

Exercise project 1 – ANN for regression

In this exercise project, regarding the use of generative AI:



Yellow – Allowed, can be used, must be reported

Artificial intelligence can be used in the creation of outputs, but the student must clearly report its use. Failure to disclose the use of AI will be interpreted as fraud. The use of AI may affect the assessment.

Note: using AI itself does not affect scoring negatively in this exercise project.

Just remember to mention everything clearly when AI has been used (including the intended purpose, reason and how can you be sure the AI's output is correct).

The aim of this exercise project is to experiment with a neural network of mainly Dense-layers for a regression dataset. **Every student has to make this exercise project with their own unique dataset.**

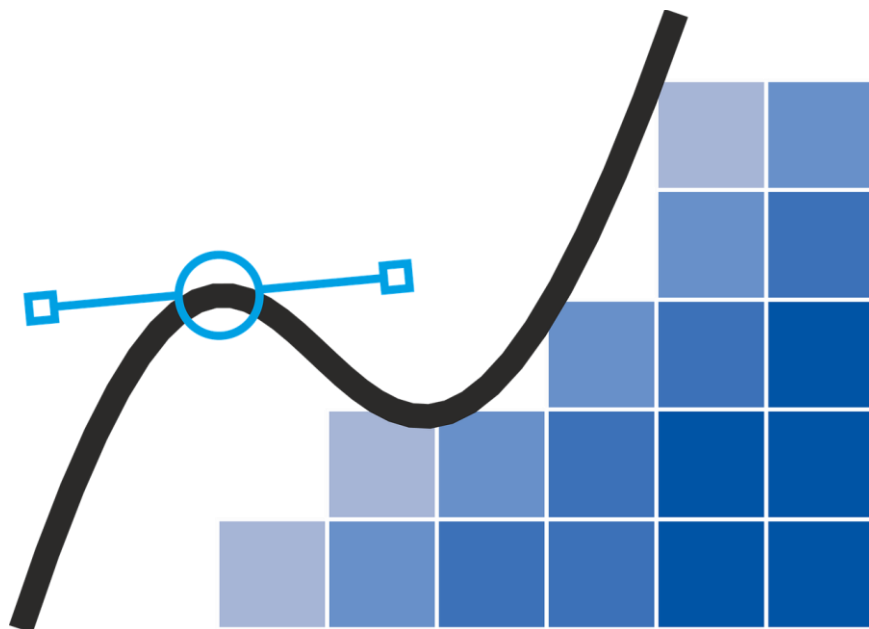
Create a Jupyter notebook for this exercise project (or multiple notebooks if you wish to do advanced tasks or multiple versions of your neural network as well).

Examples that help you with the different steps can be found in Moodle!

Remember to write your thoughts and personal analysis after each step.

Remember also: if you refer to a website in your analysis etc., also have the link in your notebook as a reference (when explaining what you found in your data analytics!).

Finally remember: Any kind of dataset/training algorithm optimizations used in the exercise project is considered extra points (see last slides of materials Part 2.)



Step 1

Find yourself a suitable dataset for this exercise project. Use either Google or Kaggle.com to find yourself a suitable dataset. If you don't have an account in Kaggle.com yet, do this first.

Requirements for the data:

- Do not use time-series data, we'll use them in later exercises (ANN regression is not suitable for time-series with conventional methods)
- At least 4 numeric columns (you can have as many as you want too)
- One of the columns have to be continuous (the target variable), something that theoretically doesn't have an upper limit (prices, measurements etc.)
 - **Other alternative:** a column that has large range of options in its values (instead of just 1, 2, 3, have something like 1-10, 1-100 or something more)
- The data has to have some kind of logical trend (linear, fluctuating, seasonal etc.)
 - If data is completely chaotic and without rhyme or reason, it won't work well with a neural network without processing the data
- Once you find a suitable dataset, **reserve your dataset in the link in Moodle, and check that it hasn't been selected by another student already.**

If multiple students have returned the same dataset => the first student will get the dataset.

Remember: every student has to have a unique dataset in the exercise.

Step 2

Clean up your dataset, handle missing values and check the balance of the data.

Remember to convert all non-numeric columns into a numeric representation (see materials and encoder-examples in Moodle).

Idea for an advanced task: if your dataset doesn't have a good balance, figure out how to get more data in order to have a better balance (for example, if data doesn't have enough examples of certain categories etc.)

For optimizations: Check also the last slides in materials part 2.

Step 3

Do the train/test –split for your cleaned dataset. You can use either 70%/30% or 80%/20% -split.

Alternatively, you can also add a validation dataset, in which the split is commonly 70% for training, 15% for testing and 15% for validation.

Step 4

Create a neural network of mostly Dense-layers and with an output layer of one node. The input layer has to match with the amount of variables you want to predict the target value with. (If you have 10 variables => input layer needs to have 10 as the input shape). You can experiment with different lengths and width with your neural network. Remember to have a complex enough neural network regarding the amount and complexity of your data.

Note: if the neural network is too simple, the neural network doesn't have enough "decision-space" or model capacity to work well with your dataset! (try to increase length or width of your neural network if this happens).

Step 5

Fit your data to the neural network model. After training, visualize with a scatter plot how well the predictions follow the true values (analyze how well your model follows a linear line). Visualize also how well the prediction distribution follows the normal distribution (see code examples in Moodle).

Also, visualize how well the training process went visually (training loss). You can also use a validation data set if you want more tools to see if the model overfits.

If your model overfits, re-train your neural network with less epochs, or alter your neural network (or even data).

Step 6

Inspect these error metrics for your model:

- MEA
- MSE
- RMSE
- R-squared
- Explained variance score

What can you say about these error metrics, how well the model performs? (R-squared: 1 is perfect, 0 is worst outcome. It measures how well the model describes the original dataset. Explained variance score does the same but regarding dataset variance).

Step 7

Try your model with some imaginary new values. Does it work as expected? Try to emphasize variables that should affect the prediction greatly (see the correlation matrix in pandas to figure out helpful variables).

Step 8

Remember to write your personal analysis after every phase you did previously within your Jupyter notebook(s). How does the code work, and where could linear regression be useful in working life? Was it easy or difficult to use? Anything else that come into mind? Any ideas for optimizations?

Final step (optional, but recommended)

Create a simple UI for your model and allow user to input test values for the regression (TKinter, PySimpleGUI, web etc.). (We'll probably experiment with this in one of the lectures)

Advanced task ideas

1. Experiment with scikit-learn's Linear Regression –algorithm (classic ML) and compare the results. Which approach works better with your data, classic ML or neural networks? (you can also think why this probably happens)
 - **Note: there are cases where classic ML methods have a better performance**
 - **If possible, experiment with two different datasets: one with strong linear data relations and one with unlinear data. Then compare scikit-learn to your neural network. What differences can you find?**
2. Instead of experimenting with the epoch-variable to optimize your neural network or to avoid overfitting, implement these features in your neural network:
 - Dropout-layers
 - EarlyStop
 - ReduceLROnPlateau
 - ModelCheckpoint

