

# Implementation of Ranknet

Ryan Greenup

February 21, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Motivation</b>	<b>2</b>
<b>3</b>	<b>Implementation</b>	<b>2</b>
3.1	Neural Networks . . . . .	2
3.1.1	The Ranknet Method . . . . .	3
3.2	Creating Data <span style="float: right;">CF9AB26</span> . . . . .	3
3.3	Creating a Neural Network <span style="float: right;">7291112</span> . . . . .	5
3.4	Train the Model with Gradient Descent <span style="float: right;">7D46636</span> . . . . .	6
3.5	Implement Ranknet <span style="float: right;">F25F376:05DF04F</span> . . . . .	8
3.6	Implement sorting <span style="float: right;">99B390A:7D46636</span> . . . . .	11
3.6.1	Visualising Model Performance . . . . .	13
3.7	Enhancing the Model . . . . .	14
<b>4</b>	<b>Further Research</b>	<b>14</b>
4.1	Improve the sorting algorithm . . . . .	14
4.2	Evaluate Model Performance . . . . .	14
4.3	Evaluate alternative machine learning models . . . . .	15
<b>5</b>	<b>Conclusion</b>	<b>15</b>
<b>6</b>	<b>Appendix</b>	<b>15</b>
6.1	Search Engines . . . . .	15
6.1.1	Fuzzy String Match . . . . .	16
6.2	Import Statements . . . . .	16
6.3	Export Data Method . . . . .	17
6.4	Version Control Repository . . . . .	17
6.4.1	Version Control Log for Walkthrough . . . . .	18
# #+TODO: TODO IN-PROGRESS WAITING DONE		

## 1 Introduction

Ranknet is an approach to *Machine-Learned Ranking* (often referred to as "*Learning to Rank*" [28]) that began development at Microsoft from 2004 onwards [7], although previous work in this area had already been undertaken as early as the 90s (see generally [13, 12, 13, 14, 51]). These earlier models did not, however, perform well compared to more modern machine learning techniques [31, §15.5].

Information retrieval is an area that demands effective ranking of queries, although straight-forward tools such as `grep`, relational databases (e.g. `sqlite`, `MariaDB`, `PostgreSQL`) and NoSQL (e.g. `CouchDB`, `MongoDB`) can be used to retrieve documents with matching characters and words, these methods do not perform well in real word tasks across large collections of documents because they do not provide any logic to rank results (see generally [47]).

## 2 Motivation

Search Engines implement more sophisticated techniques to rank results, one such example being TF-IDF weighting [32], well established search engines such as *Apache Lucene* [2] and *Xapian* [20], however, are implementing Machine-Learned Ranking in order to improve results.

This paper hopes to serve as a general introduction to the implementation of the Ranknet technique to facilitate developers of younger search engines in more modern languages (i.e. Go and Rust). This is important because these more modern languages are more accessible [19] and memory safe [38] than C/C++ respectfully, without significantly impeding performance; this will encourage contributors from more diverse backgrounds and hence improve the quality of profession-specific tooling.

For a non-comprehensive list of actively maintained search engines, see §6.1 of the appendix.

## 3 Implementation

### 3.1 Neural Networks

The Ranknet method is typically implemented using a Neural Network<sup>1</sup>, although other machine learning techniques can also be used [7, p. 1].

Neural Networks are essentially a collection of different regression models and classifiers that are fed into one another to create a non-linear classifier, a loss function is used to measure the performance of the model with respect to the parameters (e.g. RMSE<sup>2</sup> or BCE<sup>3</sup>) and the parameters are adjusted so as to reduce this error by using

---

<sup>1</sup>An early goal of this research was to evaluate the performance of different machine learning algorithms to implement the Ranknet method, as well as contrasting this with simple classification approaches, this research however is still ongoing, see §4.3

<sup>2</sup>RMSE Root Mean Square Error

<sup>3</sup>BCE Binary Cross Entropy

the *Gradient Descent Technique* (although there are other optimisation algorithms such as RMSProp and AdaGrad [33] that can be shown to perform better [4]). The specifics of Neural Networks are beyond the scope of this paper (see [17] or more generally [41]).

### 3.1.1 The Ranknet Method

The Ranknet method is concerned with a value  $p_{ij}$  that measures the probability that an observation  $i$  is ranked higher than an observation  $j$ .

A Neural Network ( $n$ ) is trained to return a value  $s_k$  from a feature vector  $\mathbf{X}_k$ :

$$n(\mathbf{X}_i) = s_i \quad \exists k$$

So as to minimise the error of:

$$p_{ij} = \frac{1}{1 + e^{\sigma \cdot (s_i - s_j)}} \quad \exists \sigma \in \mathbb{R}$$

**Version Control** The implementation in this section (§ corresponds to the `walkthrough` branch of the `git` repository used in production of this work, id values (e.g. `:08db5b0:`) will be appended to titles to denote specific changes made in that section. See §6.4 for more specific guidance.

**Code listings** The code listings provided will use a standard set of import statements (see §6.2) and so they will be omitted from the listings, for more comprehensive guidance on implementing this code refer to the [documentation page](#)<sup>4</sup> that accompanies the `git` repo.

## 3.2 Creating Data

CF9AB26

The first step is to create a simple data set and design a neural network that can classify that data set, the data set generated should have two classes of data (this could be interpreted as relevant and irrelevant documents given the features or principle components of a data set), this can be implemented using `sci kit learn` as shown below and visualized<sup>5</sup> in figure 1.<sup>6</sup>

In order to fit a Neural Network the *PyTorch* package can be used [37], this will allow the gradients of the neural network to be calculated numerically without needing to solve for the partial derivatives, hence the data will need to be in the form of tensors.

<sup>4</sup>[crmds.github.io/CRMDS-HDR-Training-2020/](https://crmds.github.io/CRMDS-HDR-Training-2020/)

<sup>5</sup>See §6.3 for the specific method definition used to export the data to a `csv`.

<sup>6</sup>Visualisations for this Report were implemented using `org-babel` [11] inside *Emacs* [44] to call *R* [43] with *GGPlot2* [49] (and *Tidyverse* [50] generally), the source code for this is available in the report manuscript available in the `git` repository available at [github.com/RyanGreenup/ranknet/blob/main/Report/Report.org](https://github.com/RyanGreenup/ranknet/blob/main/Report/Report.org)

### 3 IMPLEMENTATION

```
1 def make_data(create_plot=False, n=1000, dtype=torch.float, dev="cpu", export=""):
2     X, y = datasets.make_blobs(n, 2, 2, random_state=7)
3     # X, y = datasets.make_moons(n_samples=n, noise=0.1, random_state=0) # Moons Data
4     ↪ for later
5
6     # Save the data somewhere if necessary
7     if export != "":
8         export_data(X, y, export)
9
10    # Reshape the data to be consistent
11    y = np.reshape(y, (len(y), 1)) # Make y vertical n x 1 matrix.
12
13    # -- Split data into Training and Test Sets -----
14    data = train_test_split(X, y, test_size=0.4)
15
16    if(create_plot):
17        # Create the Scatter Plot
18        plt.scatter(X[:, 0], X[:, 1], c=y)
19        plt.title("Sample Data")
20        plt.show()
21
22    # Make sure we're working with tensors not mere numpy arrays
23    torch_data = [None]*len(data)
24    for i in range(len(data)):
25        torch_data[i] = torch.tensor(data[i], dtype=dtype, requires_grad=False)
26
27    return torch_data
28
29 # Set Torch Parameters
30 dtype = torch.float
31 dev = test_cuda()
32
33 # Generate the Data
34 X_train, X_test, y_train, y_test = make_data(
35     n=int(300/0.4), create_plot=True, dtype=dtype, dev=dev, export =
36     ↪ "/tmp/simData.csv")
```

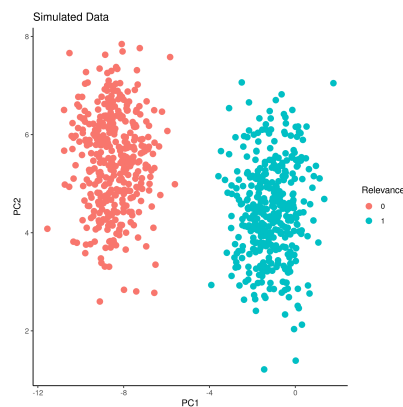


Figure 1: Generated data, output classes denote document relevance and the axis features or principle components

### 3.3 Creating a Neural Network

7291112

A Neural Network model can be designed as a class, here a 2-layer model using Sigmoid functions has been described, this design was chosen for it's relative simplicity:

```

1  class three_layer_classification_network(nn.Module):
2      def __init__(self, input_size, hidden_size, output_size, dtype=torch.float,
3          ↪ dev="cpu"):
4          super(three_layer_classification_network, self).__init__()
5          self.wi = torch.randn(input_size, hidden_size,
6                                  dtype=dtype,
7                                  requires_grad=True,
8                                  device=dev)
9          self.wo = torch.randn(hidden_size, output_size,
10                                 dtype=dtype,
11                                 requires_grad=True,
12                                 device=dev)
13
14          self.bi = torch.randn(hidden_size,
15                                  dtype=dtype,
16                                  requires_grad=True,
17                                  device=dev)
18          self.bo = torch.randn(output_size,
19                                  dtype=dtype,
20                                  requires_grad=True,
21                                  device=dev)
22
23          self.σ = torch.randn(1, dtype=dtype, requires_grad=True, device=dev)
24
25          self.losses = []      # List of running loss values
26          self.trainedQ = False # Has the model been trained yet?
27
28      def forward(self, x):
29          x = torch.matmul(x, self.wi).add(self.bi)
30          x = torch.sigmoid(x)
31          x = torch.matmul(x, self.wo).add(self.bo)
32          x = torch.sigmoid(x)
33          return x
34
35      def loss_fn(self, x, y):
36          y_pred = self.forward(x)
37          return torch.mean(torch.pow((y-y_pred), 2))
38
39      def misclassification_rate(self, x, y):
40          y_pred = (self.forward(x) > 0.5)
41          return np.average(y != y_pred)

```

A model can then be instantiated, a 2-3-1 model has, arbitrarily, been implemented in this case:<sup>7</sup>

```

1  # Set Seeds
2  torch.manual_seed(1)
3  np.random.seed(1)
4
5  # Set Torch Parameters

```

<sup>7</sup>note that the model has not yet been trained, the weights are random and the model output is not related to the data at all.

### 3 IMPLEMENTATION

---

```
6 dtype = torch.float
7 dev = test_cuda()
8
9 # Make the Data
10 X_train, X_test, y_train, y_test = make_data(
11     n=100, create_plot=True, dtype=dtype, dev=dev)
12
13 # Create a model object
14 model = three_layer_classification_network(
15     input_size=X_train.shape[1], hidden_size=2, output_size=1, dtype=dtype, dev=dev)
16
17 # Send some data through the model
18 print("\nThe Network input is:\n---\n")
19 print(X_train[7,:], "\n")
20 print("The Network Output is:\n---\n")
21 print(model.forward(X_train[7,:]).item(), "\n")
```

The Network input is:

---

tensor([-1.5129, 2.9332])

The Network Output is:

---

0.22973690927028656

#### 3.4 Train the Model with Gradient Descent

7D46636

Now that the model has been fit, a method to train the model can be implemented <sup>8</sup>:

```
1 class three_layer_classification_network(nn.Module):
2     # __init__ method goes here, see above
3     # ...
4     # ...
5
6     def train(self, x, target, η=30, iterations=2e4):
7         bar = Bar('Processing', max=iterations) # progress bar
8         for t in range(int(iterations)):
9
10             # Calculate y, forward pass
11             y_pred = self.forward(x)
12
13             # Measure the Loss
14             loss = self.loss_fn(x, target)
15
16             # print(loss.item())
17             self.losses.append(loss.item())
18
19             # Calculate the Gradients with Autograd
20             loss.backward()
21
22             with torch.no_grad():
23                 # Update the Weights with Gradient Descent
```

---

<sup>8</sup>This class definition is incomplete and serves only to show the method definition corresponding to the original class shown in §3.3

```

24         self.wi -= η * self.wi.grad; self.wi.grad = None
25         self.bi -= η * self.bi.grad; self.bi.grad = None
26         self.wo -= η * self.wo.grad; self.wo.grad = None
27         self.bo -= η * self.bo.grad; self.bo.grad = None
28         self.σ -= η * self.σ.grad; self.σ.grad = None
29         bar.next()
30         bar.finish()
31         # ; Zero out the gradients, they've been used
32
33         # Rest of the Class Definition Below ...VVV...

```

With this definition the model can hence be trained in order to produce meaningful classifications, as shown below, due to the simplicity of the data set, this model classifies the points perfectly on the testing set, the training error over time is shown in figure 2.

```

1  # Make the Data
2  X_train, X_test, y_train, y_test = make_data(
3      n=100, create_plot=True, dtype=dtype, dev=dev)
4
5  # Create a model object
6  model = three_layer_classification_network(
7      input_size=X_train.shape[1], hidden_size=2, output_size=1, dtype=dtype, dev=dev)
8
9  # Train the Model
10 model.train(X_train, y_train, η=1e-2, iterations=10000)
11
12 # Plot the Losses
13 plt.plot(model.losses)
14 plt.title("Losses at each training iteration")
15 plt.show()
16
17 print("The testing misclassification rate is:\n")
18 print(model.misclassification_rate(X_test, y_test))

```

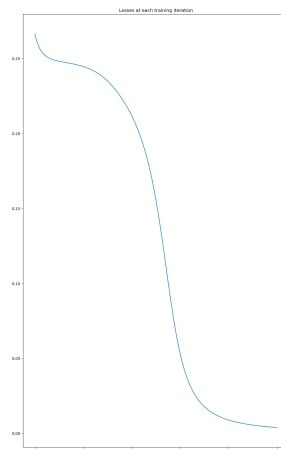


Figure 2: Training error, given by  $l(x) = \sum_{i=1}^n [(x_i - f(x_i))^2]$ , at each iteration of training

### 3.5 Implement Ranknet

F25F376:05DF04F

Now that the model can classify the data, the implementation will be modified to:

- Measure loss using a BCE function which is reported in the literature [7, 6] to perform better for Ranking problems.
- Modify the model so that it operates pairwise, such that:
  1. Two points are identified, sent through the neural network and two values returned:

$$s_i = n(\mathbf{X}_i) \quad (1)$$

$$s_j = n(\mathbf{X}_j) \quad (2)$$

The network previously created can be adapted for this and hence the method will be renamed to `forward_single` and this will represent function  $n()$  implemented in (1) and (2)

2. These values will be combined to give a single value which is intended to measure the model confidence:<sup>9</sup>

$$\hat{P}_{ij} = \text{sig}(\sigma, (s_i - s_j)), \quad \exists \sigma \in \mathbb{R} \quad (3)$$

$$= \frac{1}{1 + e^{\sigma \cdot (s_i - s_j)}} \quad (4)$$

3. The range of (4) is the interval  $\hat{P}_{ij} = [0, 1]$ , let  $\bar{P}_{ij}$  be the known probability<sup>10</sup> that  $\mathbf{X}_i \succ \mathbf{X}_j$ , the simulated data has a boolean range of  $\bar{P}_{ij} \in \{0, 1\}$ , this can be recast to  $\{-1, 0, 1\}$  and then linearly scaled to  $[0, 1]$  like so:

$$\bar{P}_{ij} \leftarrow p_i - p_j \quad (5)$$

$$\bar{P}_{ij} \leftarrow \frac{1 + \bar{P}_{ij}}{2} \quad (6)$$

These modifications only need to be made to the neural network class like so:

```

1 class three_layer_ranknet_network(nn.Module):
2     # __init__ method
3     # ...
4     # ...
5
6     def forward(self, xi, xj):
```

<sup>9</sup>This value is a measurement of the models "confidence" but could be extended to represent the "measured probability" of one item being ranked higher than an other (e.g. the probability that a person would rank one type of wine as better than the other in a random sample).

<sup>10</sup>Note the convention to that  $\triangleleft, \triangleright$  denote the ranking of two observations [6]



```

7         si = self.forward_single(xi)
8         sj = self.forward_single(xj)
9         out = 1 / (1 + torch.exp(-self.σ * (si - sj)))
10        return out
11
12    def forward_single(self, x):
13        x = torch.matmul(x, self.wi).add(self.bi)
14        x = torch.sigmoid(x)
15        x = torch.matmul(x, self.wo).add(self.bo)
16        x = torch.sigmoid(x)
17
18        return x
19
20    def loss_fn(self, xi, xj, y):
21        y_pred = self.forward(xi, xj)
22        loss = torch.mean(-y * torch.log(y_pred) -
23                          (1 - y) * torch.log(1 - y_pred))
24        return loss
25
26    def pairwise(iterable):
27        "pairwise([1,2,3,4]) --> [(1, 2), (1, 3), (1, 4), (2, 3), (2, 4), (3, 4)]"
28        s = List(iterable)
29        pair_iter = chain.from_iterable(combinations(s, r) for r in [2])
30        return pair_iter

```

The training method must be adapted to interact with these changes like so:<sup>11</sup>

```

1    class three_layer_ranknet_network(nn.Module):
2        # __init__ method
3        # ...
4        # ...
5        def train(self, x, target, η=1e-2, iterations=4e2):
6            self.trainedQ = True
7            # Create a progress bar
8            bar = Bar('Processing', max=iterations)
9            # Train for a number of iterations
10           for t in range(int(iterations)):
11               sublosses = []
12               # Loop over every pair of values
13               for pair in pairwise(range(len(x) - 1)):
14                   xi, yi = x[pair[0], :], target[pair[0]]
15                   xj, yj = x[pair[1], :], target[pair[1]]
16
17                   # encode from {0, 1} to {-1, 0, 1}
18                   y = yi - yj
19
20                   # Scale between {0,1}
21                   y = 1 / 2 * (1 + y)
22
23                   # Calculate y, forward pass
24                   y_pred = self.forward(xi, xj)
25
26                   # Measure the loss
27                   loss = self.loss_fn(xi, xj, y)
28                   sublosses.append(loss.item())
29

```

<sup>11</sup>Note the definition of the pairwise function, this was incorrectly implemented initially (f25f376) and rectified shortly after (05df04f). see §6.4.1

### 3 IMPLEMENTATION

```
30         # Calculate the Gradients with Autograd
31         loss.backward()
32
33         # Update the Weights with Gradient Descent
34         # ; Zero out the gradients, they've been used
35         with torch.no_grad():
36             self.wi -=  $\eta$  * self.wi.grad; self.wi.grad = None
37             self.bi -=  $\eta$  * self.bi.grad; self.bi.grad = None
38             self.wo -=  $\eta$  * self.wo.grad; self.wo.grad = None
39             self.bo -=  $\eta$  * self.bo.grad; self.bo.grad = None
40             self. $\sigma$  -=  $\eta$  * self. $\sigma$ .grad; self. $\sigma$ .grad = None
41
42         self.losses.append(np.average(sublosses))
43         bar.next()
44     bar.finish()
45     self.threshold_train(x, target, plot=False)
```

This can then be implemented as before, the loss function is provided at figure 3.

```
1  # Make the Data
2  X_train, X_test, y_train, y_test = make_data(
3      n=30, create_plot=True, dtype=dtype, dev=dev)
4
5  # Create a model object
6  model = three_layer_ranknet_network(
7      input_size=X_train.shape[1], hidden_size=2, output_size=1, dtype=dtype, dev=dev)
8
9  # Train the Model
10 model.train(X_train, y_train,  $\eta=1e-1$ , iterations=1e2)
11
12 # Save the Losses
13 np.savetxt(fname="/tmp/Losses.csv", X=model.losses, delimiter=',')
```

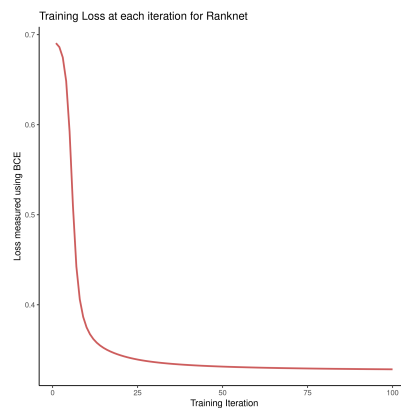


Figure 3: BCE training loss at each iteration for the Ranknet method.

### 3.6 Implement sorting

99B390A:7D46636

One of the difficulties in implementing this, however, is that it is not simple to determine whether or not the model has classified the data well,<sup>12</sup> In order to address this the model can be implemented to sort the data by ranked values and then visualised. To implement this a derivative of the *quicksort* algorithm was chosen as a sorting function [18], this was implemented by adapting code already available in the literature [24, §4.10] and online [42]:

```

1  def split(values, left, right, data, model):
2      # Define the leftmost value
3      l = (left-1)
4      # Set the right value as the pivot
5      pivot = values[right] # TODO The pivot should be random
6
7      for q in range(left, right):
8          # Only move smaller values left
9          if leq(values[q], pivot, data, model):
10             # +1 next left element
11             l = l+1
12             # Swap the current element onto the left
13             values[l], values[q] = values[q], values[l]
14
15     # Swap the pivot value into the left position from the right
16     values[l+1], values[right] = values[right], values[l+1]
17     return (l+1)
18
19
20 def qsort(values, left, right, data, model):
21     if len(values) == 1:
22         return values
23     if right > left:
24         # pi is the index of where the pivot was moved to
25         # It's position is now correct
26         pi = split(values, left, right, data, model)
27
28         # Do this again for the left and right parts
29         qsort(values, left, pi-1, data, model)
30         qsort(values, pi+1, right, data, model)
31
32 import random
33
34 def leq(a, b, data, model):
35     score = model.forward(data[a, :], data[b, :])
36     if score <= 0.5:
37         return True
38     if score > 0.5:
39         return False
40
41     if (a < b):
42         return True
43     else:
44         return False
45
46

```

<sup>12</sup>A naive misclassification method was implemented (f25f376), but it was not very insightful and so was omitted from this report.

```
47 if DEBUG:
48     for i in range(3):
49         import random
50         values = random.sample(range(9), 7)
51         n = len(values)
52         print(values)
53         qsort(values, 0, n-1, data, model)
54         print("=>", values)
```

The data can then be plotted, as in figure 4 and exported like so:<sup>13</sup>

```
1  # Main Function
2  def main():
3      # Make the Data
4      X_train, X_test, y_train, y_test = make_data(n=100,
5                                                    create_plot=True,
6                                                    dtype=dtype,
7                                                    dev=dev)
8
9      # Create a model object
10     model = three_layer_ranknet_network(input_size=X_train.shape[1],
11                                         hidden_size=2,
12                                         output_size=1,
13                                         dtype=dtype,
14                                         dev=dev)
15
16     # Train the Model
17     model.train(X_train, y_train, η=1e-1, iterations=1e2)
18
19     # Visualise the Training Error
20     plot_losses(model)
21
22     # Misclassification won't work for ranked data
23     # Instead Visualise the ranking
24     plot_ranked_data(X_test, y_test, model)
25
26
27 def plot_losses(model):
28     plt.plot(model.losses)
29     plt.title("Cost / Loss Function for Iteration of Training")
30     plt.show()
31
32
33 def plot_ranked_data(X, y, model):
34     # Create a list of values
35     n = X.shape[0]
36     order = [i for i in range(n)]
37     # Arrange that list of values based on the model
38     quicksort(values=order, left=0, right=(n - 1), data=X, model=model)
39     print(order)
40
41     ordered_data = X[order, :]
42     y_ordered = y[order]
43
44     np.savetxt("/tmp/ordered_data.csv", X=ordered_data.numpy(), delimiter=',')
45
46     p = plt.figure()
```

---

<sup>13</sup>In this case the plot has been generated by *GGPlot2* <sup>6</sup>

```

47     for i in range(len(ordered_data)):
48         plt.text(ordered_data[i, 0], ordered_data[i, 1], i)
49     plt.scatter(ordered_data[:, 0], ordered_data[:, 1], c=y_ordered)
50     plt.title("Testing Data, with ranks")
51     plt.show()
52
53
54 if __name__ == "__main__":
55     main()

```

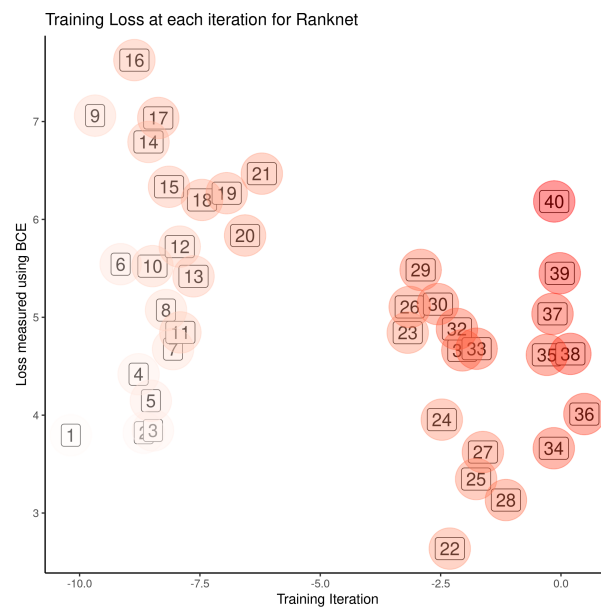


Figure 4: Using Machine Learned Ranking to order the points from most to least relevant

### 3.6.1 Visualising Model Performance

Although the results in figure 4 appear very encouraging at first, if the same sorting approach is implemented for an untrained model (i.e. a model with random weights), an uncomfortably similar, ordered pattern emerges, this is shown<sup>14</sup> in figure 5.

This pattern could be explained by the fact that the model, trained or untrained, is continuous and so the rank of nearby points could lead the emergence of an ordered pattern.

Investigating types of data where the Ranknet model will produce an ordered pattern only when trained would represent a logical next step. One approach would be to consider the rating of a product (e.g. the quality of wine, see [9]) and train a Ranknet model based only on whether or not one wine is better than the next. The order of the

<sup>14</sup>For Clarity sake the colours have been inverted.

returned results could be compared with the order of the original dataset as well as a random untrained control model in order to evaluate the efficacy of the model.

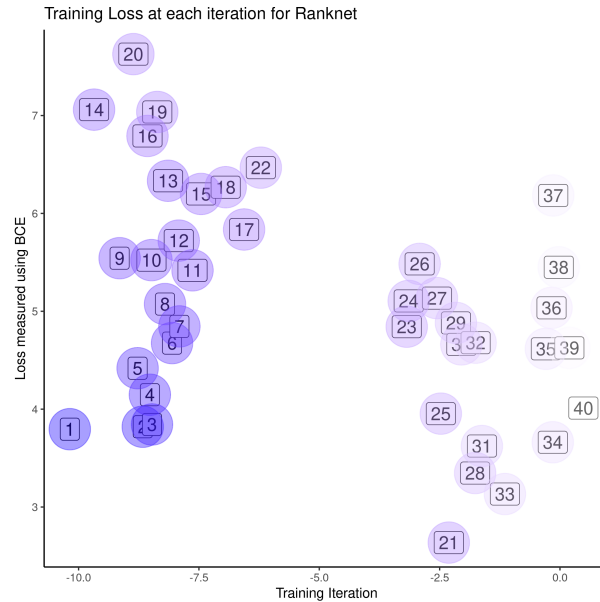


Figure 5: Using Machine Learned Ranking to order the points from most to least relevant

### 3.7 Enhancing the Model

Further enhancements (such as random batches and alternative optimisers) were made to the model and alternative datasets implemented (see `473dce`  $\rightarrow$  `d11e607`) however the results were mixed and not encouraging, rather an alternative type of data should first be considered and evaluated, as already discussed in §3.6.1.

## 4 Further Research

The implementation of this technique has proved considerably more difficult than first perceived, more work is required. In this section particular points of improvement are identified.

### 4.1 Improve the sorting algorithm

The "Quicksort" algorithm likely needs a random pivot to be efficient [45].

### 4.2 Evaluate Model Performance

It is still not clear how the performance of Ranknet compares to traditional approaches implemented by search engines (see §6.1), further study would ideally:

- Write a program to query a corpus of documents using an existing search engine.

- Or possibly just implement TF-IDF weighting in order to remove variables.
- Extend the program to implement machine learned ranking
- Measure and contrast the performance of the two models to see whether there are any significant improvements.

This could be implemented with TREC datasets [46] using a cumulated-gain cost function [22] as demonstrated in previous work [47].

It is not yet clear:

1. How this could be applied for a broad range of queries
2. How to effectively measure the "correctness" of ranked documents

### 4.3 Evaluate alternative machine learning models

An open question is whether or not alternative machine learning algorithms such as trees and SVMs could be used to fit the Ranknet model in an effective fashion, this could be ideal for performance concerns as deep neural networks can be resource intensive.

Another question is how Ranknet compares to simply using classification approaches to identify documents that might be relevant.

## 5 Conclusion

The Ranknet approach to machine-learned ranking looks promising, but it is difficult to implement and more study is needed to evaluate whether or not it brings significant advantages over traditional approaches and whether or not there are effective ways to apply the model to search problems for a broad scope of queries.

## 6 Appendix

### 6.1 Search Engines

There are many open source search engines available , a cursory review found the following popular projects:

- [Apache lucene/Solr](#) (Java) [2]
  - Implemented by [DocFetcher](#) [10]
- [Sphinx](#) (C++) [52]
- [Xapian](#) (C++) [36]
  - Implemented by [Recoll](#) [23]
- [Zettair](#) (C) [21]

More Modern Search engines include:

- [Bleve Search](#) (Go) [32]
- [Riot](#) (Go) [48]
- [LunrJS](#) (JS) [35]
- [SimSearch](#) (Rust) [29]
- [Tantivy](#) (Rust) [8]

### 6.1.1 Fuzzy String Match

Somewhat related are programs that rank string similarity, such programs don't tend to perform well on documents however (so for example these would be effective to filter document titles but would not be useful for querying documents):

- [fzf](#) [5]
- [fzy](#) [16]
- [go-fuzzyfinder](#) [26]
- [peco](#) [27]
- [Skim](#) [53]
- [Swiper](#) [25]

## 6.2 Import Statements

The following import statements were included, where used,<sup>15</sup> separate scripts were used to make the model as modular as possible, such corresponding inputs have also been listed:

```
1  # Import Packages
2  from itertools import chain
3  from itertools import combinations
4  from itertools import tee
5  from progress.bar import Bar
6  import math as m
7  import matplotlib.pyplot as plt
8  import numpy as np
9  import random
10 import sys
11 import sys
12 import torch
13 import torch
14 from torch import nn
15
16 # Sepereate Scripts Lcated below main
```

---

<sup>15</sup>Including import statements where they are not used is fine, other than complaints from a *linter* following *PEP* [34] (e.g. [autopep](#) [15]) the code will function just fine.



### 6.3 Export Data Method

```
17 from ranknet.test_cuda import test_cuda
18 from ranknet.make_data import make_data
19 from ranknet.neural_network import three_layer_ranknet_network
20 from ranknet.quicksort import quicksort
```

### 6.3 Export Data Method

The data was exported by printing the values to a text file like so:

```

1 def export_data(X, y, export):
2     try:
3         os.remove(export)
4         print("Warning, given file was over-written")
5     except:
6         pass
7
8     with open(export, "a") as f:
9         line = "x1, x2, y \n"
10        f.write(line)
11        for i in (range(X.shape[0])):
12            line = str(X[i][0]) + ", " + str(X[i][1]) + ", " + str(y[i]) + "\n"
13            f.write(line)
14    print("Data Exported")

```

## 6.4 Version Control Repository

The git repository used in production of this code is currently available on *GitHub* at [github.com/CRMDS/CRMDS-HDR-Training-2020](https://github.com/CRMDS/CRMDS-HDR-Training-2020), in order to get a local copy, execute the following commands (bash):

```
1  # Clone the repository
2  git clone https://github.com/CRMDS/CRMDS-HDR-Training-2020
3
4  # Change to the subdirectory
5  cd CRMDS-HDR-Training-2020/ranknet
6
7  # Checkout the Walkthrough branch
8  git checkout walkkthrough
9
10 # List the changes
11 git log
```

Consider the use of tools like `magit` [30] and `git-timemachine` [40] (or `GitLens` [1] and `git-temporal` [3] in VsCode) in order to effectively preview the changes at each step, alternatively a pager like `bat` [39] can also be used with something like `fzf` [5] like so:

```
1 git log | grep '^commit' | sed 's/^commit\ /\ /\ ' | \
2   fzf --preview 'git diff {}' ! | \
3   bat --color always'
```

#### 6.4.1 Version Control Log for Walkthrough

<b><i>Commit ID</i></b>	<b><i>Message</i></b>
ed5f4cf	<i>Initial Commit</i>
075acf9	<i>Walkthrough Initial Commit</i>
cf9ab26	<i>Generate data to use for classification</i>
7291112	<i>Create a Neural Network Model</i>
7d46636	<i>Implement gradient descent to train neural network</i>
f25f376	<i>Adapt Neural Network to perform Ranking</i>
42509ab	<i>Implement sorting algorithm to visualise ranking order</i>
05df04f	<i>Adapt Neural Network to perform Ranking</i>
99b390a	<i>Implement sorting algorithm to visualise ranking order</i>
473dce3	<i>Implement optimizer to replace mere gradient descent</i>
4141e92	<i>Train Model using Batches not entire dataset</i>
a2671a6	<i>Format code to make it more readable</i>
d11e607	<i>plot and only train on different ranked pairs</i>

## References

- [1] Eric Amodio. *Eamodio/Vscode-Gitlens*. Nov. 6, 2016. URL: <https://github.com/eamodio/vscode-gitlens> (visited on 02/20/2021) (cit. on p. 17).
- [2] Apache Software Foundation. *Learning To Rank | Apache Solr Reference Guide 6.6*. Solr Reference Guide. 2017. URL: [https://lucene.apache.org/solr/guide/6\\_6/learning-to-rank.html](https://lucene.apache.org/solr/guide/6_6/learning-to-rank.html) (visited on 02/20/2021) (cit. on pp. 2, 15).
- [3] Bee Wilkerson. *Git-Temporal/Git-Temporal: The Mono Repo for Components, Data Layer, and Interfaces of Git-Temporal*. Sept. 12, 2018. URL: <https://github.com/git-temporal/git-temporal> (visited on 02/20/2021) (cit. on p. 17).
- [4] Vitaly Bushaev. *Understanding RMSprop — Faster Neural Network Learning*. Medium. URL: <https://towardsdatascience.com/understanding-rmsprop-faster-neural-network-learning-62e116fcf29a> (visited on 02/16/2021) (cit. on p. 3).
- [5] Junegunn Choi. *Junegunn/Fzf*. Feb. 20, 2021. URL: <https://github.com/junegunn/fzf> (visited on 02/20/2021) (cit. on pp. 16, 17).
- [6] Christopher Burges. *From RankNet to LambdaRank to LambdaMART: An Overview (MSR-TR-2010-82)*. Jan. 1, 2010. URL: <https://www.microsoft.com/en-us/research/uploads/prod/2016/02/MSR-TR-2010-82.pdf> (cit. on p. 8).
- [7] Christopher Burges. *RankNet: A Ranking Retrospective*. Microsoft Research. July 7, 2015. URL: <https://www.microsoft.com/en-us/research/blog/ranknet-a-ranking-retrospective/> (visited on 02/15/2021) (cit. on pp. 2, 8).
- [8] Clement Renault, Marin, and Quentin de Quelen. *Meilisearch/MeiliSearch*. MeiliSearch, Feb. 20, 2021. URL: <https://github.com/meilisearch/MeiliSearch> (visited on 02/20/2021) (cit. on p. 16).
- [9] A Cortez et al. “Modeling Wine Preferences by Data Mining from Physicochemical Properties. In Decision Support Systems”. In: *Elsevier* 47.4 (2009), pp. 547–553 (cit. on p. 13).
- [10] Docfetcher Development Team. *DocFetcher - Fast Document Search*. URL: <http://docfetcher.sourceforge.net/en/index.html> (visited on 02/15/2021) (cit. on p. 15).
- [11] Carsten Dominik. *Org Mode Reference Manual*. Place of publication not identified: 12TH MEDIA SERVICES, 2018. ISBN: 978-1-68092-165-6 (cit. on p. 3).
- [12] Norbert Fuhr. “Optimum Polynomial Retrieval Functions Based on the Probability Ranking Principle”. In: *ACM Transactions on Information Systems* 7.3 (July 1, 1989), pp. 183–204. ISSN: 1046-8188. DOI: [10.1145/65943.65944](https://doi.org/10.1145/65943.65944). URL: <https://doi.org/10.1145/65943.65944> (visited on 02/15/2021) (cit. on p. 2).
- [13] Norbert Fuhr. “Probabilistic Models in Information Retrieval”. In: *The Computer Journal* 35.3 (June 1, 1992), pp. 243–255. ISSN: 0010-4620. DOI: [10.1093/comjnl/35.3.243](https://doi.org/10.1093/comjnl/35.3.243). URL: <https://doi.org/10.1093/comjnl/35.3.243> (visited on 02/15/2021) (cit. on p. 2).

- [14] Fredric C. Gey. “Inferring Probability of Relevance Using the Method of Logistic Regression”. In: *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '94. Berlin, Heidelberg: Springer-Verlag, Aug. 1, 1994, pp. 222–231. ISBN: 978-0-387-19889-7 (cit. on p. 2).
- [15] Hideo Hattori. *Autopep8: A Tool That Automatically Formats Python Code to Conform to the PEP 8 Style Guide*. Version 1.5.5. URL: <https://github.com/hhatto/autopep8> (visited on 02/20/2021) (cit. on p. 16).
- [16] John Hawthorn. *Jhawthorn/Fzy*. Feb. 20, 2021. URL: <https://github.com/jhawthorn/fzy> (visited on 02/20/2021) (cit. on p. 16).
- [17] HMKCode. *Backpropagation Step by Step*. URL: <https://hmkcode.com/ai/backpropagation-step-by-step/> (visited on 02/16/2021) (cit. on p. 3).
- [18] C. A. R. Hoare. “Algorithm 64: Quicksort”. In: *Communications of the ACM* 4.7 (July 1, 1961), p. 321. ISSN: 0001-0782. DOI: [10.1145/366622.366644](https://doi.org/10.1145/366622.366644). URL: <http://doi.org/10.1145/366622.366644> (visited on 02/21/2021) (cit. on p. 11).
- [19] Daniel Hunt. “A Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of Master of Science in Computer Science”. In: (), p. 30 (cit. on p. 2).
- [20] James Aylett. *GSoCProjectIdeas/LearningtoRankStabilisation – Xapian*. Oct. 20, 2019. URL: <https://trac.xapian.org/wiki/GSoCProjectIdeas/LearningtoRankStabilisation> (visited on 02/20/2021) (cit. on p. 2).
- [21] Dr Rolf Jansen. *Cyclaero/Zettair*. Dec. 22, 2020. URL: <https://github.com/cyclaero/zettair> (visited on 02/20/2021) (cit. on p. 15).
- [22] Kalervo Järvelin and Jaana Kekäläinen. “Cumulated Gain-Based Evaluation of IR Techniques”. In: *ACM Transactions on Information Systems* 20.4 (Oct. 1, 2002), pp. 422–446. ISSN: 1046-8188. DOI: [10.1145/582415.582418](https://doi.org/10.1145/582415.582418). URL: <http://doi.org/10.1145/582415.582418> (visited on 02/20/2021) (cit. on p. 15).
- [23] Jean-Francois Dockes. *Recoll User Manual*. URL: <https://www.lesbonscomptes.com/recoll/usermanual/usermanual.html> (visited on 02/15/2021) (cit. on p. 15).
- [24] Brian W. Kernighan and Dennis M. Ritchie. *The C Programming Language*. 2nd ed. Englewood Cliffs, N.J: Prentice Hall, 1988. 272 pp. ISBN: 978-0-13-110370-2 978-0-13-110362-7 (cit. on p. 11).
- [25] Oleh Krehel. *Abo-Abo/Swiper*. Feb. 20, 2021. URL: <https://github.com/abo-abo/swiper> (visited on 02/20/2021) (cit. on p. 16).
- [26] ktr. *Ktr0731/Go-Fuzzyfinder*. Feb. 16, 2021. URL: <https://github.com/ktr0731/go-fuzzyfinder> (visited on 02/20/2021) (cit. on p. 16).
- [27] lestrrat. *Peco/Peco*. peco, Feb. 18, 2021. URL: <https://github.com/peco/peco> (visited on 02/20/2021) (cit. on p. 16).
- [28] Tie-Yan Liu. “Learning to Rank for Information Retrieval”. In: *Foundations and Trends® in Information Retrieval* 3.3 (June 26, 2009), pp. 225–331. ISSN: 1554-0669, 1554-0677. DOI: [10.1561/15000000016](https://doi.org/10.1561/15000000016). URL: <https://www.nowpublishers.com/article/Details/INR-016> (visited on 02/15/2021) (cit. on p. 2).

## REFERENCES

---

- [29] Andy Lok. *Andylokandy/Simsearch-Rs*. Feb. 15, 2021. URL: <https://github.com/andylokandy/simsearch-rs> (visited on 02/20/2021) (cit. on p. 16).
- [30] Magit/Magit. Magit, Oct. 19, 2008. URL: <https://github.com/magit/magit> (visited on 02/20/2021) (cit. on p. 17).
- [31] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. New York: Cambridge University Press, 2008. 482 pp. ISBN: 978-0-521-86571-5 (cit. on p. 2).
- [32] Marty Schoch and Abhinav Dangeti. *Bleve Search Documentation*. URL: <http://blevesearch.com/> (visited on 02/20/2021) (cit. on pp. 2, 16).
- [33] Mahesh Chandra Mukkamala and Matthias Hein. “Variants of RMSProp and Adagrad with Logarithmic Regret Bounds”. In: *Proceedings of the 34th International Conference on Machine Learning - Volume 70. ICML’17*. Sydney, NSW, Australia: JMLR.org, Aug. 6, 2017, pp. 2545–2553 (cit. on p. 3).
- [34] Nick Coghlan. *PEP 8 – Style Guide for Python Code*. Python.org. Aug. 1, 2001. URL: <https://www.python.org/dev/peps/pep-0008/> (visited on 02/20/2021) (cit. on p. 16).
- [35] Oliver Nightingale. *Olivernn/Lunr.js*. Feb. 20, 2021. URL: <https://github.com/olivernn/lunr.js> (visited on 02/20/2021) (cit. on p. 16).
- [36] Olly Betts and Richard Boulton. *Xapian/Xapian*. Xapian, Feb. 19, 2021. URL: <https://github.com/xapian/xapian> (visited on 02/20/2021) (cit. on p. 15).
- [37] Adam Paszke et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In: *Advances in Neural Information Processing Systems 32*. Ed. by H. Wallach et al. Curran Associates, Inc., 2019, pp. 8024–8035. URL: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf> (cit. on p. 3).
- [38] Jeffrey M. Perkel. “Why Scientists Are Turning to Rust”. In: *Nature* 588.7836 (7836 Dec. 1, 2020), pp. 185–186. DOI: [10.1038/d41586-020-03382-2](https://doi.org/10.1038/d41586-020-03382-2). URL: <https://www.nature.com/articles/d41586-020-03382-2> (visited on 02/20/2021) (cit. on p. 2).
- [39] David Peter. *Sharkdp/Bat*. Apr. 21, 2018. URL: <https://github.com/sharkdp/bat> (visited on 02/20/2021) (cit. on p. 17).
- [40] Peter Stiernström and Basti. *Emacsmirror/Git-Timemachine*. Emacsmirror, June 29, 2014. URL: <https://github.com/emacsmirror/git-timemachine> (visited on 02/20/2021) (cit. on p. 17).
- [41] Philip Picton. *Neural Networks*. Basingstoke, Hampshire ; New York: Palgrave, 1994. 195 pp. ISBN: 978-0-333-94899-6 (cit. on p. 3).
- [42] *Python Program for QuickSort*. GeeksforGeeks. Jan. 7, 2014. URL: <https://www.geeksforgeeks.org/python-program-for-quicksort/> (visited on 02/15/2021) (cit. on p. 11).
- [43] R Core Team. *R: A Language and Environment for Statistical Computing*. manual. R Foundation for Statistical Computing. Vienna, Austria, 2021. URL: <https://www.R-project.org/> (cit. on p. 3).

## REFERENCES

---

- [44] Richard Stallman. *GNU Emacs Manual*. 15.ed., updated for Emacs version 21.2. Boston, Mass: Free Software Foundation, 2002. 602 pp. ISBN: 978-1-882114-85-6 (cit. on p. 3).
- [45] Tim Roughgarden, director. *Quicksort Overview*. Jan. 28, 2017. URL: [https://www.youtube.com/watch?v=ETo1cpLN7kk&list=PLEAYkSg4uSQ37A6\\_NrUnTHEKp6EkAxTma&index=25](https://www.youtube.com/watch?v=ETo1cpLN7kk&list=PLEAYkSg4uSQ37A6_NrUnTHEKp6EkAxTma&index=25) (visited on 02/15/2021) (cit. on p. 14).
- [46] US National Institute of Standards and Technology. *Text REtrieval Conference (TREC) Home Page*. URL: <https://trec.nist.gov/> (visited on 02/20/2021) (cit. on p. 15).
- [47] Vik Singh. *A Comparison of Open Source Search Engines*. Vik's Blog. July 6, 2009. URL: <https://partyondata.com/2009/07/06/a-comparison-of-open-source-search-engines-and-indexing-twitter/> (visited on 02/20/2021) (cit. on pp. 2, 15).
- [48] vz. *Go-Ego/Riot*. ego, Feb. 20, 2021. URL: <https://github.com/go-ego/riot> (visited on 02/20/2021) (cit. on p. 16).
- [49] Hadley Wickham. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN: 978-3-319-24277-4. URL: <https://ggplot2.tidyverse.org> (cit. on p. 3).
- [50] Hadley Wickham et al. "Welcome to the tidyverse". In: *Journal of Open Source Software* 4.43 (2019), p. 1686. DOI: [10.21105/joss.01686](https://doi.org/10.21105/joss.01686) (cit. on p. 3).
- [51] S. K.M. Wong and Y. Y. Yao. "Linear Structure in Information Retrieval". In: *Proceedings of the 11th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '88. New York, NY, USA: Association for Computing Machinery, May 1, 1988, pp. 219–232. ISBN: 978-2-7061-0309-4. DOI: [10.1145/62437.62452](https://doi.org/10.1145/62437.62452). URL: <https://doi.org/10.1145/62437.62452> (visited on 02/14/2021) (cit. on p. 2).
- [52] Yuri Schapov, Andrew Aksyonoff, and Ilya Kuznetsov. *Sphinxsearch/Sphinx*. Sphinx Technologies Inc, Feb. 18, 2021. URL: <https://github.com/sphinxsearch/sphinx> (visited on 02/20/2021) (cit. on p. 15).
- [53] Jinzhou Zhang. *Lotabout/Skim*. Feb. 19, 2021. URL: <https://github.com/lotabout/skim> (visited on 02/20/2021) (cit. on p. 16).