

Ranking Loss Surrogates

MSc Thesis

Abdus Salam Khazi [[Email](#)]

[Github Repository](#) [[11](#)]

Supervisors: JProf. Josif Grabocka & Sebastian Pineda

May 6, 2022

Abstract

Abstract goes here

Contents

1	Introduction	1
1.1	Hyperparameter Optimization	1
1.2	HPO Constraints	2
1.3	Objective	3
1.4	Overview	3
2	Related Work	4
2.1	Hyper Parameter Optimization	4
2.1.1	Black-box HPO	4
2.1.2	Multi-fidelity HPO,	4
2.1.3	Online HPO	4
2.2	Transfer-learning for HPO	4
2.2.1	Meta-learning of surrogates	4
2.2.2	Warm starting of initial configurations	4
2.3	Types of Surrogates for BO in HPO	4
2.3.1	GP	4
2.3.2	RF	4
2.3.3	Bayesian Neural Networks	4
2.3.4	Neural Networks	4
2.4	Types of Losses	5
2.4.1	BO with GP uses negative log likelihood	5
2.4.2	Ranking loss could be the answer in modeling hp optima better	5
2.5	Set-modeling with Neural Networks	5
2.5.1	Set-transformers	5
2.5.2	Deep Sets	5
2.6	Bayesian Optimization	7
2.6.1	Probabilistic models	8
2.6.2	Acquisition functions	10

2.7	Deep Ensembles	11
2.8	Deep Kernel Learning	12
2.9	Rank Learning	12
2.10	Deep Sets	13
3	Background	14
3.1	Benchmarking Meta-Data	14
3.1.1	Benchmark evaluation	15
3.2	Baselines	15
3.2.1	Deep Ensemble	16
3.2.2	Few Shot BO: Deep Kernel Learning	18
3.3	Evaluation Strategy	19
4	Method	20
4.1	Proposed Idea: : Ranking Loss Surrogates	20
4.2	Learning mechanism	21
4.2.1	Loss functions: Analysis and implementation	21
4.2.2	ListNet	22
4.3	Case study	24
4.3.1	Sorting with inverse mapping	24
4.3.2	Observations	25
4.3.3	Scoring function analysis	25
4.3.4	Scoring function range	25
4.4	Optimization cycle	25
4.4.1	Meta training	26
4.4.2	Fine tuning	26
4.5	Ranking Surrogate Model	26
4.5.1	Transfer Learning	26
4.5.2	Use of Deep Sets	26
4.5.3	Uncertainty implementation	26
4.6	Weighted Loss	27
4.7	Advantages and Limitations	27
4.7.1	Negative transfer learning	28
5	Research Question	30
6	Experiments and Results	31
6.1	Evaluation	31
6.1.1	Testing	31
6.1.2	Ablation	31

7 Conclusion	33
7.1 Further work	33
7.2 Conclusion	33
A More information	ii

Chapter 1

Introduction

The performance of any machine learning model is sensitive to the hyper-parameters used during the model training. Instead of using a new model type, it is more helpful to tune the hyper-parameters of an existing model to improve its performance. Learning the best hyper-parameter for an ML model is called, Hyperparameter optimization (HPO in short). This thesis studies various existing approaches to HPO and proposes a new idea for the same using the concept of ranking. The proposed idea in this thesis is called **Hyperparameter Optimization using Ranking Loss Surrogates**. The results obtained using this model are compared against the state-of-the-art results obtained using models like FSBO, RGPE, TAF, and others.

1.1 Hyperparameter Optimization

To find out the best hyper-parameter for any machine learning model m , we must first quantify a given hyper-parameter configuration \mathbf{x} by a real-valued number $v \in \mathbb{R}$. If we define that

$$\mathbf{x}_1 \succ \mathbf{x}_2 \iff v_{\mathbf{x}_1} < v_{\mathbf{x}_2}$$

then HPO can be defined mathematically by an abstract function, say, $f(\mathbf{x}) \mapsto \mathbb{R}$ as

$$\underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{S}$$

where \mathbb{S} is the hyper-parameter search space.

This function $f(\mathbf{x}) \mapsto \mathbb{R}$ is evaluated in the following chronological steps:

1. Using a given hyper-parameter configuration \mathbf{x} , we train our model m to obtain the model $m_{\mathbf{x}}^{\text{trained}}$. It consists of learning the parameters

of our model, E.g. learning the weights and biases of a Deep Neural Network. We use the training data to learn this model.

2. The validation data is passed through $m_{\mathbf{x}}^{trained}$ to obtain the required results. These results are evaluated based on an evaluation criterion 'eval'. This criterion is different for different problems, e.g. Regression, Classification, etc. The result of this evaluation is a real-value that gives a score for the configuration \mathbf{x} .

Hence the function $f(\mathbf{x}) \mapsto \mathbb{R}$ can be written as

$$\text{eval}(m_{\mathbf{x}}^{trained}(\text{Data}_{\text{val}})) \mapsto \mathbb{R}$$

Finally, the HPO problem can be defined using the following equation:

$$\underset{\mathbf{x}}{\text{argmin}} \text{ eval}(m_{\mathbf{x}}^{trained}(\text{Data}_{\text{val}})) \mapsto \mathbb{R} \quad \forall \mathbf{x} \in \mathbb{S}$$

1.2 HPO Constraints

Hyperparameter optimization is different from other optimization methods because it has different constraints. It is because of the peculiar properties of the hyper-parameter search spaces. Finding out the correct hyper-parameter setting is generally not feasible using a brute-force approach (trying out all possible combinations of hyper-parameters) because the search space itself has many dimensions, and the search space may be continuous. More specifically, some of the important constraints of this optimization problem are:

1. The evaluation of a given HPO configuration is computationally expensive.
2. It is a non-convex optimization problem.
3. The process of getting $m_{\mathbf{x}}^{trained}$ from m is stochastic hence the value $v_{\mathbf{x}}$ is noisy.
4. Some dimensions have conditional constraints. The values of some dimensions may depend on the values of others. For example, the number of neurons in layer 3 only makes sense if we have 3 or more layers in a neural network.
5. The search space is hybrid in terms of continuity. Some of the dimensions (variables) may be continuous, while others may be discrete. Using a gradient method is hence not trivial.

1.3 Objective

The aim is to study ranking loss surrogates in the context of Hyperparameter optimization.

1.4 Overview

This sections contains the overview of the paper and how the thesis report is organised.

Chapter 2

Related Work

2.1 Hyper Parameter Optimization

2.1.1 Black-box HPO

Bayesian Optimization

2.1.2 Multi-fidelity HPO,

2.1.3 Online HPO

2.2 Transfer-learning for HPO

2.2.1 Meta-learning of surrogates

2.2.2 Warm starting of initial configurations

2.3 Types of Surrogates for BO in HPO

2.3.1 GP

2.3.2 RF

2.3.3 Bayesian Neural Networks

2.3.4 Neural Networks

they don't work great, because they cant model uncertainty). Uncertainty can be modeled directly, or through ensembles.

2.4 Types of Losses

2.4.1 BO with GP uses negative log likelihood

It is not clear whether pointwise losses (regression as in GP) are the correct way to model HPO responses, because we only care for the minima regions (best performing configurations), and not for estimating all observations accurately.

2.4.2 Ranking loss could be the answer in modeling hp optima better

Pointwise

Pairwise

Listwise

2.5 Set-modeling with Neural Networks

In HPO we make predictions based on a set of observations, therefore, in this thesis, I explore the benefit of contextualizing surrogate models on the observations so far.

2.5.1 Set-transformers

2.5.2 Deep Sets

To deal with the constraints of HPO problems, researchers have used different strategies for developing HPO algorithms and models. Some of the straightforward methods include

- Manual Search
- Grid Search
- Random Search

Manual Search in the HPO search space is feasible when we have expert knowledge of the problem and the given model. The idea is to select configurations step by step by observing the results obtained so that we do not waste computation time evaluating similar configurations through intuition. This approach may be helpful for small models with lesser constraints.

However, as the HPO search space becomes very large or conditional constraints become too complex, the selection of configurations becomes more and more difficult. Hence a more systematic and automated approach is more practical.

Grid search is a more systematic approach in which we divide the search space into grid points similar to ones in a Euclidean graph. Let there be m dimensions in the search space \mathbb{S} . Let the selected number of values for each dimension be n_1, n_2, \dots, n_m . In the case of a discrete variable, the selected values will be a subset of the possible values, whereas, in the case of a continuous variable, we need to select the values based on a selected granularity. The cross-product of the selected subsets gives us the configurations to be evaluated. Hence, the number of configurations to evaluate will be $n_1 * n_2 * \dots * n_m$. The number of configurations we need to evaluate in this approach becomes a big issue for this method as the dimensions of the search spaces increase. Hence this approach becomes intractability for large search spaces.

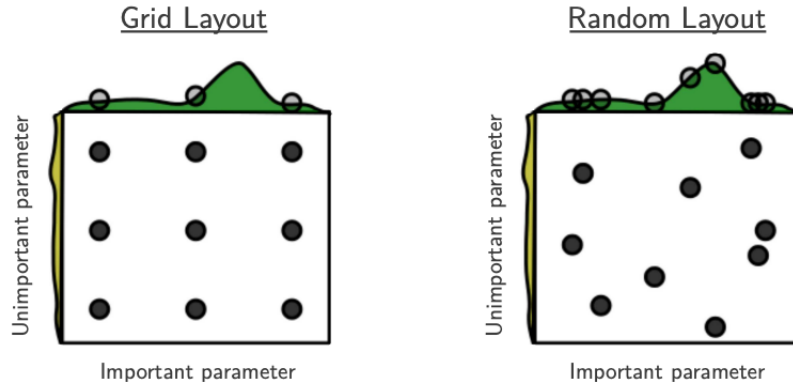


Figure 2.1: Illustrates Grid search and Random search in the case where 2 parameters are not equally important. Adapted from [2].

One issue with the Grid Search approach is that we assume that all dimensions in the HPO search space are equally important. It is not the case in many HPO problems. The Grid layout in Figure 2.1(left) shows illustrates this. For example, the learning rate in deep neural networks is much more important than many other parameters. If dimension p is the most important in the search space, then it makes sense to evaluate more values of p . Random Search helps us solve this problem. The Random layout

in Figure 2.1(right) illustrates this. Hence Random Search can be used as a trivial baseline for comparing other HPO models.

One advantage of these methods is that there are no restrictions on the HPO search spaces. Hence, they are suitable for any HPO problem at hand. On the other hand, these methods are non-probabilistic. Hence they cannot deal with noisy evaluations of the HPO configuration well. Moreover, these methods are computationally expensive. The reason is that they do not use surrogate evaluators and hence train and evaluate the whole model. Also, these search methods give us optimal HPO configurations only by chance.

In the remainder of this section, we discuss some sophisticated probabilistic models for doing HPO that use surrogates to reduce computational costs.

2.6 Bayesian Optimization

Bayesian optimization tries to solve both computational costs and noisy evaluations of our objective function (section 1). It does this by building a model of the HPO objective function. This model is called a surrogate function. Bayesian optimization uses known evaluations as its data to build the surrogate model. The data is of the form $x, f(x)$ pairs. The surrogate model is a probabilistic model. Hence, it also learns about the noise in the evaluations of the objective function.

The core procedure of the optimization process is the following:

- From known data $D = (x_1, f(x_1)), (x_2, f(x_2)), (x_3, f(x_3)), \dots$, build a probabilistic model that learns the mean and variance of the objective function
- Use the surrogate to sample the next best HPO configuration x' using a function known as acquisition function. Evaluate $f(x')$.
- Append $(x', f(x'))$ to D and repeat the process.

The above process repeats till the computational resources are finished (here time) or we find an acceptable HPO configuration. This procedure is also called SMBO (Sequential model-based Optimization). The procedure alternates between collecting data and fitting the model with the collected data [9].

Hence, there are two essential components of Bayesian optimization:

- Probabilistic model of the objective function

- The acquisition function

There are many probabilistic models such as Random Forests, Gaussian Processes, Tree parson Estimators, etc. But for the purpose of this thesis, we briefly mention Random forests and discuss in-depth the Gaussian processes.

2.6.1 Probabilistic models

Random Regression Forest

The core idea of this model is to train a Random Regression Forest, using the known data as in any SMBO procedure [9]. Random regression forests are an ensemble of regression trees. This property is used to our advantage to predict the mean and the variance. The mean of the prediction of all the trees is the mean of the surrogate model. The variance in the prediction of all trees is the variance of the surrogate model.

The advantages of this model are

- It can handle both continuous and discrete variables trivially without any modifications to the model. The data splitting during training is done using any variable be it discrete or continuous.
- It can handle conditional variables, unlike Gaussian processes, by making sure that data is not split based on a variable till it is guaranteed that no conditionality is broken by the split.

Gaussian Process Regression

Gaussian processes [20] are predictive machine learning models that work well with few data points (or data pairs). They are inherently capable of modeling uncertainty. Hence, they are used widely in problems such as hyperparameter optimization, where uncertainty estimation is essential. In this section, we briefly explain the Gaussian process regression intuitively.

Before we proceed, we need to understand normal (Gaussian) distributions. Consider a scalar random variable X that is distributed normally (a.k.a Gaussian distribution) around a mean μ with a variance of σ^2 . The following equation defines the probability density function (PDF) of X :

$$P_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Here, X represents the random variable, and x represents an instance of the variable [20]. In this case, the mean μ , variance σ^2 , and any sample x are all scalars.

If the random variable \mathbf{X} is a vector in \mathbb{R}^d where $d \in I^+$, then each component of the vector can be considered as a random variable. In this case the mean $\boldsymbol{\mu} \in \mathbb{R}^d$ whereas variance, represented by Σ , is in the $R^{d \times d}$ space. It is because the variance of all components in any valued vector random variable \mathbf{X} should contain the following two types of variance

- Variance of a vector component w.r.t itself. d diagonal values of the matrix Σ represent this variance.
- Variance of each vector component w.r.t all other components. These variances are represented by the upper/lower triangular values in the matrix Σ .

The matrix Σ , also known as the Covariance matrix, thus has all values necessary to represent the variance of any vector-valued random variable.

The probability density function of a vector valued variable $\mathbf{X} \in \mathbb{R}^d$ with a mean $\boldsymbol{\mu}$ and covariance matrix Σ is given by [14]:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{d}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

This equation defines the PDF of a multivariate normal distribution.

The core idea used in the Gaussian processes is that functions can be considered as vectors of infinite dimensions. Consider any function f that has a domain \mathbb{R} . If f is considered to be a vector in \mathbb{R}^∞ , then each point $i \in \mathbb{R}$ can be represented by a component f_i of the function f . A function, hence, is nothing but a sample from \mathbb{R}^∞ . Unfortunately, functions sampled from \mathbb{R}^∞ are too general and not useful by themselves.

The idea of Gaussian processes is to sample smooth functions from \mathbb{R}^∞ . In any smooth function f , if any point g is close to x in the domain of f , then $f(g) \approx f(x)$. It is mathematically represented by the following equation:

$$\lim_{\delta x \rightarrow 0} f_{x+\delta x} \approx f_x^+ \quad \text{and} \quad \lim_{\delta x \rightarrow 0} f_{x-\delta x} \approx f_x^-$$

$$\text{where } \delta x > 0 \text{ and } x, \delta x \in \mathbb{R}$$

The above definition is nothing but the definition of a smooth function in terms of vector notation. Moreover, nearby components of f "vary" similarly w.r.t each other. These properties can be naturally encoded using a covariance matrix. Hence, we obtain smooth functions if we sample them from a multivariate normal distribution with the required covariance matrix. The Gaussian process restricts the function sample space to a multivariate normal distribution.

The similarity between 2 points in a domain is defined by a function called **kernel** in Gaussian processes. Using this kernel function, the values in the required covariance matrix are populated. The smoothness of the sampled function f is controlled by the kernel in the GP process. Formally kernel k is defined as,

$$k(\mathbf{x}, \mathbf{x}') \mapsto \mathbb{R}$$

Here, \mathbf{x}, \mathbf{x}' belong to a domain in the most abstract sense. For example, when the input domain is a euclidean space, $\mathbf{x} \in \mathbb{R}^{\mathbb{I}^+}$.

Some well known kernels are:

- **Radial Basis Function Kernel**
- **Matern Kernel**
- **Periodic Kernel**

Finally, a Gaussian Process specifies that any new observation y^* for input \mathbf{x}^* , is jointly normally distributed with known observations \mathbf{y} (corresponding to the input \mathbf{X}) such that

$$Pr \left(\begin{bmatrix} \mathbf{y} \\ y^* \end{bmatrix} \right) = \mathcal{N}(m(\mathbf{X}), \Sigma) \quad (2.1)$$

Here, $m(\mathbf{X})$ is the mean of the vectors which is commonly taken as $\mathbf{0}$. Σ is the covariance matrix defined as

$$\Sigma = \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{bmatrix}$$

Where $\mathbf{K} = k(\mathbf{X}, \mathbf{X})$, $\mathbf{K}_* = k(\mathbf{X}, \mathbf{x}_*)$ and $\mathbf{K}_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$ for any given kernel k [20]. Due to the robustness of the GP process, we use this as one of the baselines in our thesis.

2.6.2 Acquisition functions

The acquisition functions need to do balance exploitation of information from the known data points and exploration of unknown data points in their domain. The following functions are some of the most prominent acquisition functions found in the literature [20]

- **Upper Confidence Bound (UCB)**: It returns the best possible hyperparameter configuration using a linear combination of the mean and the standard deviation.

- **Probability of Improvement:** It gives the probability with which we can get a better hyperparameter configuration than the incumbent best configuration.
- **Expected Improvement:** Given a gaussian distribution at a new input point, it finds the expectation of improvement i.e $(f(x) - f_{max})$ over the part of normal that is greater than f_{max} .

2.7 Deep Ensembles

Deep Neural Networks (DNNs) are machine learning models with very high representational capacity [8]. Due to this property, one can use them as surrogates for HPO objective functions. But the issue is that DNNs do not quantify uncertainty trivially.

In the current literature, uncertainty quantification methods using deep neural networks can be broadly classified into the following methods:

- Bayesian Neural networks [7].
- Ensemble Approach using monte carlo drop out [19].
- Ensemble approach using multiple neural networks.

In Bayesian neural network(BNN), a prior over weights and biases is specified during the intialization of the BNN. Given the data, a posterior predictive distribution is calculated for all the parameters of the network (Weights and Biases). One issue with this approach is that BNNs are very complex and difficult to train.

Monte carlo drop out is a regularization technique used during the training of neural networks. With a certain probability, connections between neurons are dropped. Using this technique one obtains possibly 2^N neural networks where N is the number of connections in the artifical neural network. We can get an ensemble of high capacity models for free. It is normally only used during training to obtain regularization.

However, if one uses monte carlo dropout during the evaluation, we can get multiple results from the same input using this approach. Given input x and output $y = \text{NN}(x)$. If we have m neural networks obtained using monte carlo dropout, we get $\{y_1, y_2 \dots y_m\}$ outputs, we can obtain the mean and variance of

$$y_{\text{mean}} = \frac{\sum y}{m} \quad y_{\text{variance}} = \frac{\sum (y - y_{\text{mean}})^2}{m - 1}$$

One way to do Hyper Parameter Optimization is by using Uncertainty prediction is the key in the paper [12].

2.8 Deep Kernel Learning

2.9 Rank Learning

Consider a set of objects $\mathbb{A} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n\}$ where each \mathbf{x} belongs to a domain \mathbb{D} . The problem of ranking is defined as finding an ordered list of objects in \mathbb{A} such that an object \mathbf{x}_i is ranked before \mathbf{x}_j if \mathbf{x}_i is more relevant/important than \mathbf{x}_j . To accomplish this objective, a ranking model needs to be learnt. In the most general case the cardinality of \mathbb{A} is not fixed. For this reason, the ranking model, say f_r , can be thought of as a process that is divided in the following steps [18]

- Obtaining a relevance score of each object in set \mathbb{A} .
- Sorting the objects based on their relevance score.

We can learn the ranking model f_r by optimizing a criteria on the output of the model i.e the sorted list of object. This criteria in the jargon of machine learning is called a loss function. Hence, it can be referred to as a *Ranking Loss*.

Since the step of sorting is non differentiable, it cannot generally be learnt during the optimization of our Ranking Loss. Hence the ranking model, f_r , boils down to a relevance scoring function. After learning f_r , one can use it to rank newly given sets of objects by first finding their relevant scores and then sorting them accordingly.

Various types of ranking losses can be used in our optimization to learn the scoring function. These can be broadly classified into the following types [5]:

- Point-wise ranking losses
- Pair-wise ranking losses
- List-wise ranking losses

In point-wise ranking loss, the loss function views the problem of ranking as that of assigning a label to each of the input data points. Hence, for learning, each instance is a single object x_i within the set \mathbb{A} . For example, in the

McRank paper [13], the authors reformulate the ranking problem as a multi-level classification problem where each data point is classified independently. They then calculate the score as the expected rank of the object based on its soft classification. Therefore, the complete scoring function comprises of a multi-level classifier and an external expectation calculation.

In pair-wise ranking loss, the loss function’s input is a pair of objects. This loss function learns to model pair-wise preferences. The function tries to separate the input data points as much as possible in the output space by minimising the pair-wise classification error [6].

In list-wise ranking losses, the loss is defined on the complete set of objects. The 2 most important list-wise loss functions are

- ListNet [3].
- ListMLE [22].

The point wise and pair wise ranking models do not view the ranking problem as a problem to rank a set of objects. This is their fundamental disadvantage. Moreover, the pairwise ranking loss may learn a scoring function that creates the issue of circular ranking. The scoring function may learn a preferences given by

$$\mathbf{x}_1 \succ \mathbf{x}_2 \quad \text{and} \quad \mathbf{x}_2 \succ \mathbf{x}_3 \quad \text{and} \quad \mathbf{x}_3 \succ \mathbf{x}_1$$

This creates an issue when creating an ordered list of objects. On the other hand, List wise loss functions are the most general ranking functions and also the most useful one for our thesis. It has been shown in [3] that list wise approaches are superior in performance to point wise and pair wise losses. We hence use the list wise approach to ranking. We discuss and analyse the loss functions ListNet and ListMLE in detail in chapter 4.

2.10 Deep Sets

Chapter 3

Background

Evaluating a novel hyper parameter surrogate is extremely expensive. For this, we would have to run the HPO technique (Bayesian optimization in our case) with our surrogate on a set of different machine learning models. This is because each ML model has a different search space. Moreover, the optimum within the HP search space may be different when the ML model is trained on different data sets. Hence we further need to do multiple HP optimizations of a model each time using a different dataset. In addition to this due to the stochastic nature of ML models as well as the surrogate models, we would have to run the HPO multiple times for each dataset.

For comparing our proposed model with other HPO models we have to do the above evaluation for each of the HPO model. As one can see, this evaluation if done right from the training of the ML model is not feasible. We use the HPO-B [15] benchmarking in this thesis. Using this benchmark, we do not need to train our ML models from scratch as the meta data also contains evaluations of multiple hyper parameters for different ML models and datasets.

In this section we first discuss about the organisation of the benchmarking meta data. Subsequently, we discuss the baselines used in our thesis.

3.1 Benchmarking Meta-Data

HPO-B is a benchmark that can be used for doing black box HPO. It can be used for both transfer models and non transfer models. The meta data consists of a list of (hyper-parameter) search spaces. These are hyper parameter search spaces of single models. The is organised in a json format with the structure illustrated in Figure 3.1.

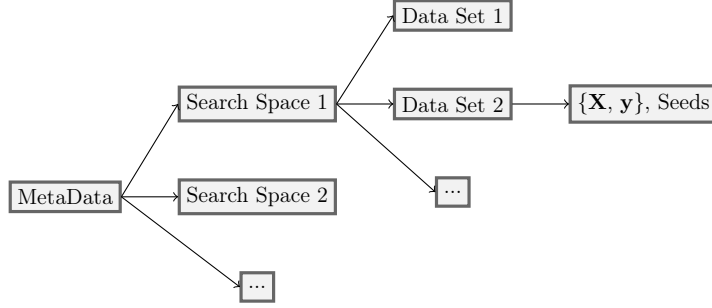


Figure 3.1: Structure of the meta data in the HPO-B benchmark

Where \mathbf{X} represents the set of configurations for a evaluated for a model in a particular data set and \mathbf{y} represents the evaluation results. The bayesian optimization used in our thesis can be started with different initial known HP configurations. One initial configuration, called a seed in the meta data, is provided by a set of values (or indices) in \mathbf{X} and \mathbf{y} . There are a total of 5 seeds provided for a search space and dataset combination.

The meta data in HPO-B comes in 3 versions namely **HPO-B-v1**, **HPO-B-v2**, and **HPO-B-v3**. Of this **HPO-B-v3** contains distilled the search spaces that have the most datasets and can be split is split into train, validation and test sets.

The meta data set in HPOB is divided into test and train data. Using meta dataset, one can learn meta learn a model using the training split. Then the model is evaluated in the testing split. This data approach is used both in the baselines and proposed idea models.

3.1.1 Benchmark evaluation

There are 2 types of HP optimization surrogates that are studied in this thesis - transfer learning surrogates and non-transfer learning surrogates. This benchmark can be used for analysing both types of models. However, in order cross compare transfer and non-transfer techniques, we only compare against HPO-B-v3 test split as recommended in the HPO-B paper.

3.2 Baselines

In this thesis, we decided to implement 2 HPO techniques from the literature - One non transfer technique and one transfer technique. The non transfer technique we implemented was Deep Ensembles. The reason we selected this

was to study how uncertainty can be modelled using deep neural networks. This was a pre-requisite to implement the uncertainty in our proposed model as we use deep neural networks as a scorer in our model.

The second baseline from the transfer learning domain was FSBO. We choose this because as this gave the state of the art results. In addition to this we used Random search and GP as standard baselines for result comparison.

Both these FSBO and DE have built in capability for uncertainty estimation. Hence we have to use an acquisition function during the fine tuning step. We use expected improvement in all models that deal with uncertainty. This is to maintain constant results across all the models. We also use this in the proposed idea due to the advantages of this acquisition function [10].

Few Shot Bayesian Optimization deals has to be run in 2 chronological steps

- Meta training using the existing meta data
- Fine tuning during the optimization cycle.

The following 2 sections discuss the implementation details of the baseline methods used in our thesis.

3.2.1 Deep Ensemble

As discussed in the literature review section, deep ensemble uses a list of artificial neural networks to find out uncertainties given data points D .

Some of the important implementation details are

- We make sure that we enforce the positivity of the variance (as mentioned in the paper) by adding $\exp(-6)$ to the values for numerical stability.
- We do not use adversarial examples because it was giving bad results.

Requirement of restarting training

We find in our experiments that the deep ensemble performs much better in our HPO cycle when we train it from scratch at every acquisition step. This is counter intuitive because an already trained model will converge to a local optima quickly. The reason for this performance however is that the model gets biased towards the points that are observed at the beginning of the optimisation cycle. Lets say we have 2 models

- DE_{restart} which always restarts training at every acquisition step.

- DE_{reuse} which uses the previously trained model for subsequent training of the optimization cycle.

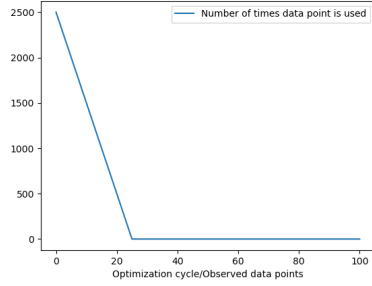


Figure 3.2: Bias at 25th optimization cycle

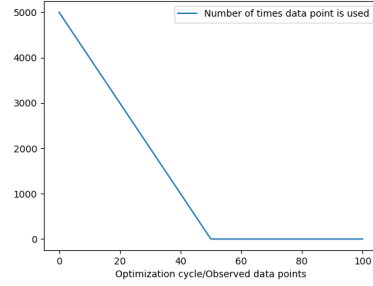


Figure 3.3: Bias at 50th optimization cycle

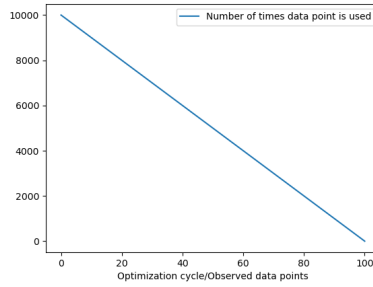


Figure 3.4: Bias at 100th optimization cycle

Let DE_{reuse} be fine tuned for 100 epochs during the optimization cycle. Let there be just 1 data point to start with. The figures 3.2, 3.3, and 3.4 show the number of times each observed data points is used for the fine tuning at the 25th, 50th and 100th optimization cycle by DE_{reuse} respectively. Generally all the data is used during fine tuning due to its scarcity at this step. Hence, the biases are scaled with the number of epochs. $DE_{restart}$ does not face this issue as the fine tuning step restarts everytime.

The heavy bias that is present in the figures is actually not intended. This is because all observations should be treated equally in any training/fine-tuning step. Hence it is better to restart the model before fine tuning step at each optimization step.

The second reason is that the model may get stuck at a stubborn local

minima at any fine tuning step n $1 \leq n \leq 100$. Coming out of this local minima may require vastly different data points than actually seen. This is not generally the case in our sequential optimization.

Optimization Procedure

Deep ensemble is not a meta transfer model. So it does not use the train data of HPO-B meta dataset. It directly optimises during the optimisation loop.

The optimization loop can be seen in the following algorithm.

The following should be written in an algorithm. First fitting the model based on observations. Predicting then the pending configuration for evaluation. Returning the index that gives us the best results.

3.2.2 Few Shot BO: Deep Kernel Learning

Some of the important implementation details are

- Scaling is ignored because it gives very bad results.
- We first used RBF kernel, then we used matern kernel for this.

Scaling should be ignored because it gives very bad results. If scaling: We used early stopping mechanism for this because we wanted to avoid

- Huge computation costs
- overfitting of the model during training.

We saved the best model in our implementation and if the validation loss went greater than lowest validation loss for some number of steps, then we stop the meta training.

Implementation issues

After completing the own implementation, we found that the results obtained were not in par with the paper. For this reason, the results are used only for the first stage. The results already present in the benchmarking data of [16] were used for comparing this result with our implementation.

Using cosine annealing

We obtained loss curves that were very bad during the training and were varying quite a bit. The reason was that the learning rate during the fine tuning was too high. Hence we used cosine annealing for this to get the a better local minima.

Show before and after curves.

The results of this are discussed in the results section.

3.3 Evaluation Strategy

Chapter 4

Method

4.1 Proposed Idea: : Ranking Loss Surrogates

The true evaluation of a hyperparameter optimisation objective function is computationally very expensive. Researchers have tried to get around this problem by proposing model based HPO algorithms. In these methods, the true objective function is modelled by a cheap surrogate function of high representational capacity. For example, the Sequential Model Based Optimization (SMBO) [1] algorithm is a very important algorithm that iterates between learning a model given a few HP configuration evaluations and using the model to propose the next candidate to evaluate.

The choice of surrogate to use in model based optimization is very crucial. Does the surrogate have enough representational capacity? Does it have the capability of representing uncertainty? How is the surrogate learnt? These are some of the questions that need to be answered before selecting a surrogate model. The selection of a surrogate model has a direct impact on the performance of the model based optimization algorithm.

In the quest to improve HPO surrogates, we propose and analyse a new type of surrogate model that is based on the concept of ranking. There are 2 components of our proposed idea

- The learning mechanism of the surrogate model.
- The surrogate model itself.

This chapter discusses in detail both these components. In this chapter, we first assume that a good surrogate model of sufficient representational capacity already exists. We use a simple Deep Neural network for this purpose. We then understand, analyse and implement the proposed learning

algorithm that uses the concept of. We then do a simple case study of inverse mapping using the learnt model. Finally we build a ranking model to improve the performance of any model based HPO algorithm. We use SMBO as a reference for our study.

4.2 Learning mechanism

In section 2.9, we discussed that list wise ranking losses were superior to the point wise and pair wise losses. In this section we analyse the prominent list loss functions - ListNet and ListMLE. Then, we discuss the implementation details of the loss function used in our model.

4.2.1 Loss functions: Analysis and implementation

Consider data in the format shown in table 4.1 is given to us.

Instance	Object Set	Ground Truth
1	$\{a_1, a_2, a_3, \dots, a_{10}\}$	$\{y_1, y_2, y_3, \dots, y_{10}\}$
2	$\{a'_1, a'_2, a'_3, \dots, a'_{15}\}$	$\{y'_1, y'_2, y'_3, \dots, y'_{15}\}$
3	$\{a''_1, a''_2, a''_3, \dots, a''_7\}$	$\{y''_1, y''_2, y''_3, \dots, y''_7\}$
...	$\{\dots\}$	$\{\dots\}$

Table 4.1: Data format used to train the scoring function using list wise ranking loss

Here let each data point a be a sample/element from the set \mathbb{A} . Each y represents the ground truth preference score of objects belonging to a set \mathbb{Y} . These preference scores of objects are relative to the objects within the input set. Let s be the scoring function to be learnt. Hence, the declaration of s is given by

$$s : \mathbb{A} \mapsto \mathbb{R}$$

As we can see from the table, one instance in our data consists of a set of objects as input and a set of corresponding ground truths to train from. To learn the function s we optimize our list wise loss function. This loss function takes as input the whole set of objects and their ground truth as one instance. If we take any set \mathbb{P} such that $\mathbb{P} \subseteq \mathbb{A}$, the declaration of the list wise loss L is hence given by

$$L : s(\mathbb{P}) \times \mathbb{Y}^{|\mathbb{P}|} \mapsto \mathbb{R} \quad (4.1)$$

Where the scoring function s applied to the \mathbb{P} gives us the set of corresponding scores of all objects in \mathbb{P} . We will consider the ground truth values to be \mathbb{R} for our analysis as this is type of value we have in our data sets.

In the next 2 sections we analyse ListNet and ListMLE, the prominent listwise loss functions in the literature.

4.2.2 ListNet

In this section we try to intuitively explain the ListNet idea proposed in [3]. Our objective is to learn the scoring function s such that it returns scores that are similar in relevance/order when compared to the ground truth scores. That is to say

$$y_3 < y_{12} < y_1 \implies s(a_3) < s(a_{12}) < s(a_1)$$

This would make the ranking of the objects equal to the ranking obtained by using the ground truth values. Note that we do not need to get the exact ground truth scores. This increases the number of acceptable functions that can be learnt by increasing the target function space. This makes it easier to learn the scoring function.

Ranking is obtained by sorting the objects based on their respective scores. Note that the sort functionality is non differentiable hence it is not a part of the ranking loss function. Our loss function needs to be constructed using the following 2 lists:

- List of scores given by scoring function s .
- List of scores given to us by ground truth.

The loss function must find some sort of a distance between the 2 given lists. It then can reduce the distance by changing the parameters of the scoring function.

In ListNet, a probabilistic approach is taken so as to account for any uncertainties in the ground truth values. Consider selecting an object from the input set with a probability

$$P = \frac{s(a)}{\sum_i s(a_i)} \quad \forall i \in \{1, 2, 3, \dots, |\mathbb{P}|\}$$

This make intuitive sense because the probability of selecting an object should be higher if it more relevant and vise-versa. Note that the score of any object by the scoring function can negative as well. Therefore the

score is passed through a strictly positive and increasing function ϕ . This changes the probability to

$$P = \frac{\phi(s(a))}{\sum_i \phi(s(a_i))} \quad \forall i \in \{1, 2, 3, \dots, |\mathbb{P}|\} \quad (4.2)$$

Equation 4.2 is also referred to as top 1 probability of an object in [3]. This is because this gives the probability of ranking the object first when we are calculating the permutation probability of given list.

The proposed way to find the distance between 2 lists in ListNet is

- Find the top 1 probabilities of each object using the scores given by the scoring function.
- Using the ground truth values, find similar top 1 probabilities.
- The cross entropy between the 2 entities gives us the "distance" between the 2 lists.

Let $P_{s(a)}$ represent the top 1 probability of an object using the scores given by the scoring function. Similarly, let P_y represent the top 1 probability using its ground truth value. The cross entropy used as a loss in ListNet is given by

$$L(\mathbf{y}, s(\mathbf{a})) = -\sum_i P_{s(a_i)} \log P_{y_i} \quad (4.3)$$

Where \mathbf{y} and $s(\mathbf{a})$ represent the ground truth values and the scores given by the scoring function.

Hyperparameter optimization using Sequential Model Based Optimization (SMBO) [1]

The loss function used is listMLE As discussed in the literature review, we use list wise ranking losses because they have been found to be more efficient in most general case.

The implementation is done by removing top element check comments The latest implementation is used due to its efficiency otherwise it is the same. (note)

Our implementation is given by the following algorithm Write the algorithm of our implementation here

Implementation given by the paper [17]. Here give the equations used by this algorithm

Main Idea: We need to use ranking losses instead of other methods for the purpose of HPO algorithms.

Applying this to HPOB... (With or without query) First learning from first search space. Even with this the loss curves are very smooth if we use $2 * \tanh(0.01 * x)$

4.3 Case study

In this case study, we study how the learnt ranking function behaves when we use different parameters to train. For this a toy example of sorting in reverse order is considered.

4.3.1 Sorting with inverse mapping

The main problem that we try to solve here is - Is it possible to train a ranking function using the Listwise maximum likelihood estimator. such that it learns inverse mapping of points on a number line. Consequently we could sort the results obtained by the scorer to obtain a descending order.

Consider a list $l = \{x_1, x_2, \dots, x_n\}$ where $x_i \in [1, 100]$. Let $s(x | \theta) \mapsto \mathbb{R}$ be our scoring function parametrised by θ . One list contains n data points sampled from $[1, 100]$. Let L_{mle} represent the maximum log likelihood loss. Using this loss function we optimize for the following equation

$$\arg \min_{\theta} L_{\text{mle}}(s(l | \theta), -k * l) \quad (4.4)$$

where $k \in \mathbb{R}$.

Note that second parameter of list wise loss function is scale invariance because it only uses the rank of the objects in the list and not the objects directly. Hence scaling the list has no effect on the loss output. This can be seen from the algorithm of implementation. (refer algorithm)

The validation data taken from 3 different ranges

- Same range as the training data $[1, 100]$
- Completely different range as seen by the scorer during training i.e $[-100, -1]$
- Hybrid range i.e $[-50, 50]$

To evaluate our scorer, we first sample the validation data from the above ranges, We then check the percentage of the lists that are correctly sorted during our testing time. The training batch is 100 and the results are averaged over 5 runs.

tabulate the results and Check the affect of list / batch size on the results
show the architecture of with only scorer

Training Epochs	List size	Learning Rate	In-range Acc.	Out-range Acc.	Hybrid Acc.
1000	3	0.0001	0.99	0.99	0.99

Table 4.2: Sorting Accuracies at test time

4.3.2 Observations

As we know that the second parameter of the loss function is not bound to any range, the score function itself can learn scores of arbitrary range. Hence the model can learn score values arbitrarily large. We must control the output domain so that the training does not lead to underflow or overflow. This is because we use the exp function for getting the strict positive increasing function.

The restriction of output was done using $\tanh()$ so that the range of values of the scorer remains within a manageable value. One drawback we found was that the learnt model was extremely sensitive to the learning rate and the number of epochs.

Then we add the following components one after another

- Deep Set
- Weight
- Uncertainty
- Different training mechanisms.

Propose using a higher list size.

4.3.3 Scoring function analysis

4.3.4 Scoring function range

4.4 Optimization cycle

Explain the training loop with an algorithm.

explain `get_batch_HPBO` implementation how it is important to have all the data sets How it is organised into a higher order tensor TO compute 2 things can be done 1. View it as a 2d tensor (Done by efficient implementation) 2. Deal with it direction (Done by us)

replace=False thing in the above function and why it is important.

4.4.1 Meta training

Give the meta training algorithm Why all tasks are taken - Because we do not want bias to crop up during the training The idea is to meta learn commonality in the search space.

4.4.2 Fine tuning

Given the Fine tuning algorithm. For reasons described in the [3.2.1](#), we restart before fine tuning at every optimization cycle step.

4.5 Ranking Surrogate Model

4.5.1 Transfer Learning

This is the transfer HPO-B

4.5.2 Use of Deep Sets

why are deep sets important. For domain conditionality meta train in source domain (HPOB training SS) meta test in target domain (HPOB test/val SS)

Architecture with deep sets.. Image

4.5.3 Uncertainty implementation

Mention about the requirement to use Module list. Show the new architecture.

The search spaces that have less data need more uncertainty The search spaces having more data need less uncertainty

Independent training

Implementation yet to be done

Training with mean and restricted output

Implementation yet to be done

4.6 Weighted Loss

General list wise loss functions are only required in the problem domains where it is essential to get the ranks of all the objects correctly. In our case where ranking functions are used as surrogates for the black box hyperparameter optimization, this constraint can be relaxed significantly. As discussed in section 4.4.2, we always restart our ranking function model before fine-tuning it. Hence, after fine-tuning we are only bothered about getting the best hyper parameter configuration from the pending configurations. This means that it is more important for our ranking function to order the the top part of the list than the bottom part. In this section we describe and propose the concept of weighted ranking loss functions.

The concept of biased ranking loss is discussed in the paper, "Top-Rank Enhanced Listwise Optimization for Statistical Machine Translation" by Chen et. al [4]. Here Chen et. al proposes a position based weighting of the ranking function such that:

$$w_j = \frac{k - j + 1}{\sum_{t=1}^k t} \quad (4.5)$$

where w_j represents weight of the object at position j in the ordered list.

However, there are also some simple weightings that a person can use. For instance inverse linear weighting given by

$$w_j = \frac{1}{j} \quad (4.6)$$

and inverse logarithmic weighting given by

$$w_j = \frac{1}{\log(j + 1)} \quad (4.7)$$

Figure 4.1 shows all 3 weighting strategies. Of the 3 weighting strategies, the inverse log weighting was found to be most suitable. This is because the this weighting is by far the least sharp. As the weighting does not decrease very rapidly, there is a possibility of parallelizing during the optimization cycle. For example, at every step the top n configurations can be selected, and all of them can be tried for optimization. The parallelism is less effective when using inverse linear weighting or weighting proposed by Chen et al.

4.7 Advantages and Limitations

Compare this loss function and method with other base lines.. Advantages of the proposed Idea: The amount of data instances for training is exponential

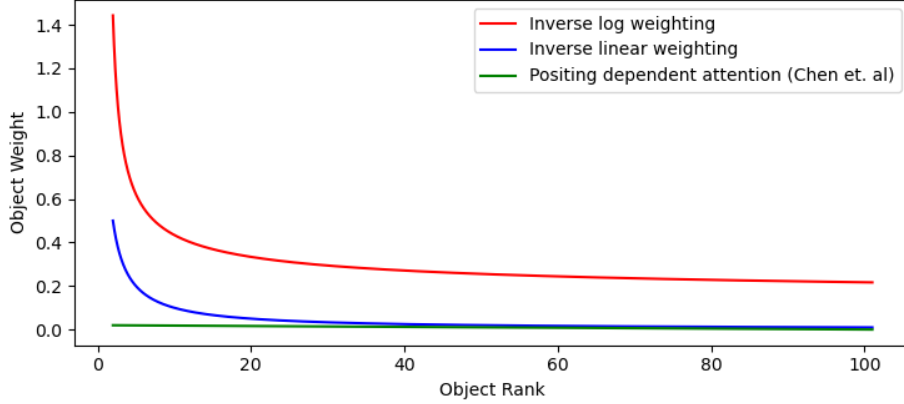


Figure 4.1: Different Weighting Functions

in number. Which is very good for a deep learning model. For example if we have 100 observation set and we use a list size of 15 to train our model, we will have $100C15$ unique instances to train.

Observed disadvantages:

- In our model, a scorer is first learnt and then using the scorer, we rank the set of objects in question. When optimising the scorer, we ignore the sorting functionality necessary to complete the process of ranking. This is because sorting is non differentiable. This means that the true evaluation of the ranked list is not optimized. This needs to be improved which is done in Pi-Rank paper.
- The learnt model is extremely sensitive to the learning rate and the number of epochs.

4.7.1 Negative transfer learning

For some search spaces, the validation errors of the ranking loss model did not reduce at all during its training. In fact the validation loss became worse. Figure 4.2 shows an illustration of this for the search space id 5527. This is a classic example of negative transfer learning [21]. Negative transfer learning occurs when the set of source tasks (here, training datasets) is very different from the target tasks (here, validation datasets).

This problem should occur only in the cases where the model is context free. In our case, the ranking loss model with deep set is context aware and results like these were unexpected. Hence the issues with negative transfer

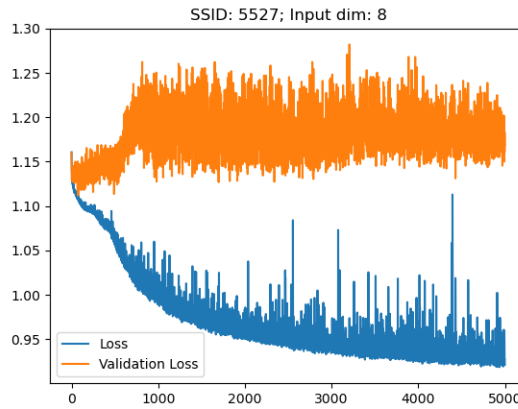


Figure 4.2: Training loss curve of ranking losses

learning remains unsolved when using our model. One remedy for getting around this problem is to fine tune the model with a larger number of epochs. Nevertheless, we cannot guarantee that the finetuning will re-learn from the small amount observed data points during the optimization cycle.

Chapter 5

Research Question

The format of this is the same as that of experiments and Results. Also known as Hypothesis

Chapter 6

Experiments and Results

Data used in the experiment Infrastructure used (That is the protocol) Results should be structured like hypothesis.

experiment and results protocol data split...

First show the results of implementation of GP ($M = 5$ and $M = 10$)

Benchmark this...

Next with DKT (Benchmark this)

Next show the best results obtained by architecture.

6.1 Evaluation

6.1.1 Testing

explain how a ranking graph works as implemented Explain the regret rank@ some location.

6.1.2 Ablation

Result tabulation of case study: sorting: 1. Within range 2. Outside range mean of 3 times should be written.

show the results of raw without deep set.

Next show different strategies used for building the ranking loss model one step at a time. First with only scorer. then with deep set. Then with raw deep set fine tuning and deep set adding uncertainty

Checking the early stop and hypothesing why is was wrong.

box plot variation of each of the scorers... for 1 or more data sets?

show results of independent training and training with output restriction

what about training independently, this requires normalization. as explained by sebastian.

Chapter 7

Conclusion

7.1 Further work

Further study required with other baselines that deal with ranking loss.

7.2 Conclusion

Bibliography

- [1] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [2] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, 13:281–305, feb 2012.
- [3] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: From pairwise approach to listwise approach. In *Proceedings of the 24th International Conference on Machine Learning, ICML '07*, page 129–136, New York, NY, USA, 2007. Association for Computing Machinery.
- [4] Huadong Chen, Shujian Huang, David Chiang, Xinyu Dai, and Jiajun Chen. Top-rank enhanced listwise optimization for statistical machine translation, 2017.
- [5] Wei Chen, Tie-yan Liu, Yanyan Lan, Zhi-ming Ma, and Hang Li. Ranking measures and loss functions in learning to rank. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009.
- [6] David Cossock and Tong Zhang. Statistical analysis of bayes optimal subset ranking. *IEEE Transactions on Information Theory*, 54(11):5140–5154, Nov 2008.
- [7] Ethan Goan and Clinton Fookes. Bayesian neural networks: An introduction and survey. In *Case Studies in Applied Bayesian Data Science*, pages 45–87. Springer International Publishing, 2020.

- [8] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [9] Frank Hutter, Holger H. Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *Proceedings of the 5th International Conference on Learning and Intelligent Optimization, LION'05*, page 507–523, Berlin, Heidelberg, 2011. Springer-Verlag.
- [10] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, Dec 1998.
- [11] Abdus Salam Khazi. [Code](#). [Online; accessed 06-Feb-2022].
- [12] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles, 2016.
- [13] Ping Li, Qiang Wu, and Christopher Burges. Mcrank: Learning to rank using multiple classification and gradient boosting. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.
- [14] Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- [15] Sebastian Pineda-Arango, Hadi S. Jomaa, Martin Wistuba, and Josif Grabocka. HPO-B: A large-scale reproducible benchmark for black-box HPO based on openml. *CoRR*, abs/2106.06257, 2021.
- [16] Sebastian Pineda-Arango, Hadi S. Jomaa, Martin Wistuba, and Josif Grabocka. HPO-B: A large-scale reproducible benchmark for black-box HPO based on openml. *Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks*, 2021.
- [17] Przemysław Pobrotyn, Tomasz Bartczak, Mikołaj Synowiec, Radosław Białobrzeski, and Jarosław Bojar. Context-aware learning to rank with self-attention, 2020.
- [18] Tao Qin, Tie-Yan Liu, Jun Xu, and Hang Li. Letor: A benchmark collection for research on learning to rank for information retrieval. *Inf. Retr.*, 13:346–374, 08 2010.

- [19] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.
- [20] Jie Wang. An intuitive tutorial to gaussian processes regression, 2020.
- [21] Karl Weiss, Taghi M. Khoshgoftaar, and DingDing Wang. A survey of transfer learning. *Journal of Big Data*, 3(1):9, May 2016.
- [22] Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. List-wise approach to learning to rank: theory and algorithm. In *ICML '08*, 2008.

Appendix A

More information

This thesis was completed in the representation learning lab of Albert-Ludwig-Universität Freiburg. (Figure [A.1](#))



Figure A.1: Logo: Albert-Ludwig-Universität Freiburg