

Improved Prediction of Deep Learners with Domain Knowledge Incorporated into the Loss Function

BIS557 Final Project

Chloe Jiang

Introduction

Deep neural networks have been wildly used in the field of machine learning, especially in areas such as image recognition, natural language understanding and game playing. One of the major tasks in constructing a deep learner is to find a way to minimize the loss function. Oftentimes the models are learnt mostly from the data, with minimal use of domain knowledge or processed input features, so as to avoid user bias being introduced into the system. Since domain knowledge often contributes to the selection of network architecture in deep learning, incorporating it into the network may potentially improve the performance of models, even if the training data is sparse or of poor quality (i.e. Missing or noisy data), which is not a rare case with real data.

There have been studies that try to include existing domain knowledge into their frameworks. Radovanovic et al. utilize a stacked generalization approach to integrate domain knowledge into a logistic regression framework and tested it on the binary classification problem of predicting 30 day hospital readmission[1]. J.-X. Wang et al. reconstruct discrepancies in a Reynolds-Averaged Navier-Stokes system using Bayesian network modelings with physics-based priors incorporated[2]. Muralidhar et al. propose a domain adapted neural network based on Bayesian approaches for integrating the domain knowledge as a prior into the loss function[3].

In this project, I plan to answer the following questions:

- (1) How can domain knowledge of a problem be built into the loss function of a deep learner for better prediction results?
- (2) When the data is of poor quality or contains noises, will domain knowledge still improve the performance of a neural network?

To address the above questions, I applied a deep neural network to noisy datasets with domain knowledge incorporated into the loss function, and compared its performance with a regular feedforward NN.

Problem and Data Description

In real life, the data we obtain for analysis may have innate knowledge with various forms. One of the most common form is a quantitative range for a particular response value Y in the system[4]. Another is the monotonic relationship between one of the predictors X and the dependent variable (i.e. $x_1 < x_2 \Rightarrow y_1 < y_2$). Due to these constraints, loss functions need to be adapted to contain such information. The generic loss function of a deep learner can be expressed with the equation:

$$\operatorname{argmin} Loss(Y, \hat{Y}) + \lambda_D Loss_D(\hat{Y}) + \lambda L_2$$

The first term is the mean squared error loss between the true values Y and the predicted results. λ_D is a parameter determining the weight of domain knowledge. $Loss_D(Y)$ represents the domain knowledge being injected into the overall loss function, and λL_2 is the L_2 regularization term for controlling model complexity.

Noisy measurements often cause the data to deviate from target values and fluctuate within a certain range. To simulate this situation, I decided to use the Bohachevsky function to generate synthetic data. Here I chose $k_1 = 1$, $k_2 = 2$, $K = 0.7$, $a_1 = 0.3$, $a_2 = 0.4$, $p_1 = 3$, and $p_2 = 4$. x_1 is a set of random draws from a normal distribution with mean = 3 and standard deviation = 1. Similarly, x_2 is generated from $N(10, 1)$. Noises are imputed by randomly selecting $p\%$ of rows in X and interchange the values of x_1 and x_2 . Based on several experiments, p is set to be 50. Both X and Y have 10,000 observations.

$$f(x_1, x_2) = k_1 x_1^2 + k_2 x_2^2 + a_1 \cos(p_1 \pi x_1) + a_2 \cos(p_2 \pi x_2) + K$$

To simulate monotonicity, I further generated another dataset X2 so that $X2 = 5X$, and Y2 was calculated using X2 and the Bohachevsky function with the same parameters. There exists a monotonic relationship such that $x_{i, 1} < x_{i, 2}$ and $y_{i, 1} < y_{i, 2}$. Then I randomly selected a subset of rows in Y1 and Y2 and switched their values, in order to create noise that violates the monotonicity. The proportion of rows being switched is 50%.

Model Construction and Implementation

1. Incorporating approximation constraints

Since the predicted values of Y are only allowed to vary within a certain range, the loss function needs to include the situations where Y exceeds the upper boundary or falls below the lower boundary. Therefore I considered using $| \hat{Y} - y_l |$ and $| \hat{Y} - y_u |$ to capture these constraints. If the predicted value \hat{Y} is within the range of $[y_u, y_l]$, then no penalty will be added to the loss function. The $Loss_D$ function with the approximation constraints is shown as the following.

$$Loss_D(\hat{Y}) = \sum ReLU(y_l - y_i) + \sum ReLU(y_i - y_u)$$

2. Incorporating monotonicity constraints

Monotonicity are violated if $x_1 < x_2$ and y_1 is not necessarily smaller or larger than y_2 . To enforce monotonicity, penalty will be added to the loss function if the monotonicity condition does not hold. Thus the $Loss_D$ functions for monotonicity constraints should be:

$$Loss_D(\hat{Y}_1, \hat{Y}_2) = \sum I[(x_1^{(1)} < x_1^{(2)}, \hat{y}^{(1)} > \hat{y}^{(2)}) ReLU(\hat{y}^{(1)} - \hat{y}^{(2)})]$$

$$Loss_D(\hat{Y}_1, \hat{Y}_2) = \sum I[(x_1^{(1)} < x_1^{(2)}, \hat{y}^{(1)} < \hat{y}^{(2)}) ReLU(\hat{y}^{(2)} - \hat{y}^{(1)})]$$

The first equation is the loss for increasing functions and the second is for decreasing functions. I is the indicator function that returns 1 if the condition is true and 0 otherwise. $x^{(i)}$ and $y^{(i)}$ are values from the i -th dataset. The above equations capture the situations where monotonicity constraints are violated and therefore add penalty to the loss function.

Scenario 1

To apply these domain knowledge formulas to the loss function of a deep learner, I constructed two feedforward neural networks using the “torch” package. One has the domain knowledge, both approximation and monotonicity, incorporated into its loss function, while the other does not. For the purpose of a fair comparison, each network has two hidden layers with the size of $\text{ncol}(X) \times \text{ncol}(X)$, and one output layer. Adam optimizer is used for both neural networks and the L_2 regularization parameter is set to be 0.01. After several experiments, learning rate used in this project is chosen to be 0.01, and the weight parameter λD for Loss_D is 1. For the approximation constraint, the domain range is set to be $[\mu - \sigma, \mu + \sigma]$, where μ and σ are the mean and standard deviation of Y before corruption. Both neural networks were trained on the same training set, which has a size of 7500 observations, and then tested on the same testing set with a size of 2500. The number of total iterations is 5000. After iterations, in sample and out of sample prediction accuracy were evaluated by the remaining average of loss.

Scenario 2

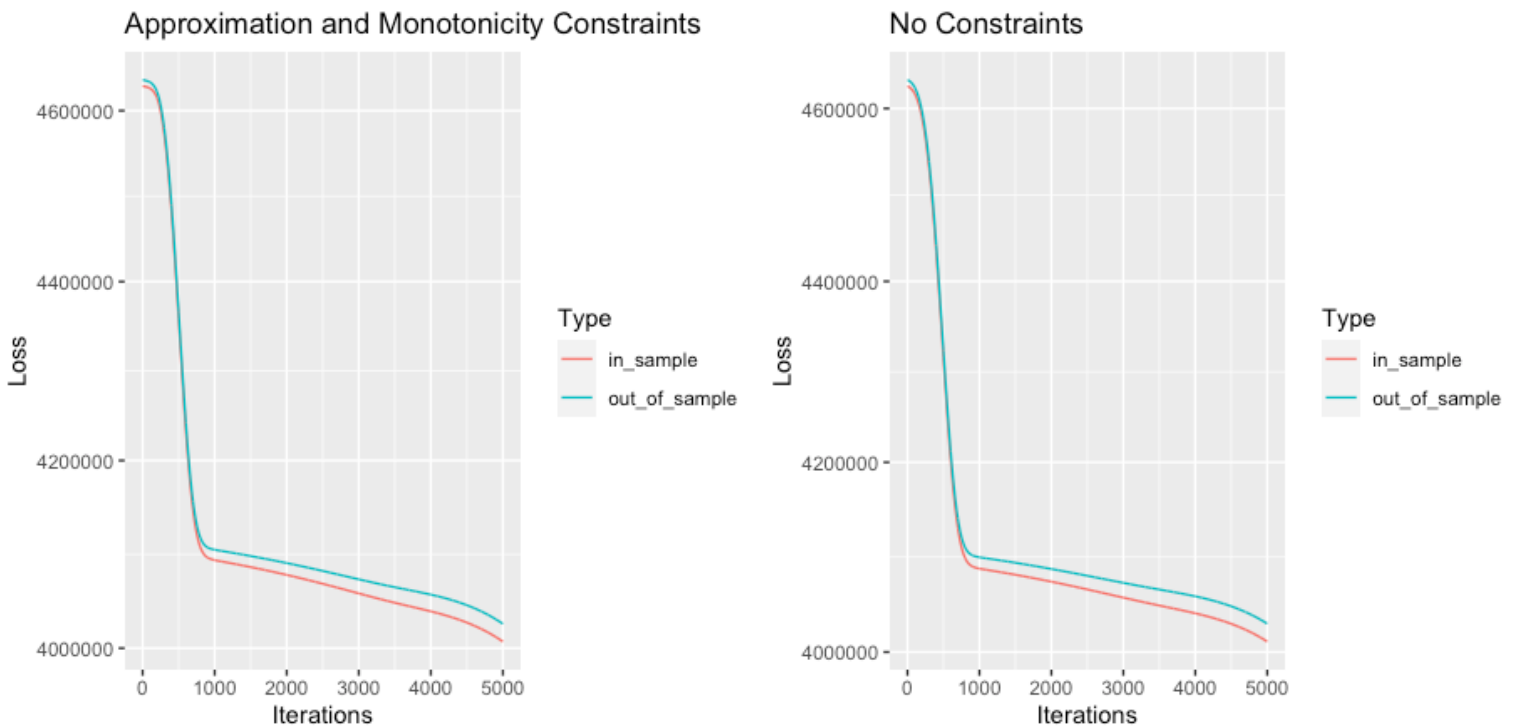
I further explored the scenario where monotonicity does not hold for all the data in X . Assume Y first increases with X_1 , and after a certain threshold it starts to decrease. To integrate such information into the loss function, I created another dataset $X_3 = 15X$, and Y_3 was generated based on X_3 using the Bohachevsky function. Threshold was chosen to be the mean of x_1 , and Y_3 will be calculated by $-10X$ if x_1 exceeds its mean. The loss function is a combination of the equations for monotonicity constraints, as the increasing part will be calculated by the first equation and the Loss_D for the rest of the data will be calculated by the second. To impute noise, I again selected a subset of

rows in Y and switched them with the same rows in Y_3 . Approximation constraints are not considered in this scenario.

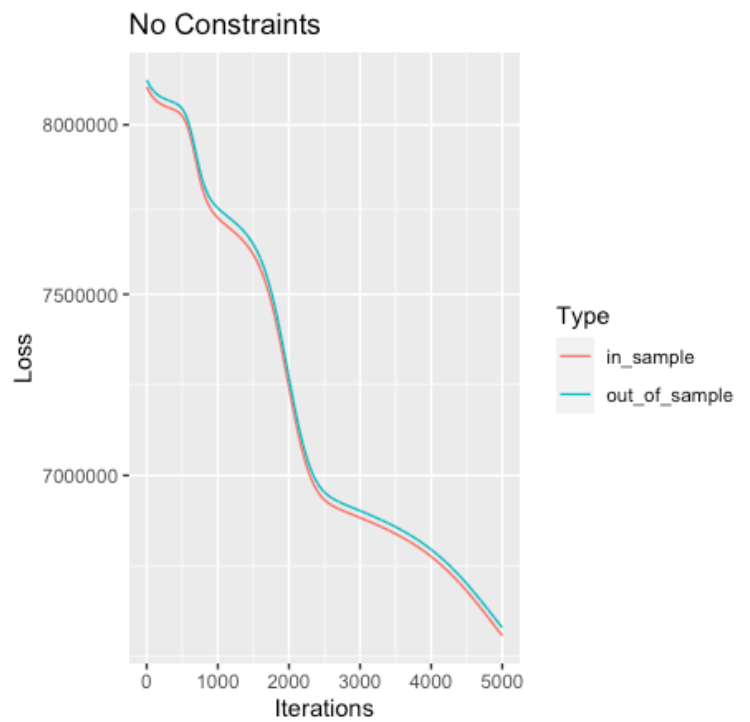
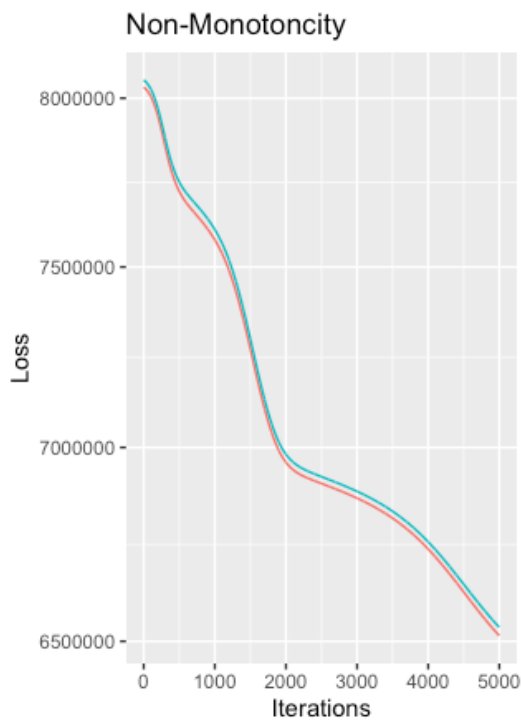
Two deep neural networks were built to compare performance of the two models with or without domain knowledge. Each network has two hidden layers and one output layer with learning rate chosen to be 0.01. L_2 regularization parameter is 0.01 and λD for Loss_D equals 1. Both networks were trained on the same training set with a size of 7500, and then tested on the testing set with 2500 observations. After 5000 iterations, in sample and out of sample prediction accuracy were evaluated by the remaining average of loss.

Results and Discussion

The trajectories of the loss functions for the two comparisons are listed below:

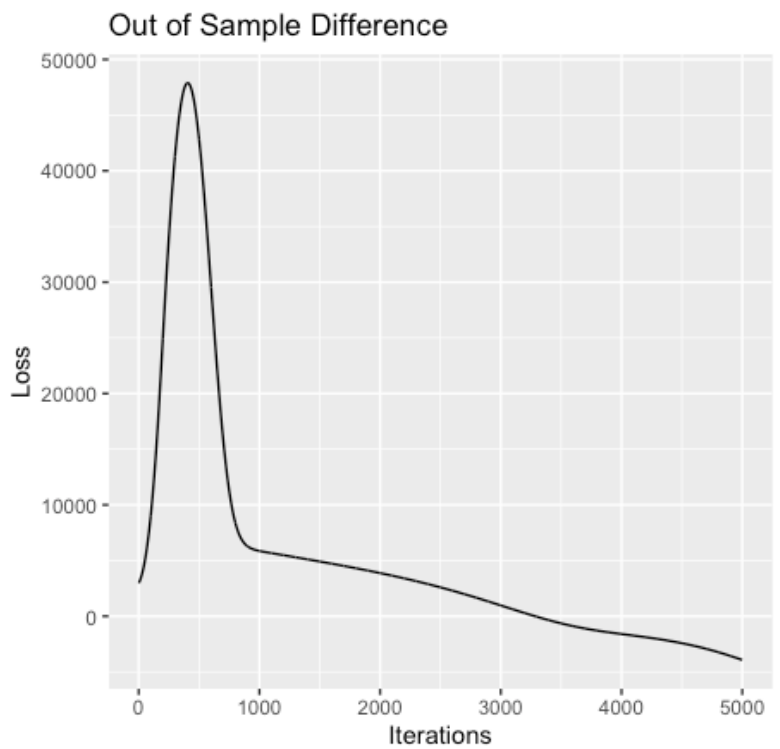
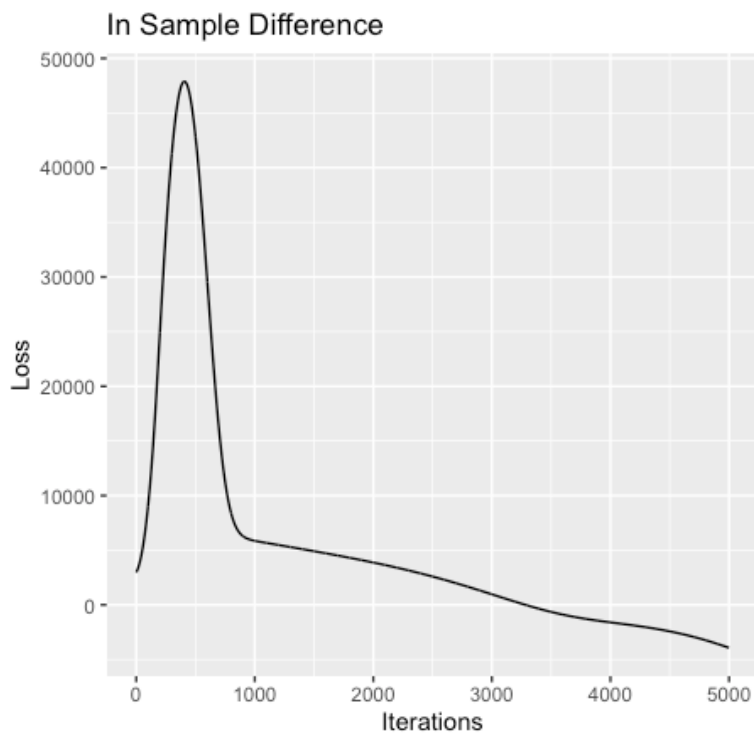


Scenario 1: Deep Neural Network with domain knowledge on range and monotonicity



Scenario 2: Deep Neural Network with domain knowledge on changed monotonicity

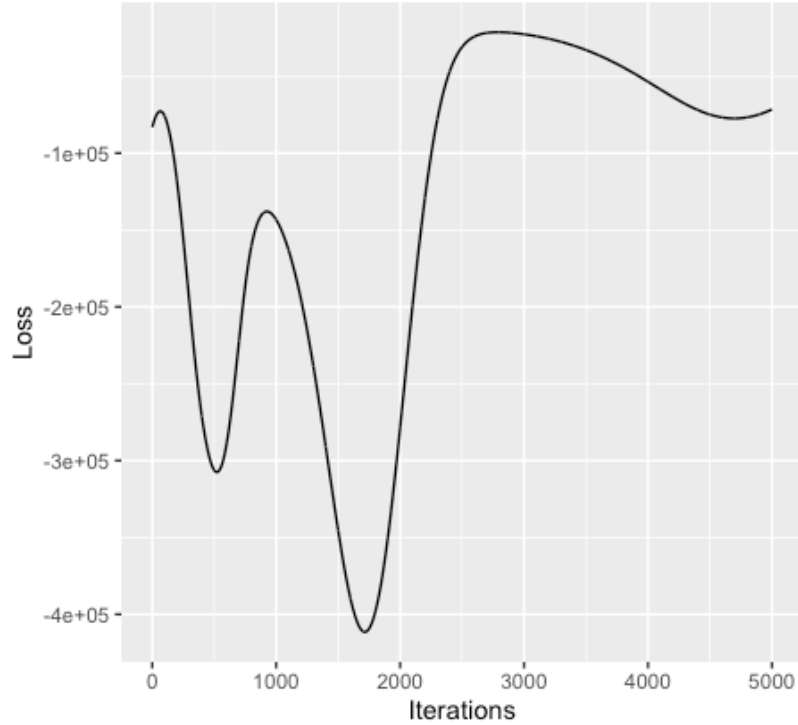
The difference between the loss functions for in sample and out of sample predictions are illustrated with the following plots:



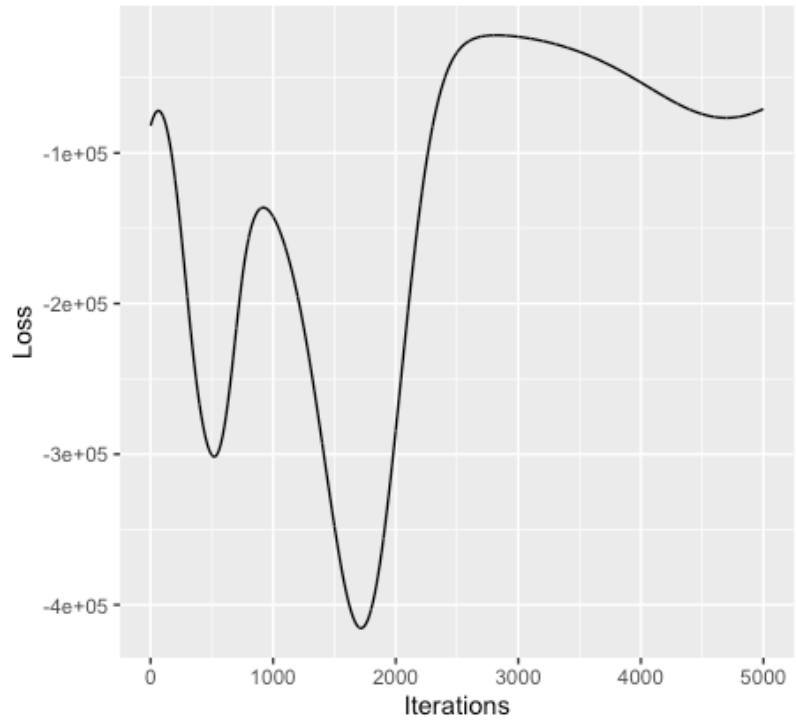
Scenario 1: Deep Neural Network with domain knowledge on range and monotonicity.

$$LOSS_{Domain} - LOSS_{FeedforwardNN}$$

In Sample Difference



Out of Sample Difference



Scenario 2: Deep Neural Network with domain knowledge on changed monotonicity.

$$LOSS_{Domain} - LOSS_{FeedforwardNN}$$

As shown in the figures, deep learners with domain knowledge incorporated into the loss function outperform regular feedforward networks, resulting in a smaller loss and a better fit for the data after iterations, even at the presence of noises within data. For scenario 1, where the domain knowledge contains information on range and monotonicity, the difference after 5000 iterations is -3927 for in sample prediction and -3881 for out of sample. Difference is much larger in scenario 2 with changed monotonicity, that is, -71542.5 for in sample and -70817 for out of sample prediction. One interesting thing is the fluctuations of the difference within the first 2000 iterations, a possible explanation could be the penalty added to the loss function.

This project constructed two deep learning neural networks with different domain knowledge incorporated directly into the loss function, and compared their predictions with a regular feedforward network under two different scenarios. The results indicate that domain knowledge integrated into the loss function does improve the performance of a deep learner, even if the data is noisy or corrupted. One of the major novelties is that this project accommodated range information and monotonicity knowledge together in the loss function, then tried to address the issue with non-monotonicity constraints, while other studies mainly focus on one aspect and have not mentioned non-monotonicity. Further improvements may be achieved by adjusting parameters of the neural network or increasing the number of hidden layers.

Reference

- [1]S. Radovanovic, B. Delibasic, M. Jovanovic, M. Vukicevic, and M. Suknovic, "Framework for integration of domain knowledge into logistic regression," *Web Intelligence, Mining and Semantics*, 2018.
- [2]J.-X. Wang, J.-L. Wu, and H. Xiao, "Physics informed machine learning approach for reconstructing reynolds stress modeling discrepancies based on dns data," *Physical Review Fluids*, Vol. 2, 2017.
- [3]N. Muralidhar, M. R. Islam, M. Marwah, A. Karpatne and N. Ramakrishnan, "Incorporating Prior Domain Knowledge into Deep Neural Networks," 2018 IEEE International Conference on Big Data, Seattle, WA, USA, 2018.
- [4]T. Arnold, M. Kane, and B. W. Lewis, *A Computational Approach to Statistical Learning*.

Appendix

Code Used in This Project

Chloe Jiang

12/18/2020

```
library(reticulate)
library(torch)

n <- 1e5
set.seed(1234)
x1 <- rnorm(n, mean = 3, sd = 1)
x2 <- rnorm(n, mean = 10, sd = 1)
X <- cbind(x1, x2)
X2 <- 5 * X

k1 <- 1
k2 <- 2
K <- 0.7
a1 <- 0.3
a2 <- 0.4
p1 <- 3
p2 <- 4

Y <- k1 * x1^2 + k2 * x2^2 + a1 * cos(p1 * pi * x1) + a2 * cos(p2 * pi * x2) + K
Y2 <- k1 * 5 * x1^2 + k2 * 5 * x2^2 + a1 * cos(p1 * pi * 5 * x1) +
  a2 * cos(p2 * pi * 5 * x2) + K

mean(Y) #212.8242
sd(Y) #40.43148
y_l <- mean(Y) - sd(Y)
y_u <- mean(Y) + sd(Y)

X3 <- 15 * X
Y3 <- k1 * 15 * x1^2 + k2 * 15 * x2^2 + a1 * cos(p1 * pi * 15 * x1) +
  a2 * cos(p2 * pi * 15 * x2) + K
threshold <- mean(X[, 1])
for (i in 1:nrow(X)){
  if (X[, 1][i] > threshold){
    Y3[i] <- k1 * (-10) * x1[i]^2 + k2 * (-10) * x2[i]^2 + a1 * cos(p1 * pi * (-10) * x1[i]) +
      a2 * cos(p2 * pi * (-10) * x2[i]) + K
  }
}
```

```

#for approximation constraints
portion <- 0.5
rows <- runif(nrow(X)) > portion
save <- X[rows, ], 1]
X[rows, ], 1] <- X[rows, ], 2]
X[rows, ], 2] <- save
Y <- k1 * x1^2 + k2 * x2^2 + a1 * cos(p1 * pi * x1) + a2 * cos(p2 * pi * x2) + K

#for monotonicity constraints
rows2 <- runif(length(Y)) > portion
save2 <- Y[rows2]
Y[rows2] <- Y2[rows2]
Y2[rows2] <- save2

#monotonicity changed after a certain threshold
rows3 <- runif(length(Y)) > portion
save3 <- Y[rows3]
Y[rows3] <- Y3[rows3]
Y3[rows3] <- save3

library(rsample)
folds <- validation_split(cbind(X, X2, Y, Y2))

torch_device(type = "cpu")

make_x_y <- function(data){
  X <- data[, 1:2]
  X2 <- data[, 3:4]
  X <- torch_tensor(X)
  X2 <- torch_tensor(X2)
  y <- matrix(data[, 5], ncol = 1)
  y <- torch_tensor(y)
  y2 <- matrix(data[, 6], ncol = 1)
  y2 <- torch_tensor(y2)
  list(X = X, y = y, X2 = X2, y2 = y2)
}

xy <- folds$splits[[1]] %>%
  analysis() %>%
  make_x_y()

x_train <- xy$X
x2_train <- xy$X2
y_train <- xy$y

```

```

y2_train <- xy$y2

xy_test <- folds$splits[[1]] %>%
  assessment() %>%
  make_x_y()

x_test <- xy_test$X
x2_test <- xy_test$X2
y_test <- xy_test$y
y2_test <- xy_test$y2

indicator <- function(condition1, condition2){
  x_dif <- as.array(condition1)
  y_dif <- as.array(condition2)
  ret <- rep(0, length(y_dif))
  for (i in 1:length(y_dif)){
    if ((x_dif[i] < 0) & (y_dif[i] > 0)){ret[i] <- 1}
    else {ret[i] <- 0}}
  return(ret)
}

#WITH DOMAIN KNOWLEDGE
model <- nn_sequential(
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), 1)
)

optimizer <- optim_adam(model$parameters, weight_decay = 0.01)
learning_rate <- 0.01
itr <- 5000
lambdaD <- 1

ll <- rep(NA_real_, itr)
opa <- rep(NA_real_, itr)

for (i in seq_len(itr)){
  y_pred <- model(x_train)
  y2_pred <- model(x2_train)
  #approximation constraints
  lossD1 <- sum(max(y_l - y_pred, 0) + max(y_pred - y_u, 0))
  #monotonicity constraints
  lossD2 <- sum((indicator((x_train - x2_train)[, 1], (y_pred - y2_pred))) *
    max(y_pred - y2_pred, 0))

```

```

loss <- ((y_train$sub(y_pred)$pow(2)$sum())$add(y2_train$sub(y2_pred)$pow(2)$sum())$
  add(lambdaD * lossD1)$add(lambdaD * lossD2))$div(nrow(y_train))

optimizer$zero_grad()
loss$backward()
optimizer$step()

ll[i] <- loss$item()
y_out_pred <- model(x_test)
y2_out_pred <- model(x2_test)
lossD1_out <- sum(max(y_l - y_out_pred, 0) + max(y_out_pred - y_u, 0))
lossD2_out <- sum((indicator((x_test - x2_test)[, 1], (y_out_pred - y2_out_pred))) *
  max(y_out_pred - y2_out_pred, 0))
loss_out <- ((y_test$sub(y_out_pred)$pow(2)$sum())$add(y2_test$sub(y2_out_pred)$pow(2)$
  sum())$add(lambdaD * lossD1_out)$add(lambdaD * lossD2_out))$div(nrow(y_test))
opa[i] <- loss_out$item()
}

library(tibble)
results1 <- tibble(in_sample = ll, out_of_sample = opa)
results1$iteration <- seq_len(nrow(results1))

library(ggplot2)
library(tidyr)

results1 <- results1 %>%
  gather(key = "Type", value = "Value", -iteration)

ggplot(results1, aes(x = iteration, y = Value, color = Type, group = Type)) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "Approximation and Monotonicity Constraints")

##WITHOUT CONSTRAINTS

model <- nn_sequential(
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), 1)
)

optimizer <- optim_adam(model$parameters, weight_decay = 0.01)
learning_rate <- 0.01
itr <- 5000
lambdaD <- 1

ll <- rep(NA_real_, itr)
opa <- rep(NA_real_, itr)

```

```

for (i in seq_len(itr)){
  y_pred <- model(x_train)
  y2_pred <- model(x2_train)

  loss <- ((y_train$sub(y_pred)$pow(2)$sum())$add(y2_train$sub(y2_pred)$pow(2)$
                                                    sum()))$div(nrow(y_train))

  optimizer$zero_grad()
  loss$backward()
  optimizer$step()

  ll[i] <- loss$item()
  y_out_pred <- model(x_test)
  y2_out_pred <- model(x2_test)
  loss_out <- ((y_test$sub(y_out_pred)$pow(2)$sum())$add(y2_test$sub(y2_out_pred)$pow(2)$
                                                         sum()))$div(nrow(y_test))

  opa[i] <- loss_out$item()
}

library(tibble)
results2 <- tibble(in_sample = ll, out_of_sample = opa)
results2$iteration <- seq_len(nrow(results2))

library(ggplot2)
library(tidyr)

results2 <- results2 %>%
  gather(key = "Type", value = "Value", -iteration)

ggplot(results2, aes(x = iteration, y = Value, color = Type, group = Type)) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "No Constraints")

#WHAT IF MONOTONICITY CHANGED AFTER A CERTAIN THRESHOLD

#INCREASE THEN DECREASE

folds <- validation_split(cbind(X, X3, Y, Y3))

make_x_y2 <- function(data){
  X <- data[, 1:2]
  X3 <- data[, 3:4]
  X <- torch_tensor(X)

```

```

X3 <- torch_tensor(X3)
y <- matrix(data[, 5], ncol = 1)
y <- torch_tensor(y)
y3 <- matrix(data[, 6], ncol = 1)
y3 <- torch_tensor(y3)
list(X = X, y = y, X3 = X3, y3 = y3)
}

xy2 <- folds$splits[[1]] %>%
  analysis() %>%
  make_x_y2()

x_train3 <- xy2$X
x3_train3 <- xy2$X3
y_train3 <- xy2$y
y3_train3 <- xy2$y3

xy_test2 <- folds$splits[[1]] %>%
  assessment() %>%
  make_x_y2()

x_test3 <- xy_test2$X
x3_test3 <- xy_test2$X3
y_test3 <- xy_test2$y
y3_test3 <- xy_test2$y3

indicator2 <- function(condition1, condition2){
  x_dif <- as.array(condition1)
  y_dif <- as.array(condition2)
  ret <- rep(0, length(y_dif))
  for (i in 1:length(y_dif)){
    if ((x_dif[i] < 0) & (y_dif[i] < 0)){ret[i] <- 1}
    else {ret[i] <- 0}}
  return(ret)
}

model <- nn_sequential(
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), 1)
)

optimizer <- optim_adam(model$parameters, weight_decay = 0.01)
learning_rate <- 0.01

```

```

itr <- 5000
lambdaD <- 1

ll <- rep(NA_real_, itr)
opa <- rep(NA_real_, itr)

for (i in seq_len(itr)){
  y_pred <- model(x_train3)
  y3_pred <- model(x3_train3)
  #monotonicity constraints
  x_mat3 <- as_array(x_train3)
  index <- row(x_mat3)[which(x_mat3[, 1] < threshold)]
  index2 <- row(x_mat3)[which(x_mat3[, 1] >= threshold)]

  increase <- (indicator((x_train3[index, ] - x3_train3[index, ])[, 1], (y_pred[index] -
    y3_pred[index]))) * max(y_pred[index] - y3_pred[index], 0)
  decrease <- (indicator2((x_train3[index2, ] - x3_train3[index2, ])[, 1], (y_pred[index2] -
    y3_pred[index2]))) * max(y3_pred[index2] - y_pred[index2], 0)

  lossD3 <- sum(increase) + sum(decrease)

  loss <- ((y_train3$sub(y_pred)$pow(2)$sum())$add(y3_train3$sub(y3_pred)$pow(2)$
    sum())$add(lambdaD * lossD3))$div(nrow(y_train3))

  optimizer$zero_grad()
  loss$backward()
  optimizer$step()

  ll[i] <- loss$item()
  y_out_pred <- model(x_test3)
  y3_out_pred <- model(x3_test3)

  x_mat33 <- as_array(x_test3)
  ind <- row(x_mat33)[which(x_mat33[, 1] < threshold)]
  ind2 <- row(x_mat33)[which(x_mat33[, 1] >= threshold)]

  increase2 <- (indicator((x_test3[ind, ] - x3_test3[ind, ])[, 1], (y_out_pred[ind] -
    y3_out_pred[ind]))) * max(y_out_pred[ind] - y3_out_pred[ind], 0)
  decrease2 <- (indicator2((x_test3[ind2, ] - x3_test3[ind2, ])[, 1], (y_out_pred[ind2] -
    y3_out_pred[ind2]))) * max(y3_out_pred[ind2] - y_out_pred[ind2], 0)

  lossD3_out <- sum(increase2) + sum(decrease2)

  loss_out <- ((y_test3$sub(y_out_pred)$pow(2)$sum())$add(y3_test3$sub(y3_out_pred)$
    pow(2)$sum())$add(lambdaD * lossD3_out))$div(nrow(y_test3))

```

```

    opa[i] <- loss_out$item()
  }

library(tibble)
results3 <- tibble(in_sample = ll, out_of_sample = opa)
results3$iteration <- seq_len(nrow(results3))

library(ggplot2)
library(tidyr)

results3 <- results3 %>%
  gather(key = "Type", value = "Value", -iteration)

ggplot(results3, aes(x = iteration, y = Value, color = Type, group = Type)) + geom_line() +
  scale_y_log10() + labs(x = "Iterations", y = "Loss", title = "Non-Monotonicity")

#WITHOUT DOMAIN KNOWLEDGE

model <- nn_sequential(
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), ncol(X)),
  nn_linear(ncol(X), 1)
)

optimizer <- optim_adam(model$parameters, weight_decay = 0.01)
learning_rate <- 0.01
itr <- 5000
lambdaD <- 1

ll <- rep(NA_real_, itr)
opa <- rep(NA_real_, itr)

for (i in seq_len(itr)){
  y_pred <- model(x_train3)
  y3_pred <- model(x3_train3)

  loss <- ((y_train3$sub(y_pred)$pow(2)$sum())$add(y3_train3$sub(y3_pred)$
    pow(2)$sum()))$div(nrow(y_train3))

  optimizer$zero_grad()
  loss$backward()
  optimizer$step()

```



```

ll[i] <- loss$item()
y_out_pred <- model(x_test3)
y3_out_pred <- model(x3_test3)

loss_out <- ((y_test3$y - y_out_pred)^2)$sum()$add(y3_test3$y - y3_out_pred)^2$sum())$div(nrow(y_test3))
opa[i] <- loss_out$item()
}

library(tibble)
results4 <- tibble(in_sample = ll, out_of_sample = opa)
results4$iteration <- seq_len(nrow(results4))

library(ggplot2)
library(tidyr)

results4 <- results4 %>%
  gather(key = "Type", value = "Value", -iteration)

ggplot(results4, aes(x = iteration, y = Value, color = Type, group = Type)) + geom_line() +
  scale_y_log10() + labs(x = "Iterations", y = "Loss", title = "No Constraints")

results_in <- data.frame(1:itr, (results1[, 3] - results2[, 3])[1:itr, ])

results0_in <- data.frame(1:itr, (results3[, 3] - results4[, 3])[1:itr, ])
results_out <- data.frame(1:itr, (results1[, 3] - results2[, 3])[5001:10000, ])
results0_out <- data.frame(1:itr, (results3[, 3] - results4[, 3])[5001:10000, ])

ggplot(results_in, aes(x = results_in[, 1], y = results_in[, 2])) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "In Sample Difference")

ggplot(results_out, aes(x = results_out[, 1], y = results_in[, 2])) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "Out of Sample Difference")

ggplot(results0_in, aes(x = results0_in[, 1], y = results0_in[, 2])) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "In Sample Difference")

ggplot(results0_out, aes(x = results0_out[, 1], y = results0_out[, 2])) + geom_line() +
  labs(x = "Iterations", y = "Loss", title = "Out of Sample Difference")

```