



## NeuralNetTools: Visualization and Analysis Tools for Neural Networks

Marcus W. Beck

Oak Ridge Institute for Science and Education  
US Environmental Protection Agency

---

### Abstract

Functions within this package can be used for the interpretation of neural network models created in R, including functions to plot a neural network interpretation diagram, evaluation of variable importance, and a sensitivity analysis of input variables.

*Keywords:* neural networks, plotnet, sensitivity, variable importance, R.

---

## 1. Introduction

The increasing quantity of information and computational capacity to address relevant research questions has contributed to the growth of data science as a legitimate field of study. Data science is a relatively new paradigm of analysis that focuses on the synthesis of unstructured information from multiple sources to identify patterns or trends ‘born from the data’ (Kelling, Hochachka, Fink, Riedewald, Caruana, Ballard, and Hooker 2009). A central theme is the focus on data exploration and prediction as compared to hypothesis-testing using domain-specific methods for scientific exploration (Kell and Oliver 2003). Demand for quantitative toolsets to address challenges in data-rich environments has increased drastically with the advancement of techniques for rapid acquisition of data. Fields of research characterized by high-throughput data have a strong foundation in computationally-intensive methods of analysis (e.g., Saeys, Inza, and naga (2007)). Moreover, disciplines that have historically been limited by data quantity, such as ecological studies across broad temporal and spatial scales, have also realized the importance of data intensive approaches of analysis given the improved ability to acquire information (e.g., Swanson, Kosmala, Lintott, Simpson, Smith, and Packer (2015)). Regardless of the discipline, quantitative methods that explicitly focus on inductive reasoning can serve a complementary role to conventional, hypothesis-driven approaches to scientific discovery (Kell and Oliver 2003).

Statistical methods that have been used to support data exploration for inductive analysis are numerous (Jain, Duin, and Mao 2000). A common theme among these methods is the use of machine-learning algorithms where the primary objective is to identify emergent patterns in the data with minimal human intervention. Neural networks, in particular, are designed to mimic the neuronal structure of the human brain by ‘learning’ inherent data structures through adaptive algorithms (Rumelhart, Hinton, and Williams 1986; Ripley 1996). Although the conceptual model was introduced several decades ago (McCulloch and Pitts 1943), neural networks have had a central role in data intensive science. The most popular form of neural network is the multilayer perceptron (mlp) trained using the backpropagation algorithm (Rumelhart *et al.* 1986). This model is typically used to predict the response of one or more variables given co-occurrence of one to many explanatory variables. The hallmark feature of the mlp is the characterization of relationships using an arbitrary number of parameters (i.e., the hidden layer) that are chosen through an iterative training process with the backpropagation algorithm. Conceptually, the mlp is nothing more than a hyper-parameterized non-linear model that can fit a smooth function to any dataset with almost non-existent residual error (Hornik 1991).

An arbitrarily large number of parameters to fit a neural network provides obvious predictive advantages, but conversely complicates the extraction of critical model information. Information such as variable importance or model sensitivity are necessary aspects of exploratory data analysis that are not easily obtained from a neural network. As such, a common criticism is that neural networks are ‘black-boxes’ that offer minimal insight into relationships among variables (e.g., Paruelo and Tomasel 1997). Olden and Jackson (2002) provide a rebuttal to this concern by describing methods to extract information from neural networks, most of which were previously available but not commonly used. For example, Olden and Jackson (2002) describe neural interpretation diagrams for plotting (Özesmi and Özesmi 1999), the Garson algorithm for variable importance (Garson 1991), and the Profile method for sensitivity analysis (Lek, Delacoste, Baran, Dimopoulos, Lauga, and Aulagnier 1996). These quantitative tools ‘illuminate the black box’ by disaggregating the network parameters to characterize relationships between variables that are described by the model. In essence, mlp neural networks were developed for prediction but methods described in (Olden and Jackson 2002) leverage these models to describe data signals. Increasing the accessibility of these diagnostic tools will have value for exploratory analysis in data science.

This article describes the **NeuralNetTools** package for R that was developed to improve the breadth and quality of information obtained from the mlp neural network. Functions provided by the package are those previously described in (Olden and Jackson 2002) but have not been available in an open-source programming environment. The reach of the package is all-inclusive such that generic functions were developed using S3 methods for all neural network object classes available in R. The objectives of this article are to 1) provide an overview of the statistical foundation the mlp network, 2) briefly describe similarities and differences between existing neural network packages in R, and 3) describe the theory and application of the primary functions in the **NeuralNetTools** package. The package is currently available on CRAN, whereas the development version is maintained as a GitHub repository.

## 2. Theoretical foundation and existing R packages

The following is a brief description of the basic structure and methods for creating a network

to provide sufficient context for the **NeuralNetTools** package. A detailed discussion of the mlp neural network is beyond the scope of this manuscript and interested readers should consult [Rumelhart \*et al.\* \(1986\)](#); [Ripley \(1996\)](#) and references therein for a more thorough review. An intriguing characteristic of the neural network approach to predictive modelling is the absence of assumptions regarding the distributional characteristics of the response variables. Unlike conventional approaches, neural networks have been popularized by their ability to model response variables with arbitrary distributions and can describe relationships in datasets with noisy or imprecise information. The typical mlp network is composed of multiple layers that define the transfer of information between input and response layers. Information travels in one direction where a set of values for one to many variables in the input layer propagates through one or more hidden layers to the resulting layer of the response variables. ‘Hidden’ layers between the input and response layers are key components of a neural network that mediate the transfer of information. Just as the input and response layers are composed of variables or ‘nodes’, each hidden layer is composed of nodes with weighted connections that define the strength of information flow between layers. ‘Bias’ layers connected to hidden and response layers may also be used that are analagous to intercept terms in a standard regression model.

Training a neural network model requires identifying the ‘optimal’ weights that define the connections between the model layers. The optimal weights are those that minimize prediction error for a test dataset that is independent of the training dataset. Training is commonly achieved using the backpropagation algorithm described in detail in ([Rumelhart \*et al.\* 1986](#)). This algorithm identifies the optimal weighting scheme between layers through an iterative process where weights are gradually changed through a forward- and backward-propagation process ([Rumelhart \*et al.\* 1986](#); [Lek and Guégan 2000](#)). The algorithm begins by assigning an arbitrary weighting scheme to the connections in the network, followed by estimating the output in the response variable through the forward-propagation of information through the network, and finally calculating the difference between the predicted and actual value of the response. The weights are then changed through a back-propagation step that begins by changing weights in the output layer and then the remaining hidden layers. The process is repeated until the chosen error function is minimized. Formulaically, a generic mlp neural network can be represented as ([Ripley 1996](#)):

$$y_k = f_k \left( \sum_{j=1}^k w_{jk} f_j \left( \sum_{i=1}^j w_{ij} x_i \right) \right) \quad (1)$$

where the estimated value of the response variable  $y_k$  is a summation of the product between multiple input variables  $x$  and multiple hidden nodes  $j$ , as mediated by the respective weights  $w$  and activation functions  $f_j$  and  $f_k$  for each hidden and output node as the information progresses through the network.

Packages available in R to create mlp neural networks include **neuralnet**, **nnet**, and **RSNNS**. An additional package, **FCNN4R**, provides an interface between **R** and the **FCNN** C++ library. At the time of writing, **FCNN4R** is a relatively new package that represents less than 1% of the total downloads of all neural network packages on CRAN. As such, **NeuralNetTools** does not include methods for the **FCNN4R** package, although future versions may do so depending on the popularity of **FCNN4R**. In general, the remaining packages use similar methods to create mlp networks such that the resulting models should be comparable.

Differences between the packages primarily include default values of the arguments and the ability to handle multiple response variables or multiple hidden layers in the same model. The `nnet` function can take separate (or combined) input or response inputs as each as a separate `data.frame` or as a `formula`, the `neuralnet` function can only use a `formula` as input, and the `mlp` function (**RSNNS**) can only take a `data.frame` as combined or separate variables as input. The `neuralnet` function is not capable of modelling multiple response variables, unless the response is a categorical variable where one node for each outcome is used in the response. Additionally, the default output for the `neuralnet` function is linear, whereas the opposite is true for the other two functions. Other differences that influence how functions in **NeuralNetTools** handle object classes from the different packages are described below.

### 3. Package structure

#### 3.1. Visualizing neural networks

The number of existing functions in **R** to view neural networks is minimal. Such tools have practical use for visualizing network architecture and connections between layers that mediate variable importance. To our knowledge, only the `neuralnet` and **FCNN4R** packages provide plotting methods for `mlp` networks. Although useful for viewing the basic structure, the output is visually minimal and does not include options for customization (verify).

This function plots a neural network as a neural interpretation diagram as in [Özesmi and Özesmi \(1999\)](#). Options to plot without color-coding or shading of weights are also provided. The default settings plot positive weights between layers as black lines and negative weights as grey lines. Line thickness is in proportion to relative magnitude of each weight. The first layer includes only input variables with nodes labelled arbitrarily as I1 through In for n input variables. One through many hidden layers are plotted with each node in each layer labelled as H1 through Hn. The output layer is plotted last with nodes labeled as O1 through On. Bias nodes connected to the hidden and output layers are also shown. Neural networks created using `mlp` do not show bias layers.

A primary network and a skip layer network can be plotted for `nnet` models with a skip layer connection. The default is to plot the primary network, whereas the skip layer network can be viewed with `skip = TRUE`. If `nid = TRUE`, the line widths for both the primary and skip layer plots are relative to all weights. Viewing both plots is recommended to see which network has larger relative weights. Plotting a network with only a skip layer (i.e., no hidden layer, `size = 0`) will include bias connections to the output layer, whereas these are not included in the plot of the skip layer if size is greater than zero.

Pruned networks in **RSNNS**.

#### 3.2. Evaluating variable importance

The `garson` function uses Garson's algorithm to evaluate relative variable importance. This function identifies the relative importance of explanatory variables for a single response variable by deconstructing the model weights. The importance of each variable can be determined by identifying all weighted connections between the layers in the network. That is, all weights connecting the specific input node that pass through the hidden layer to the response variable

are identified. This is repeated for all other explanatory variables until a list of all weights that are specific to each input variable is obtained. The connections are tallied for each input node and scaled relative to all other inputs. A single value is obtained for each explanatory variable that describes the relationship with the response variable in the model. The results indicate relative importance as the absolute magnitude from zero to one. The function cannot be used to evaluate the direction of the response. Only neural networks with one hidden layer and one output node can be evaluated.

The `olden` function is an alternative and more flexible approach to evaluate variable importance. The function calculates importance as the product of the raw input-hidden and hidden-output connection weights between each input and output neuron and sums the product across all hidden neurons. An advantage of this approach is the relative contributions of each connection weight are maintained in terms of both magnitude and sign as compared to Garson's algorithm which only considers the absolute magnitude. For example, connection weights that change sign (e.g., positive to negative) between the input-hidden to hidden-output layers would have a cancelling effect whereas Garson's algorithm may provide misleading results based on the absolute magnitude. An additional advantage is that Olden's algorithm is capable of evaluating neural networks with multiple hidden layers and response variables. The importance values assigned to each variable are in units that are based directly on the summed product of the connection weights. The actual values should only be interpreted based on relative sign and magnitude between explanatory variables. Comparisons between different models should not be made.

Issues with different indications of variable importance as a model is refit...

### 3.3. Sensitivity analysis

The Lek profile method is described briefly in [Lek \*et al.\* \(1996\)](#) and in more detail in [Gevrey, Dimopoulos, and Lek \(2003\)](#). The profile method is fairly generic and can be extended to any statistical model in R with a `predict` method. However, it is one of few methods used to evaluate sensitivity in neural networks.

The profile method can be used to evaluate the effect of explanatory variables by returning a plot of the predicted response across the range of values for each separate variable. The original profile method evaluated the effects of each variable while holding the remaining explanatory variables at different quantiles (e.g., minimum, 20th percentile, maximum). This is implemented in the function by creating a matrix of values for explanatory variables where the number of rows is the number of observations and the number of columns is the number of explanatory variables. All explanatory variables are held at their mean (or other constant value) while the variable of interest is sequenced from its minimum to maximum value across the range of observations. This matrix (or data frame) is then used to predict values of the response variable from a fitted model object. This is repeated for each explanatory variable to obtain all response curves. Values passed to `split_vals` must range from zero to one to define the quantiles for holding unevaluated explanatory variables.

An alternative implementation of the profile method is to group the unevaluated explanatory variables using groupings defined by the statistical properties of the data. Covariance among predictors may present unlikely scenarios if holding all unevaluated variables at the same level. To address this issue, the function provides an option to hold unevaluated variable at mean values defined by natural clusters in the data. `kmeans` clustering is used on the

input data.frame of explanatory variables if the argument passed to `split_vals` is an integer value greater than one. The centers of the clusters are then used as constant values for the unevaluated variables. An arbitrary grouping scheme can also be passed to `split_vals` as a data.frame where the user can specify exact values for holding each value constant (see the examples). Examples in Beck, Wilson, Vondracek, and Hatch (2014) show this...

For all plots, the legend with the 'splits' label indicates the colors that correspond to each group. The groups describe the values at which unevaluated explanatory variables were held constant, either as specific quantiles, group assignments based on clustering, or in the arbitrary grouping defined by the user. The constant values of each explanatory variable for each split can be viewed as a barplot by using `split_show = TRUE`.

Note that there is no predict method for neuralnet objects from the nn package. The lekprofile method for nn objects uses the nnet package to recreate the input model, which is then used for the sensitivity predictions. This approach only works for networks with one hidden layer.

## 4. Future development

## 5. Conclusions

A cautionary note about the 'optimal network' and reproducibility of results. Note that most literature sources suggest variable standardization prior to using a model yet none of the existing packages in R provide this functionality as a default option (verify).

## 6. Acknowledgments

## References

- Beck MW, Wilson BN, Vondracek B, Hatch LK (2014). "Application of neural networks to quantify the utility of indices of biotic integrity for biological monitoring." *Ecological Indicators*, **45**, 195–208.
- Garson GD (1991). "Interpreting neural network connection weights." *Artificial Intelligence Expert*, **6**(4), 46–51.
- Gevrey M, Dimopoulos I, Lek S (2003). "Review and comparison of methods to study the contribution of variables in artificial neural network models." *Ecological Modelling*, **160**(3), 249–264.
- Hornik K (1991). "Approximation capabilities of multilayer feedforward networks." *Neural Networks*, **4**(2), 251–257.
- Jain AK, Duin RPW, Mao JC (2000). "Statistical pattern recognition: A review." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **22**(1), 4–37.

- Kell DB, Oliver SG (2003). “Here is the evidence, now what is the hypothesis? The complementary roles of inductive and hypothesis-driven science in the post-genomic era.” *BioEssays*, **26**(1), 99–105.
- Kelling S, Hochachka WM, Fink D, Riedewald M, Caruana R, Ballard G, Hooker G (2009). “Data-intensive science: A new paradigm for biodiversity studies.” *BioScience*, **59**(7), 613–620.
- Lek S, Delacoste M, Baran P, Dimopoulos I, Lauga J, Aulagnier S (1996). “Application of neural networks to modelling nonlinear relationships in ecology.” *Ecological Modelling*, **90**(1), 39–52.
- Lek S, Guégan JF (2000). *Artificial Neuronal Networks: Application to Ecology and Evolution*. Springer-Verlag, Berlin, Germany.
- McCulloch WS, Pitts W (1943). “A logical calculus of the ideas imminent in nervous activity.” *Bulletin of Mathematical Biophysics*, **5**, 115–133.
- Olden JD, Jackson DA (2002). “Illuminating the “black box”: A randomization approach for understanding variable contributions in artificial neural networks.” *Ecological Modelling*, **154**(1-2), 135–150.
- Özesmi SL, Özesmi U (1999). “An artificial neural network approach to spatial habitat modelling with interspecific interaction.” *Ecological Modelling*, **116**(1), 15–31.
- Paruelo JM, Tomasel F (1997). “Prediction of functional characteristics of ecosystems: A comparison of artificial neural networks and regression models.” *Ecological Modelling*, **98**(2-3), 173–186.
- Ripley BD (1996). *Pattern Recognition and Neural Networks*. Cambridge University Press, Cambridge, United Kingdom.
- Rumelhart DE, Hinton GE, Williams RJ (1986). “Learning representations by back-propagating errors.” *Nature*, **323**(6088), 533–536.
- Saeys Y, Inza I, Laga PL (2007). “A review of feature selection techniques in bioinformatics.” *Bioinformatics*, **23**(19), 2507–2517.
- Swanson A, Kosmala M, Lintott C, Simpson R, Smith A, Packer C (2015). “Snapshot Serengeti: High-frequency annotated camera trap images of 40 mammalian species in African savanna.” *Scientific Data*, **2**, 150026.

**Affiliation:**

Marcus W. Beck

Oak Ridge Institute for Science and Education

US Environmental Protection Agency

National Health and Environmental Effects Research Laboratory

Gulf Ecology Division, 1 Sabine Island Drive

Gulf Breeze, Florida, 32561, USA

E-mail: [beck.marcus@epa.gov](mailto:beck.marcus@epa.gov)