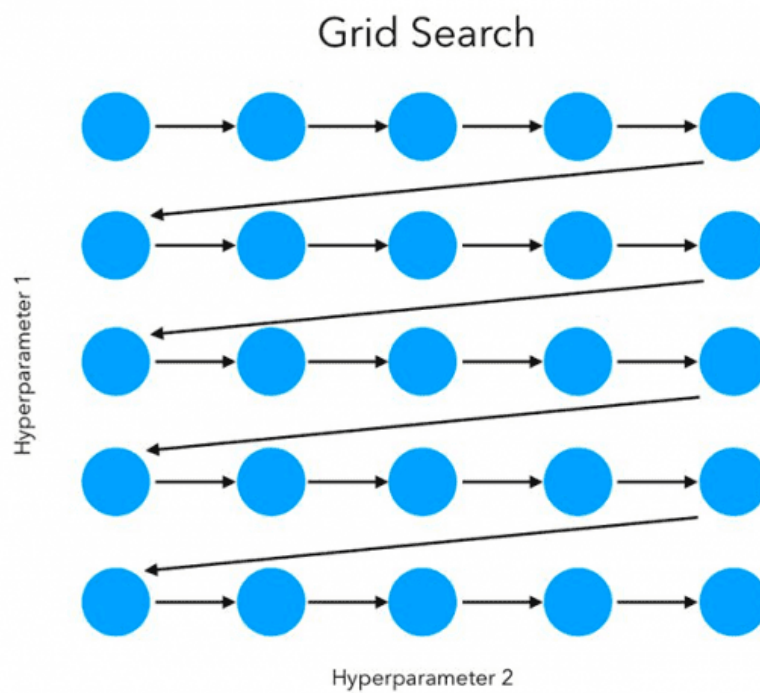


O3

SWMAL-01 GROUP 13



Group members

Christoffer Regnar Mølck - 202009347
Martin Teit Bruun Andersen - 202109592
Morten Kierkegaard Erichsen - 201807571

17. april 2024

1 07 CNN

1.1 What are CNNs?

Convolutional neural network's (CNN's) were invented in the 80's by Kunihiko Fukushima. His findings are inspired by Hubel and Wiesel's work on cats visual cortices neurons that had small individual visual fields. This inspiration was the beginning for feature extraction, pooling layers, recognition and classification. Kunihiko Fukushima's full article on Neocognitron can be read here: <https://www.rcin.org/bruno/public/papers/Fukushima1980.pdf>

Today CNN's are most commonly known for its use in image analysis. However, it is also commonly used for analysing stock prices and audio data.

Looking at image analysis, the CNN works by taking in a series of layers. Each of these layers looks at different features in the image. These features come in different sizes. One filter could be looking at vertical lines, while another could be finding eye features. However, the whole thing starts by a convolution operation. This operation is like a sliding filter of a given size. The size is normally a 3x3 square, however as mentioned it is user decided. The filter is called a kernel. Inside the kernel it reads the individual pixels within its square. And depending on the RGB color value it returns a number input (from 0-255) in the kernel square. Each complete kernel run over the input layer and creates a feature map. If an example were to be made on the MNIST numbers, one of these created feature maps could be identifying black values in the input layer. From the convolution, each layer is pooled to a smaller layer size. Pooling layers downsample the convolution created feature maps. And on and on... However, for each convolution and pooling one has to be aware of overfitting the CNN model.

In the CNN code below, generated by Google Gemini, we have inserted a visual learning curve, that visualizes how the loss is throughout the epochs. Our CNN's input layer is a 28 by 28 by 1. Through our convolution we then take and expand the last layer, creating 32 channels. Each of these channels is then pooled through a 2 by 2 filter reducing the convolution layer from 26 by 26 by 32 to 13 by 13 by 32. Our second convolution layer then takes on our pooled 32 layers and expands it to 64 layers with a smaller size. Before the data is flattened our second convolution layer is then pooled to its last size. When the data is flattened we see, from the `model.summary()`, that the value is 1600. This value represents the total amount pixels our last pool contained ($5*5*64 = 1600$).

The flattened data is then held on to the first dense layer of 128 feature recognizers, meaning that this dense layer has 128 neurons/categories. Afterwards this layer is then broken down to our input size (mean we have 10 different input images).

1.2 Conclusion

As a conclusion to the assignment the group tried out multiple amount of convolution layers and decided on including the one returning the best score for us. High score is commented in the bottom of the code. We found that for each convolution layer we added the worse the CNN performed leaving us with two since one was to little. We found that increasing the amount of feature maps in the first "Conv2D" would lead to increased step time when running the evaluation of the model. This makes sense due to the increased amount of features the model would look for.

1.3 Our model

```
1 from tensorflow import keras
2 from tensorflow.keras import layers
3 import matplotlib.pyplot as plt
4
5 # Define the model (replace with your final CNN architecture)
6 model = keras.Sequential([
7     # First convolutional layer
8     layers.Conv2D(32, kernel_size=(3, 3), activation="relu", input_shape=(28, 28,
9         ↪ 1)),
10     layers.MaxPooling2D(pool_size=(2, 2)),
11
12     # Second convolutional layer
13     layers.Conv2D(64, kernel_size=(3, 3), activation="relu"),
14     layers.MaxPooling2D(pool_size=(2, 2)),
15
16     layers.Flatten(),
17     layers.Dense(128, activation="relu"),
18     layers.Dense(10, activation="softmax")
19 ])
20 print(model.summary())
21
22 # Compile the model
23 model.compile(optimizer="adam", loss="categorical_crossentropy",
24     ↪ metrics=["accuracy"])
25
26 # Load the MNIST dataset
27 (x_train, y_train), (x_test, y_test) = keras.datasets.mnist.load_data()
28
29 # Preprocess the data
30 x_train = x_train.astype("float32") / 255.0
```

```

30 x_test = x_test.astype("float32") / 255.0
31 x_train = x_train.reshape(-1, 28, 28, 1)
32 x_test = x_test.reshape(-1, 28, 28, 1)
33 y_train = keras.utils.to_categorical(y_train, num_classes=10)
34 y_test = keras.utils.to_categorical(y_test, num_classes=10)
35
36 # Store training history
37 history = model.fit(x_train, y_train, epochs=10, validation_data=(x_test,
    ↪ y_test))
38
39 # Plot the learning curve
40 plt.plot(history.history['loss'], label='Training Loss')
41 plt.plot(history.history['val_loss'], label='Validation Loss')
42 plt.title('Training and Validation Loss')
43 plt.ylabel('Loss')
44 plt.xlabel('Epoch')
45 plt.legend(loc='upper right')
46 plt.show()
47
48 # Evaluate the model on test data
49 loss, accuracy = model.evaluate(x_test, y_test)
50 print("Test accuracy:", accuracy)
51
52 # Highest score run 99.33

```

```

... Model: "sequential_25"
...

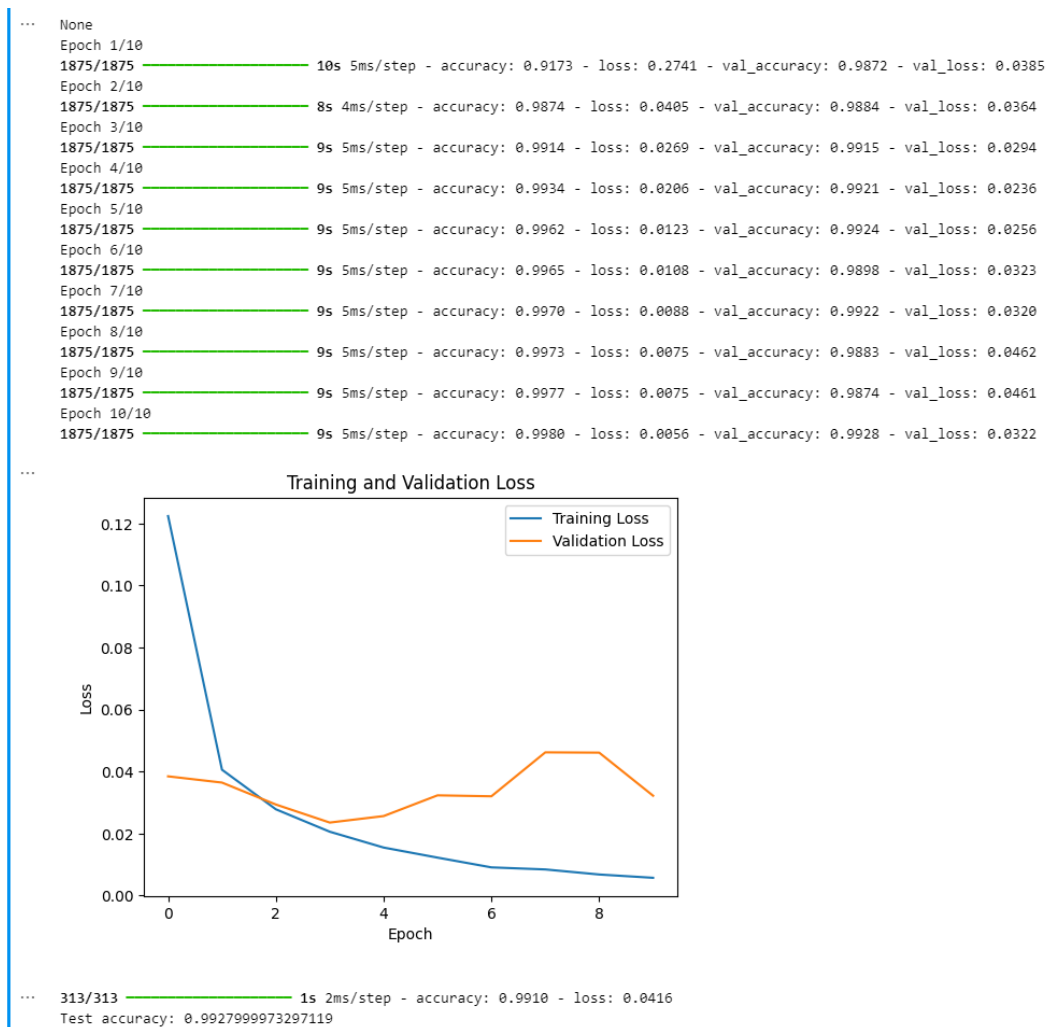
```

Layer (type)	Output Shape	Param #
conv2d_66 (Conv2D)	(None, 26, 26, 32)	320
max_pooling2d_66 (MaxPooling2D)	(None, 13, 13, 32)	0
conv2d_67 (Conv2D)	(None, 11, 11, 64)	18,496
max_pooling2d_67 (MaxPooling2D)	(None, 5, 5, 64)	0
flatten_25 (Flatten)	(None, 1600)	0
dense_52 (Dense)	(None, 128)	204,928
dense_53 (Dense)	(None, 10)	1,290

```

...
Total params: 225,034 (879.04 KB)
...
Trainable params: 225,034 (879.04 KB)
...
Non-trainable params: 0 (0.00 B)

```



Figur 1.1: CNN model output

2 GENERALIZATION ERROR

2.1 Qa) On Generalization Error

Training error is the error we get when we compare the training data to the model prediction, while generalization error is what we get when we compare the model to previously unseen data.

As a model can be very well fitted to data it has seen before, but poorly regularized. There are two types, under and over fitting. Under fitting, is when the model fits poorly to both the training and unseen data. While over fitting is when the model fits very well to the training data but poorly to any unseen data. Ideally the model should be general enough to fit both somewhat equally well. This point is called optimal capacity.

The generalization gap is the gap between the training error and the generalization error. Capacity is simply a term to express a models complexity or its capacity to handle complex data relations.

2.2 Qb A MSE-Epoch/Error Plot

The code performs polynomial regression using SGD, demonstrating how the model learns over time by plotting the RMSE against epochs. An epoch represents one complete pass through the entire training dataset, allowing the model to update its weights.

`mse_train` and `mse_val` refer to the mean squared error on the training and validation datasets, respectively, offering a measure of how well the model predicts compared to the actual values, with lower values indicating better performance.

The generation of the array is good for visualizing the model's learning progress, highlighting the balance between fitting the training data well (low `mse_train`) and generalizing effectively to new data (low `mse_val`). Tracking these metrics helps in identifying the best model configuration that avoids overfitting while maintaining good predictive accuracy.

The code has been changed slightly, iteration count was changed from 0 to 1000 and the `-float(inf)` to 0.0 as it otherwise gave an error and did not run.

(Initial suggestion by ChatGPT, text changed significantly)

2.3 Qc) Early Stopping

I would record the current validation score, if it does not increase in the next 5 iterations i would break out of the loop. so:

```
1  if mse_val > mse_val_highest
2      i++
3      if i == 5
4          break
5  else i = 0
```

2.4 Qd) Explain the Polynomial RMSE-Capacity plot

The x-axis shows the model complexity, indicated by the polynomial degree. The y-axis is the prediction error's standard deviation. in other words how much the data deviates from the prediction.

The training RMSE declines as model capacity increases because a higher degree polynomial can better fit the training data (it has more complexity). However, the validation RMSE decreases only to a point before it starts to increase, suggesting the model is over fitting. it's learning noise so the generalization becomes much worse.

Increasing the polynomial degree to 10 leads to higher validation RMSE, indicating more severe over fitting. The optimal model strikes a balance, fitting the training data well without capturing excessive noise, and thus generalizes better to new data.

3 GRIDSEARCH

3.1 Qa Explain GridSearchCV

Cell to starts by loading the data using the load LoadAndSetupData function defined in cell one. This just loads a specific dataset depending on a passed in string.

The model is then instantiated and the parameters we want to grid search over are created. In this case we want to search through the kernal type, and the C values to find the best model.

the grid search is then run on the new model using the tuning parameters to tell it waht to search over.

The result is reported using the function full report that is facade that runs a few other functions to create a report on the result of the gridsearch. The report can then be used to create a new model with the found optimal parameters.

3.2 Qb Hyperparameter Grid Search using an SDG classifier

I repeat a lot of the code from Qa but using the SDG model instead, and using a few more randomly selected parameters for tuning.

```
1 model = SGDClassifier() #Instantiate mdoel with default params.
2
3 tuning_parameters_SGD = {
4     'loss': ['hinge', 'modified_huber', 'squared_hinge', 'squared_error'],
5     'penalty': ['l2', 'l1', 'elasticnet'],
6     'alpha': [0.0001, 0.001, 0.01, 0.1],
7     'learning_rate': ['constant', 'optimal', 'invscaling', 'adaptive'],
8     'eta0': [0.0001, 0.001, 0.01, 0.1],
9     'max_iter': [1000, 2000, 3000],
10    'tol': [1e-3, 1e-4, 1e-5]
11 }
12
13 CV = 5
14 VERBOSE = 0
15
16 # Run GridSearchCV for the model
17 grid_tuned_SGD = GridSearchCV(model,
18                               tuning_parameters_SGD,
```



```

19         cv=CV,
20         scoring='f1_micro',
21         verbose=VERBOSE,
22         n_jobs=-1)
23
24 start = time()
25 grid_tuned_SGD.fit(X_train, y_train)
26 t = time() - start
27
28 # Report result
29 b0, m0 = FullReport(grid_tuned_SGD, X_test, y_test, t)
30 print('OK(grid-search)')
```

this takes a few minute to run, but nothing too bad. Outputs:

```

Detailed classification report:
  The model is trained on the full development set.
  The scores are computed on the full evaluation set.

      precision recall f1-score support

0 1.00 1.00 1.00 16
1 0.94 0.89 0.91 18
2 0.83 0.91 0.87 11

accuracy 0.93 45
macro avg 0.92 0.93 0.93 45
weighted avg 0.94 0.93 0.93 45

CTOR for best model: SGDClassifier(alpha=0.1, eta0=0.0001, loss='modified_huber', max_iter=3000,
    penalty='l1', tol=0.0001)

best: dat=iris, score=1.00000,
    ↳ model=SGDClassifier(alpha=0.1,eta0=0.0001,learning_rate='optimal',loss='modified_huber',
    ↳ max_iter=3000,penalty='l1',tol=0.0001)
```

3.3 Qc Hyperparameter Random Search using an SDG classifier

So the randomsearch essentially uses a grid as well, however it just picks random points on that grid and test them. n_iter defines how many tests it runs. so an n_iter of 20 runs 20 tests and it returns the best one of those 20. This is obviously less accurate then gridsearch that brute forces every combination, however in my small test here, the scores ended up being the same as the hyper parameters did not have a significant effect on the scoring. So in this case, a shorter more random search works as well as grid search does. This is however rarely true when applied to real data with more complex models.

3.4 Qd MNIST Search Quest II

We have chosen to use HistGradientBoostingClassifier as our model, as it can be very accurate within its dataset and is well suited for training on large datasets.

This type of model which uses boosted trees also does not need extra data preprocessing which is convenient. This type of model however performs very poorly outside the extremes of its training data.

We started running a gridsearch, but due to the complexity of the dataset it took far too long to complete. So we changed the search to random with 20 iterations hoping that would be sufficient to find a good model.

sadly we are running out of time and had to do only a single iteration as running the boosted trees takes quite a while. I am sure we could have gotten a much better result with some tweaking. Still a score of 0.95 is quite good for a first try.

we ended up with a score:

```
best: dat=mnist, score=0.95196,  
      ↳ model=HistGradientBoostingClassifier(l2_regularization=0.1,learning_rate=0.01,  
      ↳ loss='log_loss',max_bins=255,max_iter=200,max_leaf_nodes=31)
```

4 REGULIZERS

4.1 Qa) The Penalty Factor

Code:

```
1 def Omega(w):  
2 return np.linalg.norm(w[1:])**2
```

Explanation: The penalty factor is low when the weights are small (numerically), and large when the weights are large as well (also numerically).

4.2 Qb) Explain the Ridge Plot

Explanation The α in the 'Ridge' model is a constant that multiplies the objective function:

$$\|y - Xw\|_2^2 + \alpha \cdot \|w\|_2^2$$

in an attempt to minimize it. In the case of $\alpha = 0$, it looks as if there's a lot of overfitting, since the model being used is linear. The same with $\alpha = 1e - 5$, a bit less overfitting, but still significant. With $\alpha = 1$, however, the fit looks to be underfitted, since J actually gets penalized by the entire $\|w\|_2^2$ value. To determine whether or not $\alpha = 1$ indeed is a good fit, one would have to check the model score.

With $\alpha = 0$, the fitting is the same as the regular least squares method, since the penalty term simply becomes 0. The SKLearn actually advises not to use the Ridge function with $\alpha = 0$, for "Numerical reasons". Source: [.](#)

4.3 Qc Explain the Ridge, Lasso and ElasticNet Regularization Methods

Explanation

First, looking at the formulas, from the Scikit-learn documentation:

Ridge:

$$\|y - Xw\|_2^2 + \alpha \cdot \|w\|_2^2$$

Lasso:

$$(1/(2 \cdot n_{samples})) \cdot \|y - Xw\|_2^2 + \alpha \cdot \|w\|_1$$

ElasticNet:

$$1/(2n_{samples}) \cdot \|y - Xw\|_2^2 + \alpha \cdot l1_{ratio} \cdot \|w\|_1 + 0.5 \cdot \alpha \cdot (1 - l1_{ratio}) \cdot \|w\|_2^2$$

It is quite apparant that the main difference are as follows:

(note: coefficients = weights)

Ridge uses the L2 norm, and no constant depending on the number of samples

Lasso uses the L1 norm, and is dependant on the number of samples

Elasticnet is also dependant on the number of samples, while using both the L2 and L1 norms. Furthermore, it specifies a(n) $L1_{ratio}$, that gives the norm L2 and L1 their own form of weight. It is also clear, that if one sets the $L1_{ratio} = 1$, one would get the Lasso regulizer.

The α value in Ridge controls how much the coefficients "shrink", with a larger α shrinking the coefficients closer to zero. With a larger α , the model becomes better at handling co-linearity, since the coefficients become smaller and smaller.

Lasso estimates "sparse coefficients", which are coefficients that are non-zero, which is useful since that reduces the number of features on which a solution is dependant. In other words, it "prefers" solutions with fewer non-zero coefficients, meaning that under some circumstances, it will be able to return the exact values of the coefficients on which the model is dependant.

Being dependant on both L2 and L1, ElasticNet allows for a combination-ish of Lasso and Ridge, since it will then include some of the sparse, non-zero coefficients, while also maintaining the regularization properties of Ridge (second term in the formula).

4.4 Qd) Regularization and Overfitting

Regularization is useful in regards to reducing a models tendency to overfit, since the model gets penalized by some value, making it impossible for the model to fit too closely to the data. If theres no penalty, the model wont be able to generalize for unseen data, resulting in overfitting.

the tug-of-war is kind of a balancing act between the MSE and the regularization Ω . The MSE aims to minimize the error between what the model predicts, and what the actual target values are, pushing for a better fit of the training data. The regularization term Ω "wants" to keep the model parameters small, putting shackles on the model, preventing it from overfitting (and reduces its complexity).

When $\alpha = 1$ the regularization is much stronger, which can be beneficial for small training sets, since fewer data points means less certainty (large data good, small data bad, ish). But, as we saw in the Ridge model in Qa, with $\alpha = 1$, the model doesnt really capture any behaviour in the middle section of the data, making it underfitted. Therefore its important to tune the α value carefully.

Optional answer: Scaling would be of help, since scaling the data down to some range, eg. $[-1:1]$, typically helps the model converge quicker, preventing the regularization term from being dominated by larger scales, eg $[-1e5:1e5]$. By standardizing the data, $\mu = 0$ and $\sigma = 1$, the data becomes more comparable, and as with scaling, prevents the regularization term from being dominated by larger values. By normalizing the data, each data-point would contribute equally to the regularization term, meaning no outliers (eg. faulty data) would contribute more heavily when the model converges.

4.5 Qe) Regularization Methods for Neural Networks

The explanation part is covered in the previous questions.

The following code is written by ChatGPT, which loads a housing dataset from Boston, and trains a model using Ridge, Lasso, ElasticNet and no regularization:

Code:

```

1 from sklearn.datasets import load_boston
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.pipeline import Pipeline
5 from sklearn.neural_network import MLPRegressor
6 from sklearn.linear_model import Ridge, Lasso, ElasticNet
7 from sklearn.metrics import mean_squared_error
8
9 # Load the Boston Housing dataset
10 data = load_boston()
11 X, y = data.data, data.target
12
13 # Split the data into train and test sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
15     ↪ random_state=42)
16
17 # Define pipelines for preprocessing and model training
18 pipelines = {
19     "Ridge": Pipeline([
20         ("scaler", StandardScaler()),
21         ("regressor", Ridge(alpha=0.5)) # Ridge (L2) regularization
22     ]),
23     "Lasso": Pipeline([
24         ("scaler", StandardScaler()),
25         ("regressor", Lasso(alpha=0.5)) # Lasso (L1) regularization
26     ]),
27     "ElasticNet": Pipeline([
28         ("scaler", StandardScaler()),
29         ("regressor", ElasticNet(alpha=0.5, l1_ratio=0.5)) # ElasticNet
30     ↪ regularization
31     ]),
32     "No Regularization": Pipeline([
33         ("scaler", StandardScaler()),
34         ("regressor", MLPRegressor(hidden_layer_sizes=(100,), activation='relu',
35     ↪ solver='adam', max_iter=1000)) # Neural network without
36     ↪ regularization
37     ])
38 }
```

```

34 }
35
36 import matplotlib.pyplot as plt
37 import numpy as np
38 from sklearn.metrics import r2_score
39
40 # Define function to plot predicted vs. actual prices
41 def plot_predicted_vs_actual(y_test, y_pred, title, r2_score):
42     plt.figure(figsize=(8, 6))
43     plt.scatter(y_test, y_pred, alpha=0.5)
44     plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--',
45             ↪ lw=2)
46     plt.xlabel('Actual Prices')
47     plt.ylabel('Predicted Prices')
48     plt.title(title)
49     plt.text(0.05, 0.9, f'R-squared: {r2_score:.2f}',
50             ↪ transform=plt.gca().transAxes)
51     plt.grid(True)
52     plt.show()
53
54 # Define function to plot coefficients of regularized models
55 def plot_coefficients(coefficients, feature_names, title):
56     plt.figure(figsize=(10, 6))
57     plt.barh(np.arange(len(feature_names)), coefficients,
58             ↪ tick_label=feature_names)
59     plt.xlabel('Coefficient Value')
60     plt.ylabel('Feature')
61     plt.title(title)
62     plt.grid(True)
63     plt.show()
64
65 # Train and evaluate models
66 for name, pipeline in pipelines.items():
67     pipeline.fit(X_train, y_train)
68     y_pred = pipeline.predict(X_test)
69     mse = mean_squared_error(y_test, y_pred)
70     r2 = r2_score(y_test, y_pred)
71     print(f"{name} MSE: {mse:.2f}, R-squared: {r2:.2f}")
72
73     # Plot predicted vs. actual prices
74     plot_predicted_vs_actual(y_test, y_pred, f"Predicted vs. Actual Prices -
75     ↪ {name}", r2)
76
77     # If it's a regularized model, plot coefficients
78     if name != "No Regularization":
79         feature_names = data.feature_names
80         if hasattr(pipeline.named_steps['regressor'], 'coef_'):

```

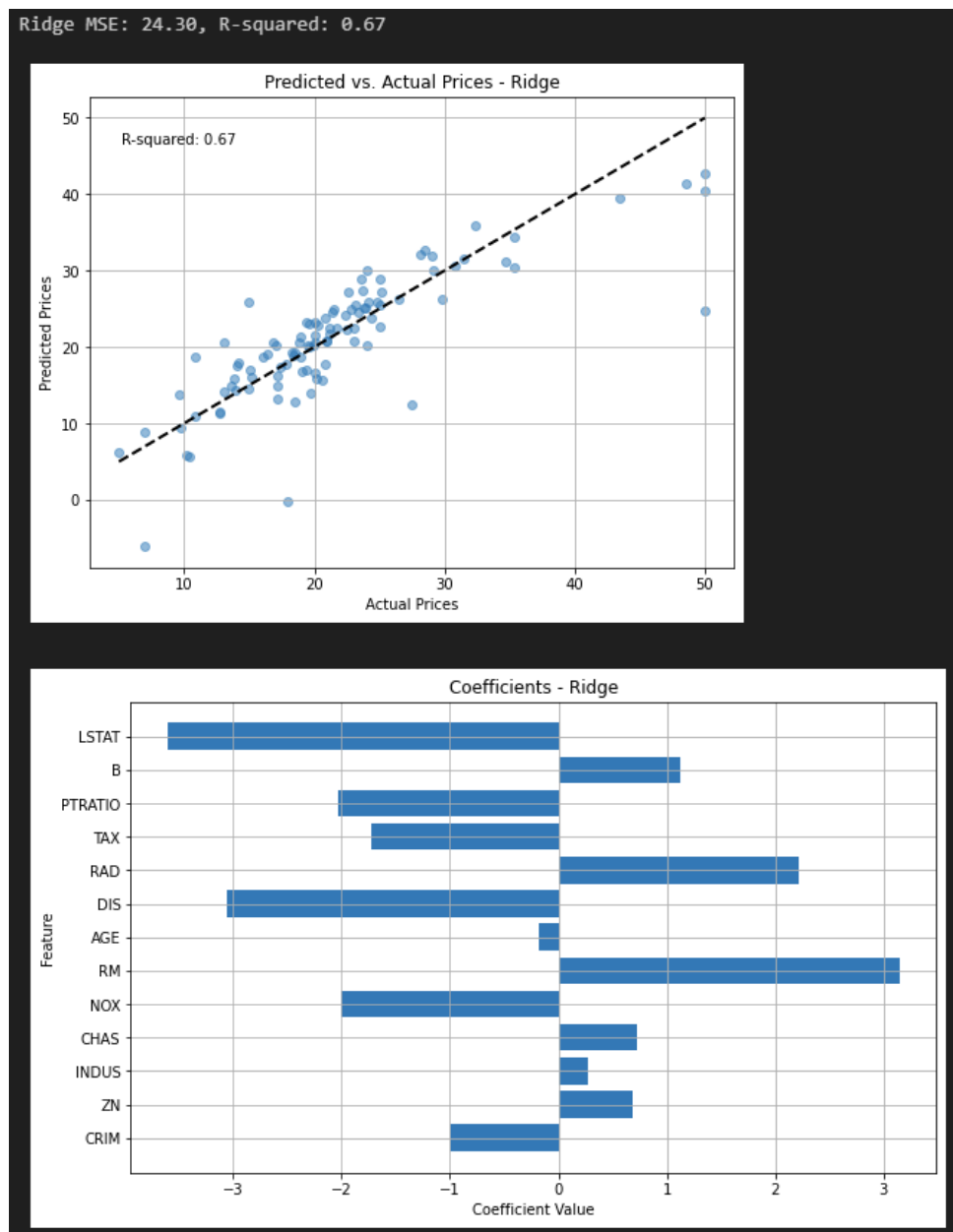
4.5. QE) REGULARIZATION METHODS FOR NEURAL NETWORKS Group 13

```
77         coefficients = pipeline.named_steps['regressor'].coef_  
78     elif hasattr(pipeline.named_steps['regressor'], 'feature_importances_'):  
79         coefficients = pipeline.named_steps['regressor'].feature_importances_  
80     plot_coefficients(coefficients, feature_names, f"Coefficients - {name}")
```


4.5. QE) REGULARIZATION METHODS FOR NEURAL NETWORKS Group 13

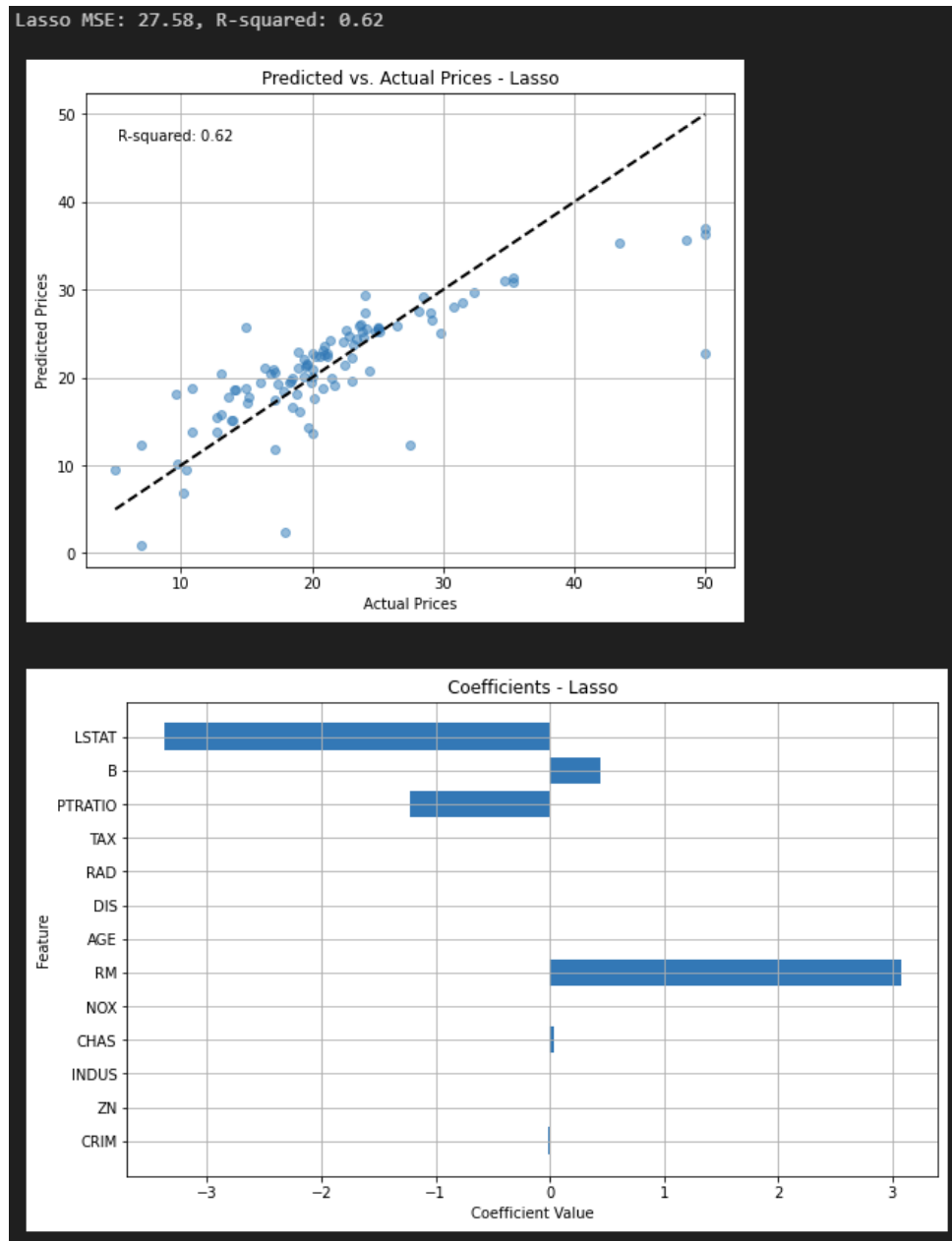
Which outputs different plots with R^2 and MSE values, as well as a graphic representation of the values of the different coefficients:

4.5.1 Ridge model



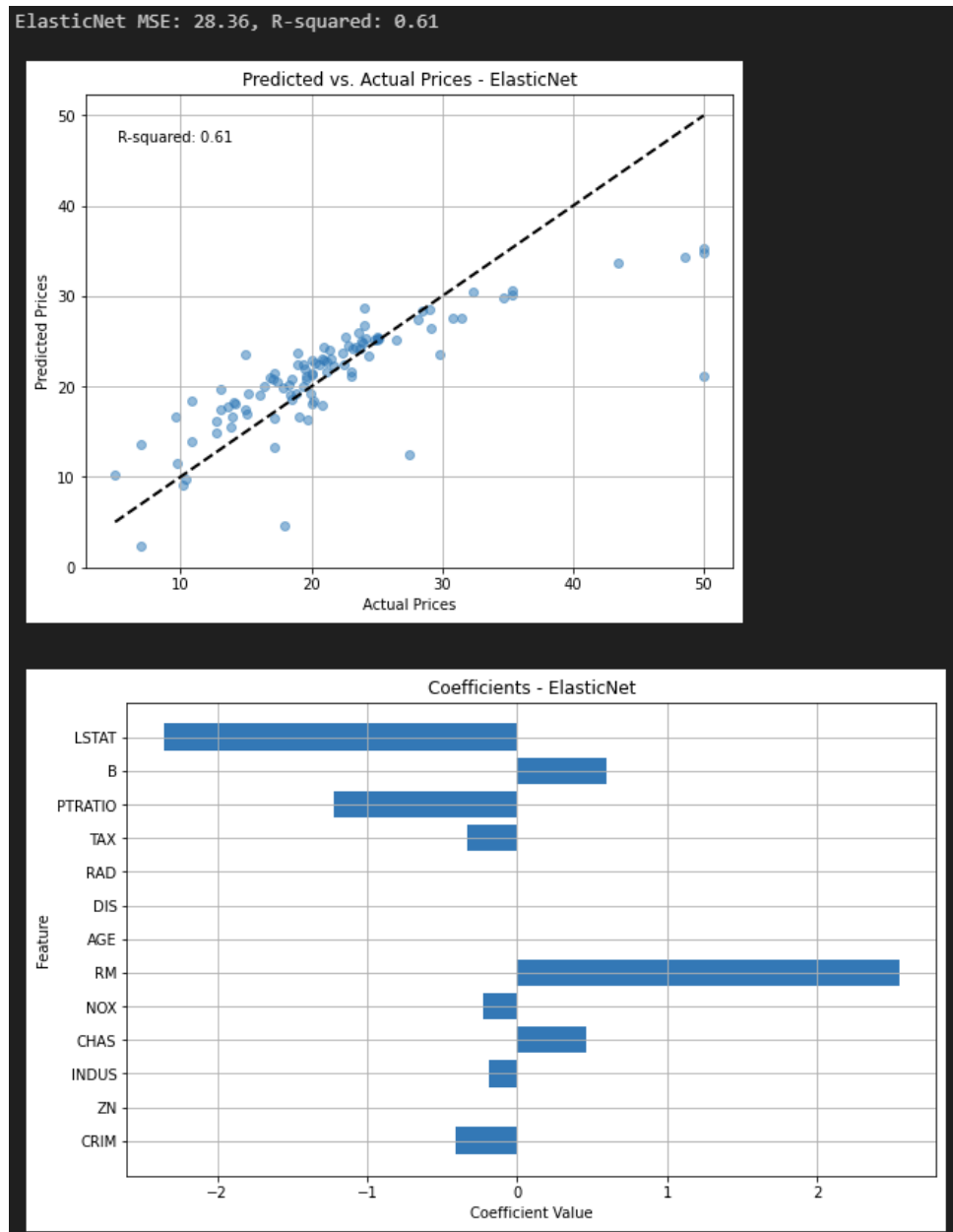
Figur 4.1: Ridge model plots

4.5.2 Lasso model



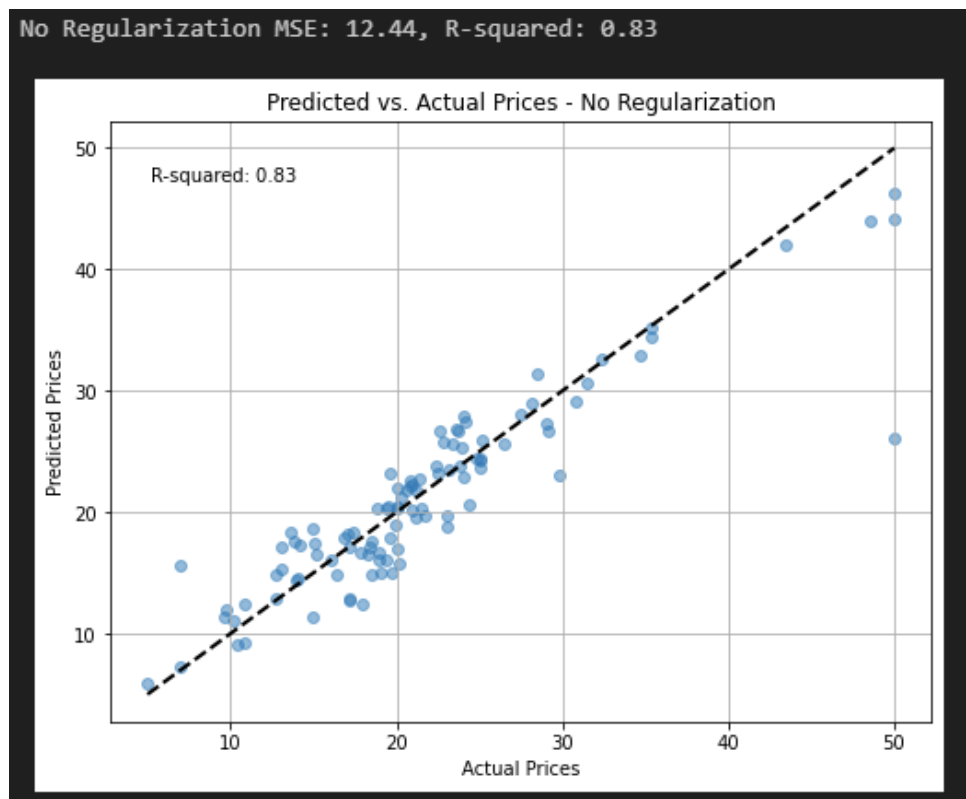
Figur 4.2: Lasso model plots

4.5.3 ElasticNet



Figur 4.3: ElasticNet model plots

4.5.4 No regularization, regular MLPRegressor



Figur 4.4: Regular MLP regression

Explanation: We see that with the same value of $\alpha = 0.5$, the models vary differently - mostly Ridge when comparing the regularization methods. This does make sense, since it punishes the least, compared to Lasso and ElasticNet (also apparent in its lower MSE value). The regular MLPRegressor also shows the highest R^2 value, which it also should, since there is no punishment included in that model.