# Deep Learning
# Parameter Optimization

Nagiza F. Samatova, samatova@csc.ncsu.edu
**Professor, Department of Computer Science**
**North Carolina State University**

**NC STATE** UNIVERSITY
Department of Computer Science

# Outline

- Model Hyperparameters
- Parameter Optimization Problem
- Parameter Optimization Strategies
    - Trial and error
    - Grid Search
    - Random Search
    - Bayesian Parameter Optimization
- Grid Search
- Random Search

# Model Parameters and Hyperparameters

- **Model Parameters**:
  - Parameters that are learned during model training
  - Examples: weights, intercept/bias parameters
- **Hyperparameters**:
  - Parameters that **must be set** at the offset and are **not optimized during model training**
- **Examples of Hyperparameters**:
  - **The choice of a model**: deep feedforward neural network, random forest, support vector machines, etc.
  - **Model architecture**: the number of hidden layers, the number of nodes in each hidden layer, types of activation functions
  - **Model Convergence and Stopping Criteria**: the number of training epochs (iterations), threshold on the difference in optimization function between two subsequent iterations, learning rate
  - **Form of regularization to avoid model overfitting**: $L_1$- or $L_2$-norm regularizer, the number of dropout neurons,

3

# Deep NN Models: Hyperparameters

```
38  system.time(ex1 <- h2o.deeplearning(
39    x = xnames,
40    y = "label",
41    training_frame= h2odigits.train,
42    validation_frame = h2odigits.test,
43    activation = "RectifierWithDropout",
44    hidden = c(100),
45    epochs = 10,
46    adaptive_rate = FALSE,
47    rate = .001,          <---
48    input_dropout_ratio = 0,
49    hidden_dropout_ratios = c(.2)
50  ))
```

```
Validation Set Metrics:
=====================

Extract validation frame with `h2o.getFrame(
MSE: (Extract with `h2o.mse`) 0.032
RMSE: (Extract with `h2o.rmse`) 0.18
Logloss: (Extract with `h2o.logloss`) 0.21
Mean Per-Class Error: 0.038
Confusion Matrix: Extract with `h2o.confusi
```

Although ex1 took longer to train, it performs substantially better on the test data than does ex2.

```
52  system.time(ex2 <- h2o.deeplearning(
53    x = xnames,
54    y = "label",
55    training_frame= h2odigits.train,
56    validation_frame = h2odigits.test,
57    activation = "RectifierWithDropout",
58    hidden = c(100),
59    epochs = 10,
60    adaptive_rate = FALSE,
61    rate = .01,          <---
62    input_dropout_ratio = 0,
63    hidden_dropout_ratios = c(.2)
64  ))
```

```
Validation Set Metrics:
=====================

Extract validation frame with `h2o.getFram
MSE: (Extract with `h2o.mse`) 0.08
RMSE: (Extract with `h2o.rmse`) 0.28
Logloss: (Extract with `h2o.logloss`) 2.5
Mean Per-Class Error: 0.082
Confusion Matrix: Extract with `h2o.confus
```

# Overarching Goal: Parameter Optimization

- **Motivation:** Some of the key reasons for poor model performance:
  - lack of the variables required for good prediction
  - not enough data to support training a complex enough model
  - **poorly tuned and optimized hyperparameters**

- **Goal:** **To choose the best combination of and best values for the hypermeters of the machine learning model**.
  - **Assumption:** Better hyperparameters can often improve the accuracy of a model.

The values chosen for the hyperparameters can have a dramatic impact on the accuracy and training speed of a model.

# Parameter Optimization Strategies

- **Trial and error**
- **Grid Search:**
  - pros: great if there are only a **few values** for a **few parameters**
  - cons: combinatorial (brute-force exhaustive enumeration of all possible combinations)
- **Random Search:**
  - pros: searching via random sampling; no need to pre-specify all the values to try and create all possible combinations
  - cons: computationally demanding for large sample sizes
- **Bayesian Optimization (e.g., using Spearmint Python library):**
  - iteratively adjusts a number of parameters so as to minimize some objective in as few runs as possible
  - https://github.com/HIPS/Spearmint

# Grid Search for Optimal Parameters

- **Basic Idea:**
  - several values for hyperparameters are specified and
  - all possible combinations of these values are tried
- To create all possible combinations in R:
  - Use **expandGrid()** function in **gridExtra** package

DeepNeuralNetwork.ParameterOptimization.R

```
14  expand.grid(
15      layers = c(1, 2, 4),
16      epochs = c(50, 100),
17      l1 = c(.001, .01, .05))
```

DeepNeuralNetwork.R

```
40  digits.m1 <- train(digits.X, digits.y,
41                  method = "nnet",
42                  tuneGrid = expand.grid(
43                      .size = c(5),
44                      .decay = 0.1),
45                  trControl = trainControl(method = "n
46                  MaxNWts = 10000,
47                  maxit = 100)
```

all possible combinations

|    | layers | epochs | l1    |
|----|--------|--------|-------|
| 1  | 1      | 50     | 0.001 |
| 2  | 2      | 50     | 0.001 |
| 3  | 4      | 50     | 0.001 |
| 4  | 1      | 100    | 0.001 |
| 5  | 2      | 100    | 0.001 |
| 6  | 4      | 100    | 0.001 |
| 7  | 1      | 50     | 0.010 |
| 8  | 2      | 50     | 0.010 |
| 9  | 4      | 50     | 0.010 |
| 10 | 1      | 100    | 0.010 |
| 11 | 2      | 100    | 0.010 |
| 12 | 4      | 100    | 0.010 |
| 13 | 1      | 50     | 0.050 |
| 14 | 2      | 50     | 0.050 |
| 15 | 4      | 50     | 0.050 |
| 16 | 1      | 100    | 0.050 |
| 17 | 2      | 100    | 0.050 |
| 18 | 4      | 100    | 0.050 |

# Random Search for Optimal Parameters

- **What to specify for random sampling:**
  - the values to randomly sample or
  - distributions to randomly draw from.
  - some limits, e.g., # of hidden layers in a range from 1 to 10
- **Examples**:
  - The **number of dropout neurons** will be drawn from beta distribution (20% dropouts in input layer and 50% on hidden layer): **dbeta()**
  - The depth or **number of layers** is sampled from 1 to 5
  - Adaptive learning rates ($\rho$ and $\epsilon$) are drawn from a uniform distribution: **runif()**
- **How to do random sampling:**
  - write a function that takes a seed and
  - then randomly samples a number of hyperparameters,
  - stores the sampled parameters,
  - runs the model, and
  - returns the results

8

# Step-1: Set Up Random Search Parameters

## Setting Up Model Hyperparameters

```
 92  run <- function(seed,
 93                    name = paste0("m_", seed),
 94                    run = TRUE) {
 95    set.seed(seed)
 96
 97    p <- list(
 98      Name = name,
 99      seed = seed,
100      depth = sample(1:5, 1),
101      l1 = runif(1, 0, .01),
102      l2 = runif(1, 0, .01),
103      input_dropout = rbeta(1, 1, 12),
104      rho = runif(1, .9, .999),
105      epsilon = runif(1, 1e-10, 1e-4))
106
107    p$neurons <- sample(20:600, p$depth, TRUE)
108    p$hidden_dropout <- rbeta(p$depth, 1.5, 1)/2
```

# Step-2: Build Models for these Parameters

**Building Models for Different Model Hyperparameters**

```
110 -    if (run) {
111         model <- h2o.deeplearning(
112             x = colnames(use.train.x),
113             y = "Outcome",
114             training_frame = h2oactivity.train,
115             activation = "RectifierWithDropout",
116             hidden = p$neurons,
117             epochs = 100,
118             loss = "CrossEntropy",
119             input_dropout_ratio = p$input_dropout,
120             hidden_dropout_ratios = p$hidden_dropout,
121             l1 = p$l1,
122             l2 = p$l2,
123             rho = p$rho,
124             epsilon = p$epsilon,
125             export_weights_and_biases = TRUE,
126             model_id = p$Name
127         )
```

# Step-3: Measure Models' Performance

**Measure Model Performance for Different Model Hyperparameters**

```
131    ## Measure performance on training data
132    p$MSE <- h2o.mse(model)
133    p$R2 <- h2o.r2(model)
134    p$Logloss <- h2o.logloss(model)
135    p$CM <- h2o.confusionMatrix(model)
136
137    ## Measure performance on testing data
138    perf <- h2o.performance(model, h2oactivity.test)
139    p$T.MSE <- h2o.mse(perf)
140    p$T.R2 <- h2o.r2(perf)
141    p$T.Logloss <- h2o.logloss(perf)
142    p$T.CM <- h2o.confusionMatrix(perf)
```

# Step-4: Return Parameters, Model & Results

**Return Model Results for Different Model Hyperparameters**

```
148    return(list(
149       Params = p,
150       Model = model))
```

# Step-5: Generate Results for Different Seeds
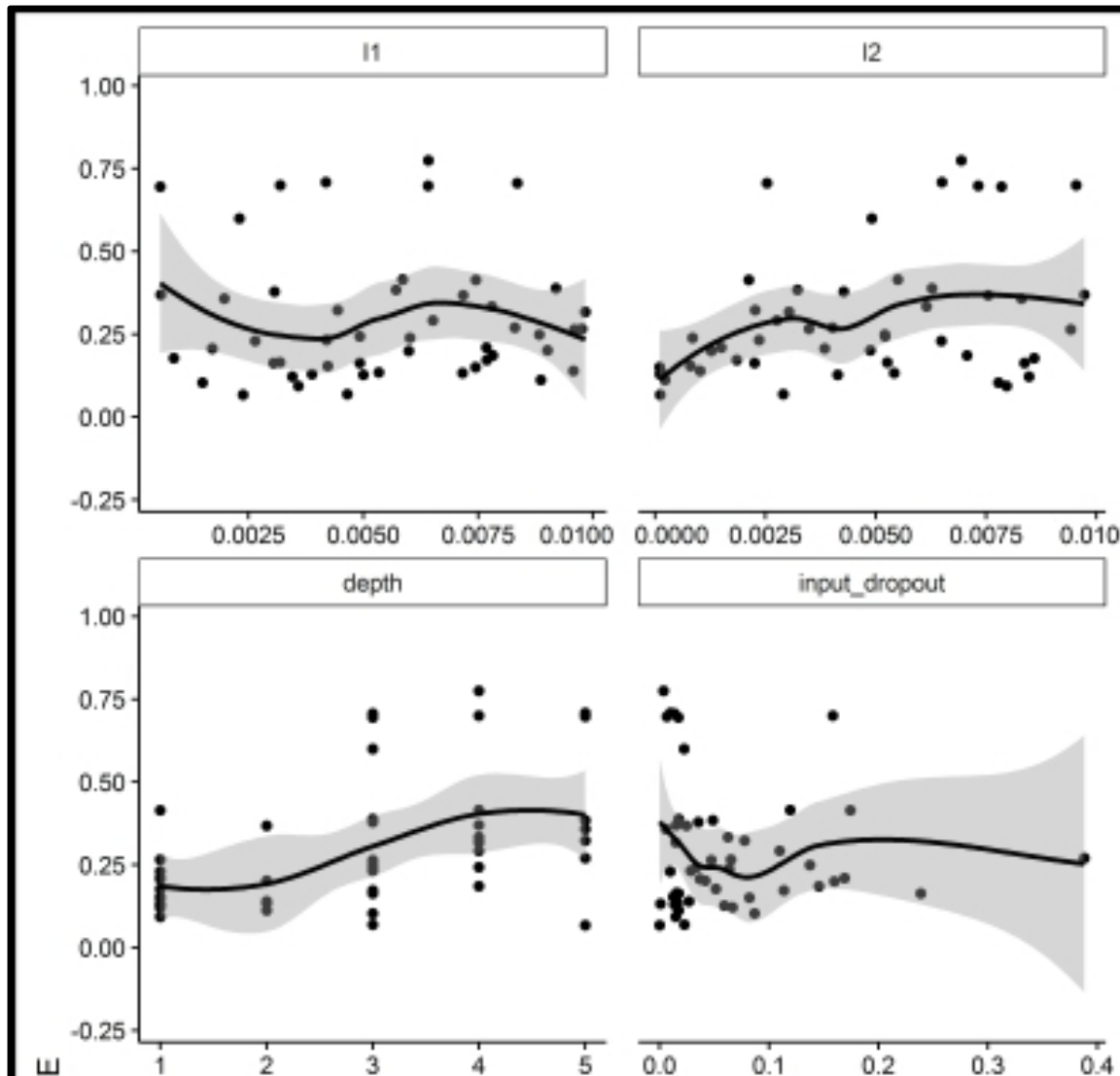
To make results reproducible, use different seeds:

```
188  # To make the parameters reproducible,
189  # specify a list of random seeds to loop through to run the models
190  use.seeds <- c(403L, 10L)
191  # Step-5: To run the models simply by looping through the seeds:
192  model.res <- lapply(use.seeds, run)
```

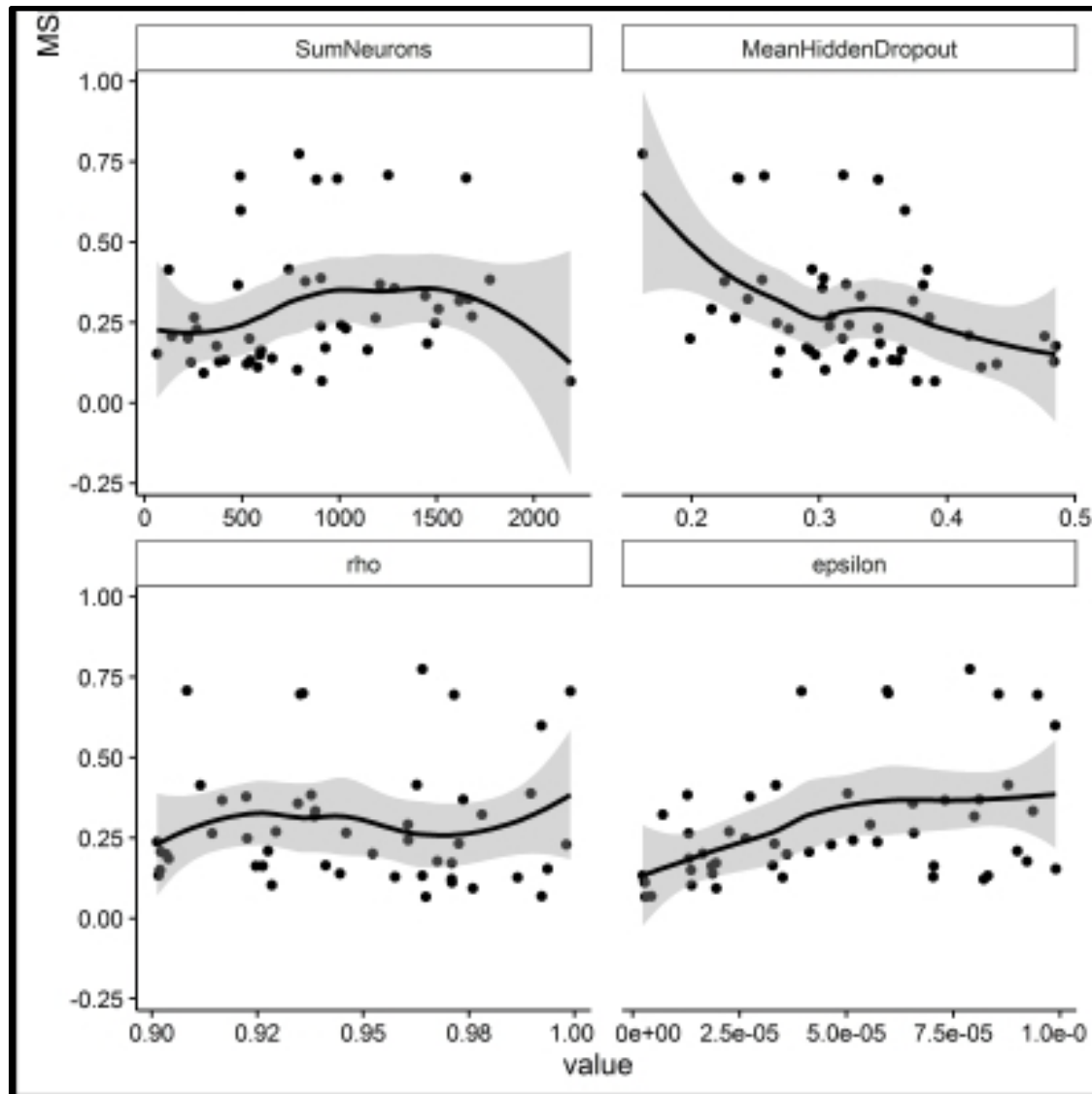# Step-6: Examine MSE for Diff. Parameters

Plot the mean squared error (MSE) against the different parameters

```
196 model.res.dat <- do.call(rbind,
197         lapply(model.res,
198             function(x) with(x$Params,
199                 data.frame(l1 = l1, l2 = l2,
200                 depth = depth,
201                 input_dropout = input_dropout,
202                 SumNeurons = sum(neurons),
203                 MeanHiddenDropout = mean(hidden_dropout),
204                 rho = rho,
205                 epsilon = epsilon,
206                 MSE = T.MSE))))
207
208
209 p.perf <- ggplot(melt(model.res.dat, id.vars = c("MSE")), aes(value, MSE)) +
210     geom_point() +
211     stat_smooth(colour = "black") +
212     facet_wrap(~ variable, scales = "free_x", ncol = 2) +
213     theme_classic()
214
215 print(p.perf)
```

# Step-6: MSE Plots for Different Seeds

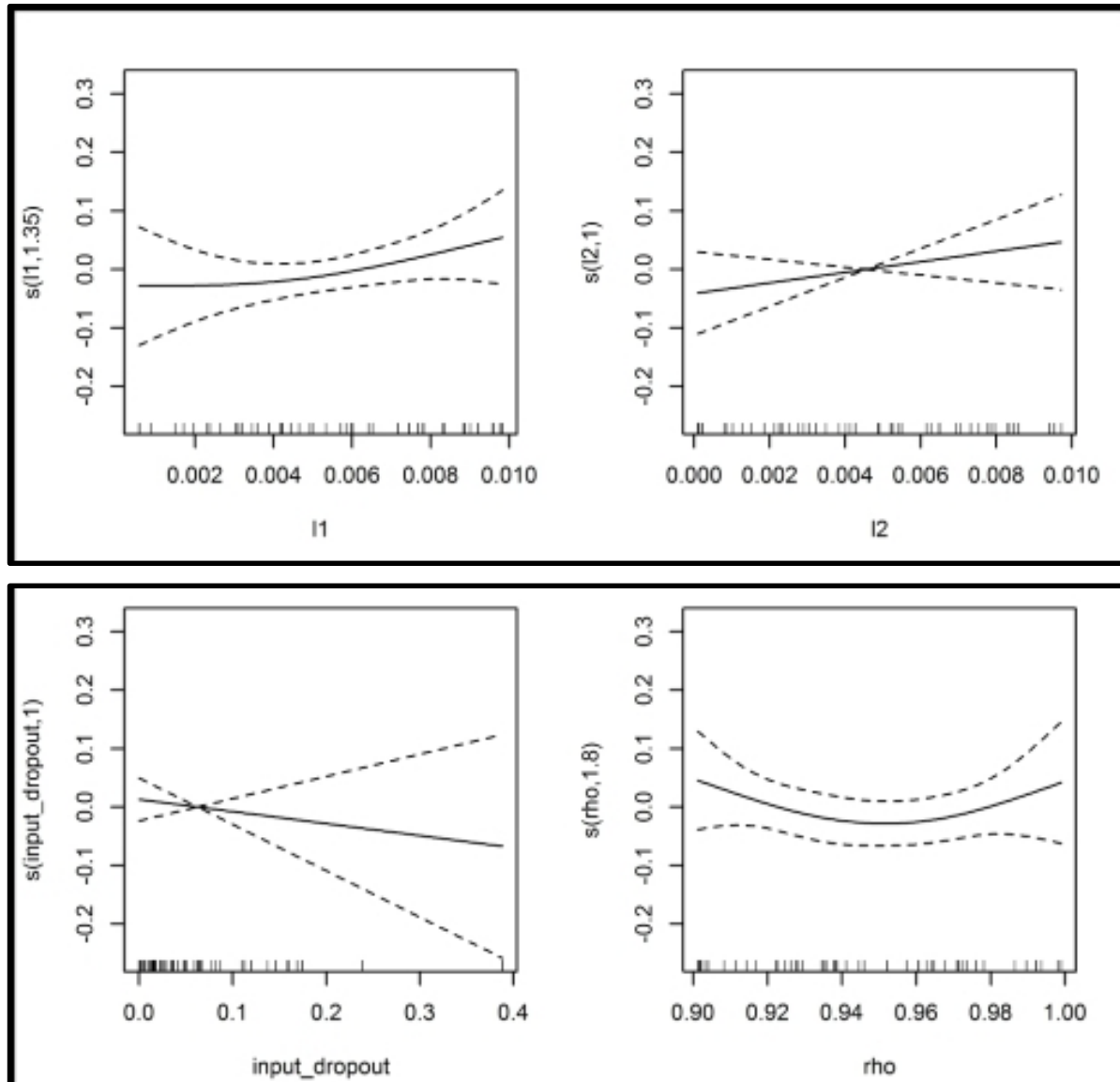# Step-6: MSE Plots for Different Seeds

# Step-7: Multivariate Parameter Relationships

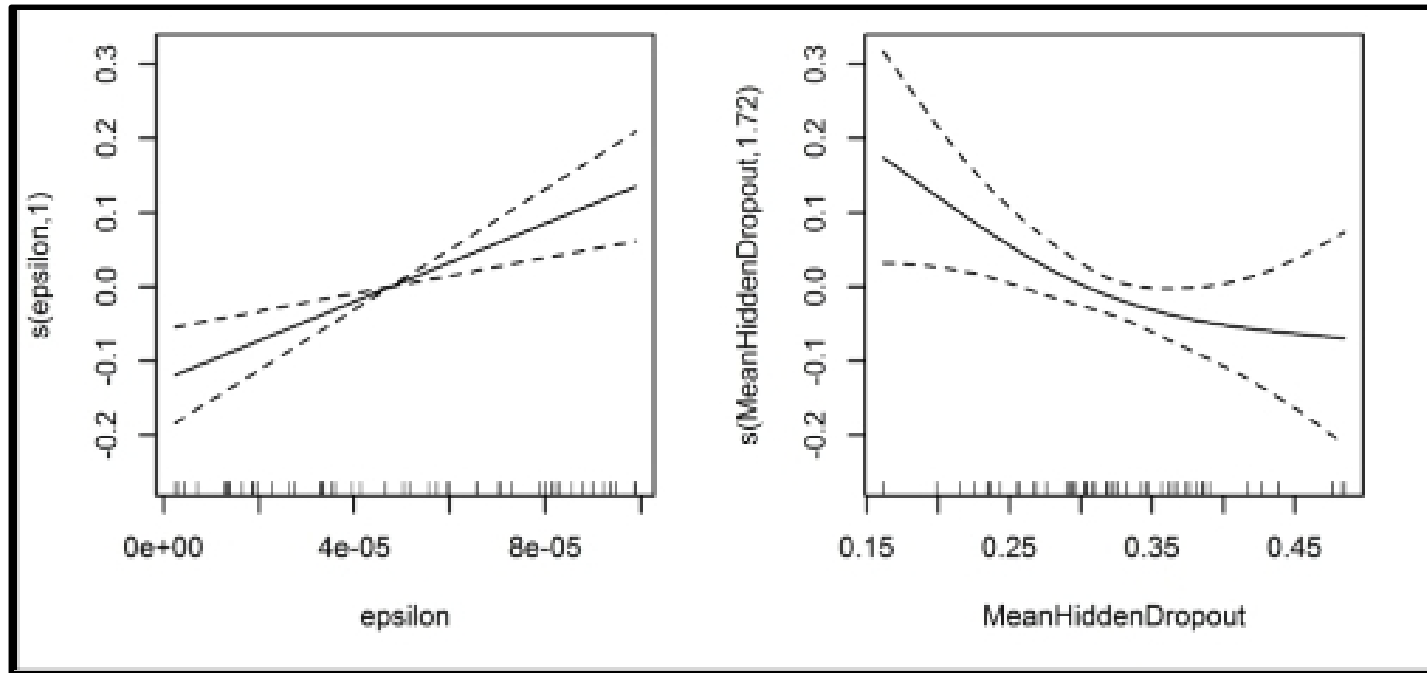Examine multivariate parameter relationships w.r.t. MSE

```
227  summary(m.gam <- gam(MSE ~ s(l1, k = 4) +
228                                s(l2, k = 4) +
229                                s(input_dropout) +
230                                s(rho, k = 4) +
231                                s(epsilon, k = 4) +
232                                s(MeanHiddenDropout, k = 4) +
233                                te(depth, SumNeurons, k = 4),
234                          data = model.res.dat))
235· for (i in 1:6) {
236    plot(m.gam, select = i)
237  }
238  plot(m.gam, select = 7)
```

# Step-7: MSE Multivariate Relation Plots

# Step-7: MSE Multivariate Relation Plots

# Step-8: Choose Optimal Model Parameters

```r
247  model.optimized <- h2o.deeplearning(
248    x = colnames(use.train.x),
249    y = "Outcome",
250    training_frame = h2oactivity.train,
251    activation = "RectifierWithDropout",
252    hidden = c(500, 500, 500),
253    epochs = 100,
254    loss = "CrossEntropy",
255    input_dropout_ratio = .08,
256    hidden_dropout_ratios = c(.50, .50, .50),
257    l1 = .002,
258    l2 = 0,
259    rho = .95,
260    epsilon = 1e-10,
261    export_weights_and_biases = TRUE,
262    model_id = "optimized_model"
263  )
264
265  h2o.performance(model.optimized, h2oactivity.test)
266
267  model.res.dat[which.min(model.res.dat$MSE), ]
```

# References

- **Chapter 19, Design and Analysis of Machine Learning Experiments**
- **Bengio, Y.** (2012), **Section 3**, **Hyper-Parameters** in Practical Recommendations for Gradient-Based Training of Deep Architectures. In Neural Networks: Tricks of the Trade (pp. 437-478). Springer Berlin Heidelberg. (Also on the arXiv: http://arxiv.org/pdf/1206.5533.pdf)
- **Zeiler, M. D.** (2012),  **ADADELTA**: An Adaptive Learning Rate Method. arXiv:1212.5701
- Bayesian optimization of hyperparameters, available online:
    - https://github.com/HIPS/Spearmint
    - see references on github for Spearmint library
- Joshua F. Wiley , R Deep Learning Essentials