

Stat 666 Term Project: Relative Predictive Performance for Multivariate Regression and Artificial Neural Networks

Josh Christensen

April 17, 2023

1 Introduction

In this project we seek to understand and compare the robustness of predictions from multivariate regression and neural networks. We explore this phenomenon in two ways. The first is through violation of the assumption of constant variance that is part of the typical multivariate regression model. The second is through the introduction of non-linearity in approximately linear relationships. We perform two simulation studies to explore these ideas and then apply both methods to a real data set of imported cars. This data set exhibits changes in the covariance structure of the response at different levels of the covariates.

We investigate this research question for two main reasons. The first is that linear models form the backbone of many statistical models and are often the default against which other models are compared. This makes their robustness to assumption violations an interesting topic since the use of these models is often the first choice for analysts due to their simplicity, computational efficiency and interpretability.

The second reason for this project relates to neural networks. We use neural networks as the comparison case in this project because they have become commonplace in the machine learning community and are frequently touted as improved alternatives to the linear model due to their ability to automatically incorporate non-linear effects, interactions and other complexities that must be explicitly included in a linear model. While practitioners recognize that no model is perfect, neural networks are sometimes presented as a panacea for problems where prediction is the primary concern. Due to the unique positions that linear regression and neural networks occupy in the landscape of quantitative predictions we hope this project will be able to explore some of the merits of each method and provide insight on which to choose for a given analysis.

We acknowledge that a full comparison between these two models would have to be far more detailed than the prediction comparisons we share here. Indeed the interpretability and uncertainty quantification that come with multivariate regression are often reason enough to forego the black box of a neural network. On the other hand, many forms of unstructured data would effectively preclude the use of linear regression. Since these considerations are on a qualitative level we will not address them directly in this project. Instead, we will focus on providing quantitative comparisons of predictive performance to inform analysts in cases where point prediction is the primary goal.

The rest of the project proceeds as follows. First, we include an introduction to the data as well as brief explanations of multivariate regression and simple neural networks. We then describe and present results from our two simulation studies. Next we apply each method to our chosen data set. Finally, we will provide a brief discussion on the practical insights gained from the project and conclude.

1.1 Data Introduction

The data set to which we will apply the two methods comes from the University of California Irvine (UCI) Machine Learning Repository (Dua and Graff, 2019). It contains 26 variables related to the style and function of 205 different cars produced around the world. We eliminated a variable related to insurance company losses for each car due to high levels of missing values (41 of 205). We then removed all rows with missing values for any variable which left us with 203 observations of 25 variables. We used city miles per gallon (MPG), highway MPG and horsepower as our trivariate response vector. We are interested in finding

relationships to these response variables since they are among the most important functional measures of a car's performance.

2 Multivariate Multiple Regression

Multivariate multiple linear regression is a generalized case of the more common multiple linear regression. In multiple linear regression we have one response variable and multiple explanatory variables. In the case of multivariate regression we have multiple response variables which we wish to predict using one or more explanatory variables. Multivariate regression includes a host of related methods for quantifying uncertainty in parameter estimates and predictions. We will focus only on how to use multivariate regression to produce point estimates for the prediction of new data points.

The most common univariate multiple regression model can be described succinctly as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where \mathbf{X} represents the matrix of covariates, \mathbf{y} represents the response vector, $\boldsymbol{\beta}$ represents the regression coefficients and $\boldsymbol{\epsilon} \sim \mathcal{N}_n(0, \sigma^2 I_n)$. To predict new values of the response from this model we need only the values for our explanatory variables in the \mathbf{X} matrix and an estimate for $\boldsymbol{\beta}$ since

$$\begin{aligned} E(y_i) &= E(\mathbf{x}_i\boldsymbol{\beta} + \epsilon_i) \\ &= \mathbf{x}_i\boldsymbol{\beta} + E(\epsilon_i) \\ &= \mathbf{x}_i\boldsymbol{\beta}. \end{aligned}$$

The most common estimator used for $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

which can be obtained via minimization of the squared error loss function (generally) or via maximum likelihood estimation (using the assumption of normality). The Gauss-Markov theorem indicates that for this model $\hat{\boldsymbol{\beta}}$ is the linear unbiased estimator with minimum variance. Derivations of this estimator can be found in a majority of textbooks dealing with introductory statistics and linear models, and we recommend Rencher and Schaalje (2008) to the interested reader. Given the results above we then make predictions for a new observation using $\hat{\mathbf{y}} = \mathbf{x}\hat{\boldsymbol{\beta}}$.

The estimator for the multivariate regression is closely related to $\hat{\boldsymbol{\beta}}$. For the multivariate regression model we have

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \boldsymbol{\Xi}$$

where \mathbf{Y} represents the matrix of responses, \mathbf{X} represents the matrix of covariates, \mathbf{B} represents the matrix of regression coefficients and $\boldsymbol{\Xi}$ is the matrix of error terms. We assume that the error terms for any given response vector (values of all response variables for a single observation) are uncorrelated with all others, i.e., $\text{cov}(\mathbf{y}_i, \mathbf{y}_j) = \mathbf{0}$. We also assume a common covariance matrix for the response variables across observations, meaning that for all $i \in \{1, \dots, n\}$ we have $\text{cov}(\mathbf{y}_i, \mathbf{y}_i) = \boldsymbol{\Sigma}$. Finally we once again assume that the errors have mean zero. This means that we once again only need covariates and an estimate for \mathbf{B} in order to make predictions. We will use the least-squares estimator

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

for our analysis. This estimator has many properties closely related to those of $\hat{\boldsymbol{\beta}}$ but their details are more complex. For details on why this estimator is the default choice, how it is derived and some its properties we recommend chapter 10 of Rencher and Christensen (2012). Once we have this estimator we generate predictions using $\hat{\mathbf{y}} = \mathbf{x}\hat{\mathbf{B}}$.

3 Neural Networks

We now turn our attention to neural networks. Neural networks have received a great deal of attention from quantitative researchers in recent years and consequently have seen a host of innovations in their estimation,

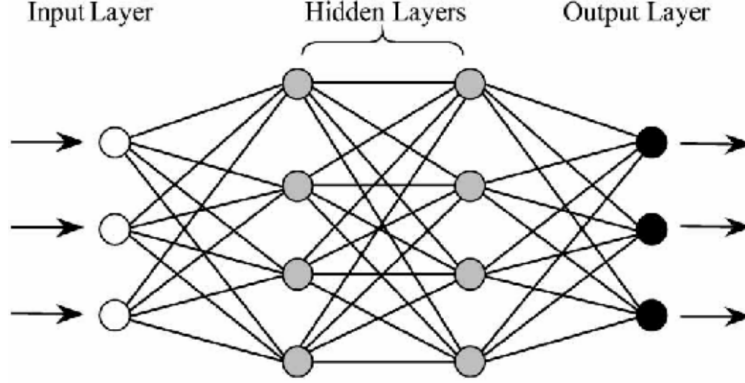


Figure 1: Modified figure 3.7 from Lesschen et al. (2005).

architecture and applications. We will focus only on feedforward neural networks for simplicity and because for simple regression problems more complex architectures are usually not needed. Feedforward neural networks are composed of layers of nodes, each of which passes information to nodes in the subsequent layer. Each neural network begins with an input layer composed of the information we are using to predict the outcome, and ends with an output layer that consists of the networks predictions of the outcomes. A basic architecture can be seen in fig. 1.

Each hidden node in a neural network can be thought of as its own multiple regression in which the explanatory variables are the outputs from preceding nodes and the response is the information passed to future nodes. The output is then modified by a secondary function, the activation function, which scales the output in some non-linear way before relaying it to future nodes. In our case we will use the rectified linear activation function which is a function that moves all negative values to zero and does not change positive values. The activation function will allow us to capture non-linear effects in our response. The scaling values for the input values of each node are called weights and the constant added to the scaled inputs is called the bias.

In supervised regression settings like the ones considered in this project neural networks are trained by evaluating predicted output from the network against known values from a training data set. A loss function is used to quantify how far off the predictions are from the true values. The loss function is then minimized via an optimization routine. In most cases this consists of backpropagation and some form of stochastic gradient descent. Predictions are obtained by giving the network new input values and then feeding those values through the network using the weights and biases found through training to produce predictions from the output nodes.

4 Simulation Studies

The first simulation study was designed to test the robustness of the two models to heteroskedasticity in a multivariate response. In each iteration of each run of the simulation we created a data set with ten covariates and three response variables. The responses consisted of two subpopulations with different correlation structures that were treated as a single population. Marginal variances were kept constant across the two subpopulations and across all runs. The covariate values were randomly generated from a standard normal distributions. The beta coefficients were $\beta_1 = (1, 2, \dots, 10)$, $\beta_2 = (.1, .2, \dots, 1)$ and $\beta_3 = (-.5, -.4, \dots, 0, .1, \dots, .4)$ which correspond to the three variables in the trivariate response respectively. The marginal variances for the three response variables were 1, 100 and 25 respectively. For each run we then selected two different sets of correlations to correspond to the two subpopulations. The first run was a control run where correlations for both subpopulations were set to 0, which generated homoskedastic data sets with no covariance between responses. For the second run, one subpopulation had no correlation while the other had $\rho = (0.9, -0.1, 0.3)$ which represents the correlation between variables one and two, one and three and two and three respectively. For the last run the two subpopulations had correlation structures of

$\boldsymbol{\rho}_1 = (0.8, -0.1, -0.5)$ and $\boldsymbol{\rho}_2 = (-0.8, 0.1, 0.5) = -\boldsymbol{\rho}_1$. For each parameter combination 100 data sets of 100 points were generated. For each data set of 100 points we then fit both models to 80 data points and evaluated mean absolute error on the remaining 20 data points. The results from this simulation study are shown in table 1. Note that the two methods performed comparably well in all three runs of simulation one.

	Run 1	Run 2	Run 3
Multivariate Regression	4.282	4.275	4.271
Neural Network	4.256	4.256	4.249

Table 1: Mean absolute error results from the first simulation study.

Our original hypothesis was that the multivariate regression predictions would be biased by the different behavior in the distinct subpopulations. However, the point predictions appear to be robust to violation of this assumption, provided that the underlying mean function is linear in nature. For the neural networks in this simulation we used one hidden layer with two nodes. This neural network very close to a linear regression due to the small number of hidden nodes, and was chosen from many different architectures by evaluating performance on training and validation sets from the control example. This makes sense because the underlying function was linear and the fewer hidden nodes and hidden layers we have the closer the system of the network is to being linear.

Our second simulation introduced nonlinearity and sought to explore the competing ideas of a multivariate regression being unable to fit obvious nonlinear trends, while a neural network may be prone to overfitting noise. In this simulation we had one covariate, whose values were uniformly generated between 0 and 10. Each data set once again contained 100 data points and we generated 100 data sets per simulation. We once again fit 80 points and used 20 points as a test set to evaluate mean absolute error. Every β coefficient was equal to 2 to create a clear trend and a term equal to $\sin(8x/5)$ was added to the mean function to create a cyclical wave around the trend. We then added Gaussian noise to the 4-variate response with a constant marginal variance of σ^2 and constant correlation equal to ρ . An example of one dimension of a generated response can be seen in fig. 2. A table of the average mean absolute errors from the 100 runs at

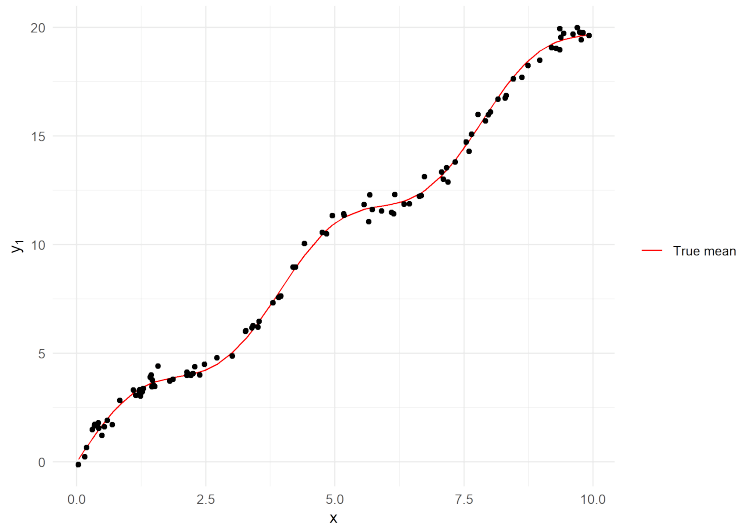


Figure 2: Demonstration of the data generating process for simulation two.

each parameter combination can be found in table 2.

First we note that it appears that neither method was affected by the correlation between the response variables. This is once again consistent with our observation from simulation one that the form the underlying mean function is far more important to the performance of our model predictions than the correlation structure. We also observe that at low observation variances where the nonlinear trends can be easily picked out, the neural network was able to outperform linear regression by fitting the sinusoidal function. However,

	σ^2							
	Multivariate Regression				Neural Network			
	0.1	1	3	8	0.1	1	3	8
$\rho = 0.1$	0.654	0.992	1.499	2.372	0.554	0.964	1.500	2.387
$\rho = 0.3$	0.658	0.985	1.501	2.333	0.545	0.948	1.505	2.366
$\rho = 0.5$	0.654	0.994	1.501	2.362	0.563	0.965	1.500	2.405
$\rho = 0.7$	0.658	1.009	1.535	2.410	0.545	0.964	1.547	2.428
$\rho = 0.9$	0.645	0.979	1.570	2.310	0.548	0.946	1.566	2.360

Table 2: Mean absolute error results from the second simulation study.

as noise increased and the underlying linear trend became the only distinguishable feature of the function, multivariate regression performed as well or better than the neural networks. At even higher levels of noise we may start to see the linear regression outperform the neural networks if we do not carefully retune the architecture and training of the neural network to avoid overfitting. For this simulation we used two hidden layers with 20 nodes each.

5 Results from the Car Data

In this section we briefly summarize the results of our application of both methods to a real data set. This exercise revealed many practical considerations relevant to the two methods. Originally we wished to compare predictions using a train data set of consisting of 80% of the data and a test set of 20% of the data. We decided on a neural network architecture of two hidden layers with 40 nodes each after comparing training and validation set performance for many different architectures. This highlights one of the practical considerations of the neural network: there are many tuning parameters and if tuning is not performed predictions may be unreliable. In addition to the architecture of the network one must also decide on an optimization algorithm, learning rate, number of training passes, etc. While these decisions can be informed by quantitative metrics many of them are ultimately left to the discretion of the modeler. We obtained a mean absolute error of 3.86 across all response values in the test set.

We then attempted to apply the same methodology to the multivariate regression, but encountered a problem where some factor values were present in the test set, but not the training set. Predictions from multivariate regression require estimates for all relevant beta values. With no data for some values we could not predict. Where the neural net could simply train nodes not to fire if no information was given, the multivariate regression failed to produce results. This could be a blessing and a curse as we can still produce predictions, but those predictions may not represent the information we think they do.

We ultimately changed our train/test split to 90% and 10% for the multivariate regression and sampled the training set such that all levels of the factors were present in the training set. We then obtained a mean absolute error of 35.57. Notice that this is an order of magnitude larger than the mean absolute error of the neural network. Based on our simulation results we expect this is a result of nonlinear relationships among the variables being detected by the neural network, but overlooked by the linear model. With the large number of variables this is a very likely explanation. Based on the simulation studies the changing correlation structure across levels of the covariates for which the data set was originally selected is likely less of a factor in the predictive difference observed.

6 Discussion

Based on the quantitative results from our simulations and our cars data set we now provide some interpretation of the results and some practical suggestions.

In comparing the predictive power of neural networks to multivariate regression for multivariate responses we discovered that the correlation structure of the response variables played little role in the models ability to produce accurate point predictions. Indeed it appears that a choice by a modeler to completely ignore underlying correlation among responses will have little to no effect on point estimate accuracy.

Since the multivariate regression is robust to violations of this assumption it is an excellent choice for any predictive problem in which it is reasonable to assume that the relationship between explanatory and response variables is linear in nature. This would be easy to investigate in smaller dimensional systems, but becomes difficult once we reach the number of variables we see in the cars data. If this assumption is reasonable one can greatly reduce the programming, computation and modeling time required compared to the neural network for predictions of similar accuracy.

On the other hand, if there is doubt as to the linearity of the underlying system and it would be difficult to produce informed hypotheses about the true nature of the underlying relationships we suggest the use of a neural network. Neural networks are much more involved to program, take more time to estimate and take time and care to tune. However, their extreme flexibility makes them an appealing choice in cases where it would be difficult or impossible to investigate the relationships between all of the variables at play. This is evident by the neural network's outperformance on the cars data where we had a much higher dimensional problem than our simulations and where the underlying relationships may not have been even approximately linear.

7 Conclusion

In this project we compared predictive capabilities of multivariate regression and neural networks for scenarios that included non-homogeneous covariance structures in the response variables and differing levels of apparent non-linearity in a multivariate response. We found that point predictions from both methods are robust to heteroskedasticity. The neural network outperforms the linear model when non-linear trends are clear, but as noise increases the neural network's advantage disappears and the predictions from the two methods become comparable.

In future work we hope to expand our second simulation study to a wider variety of non-linear settings and observe whether extremely high levels of noise cause the neural network to overfit the data. Also of future interest is a sensitivity analysis for the various tuning parameters of the neural network, including learning rate, optimization routine, number of hidden layers, number of hidden nodes, etc. The selection of these values is non-trivial and can result in drastic differences in a network's predictive performance.

References

- Dua, D. and Graff, C. (2019). UCI machine learning repository.
- Lesschen, J. P., Verburg, P., and Staal, S. (2005). Statistical methods for analysing the spatial dimension of changes in land use and farming systems.
- Rencher, A. C. and Christensen, W. F. (2012). *Methods of multivariate analysis*. Wiley, Hoboken, New Jersey, third edition. edition.
- Rencher, A. C. and Schaalje, G. B. (2008). *Linear models in statistics*. John Wiley & Sons.

Appendix: Analysis Code

Simulation 1 R code (multivariate regression)

```
1 # Simulate data frame
2 run_simulation <- function(n_sim = 100, sigmas) {
3   # Extract sigmas
4   sigma_1 <- sigmas[[1]]
5   sigma_2 <- sigmas[[2]]
6
7   # Set sample size for each simulated dataset
8   n <- 100
9
10  # set the betas for the simulation
11  beta <- matrix(c(1:10,
12                  seq(0.1,1,by = 0.1),
13                  seq(-0.5, 0.4, by = 0.1)),
14                nrow = 10, ncol = 3)
15
16  # Create storage
17  linear_model_mae <- numeric(n_sim)
18
19  # Storage for the datasets so they can be fit with the neural net later
20  generated_datasets <- vector("list", n_sim)
21
22  for(i in 1:n_sim) {
23    # Generate the covariate matrices
24    x_1 <- matrix(rnorm(10 * n),
25                  nrow = n, ncol = 10)
26    x_2 <- matrix(rnorm(10 * n),
27                  nrow = n, ncol = 10)
28
29    # Generate the first dataset
30    y_1 <- x_1 %*% beta + t(replicate(n, crossprod(chol(sigma_1), rnorm(3)), simplify = "
31      matrix"))
32
33    # Generate the second dataset
34    y_2 <- x_2 %*% beta + t(replicate(n, crossprod(chol(sigma_2), rnorm(3)), simplify = "
35      matrix"))
36
37    # Aggregate
38    y <- rbind(y_1, y_2)
39    x <- rbind(x_1, x_2)
40
41    # Split train and test
42    test_sample_size <- round(n * .2)
43    test_index <- sample.int(n * 2, test_sample_size)
44    y_train <- y[-test_index,]
45    y_test <- y[test_index,]
46    x_train <- x[-test_index,]
47    x_test <- x[test_index,]
48
49    # Fit the linear model
50    betahat <- coef(lm(y_train ~ x_train))
51
52    # Evaluate MAE for predictions
53    linear_model_mae[i] <- mean(abs(y_test - cbind(1,x_test) %*% betahat))
54
55    # Store the dataset
56    generated_datasets[[i]] <- list(x_train = x_train,
57                                   y_train = y_train,
58                                   x_test = x_test,
59                                   y_test = y_test)
60  }
61  return(
62    list(
```



```

62     linear_model_mae = linear_model_mae,
63     generated_datasets = generated_datasets
64 )
65 )
66 }
67
68 make_sigmas <- function(cor1, cor2){
69     sigma_1 <- matrix(c(1,cor1[1],cor1[2],
70                        cor1[1],1,cor1[3],
71                        cor1[2],cor1[3],1), nrow = 3, ncol = 3, byrow = TRUE)
72
73     sigma_2 <- matrix(c(1,cor2[1],cor2[2],
74                        cor2[1],1,cor2[3],
75                        cor2[2],cor2[3],1), nrow = 3, ncol = 3, byrow = TRUE)
76
77     stdevs <- c(1,10,5)
78
79     scaling_matrix <- matrix(stdevs, nrow = 3, ncol = 3)
80     scaling_matrix <- scaling_matrix * t(scaling_matrix)
81     scaling_matrix
82
83     sigma_1 <- sigma_1 * scaling_matrix
84     sigma_2 <- sigma_2 * scaling_matrix
85
86     return(
87         list(
88             sigma_1,
89             sigma_2
90         )
91     )
92 }
93
94 control_results <- run_simulation(sigmars = make_sigmas(rep(0,3), rep(0,3)))
95 sim_results_ind_mixed <- run_simulation(sigmars = make_sigmas(rep(0,3), c(.9,-.1,.3)))
96 sim_results_swap <- run_simulation(sigmars = make_sigmas(c(.8,-.1,-.5), c(-.8,.1,.5)))
97
98 c(
99     mean(control_results$linear_model_mae),
100     mean(sim_results_ind_mixed$linear_model_mae),
101     mean(sim_results_swap$linear_model_mae)
102 )

```

Simulation 1 Python code (neural networks)

```

1 # Re run the same simulations that we ran in R with a neural net
2 # Import packages
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6 import tensorflow as tf
7 from tensorflow import keras
8 from tensorflow.keras import layers
9
10 # Make NumPy printouts easier to read.
11 np.set_printoptions(precision = 3, suppress = True)
12
13 print(tf.__version__)
14
15 control_sets = r.control_results['generated_datasets']
16 mixed_sets = r.sim_results_ind_mixed['generated_datasets']
17 swap_sets = r.sim_results_swap['generated_datasets']
18
19 # This is where I could begin a function that would take first_set
20 def fit_nn(generated_set):
21     normalizer = tf.keras.layers.Normalization(axis=-1)
22

```

```

23 normalizer.adapt(np.array(generated_set['x_train']))
24
25 model = keras.Sequential([
26     normalizer,
27     layers.Dense(2, activation='relu'),
28     # layers.Dense(2, activation='relu'),
29     layers.Dense(3)
30 ])
31
32 model.compile(loss='mean_absolute_error',
33               optimizer=tf.keras.optimizers.Adam(0.01))
34
35 dnn_model = model
36
37 history = dnn_model.fit(
38     generated_set['x_train'],
39     generated_set['y_train'],
40     # validation_split=0.2,
41     verbose=0, epochs = 300)
42
43 return dnn_model.evaluate(generated_set['x_test'], generated_set['y_test'], verbose=0)
44
45 fit_nn(control_sets[0])
46
47 def nn_sim(dataset_list):
48     n_sim = len(dataset_list)
49     nn_mae = [None] * n_sim
50     for i in range(n_sim):
51         nn_mae[i] = fit_nn(dataset_list[i])
52
53     return nn_mae
54
55 control_results = nn_sim(control_sets)
56 mixed_results = nn_sim(mixed_sets)
57 swap_results = nn_sim(swap_sets)
58
59 print([np.mean(control_results), np.mean(mixed_results), np.mean(swap_results)])

```

Simulation 2 R code (multivariate regression)

```

1 # Simulate data frame
2 run_simulation <- function(n_sim = 100, rho = 0.5, sigma2 = 1) {
3     # Create sigma matrix
4     sigma <- (matrix(rho, nrow = 4, ncol = 4) + diag(1 - rho, 4)) * sigma2
5
6     # Set sample size for each simulated dataset
7     n <- 100
8
9     # set the betas for the simulation
10    beta <- matrix(1,
11                  nrow = 1, ncol = 4)
12
13    # Create storage
14    linear_model_mae <- numeric(n_sim)
15
16    # Storage for the datasets so they can be fit with the neural net later
17    generated_datasets <- vector("list", n_sim)
18
19    for(i in 1:n_sim) {
20        # Generate the covariate vector
21        x <- runif(n, 0, 10)
22
23        # Generate the y sequences
24        f <- function(x) x + sin(x * 1.6)
25        cyclical_component <- matrix(rep(f(x), 4), nrow = n, ncol = 4)

```

```

26 y <- x %%% beta + cyclical_component + t(replicate(n, crossprod(chol(sigma), rnorm(4)),
simplify = "matrix"))
27
28 # Split train and test
29 test_sample_size <- round(n * .2)
30 test_index <- sample.int(n, test_sample_size)
31 y_train <- y[-test_index,]
32 y_test <- y[test_index,]
33 x_train <- x[-test_index]
34 x_test <- x[test_index]
35
36 # Fit the linear model
37 betahat <- coef(lm(y_train ~ x_train))
38
39 # Evaluate MAE for predictions
40 linear_model_mae[i] <- mean(abs(y_test - cbind(1,x_test) %%% betahat))
41
42 # Store the dataset
43 generated_datasets[[i]] <- list(x_train = x_train,
44                                y_train = y_train,
45                                x_test = x_test,
46                                y_test = y_test)
47 }
48
49 return(
50   list(
51     linear_model_mae = linear_model_mae,
52     generated_datasets = generated_datasets
53   )
54 )
55 }
56
57 sims <- list(
58   sim_1 = run_simulation(rho = 0),
59   sim_2 = run_simulation(rho = .3),
60   sim_3 = run_simulation(rho = .5),
61   sim_4 = run_simulation(rho = .7),
62   sim_5 = run_simulation(rho = .9),
63   sim_6 = run_simulation(rho = 0, sigma2 = 3),
64   sim_7 = run_simulation(rho = .3, sigma2 = 3),
65   sim_8 = run_simulation(rho = .5, sigma2 = 3),
66   sim_9 = run_simulation(rho = .7, sigma2 = 3),
67   sim_10 = run_simulation(rho = .9, sigma2 = 3),
68   sim_11 = run_simulation(rho = 0, sigma2 = 0.1),
69   sim_12 = run_simulation(rho = .3, sigma2 = 0.1),
70   sim_13 = run_simulation(rho = .5, sigma2 = 0.1),
71   sim_14 = run_simulation(rho = .7, sigma2 = 0.1),
72   sim_15 = run_simulation(rho = .9, sigma2 = 0.1),
73   sim_16 = run_simulation(rho = 0, sigma2 = 8),
74   sim_17 = run_simulation(rho = .3, sigma2 = 8),
75   sim_18 = run_simulation(rho = .5, sigma2 = 8),
76   sim_19 = run_simulation(rho = .7, sigma2 = 8),
77   sim_20 = run_simulation(rho = .9, sigma2 = 8)
78 )
79
80 unlist(lapply(lapply(sims, \(x) x$linear_model_mae), mean))

```

Simulation 2 Python code (neural networks)

```

1 # Re run the same simulations that we ran in R with a neural net
2 # Import packages
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6 import tensorflow as tf
7 from tensorflow import keras

```

```

8 from tensorflow.keras import layers
9
10 # Make NumPy printouts easier to read.
11 np.set_printoptions(precision = 3, suppress = True)
12
13 print(tf.__version__)
14
15 sims = r.sims
16
17 # This is where I could begin a function that would take first_set
18 def fit_nn(generated_set):
19     normalizer = tf.keras.layers.Normalization(input_shape=[1,], axis=None)
20
21     normalizer.adapt(np.array(generated_set['x_train']))
22
23     model = keras.Sequential([
24         normalizer,
25         layers.Dense(20, activation='relu'),
26         layers.Dense(20, activation='relu'),
27         layers.Dense(4)
28     ])
29
30     model.compile(loss='mean_absolute_error',
31                  optimizer=tf.keras.optimizers.Adam(0.01))
32
33     dnn_model = model
34
35     history = dnn_model.fit(
36         np.array(r.x), # generated_set['x_train']
37         np.array(r.y), # generated_set['y_train']
38         validation_split=0.2,
39         verbose = 0, epochs = 100)
40
41     return dnn_model.evaluate(np.array(r.x), np.array(r.y), verbose=0)
42
43 fit_nn(sims['sim_1']['generated_datasets'][0])
44
45 def nn_sim(dataset_list):
46     n_sim = len(dataset_list)
47     nn_mae = [None] * n_sim
48     for i in range(n_sim):
49         nn_mae[i] = fit_nn(dataset_list[i])
50
51     return nn_mae
52
53 results_1 = nn_sim(sims['sim_1']['generated_datasets'])
54 results_2 = nn_sim(sims['sim_2']['generated_datasets'])
55 results_3 = nn_sim(sims['sim_3']['generated_datasets'])
56 results_4 = nn_sim(sims['sim_4']['generated_datasets'])
57 results_5 = nn_sim(sims['sim_5']['generated_datasets'])
58 results_6 = nn_sim(sims['sim_6']['generated_datasets'])
59 results_7 = nn_sim(sims['sim_7']['generated_datasets'])
60 results_8 = nn_sim(sims['sim_8']['generated_datasets'])
61 results_9 = nn_sim(sims['sim_9']['generated_datasets'])
62 results_10 = nn_sim(sims['sim_10']['generated_datasets'])
63 results_11 = nn_sim(sims['sim_11']['generated_datasets'])
64 results_12 = nn_sim(sims['sim_12']['generated_datasets'])
65 results_13 = nn_sim(sims['sim_13']['generated_datasets'])
66 results_14 = nn_sim(sims['sim_14']['generated_datasets'])
67 results_15 = nn_sim(sims['sim_15']['generated_datasets'])
68 results_16 = nn_sim(sims['sim_16']['generated_datasets'])
69 results_17 = nn_sim(sims['sim_17']['generated_datasets'])
70 results_18 = nn_sim(sims['sim_18']['generated_datasets'])
71 results_19 = nn_sim(sims['sim_19']['generated_datasets'])
72 results_20 = nn_sim(sims['sim_20']['generated_datasets'])
73
74 print([np.mean(results_1), np.mean(results_2), np.mean(results_3), np.mean(results_4),
75        np.mean(results_5), np.mean(results_6), np.mean(results_7), np.mean(results_8),

```

```

76 np.mean(results_9), np.mean(results_10), np.mean(results_11), np.mean(results_12),
77 np.mean(results_13), np.mean(results_14), np.mean(results_15), np.mean(results_16),
78 np.mean(results_17), np.mean(results_18), np.mean(results_19), np.mean(results_20)])

```

Car Data R code (multivariate regression)

```

1 # Set up the multivariate regression on the real data
2
3 # Load necessary libraries
4 library(dplyr)
5 library(magrittr)
6
7 # Import the data
8 url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data'
9 column_names = c('symboling', 'normalized-losses', 'make', 'fuel-type', 'aspiration',
10                  'num-of-doors', 'body-style', 'drive-wheels', 'engine-location',
11                  'wheel-base', 'length', 'width', 'height', 'curb-weight', 'engine-type',
12                  'num-of-cylinders', 'engine-size', 'fuel-system', 'bore', 'stroke',
13                  'compression-ratio', 'horsepower', 'peak-rpm', 'city-mpg', 'highway-mpg',
14                  'price')
15
16 dataset <- read.csv(url, header = FALSE, col.names = column_names, na.strings = "?")
17
18 dataset <- dataset[,-2]
19 ncol(dataset)
20
21 sum(is.na(dataset))
22
23 dataset <- na.omit(dataset)
24
25 dataset %<>% mutate(num.of.cylinders = case_when(num.of.cylinders == "three" ~ 3,
26                                                num.of.cylinders == "four" ~ 4,
27                                                num.of.cylinders == "five" ~ 5,
28                                                num.of.cylinders == "six" ~ 6,
29                                                num.of.cylinders == "eight" ~ 8,
30                                                num.of.cylinders == "twelve" ~ 12))
31
32 train_index <- sample.int(nrow(dataset), size = round(nrow(dataset) * .5))
33 train <- dataset[train_index,]
34 test <- dataset[-train_index,]
35
36 regression_fit <- lm(cbind(city.mpg, highway.mpg, horsepower) ~ ., data = train)
37
38 predictions <- predict(regression_fit, test)
39
40 mean(abs(as.matrix(predictions[,3] - dataset[c("city.mpg", "highway.mpg", "horsepower")]))))

```

Car Data Python code (neural network)

```

1 # Import packages
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import pandas as pd
5 import tensorflow as tf
6 from tensorflow import keras
7 from tensorflow.keras import layers
8
9 # Make NumPy printouts easier to read.
10 np.set_printoptions(precision = 3, suppress = True)
11
12 print(tf.__version__)
13
14 url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data'
15 column_names = ['symboling', 'normalized-losses', 'make', 'fuel-type', 'aspiration',

```

```

16         'num-of-doors', 'body-style', 'drive-wheels', 'engine-location',
17         'wheel-base', 'length', 'width', 'height', 'curb-weight', 'engine-type',
18         'num-of-cylinders', 'engine-size', 'fuel-system', 'bore', 'stroke',
19         'compression-ratio', 'horsepower', 'peak-rpm', 'city-mpg', 'highway-mpg',
20         'price']
21
22 raw_dataset = pd.read_csv(url, names=column_names,
23                             na_values='?', comment='\t',
24                             sep=',', skipinitialspace=True)
25
26 dataset = raw_dataset.copy()
27 dataset = dataset.drop("normalized-losses", axis = 1)
28 len(dataset.columns)
29
30 dataset.isna().sum()
31
32 dataset = dataset.dropna()
33
34 dataset.loc[:, 'num-of-cylinders'] = dataset['num-of-cylinders'].map({'three': 3, 'four': 4,
35                               'five': 5, 'six': 6, 'eight': 8, 'twelve': 12})
36
37 dataset = pd.get_dummies(dataset, columns=['make', 'fuel-type', 'aspiration', 'num-of-doors',
38     'body-style', 'drive-wheels', 'engine-location', 'engine-type', 'fuel-system'],
39     prefix='', prefix_sep='')
40
41 train_dataset = dataset.sample(frac=0.8, random_state=0)
42 test_dataset = dataset.drop(train_dataset.index)
43
44 train_dataset.describe().transpose()
45
46 train_features = train_dataset.drop(['city-mpg', 'highway-mpg', 'horsepower'], axis = 1)
47 test_features = test_dataset.drop(['city-mpg', 'highway-mpg', 'horsepower'], axis = 1)
48
49 train_labels = train_dataset[['city-mpg', 'highway-mpg', 'horsepower']]
50 test_labels = test_dataset[['city-mpg', 'highway-mpg', 'horsepower']]
51
52 train_dataset.describe().transpose()[['mean', 'std']]
53
54 normalizer = tf.keras.layers.Normalization(axis=-1)
55
56 normalizer.adapt(np.array(train_features))
57
58 print(normalizer.mean.numpy())
59
60 first = np.array(train_features[:1])
61
62 with np.printoptions(precision=2, suppress=True):
63     print('First example:', first)
64     print()
65     print('Normalized:', normalizer(first).numpy())
66
67 def plot_loss(history):
68     plt.figure()
69     plt.plot(history.history['loss'], label='loss')
70     plt.plot(history.history['val_loss'], label='val_loss')
71     plt.ylim([0, 10])
72     plt.xlabel('Epoch')
73     plt.ylabel('Error')
74     plt.legend()
75     plt.grid(True)
76     plt.show()
77     plt.close()
78
79 test_results = {}
80
81 def build_and_compile_model(norm):
82     model = keras.Sequential([
83         norm,

```

```

82     layers.Dense(40, activation='relu'),
83     layers.Dense(40, activation='relu'),
84     layers.Dense(3)
85 ])
86
87 model.compile(loss='mean_absolute_error',
88               optimizer=tf.keras.optimizers.Adam(0.001))
89 return model
90
91 dnn_model = build_and_compile_model(normalizer)
92 dnn_model.summary()
93
94 # %%time
95 history = dnn_model.fit(
96     train_features,
97     train_labels,
98     validation_split=0.2,
99     verbose=0, epochs=1500)
100
101 plot_loss(history)
102
103 test_results['dnn_model'] = dnn_model.evaluate(test_features, test_labels, verbose=0)
104
105 pd.DataFrame(test_results, index=['Mean absolute error']).T
106
107 test_predictions = dnn_model.predict(test_features)
108
109 error = test_predictions - test_labels
110 plt.hist(error, bins=25)
111 plt.xlabel('Prediction Error')
112 _ = plt.ylabel('Count')
113 plt.show()
114 plt.close()
115
116 dnn_model.save('term_project/project_dnn_model')
117
118 reloaded = tf.keras.models.load_model('term_project/project_dnn_model')
119
120 test_results['reloaded'] = reloaded.evaluate(
121     test_features, test_labels, verbose=0)
122
123 pd.DataFrame(test_results, index=['Mean absolute error']).T
124
125 # Now fit to all of the data
126
127 dataset_features = dataset.drop(['city-mpg', 'highway-mpg', 'horsepower'], axis = 1)
128 dataset_labels = dataset[['city-mpg', 'highway-mpg', 'horsepower']]
129
130 full_normalizer = tf.keras.layers.Normalization(axis=-1)
131
132 full_normalizer.adapt(np.array(dataset_features))
133
134 print(full_normalizer.mean.numpy())
135
136 first = np.array(dataset_features[:1])
137
138 with np.printoptions(precision=2, suppress=True):
139     print('First example:', first)
140     print()
141     print('Normalized:', full_normalizer(first).numpy())
142
143 full_dnn_model = build_and_compile_model(full_normalizer)
144 full_dnn_model.summary()
145
146 # %%time
147 history = full_dnn_model.fit(
148     train_features,
149     train_labels,

```

```
150     validation_split = 0.2,
151     verbose = 0, epochs = 1500)
152
153 plot_loss(history)
154
155 test_results['full_dnn_model'] = full_dnn_model.evaluate(dataset_features, dataset_labels,
156                                                         verbose=0)
157
158 test_predictions = full_dnn_model.predict(test_features)
159
160 error = test_predictions - test_labels
161 plt.hist(error, bins=25)
162 plt.xlabel('Prediction Error')
163 _ = plt.ylabel('Count')
164 plt.show()
165 plt.close()
```