# NEW YORK TAXI FARE CHALLENGE

**Till Bungert**

## ABSTRACT

## 1 INTRODUCTION

## 2 RELATED WORK

## 3 METHODS

### 3.1 ENTITY EMBEDDINGS

We map categorical variables with $C$ categories represented by indices $c \in \{0, 1, \ldots, C-1\}$ to real-numbered vectors $\mathbf{x}_c \in \mathbb{R}^n$

$$Embedding : \{0, 1, \ldots, C-1\} \to \mathbb{R}^n, c \mapsto Embedding(c) = \mathbf{x}_c. \tag{1}$$

These embedding layers are implemented as lookup tables. The vector associated with each index is a parameter of the model and is learnd jointly with the rest of the model.

If the input to our model is a mixture of continuous and categorical variables as is the case here, we learn one embedding layer for each of the categorical variables and concatenate the vector components of each embedding output together with the continuous variables to one vector. This concatenated vector then serves as the input to the rest of the model.
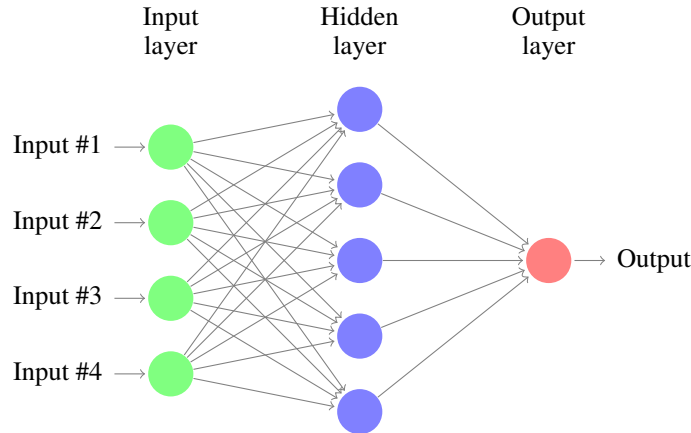
### 3.2 NEURAL NETWORKS



**Figure 1:** 2-Layer neural network

Feed-forward neural networks, sometimes called multilayer perceptrons, are one of many machine learning models designed to approximate some function $f^*$. They define a mapping $y = f(x, \theta)$ where $\theta$ is learned to result in the best approximation.

The name feed-forward neural network stems from the fact that they consist of intermediary functions called layers, that are chained together. The length of the chain of these intermediary functions

gives the depth of the network. As $f^*$, $f$ and all intermediary functions are vertor valued, the dimensionality of the vector gives the width of the layer. If a layer is not the input or output layer it is called hidden. By depicting each vector component as a node, neural networks can be described by directed acyclic graphs as in Figure 1 [2].

In most cases, each layer consists of a linear function $y = w^T x$, where $w$ are called the weights of that layer and are part of $\theta$. As a chain of linear functions is still a linear function, we need something to make the neural network nonlinear to learn general functions. To acomplish this, each layer is associated with a nonlinear function called the activation function, that is applied to the result of the linear function. Popular choices for activation functions are the rectified linear unit (ReLU)

$$ReLU(x) = \max(0, x) \tag{2}$$

and the softmax [2].

## 3.3  GAUSSIAN PROCESSES

Gaussian processes are nonparametric [1].

## 3.4  TREE-BASED METHODS

# 4  EXPERIMENTS

## 4.1  DATA SET

## 4.2  MODEL SETUPS

### 4.2.1  SIMPLE NEURAL NETWORK

### 4.2.2  DEEP NEURAL NETWORK

### 4.2.3  GAUSSIAN PROCESSES

### 4.2.4  GRADIENT BOOSTING

## 4.3  RESULTS

# 5  CONCLUSION

## REFERENCES

[1]  D. Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.

[2]  Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. `http : / / www . deeplearningbook.org`. MIT Press, 2016.