which you are a data scientist working for some organization. The context must be meaningful such that it could or does exist in real life.

## Task 1 – The Data

Load and present a dataset for the subsequent experiments.

- Explain the dataset itself (e.g., what do the features represent, units, . . . ).

- Explain how the dataset is suitable for the project proposal in you pitch.

- Specifically keep the conditions for selecting datasets in mind!

- When presenting the data, make sure to present ALL columns, e.g., use the pandas option `pd.options.display.max_columns = None`.

## Task 2 – IDA

Conduct an initial data analysis.

- Present relevant quantities of the data that inform the reader about the dataset or that might indicate the need for preprocessing.

## Task 3 – Preprocessing

Bring the dataset into the form that you need for the experiments.

- Conduct necessary transformations according to your findings in Task 2 until the dataset is suitable.

- If you change data, do not forget to present updated relevant quantities of the result.

## Task 4 – EDA

Use the preprocessed dataset from Task 3 to compute interesting statistics, distributions, and feature relations that are relevant to the pitched task. Limit the output to **3 particularly interesting aspects**!

## Task 5 – Baselines

As a comparison for the upcoming tasks, compute two baselines:

- Pick two evaluation metrics to evaluate the final model and explain your choice.

- Use a very simple **suitable heuristic** (not a trained model) as a baseline.

- Run and evaluate **one classical machine learning** algorithm. In this exam, no hyperparameter tuning is required for the baselines. Use default values of the implementation.[1]

## Task 6 – Deep Learning Experiments

This is the largest and most important task in the portfolio. Build, train, optimize, and evaluate a suitable deep learning model for the pitched task.

- Demonstrate, describe, and interpret all necessary steps for conducting the deep learning experiment for the project!

- Conduct a reproducible experiment!

- Select, implement, and explain exactly one architecture of a deep learning model. Explain your choice for the usecase at hand. The choice should be reasonable in light of the pitched project and use methodology that seems the most promising for the task. For that purpose choose from the layers and cells discussed in the lectures.[2]

- Use at least one dropout layer in the model.

- When training the model

  - use a `DataLoader` to load the training data in three or more batches.

  - use an early stopping mechanism.

- Optimize the model by finding a good choice for the dropout probability in the dropout layer(s). If you use more than one dropout layers, use the same parameter for all such layers. Visualize your results.

- Select reasonable values for all other learning and model hyperparameters (no further optimization).

- Select one final choice for the dropout parameter and train a final model. Track the learning progress using `Tensorboard` and include and interpret a suitable visual in the notebook.

- Evaluate your final model using the above discussed metrics and present a comparison to the results of your baselines in a table.

---

[1]In a real experiment, of course proper hyperparameter tuning would be part of the experiments.

[2]Since for the exam, you should choose only one model (without trying different approaches), you must find an architecture that looks promising, given the task at hand. You are of course free to try various different models. However in the submitted version, please have exactly one model.

## Task 7 – Conclusions and Future Work

Please address the following points separately and in that order.

1. Summarize and interpret the achieved results.

2. Recommend a course of action for the organization in your story based on the results of Task 6.

3. Reflect on limitations and possible pitfalls of using these results.

4. Propose three ideas for future work (a short sketch or enumeration of ideas is sufficient, no further experiments). The ideas should not be too general (e.g., "try further algorithms") but be specific to the project (e.g., "try Algorithm X, as because of Property Y, it might work specifically well on this dataset").

5. Critically discuss the employed methodology (your choices as well as the choices given in these tasks). What could or even should have been done differently?

6. Critically reflect the original task you pitched. In hindsight, were the goals realistic? What could have been changed at the time of the pitch?