
Оптимизация количества деревьев в ансамбле градиентного бустинга с использованием стратегий выбора агрегирующих функций

A Preprint

Якушевич Антон Сергеевич
ВМК
МГУ
Москва
s02220301@gse.cs.msu.ru

Сенько Олег Валентинович
ВМК
МГУ
Москва
senkoov@mail.ru

Abstract

В работе исследуется задача оптимизации ансамблей градиентного бустинга за счёт выбора обобщённой агрегирующей функции. Исследование проводится для уменьшения количества деревьев в ансамбле с целью повышения интерпретируемости модели при сохранении её прогностической точности. Для этого предлагается использовать усечённое разложение агрегирующей функции в ряд Тейлора и оптимизировать её коэффициенты с помощью градиентного спуска, что позволяет адаптивно настраивать структуру ансамбля и улучшать баланс между точностью и интерпретируемостью. Экспериментальные результаты показывают, что оптимизация коэффициентов ряда Тейлора приводит к уменьшению числа деревьев без потери качества, а в ряде случаев даже к улучшению точности.

Keywords First keyword · Second keyword · More

1 Введение

Градиентный бустинг остаётся одним из наиболее эффективных методов машинного обучения, применяемых в задачах регрессии и классификации. Однако высокая точность ансамбля достигается ценой увеличения числа деревьев, что усложняет интерпретацию модели и приводит к увеличению вычислительных затрат. Поэтому важной задачей является разработка подходов, позволяющих уменьшить размер ансамбля без ухудшения качества прогнозов.

Одним из направлений снижения сложности ансамблей является модификация агрегирующей функции. Наиболее простой подход заключается во введении весов для отдельных деревьев, которые подбираются либо вручную, либо оптимизируются на валидационных данных [Domke, 2012]. Более сложные методы рассматривают автоматическую настройку весов или коэффициентов, что позволяет повысить адаптивность ансамбля [Luca Franceschi, 2017]. Другая линия исследований связана с использованием мета-моделей для агрегирования, таких как нейронные сети, которые применяются для того, чтобы обучаться более сложным функциям взвешивания базовых предсказаний [Anonymus, 2022]. Подобные подходы расширяют пространство возможных агрегирующих функций и потенциально позволяют достичь более высокой точности. В смежных направлениях рассматривались и дифференцируемые методы оптимизации гиперпараметров [Dougal Maclaurin, 2015, Jonathan Lorraine, 2019, Amirreza Shaban, 2019], где подбираются параметры регуляризации и обучения, что демонстрирует применимость автоматического дифференцирования к задачам настройки сложных моделей.

Несмотря на успехи, существующие подходы имеют ряд ограничений. Во-первых, во многих работах агрегирующая функция фиксируется заранее (например, простое суммирование или линейная

комбинация), что ограничивает адаптивность ансамбля [Domke, 2012]. Во-вторых, даже в случаях, когда веса подбираются автоматически, процесс оптимизации часто является дорогостоящим и требует перебора гиперпараметров или применения эвристик [Luca Franceschi, 2017]. В-третьих, подходы, где агрегирующая функция параметризуется нейронной сетью [Anonymous, 2022], сталкиваются с ограничением: обучение функции проводится поверх уже зафиксированного набора деревьев, что не позволяет агрегирующей функции влиять на сам процесс построения ансамбля. Это делает такие методы менее гибкими и усложняет их использование на практике.

В данной работе агрегирующая функция ансамбля представляется в виде усечённого разложения в ряд Тейлора, где коэффициенты ряда рассматриваются как параметры, подлежащие обучению. Для оптимизации этих коэффициентов используется следующая идея. Сначала строится оператор, который принимает на вход набор коэффициентов, затем обучает ансамбль деревьев с фиксированными гиперпараметрами, и, наконец, вычисляет значение функции потерь на валидационных данных. Возвращаемое оператором значение служит целевой функцией для обновления коэффициентов. Таким образом, коэффициенты агрегирующей функции оптимизируются напрямую с помощью градиентного спуска. Благодаря этому достигается совместная настройка структуры ансамбля и параметров агрегирования, что даёт дополнительную гибкость по сравнению с фиксированными или отдельно обучаемыми функциями.

Предложенный метод открывает новый способ интеграции оптимизации агрегирующих функций в процесс обучения ансамблей. В отличие от подходов, где агрегирующая функция задаётся заранее или обучается постфактум, наша методика позволяет адаптивно корректировать её форму одновременно с построением деревьев. Экспериментальные результаты показывают, что оптимизация коэффициентов ряда Тейлора приводит к уменьшению числа деревьев без потери качества, а в ряде случаев даже к улучшению точности. Тем самым достигается более выгодный баланс между точностью, интерпретируемостью и вычислительной эффективностью. Вклад работы заключается в том, что она расширяет область применения методов оптимизации, показывая, что обучение агрегирующих функций может быть встроено непосредственно в процесс градиентного бустинга, а не рассматриваться как внешний этап.

2 Литературный обзор

Ансамблевые методы давно зарекомендовали себя как один из наиболее эффективных классов алгоритмов машинного обучения [Breiman, 1996, Yoav Freund, 1997]. Среди них особое место занимает градиентный бустинг [Friedman, 2001], который благодаря высокой точности и универсальности получил широкое распространение как в академических исследованиях, так и в промышленных приложениях [Tianqi Chen, 2016, Guolin Ke, 2017]. Однако рост числа деревьев в ансамбле приводит к ухудшению интерпретируемости и повышению вычислительных затрат, что стимулировало поиск методов сокращения сложности без потери качества.

Одним из направлений развития бустинга стали методы регуляризации и усечения ансамблей. Ряд работ показал эффективность ограничений на глубину деревьев и скорости обучения [Tong Zhang, 2005, Peter Bühlmann, 2007], а также введения различных форм регуляризации [Llew Mason, 2000, Friedman, 2002]. Другой класс подходов связан с сокращением ансамбля путём отбора наиболее информативных моделей [Dragos D Margineantu, 1997, Rich Caruana, 2004], что позволяет уменьшить число деревьев при минимальных потерях в точности.

Ключевым элементом ансамблевых методов является агрегирующая функция, которая объединяет предсказания отдельных базовых алгоритмов в итоговый результат. В классических вариантах бустинга используется простая сумма или усреднение выходов деревьев с фиксированными весами [Trevor Hastie, 2009]. Однако такая схема имеет ограниченную гибкость и не всегда позволяет учесть различный вклад деревьев в итоговую точность. Для повышения адаптивности были предложены методы взвешивания, где каждому дереву приписывается коэффициент, подбираемый либо с помощью оптимизации на валидационных данных, либо через регуляризацию [Domke, 2012, Luca Franceschi, 2017]. Подобные подходы позволяют усилить значимость наиболее полезных моделей и, наоборот, снизить влияние переобученных деревьев.

В смежных исследованиях развивались методы *stacking*, где агрегирующая функция задавалась в виде отдельной модели, обучаемой на предсказаниях базовых алгоритмов [Wolpert, 1992, K. M. Ting, 1999, R. Maclin, 1999]. В простейшем случае это линейная или логистическая регрессия, а в более сложных вариантах — регуляризованные модели или мета-классификаторы, способные учитывать взаимосвязи

между деревьями. Такие методы обеспечивали более богатое пространство комбинаций по сравнению с фиксированными весами, но усложняли интерпретацию и повышали риск переобучения.

В последние годы начали активно применяться подходы, где агрегирующая функция параметризуется нейронными сетями [Corinna Cortes, 2017, Anonymous, 2022]. Нейросетевой агрегатор способен моделировать нелинейные зависимости между выходами деревьев и строить адаптивные правила комбинирования, выходящие за рамки линейных комбинаций. Это открывает возможность учёта сложных взаимодействий внутри ансамбля. Однако ключевым ограничением таких методов является то, что обучение агрегирующей функции проводится поверх уже построенного ансамбля деревьев: деревья фиксируются, и нейросеть лишь учится их комбинировать. В результате агрегатор не влияет на сам процесс формирования ансамбля, что снижает потенциал такого подхода.

Параллельно активно развивалось направление дифференцируемой оптимизации гиперпараметров. Ключевая идея была предложена в работах Dougal Maclaurin [2015] и Domke [2012], где рассматривалось дифференцирование процесса обучения для вычисления градиентов валидационной ошибки. Дальнейшие исследования сосредоточились на масштабировании этого подхода [Luca Franceschi, 2017, Jonathan Lorraine, 2019, Amirreza Shaban, 2019], использовании неявного дифференцирования и аппроксимаций обратного Гессияна [Pedregosa, 2016], а также применении автоматического дифференцирования в широком классе моделей [Atilim Gunes Baydin, 2018]. Эти работы подтвердили возможность точной и вычислительно эффективной оптимизации параметров, выходящих за рамки традиционных методов перебора и байесовской оптимизации [Jasper Snoek, 2012, James Bergstra, 2011].

Несмотря на прогресс, существующие методы имеют ограничения. Весовые и линейные комбинации деревьев не всегда обеспечивают достаточную гибкость, а нейросетевые агрегаторы обучаются поверх фиксированного ансамбля, не влияя на сам процесс построения деревьев. Дифференцируемые методы гипероптимизации сосредоточены преимущественно на настройке параметров обучения и регуляризации, а не на структурных характеристиках агрегирования. Эти ограничения формируют исследовательский разрыв, восполняемый настоящей работой, где агрегирующая функция представляется в виде усечённого разложения в ряд Тейлора, а её коэффициенты оптимизируются напрямую градиентным спуском в ходе построения ансамбля.

3 Постановка задачи

3.1 Данные

В работе рассматриваются наборы данных для задач регрессии:

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

где $x_i \in \mathbb{R}^d$ - вектор признаков, а $y_i \in \mathbb{R}$ - целевая переменная.

В исследовании используются 20 наборов данных из различных регрессионных задач

3.2 Обобщение градиентного бустинга

Обобщим метод градиентного бустинга. Пусть задана последовательность функций $\{F_N\}_{N=1}^{+\infty}$, где $F_N : \mathbb{R}^N \rightarrow \mathbb{R}$. На $(N + 1)$ -ом шаге обучения вместо стандартной суммы $\sum_{i=1}^N T_i(x)$ будем использовать обобщённую функцию ансамбля $F_N(T_1(x), T_2(x), \dots, T_N(x))$, где $\{T_i\}_{i=1}^N$ - базовые модели (решающие деревья), обученные на предыдущих шагах.

Задача состоит в нахождении такой последовательности функций F_N , которая при фиксированном числе базовых моделей N обеспечивала бы лучшее качество приближения целевой функции $f(x)$ (относительно заданной функции потерь \mathcal{L}) по сравнению с классическим градиентным бустингом, использующим сумму в качестве агрегирующей функции.

3.3 Параметризация агрегирующей функции

Зафиксируем число N и рассмотрим функцию от N переменных $F_N(T_1(x), T_2(x), \dots, T_N(x))$. Разложим её в ряд Тейлора в окрестности точки (a_1, a_2, \dots, a_N) :

$$F_N(T_1(x), T_2(x), \dots, T_N(x)) = \sum_{k=0}^{+\infty} \left[\sum_{k_1=0}^{k_1+k_2+\dots+k_N=k} \sum_{k_2=0} \dots \sum_{k_N=0} \frac{1}{k_1!k_2! \dots k_N!} \frac{\partial^k F_N(a_1, a_2, \dots, a_N)}{\partial T_1^{k_1} \partial T_2^{k_2} \dots \partial T_N^{k_N}} \prod_{i=1}^N (T_i(x) - a_i)^{k_i} \right]$$

Тогда вместо классической линейной комбинации деревьев в решающем лесу предлагается использовать обобщённую функцию ансамбля, которая является усечённым вариантом ряда Тейлора:

$$\hat{F}_N(T_1(x), T_2(x), \dots, T_N(x)) = \sum_{k=1}^M \left[\sum_{k_1=0}^{\overbrace{k_1+k_2+\dots+k_N=k}} \sum_{k_2=0} \cdots \sum_{k_N=0} b_{k_1, k_2, \dots, k_N} \prod_{i=1}^N (T_i(x) - a_i)^{k_i} \right]$$

где M , b_{k_1, k_2, \dots, k_N} , a_i – гиперпараметры

При этом выбор значений b_{k_1, k_2, \dots, k_N} и a_i позволяет моделировать различные виды функций F_N при $M \rightarrow +\infty$.

Если количество используемых переменных меньше общего числа $L < N$ (например при обучении $(L+1)$ -го дерева), то функция принимает вид:

$$\hat{F}_L(T_1(x), T_2(x), \dots, T_L(x)) = \sum_{k=1}^M \left[\sum_{k_1=0}^{\overbrace{k_1+k_2+\dots+k_L=k}} \sum_{k_2=0} \cdots \sum_{k_L=0} b_{k_1, k_2, \dots, k_L, 0, \dots, 0} \prod_{i=1}^L (T_i(x) - a_i)^{k_i} \right]$$

Здесь нулевые значения индексов k_{L+1}, \dots, k_N соответствуют игнорированию лишних переменных. Таким образом, при добавлении нового алгоритма в ансамбль старые слагаемые линейной комбинации остаются неизменными, а новые слагаемые просто дополняют разложение. Такой подход обеспечивает стабильность модели и позволяет эффективно наращивать ансамбль без пересчета уже найденных коэффициентов.

Для упрощения вычислений и повышения устойчивости модели будем нормализовать целевую переменную с помощью z-нормализации. Тогда центр разложения можно выбрать равным нулю: $(a_1, a_2, \dots, a_N) = \bar{0}$, а в качестве M взять небольшое число.

Везде далее для упрощения восприятия под F_N будет подразумеваться \hat{F}_N .

3.4 Стратегии обучения нового алгоритма

Введём обозначение $F_N(x) = \sum_{k=1}^M \sum_{k_1=0}^{\overbrace{k_1+k_2+\dots+k_N=k}} \sum_{k_2=0} \cdots \sum_{k_N=0} b_{k_1, k_2, \dots, k_N} T_1^{k_1}(x) T_2^{k_2}(x) \dots T_N^{k_N}(x)$. Для повышения

точности приближения при добавлении нового алгоритма $T_{N+1}(x)$ вводится параметр скорости обучения (learning rate), который вычисляется как $\alpha = \frac{\delta}{N^\theta}$, где N – порядковый номер добавляемого алгоритма в ансамбле, а δ и θ – гиперпараметры. Выбор вектора, на котором будет обучаться новый алгоритм, может осуществляться по одной из трёх стратегий:

1. Нахождение точного корня:

При добавлении нового алгоритма $T_{N+1}(x)$ в ансамбль требуется, чтобы функция $F_{N+1}(x)$ аппроксимировала целевую переменную y для каждого объекта x из обучающей выборки (т.е. $F_{N+1}(x) \approx y$). Если подставить известные значения $T_1(x), T_2(x), \dots, T_N(x)$ в выражение для $F_{N+1}(x)$, то задача сводится к решению полиномиального уравнения степени M относительно неизвестного $T_{N+1}(x)$. Для каждого объекта x_i из обучающей выборки необходимо решить уравнение вида $\sum_{k=0}^M C_k(x_i) \cdot z_i^k = y_i$, где $C_k(x_i)$ – функции, однозначно задаваемые $T_1(x_i), \dots, T_N(x_i)$, а z_i представляет собой искомое значение $T_{N+1}(x_i)$. Однако данное уравнение не всегда имеет точное решение; в таких случаях в качестве z_i выбирается точка, в которой полином принимает значение, минимальное по модулю.

После нахождения значений z_i для всех объектов обучается новый алгоритм T_{N+1} на векторе $\alpha \cdot (z_1, z_2, \dots, z_N)$. Отметим, что первый алгоритм всегда обучается именно по этой стратегии.

2. Градиент:

Пусть имеется функция потерь $\mathcal{L}(y, F(x))$, тогда можно записать, что $\mathcal{L}(y, F_{N+1}(x)) = \mathcal{L}(y, F_N(x) + \sum_{k=1}^M C_k(x) T_{N+1}^k(x))$, где $C_k(x)$ – функции, однозначно задаваемые $T_1(x), \dots, T_N(x)$. Разложим \mathcal{L} в ряд Тейлора в точке $T_{N+1}(x) = 0$ (для удобства пере-

именуем $T_{N+1}(x)$ в z):

$$\begin{aligned} \mathcal{L}(y, F_N(x) + \sum_{k=1}^M C_k(x)z^k) &= \mathcal{L}(y, F_N(x)) + z \cdot \frac{\partial \mathcal{L}(y, F_N(x) + \sum_{k=1}^M C_k(x)z^k)}{\partial z} \Big|_{z=0} = \\ &= \mathcal{L}(y, F_N(x)) + z \cdot C_1(x) \cdot \frac{\partial \mathcal{L}(y, F)}{\partial F} \Big|_{F=F_N(x)} \end{aligned}$$

$$\text{где } C_1 = \sum_{k=1}^M \sum_{k_1=0}^{k_1+k_2+\dots+k_N=k-1} \dots \sum_{k_N=0} b_{k_1, k_2, \dots, k_N} T_1^{k_1}(x) T_2^{k_2}(x) \dots T_N^{k_N}(x)$$

Возвращаясь к T_{N+1} :

$$T_{N+1}(x_i) = \arg \min_T \mathcal{L}(y_i, F_N(x_i) + \sum_{k=1}^M C_k(x_i)T^k(x_i)) = \arg \min_T \left(T(x_i) \cdot C_1(x_i) \cdot \frac{\partial \mathcal{L}(y, F)}{\partial F} \Big|_{F=F_N(x_i)} \right)$$

$$\text{Отсюда получаем выражение для нового алгоритма } T_{N+1}(x_i) = -\alpha \cdot C_1(x_i) \cdot \frac{\partial \mathcal{L}(y, F)}{\partial F} \Big|_{F=F_N(x_i)},$$

где $\alpha \in (0, 1]$

Пример для MSE: $\mathcal{L}(y, F) = \frac{1}{2}(y - F(x))^2$

$$\frac{\partial \mathcal{L}(y, F)}{\partial F} = F(x) - y \Rightarrow T_{N+1}(x_i) = \alpha \cdot C_1(x_i) \cdot (y_i - F_N(x_i))$$

3.5 Стратегия оптимизации коэффициентов агрегирующей функции

При фиксированном наборе данных \mathcal{D} можно определить оператор $A_{\mathcal{D}}(b) : \mathbb{R}^P \rightarrow \mathbb{R}$, который для заданного вектора коэффициентов b обучает градиентный бустинг и возвращает значение функции потерь для полученной модели.

Допустим, что этот оператор обладает непрерывными частными производными во всех точках. В таком случае, для вычисления этих производных можно применить следующий подход:

$$\frac{\partial A_{\mathcal{D}}}{\partial b_i} \Big|_{b=b_0} \approx \frac{A_{\mathcal{D}}(b_0 + \epsilon \cdot e_i) - A_{\mathcal{D}}(b_0)}{\epsilon}$$

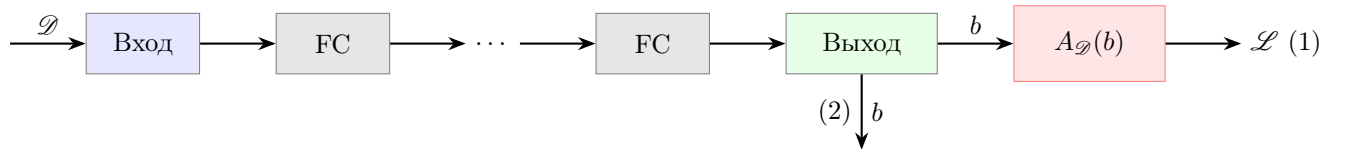
где e_i - это единичный вектор, направленный вдоль i -й координатной оси.

Таким образом, данный оператор может быть оптимизирован методом градиентного спуска:

$$b_i^{k+1} = b_i^k - \alpha \cdot \frac{\partial A_{\mathcal{D}}}{\partial b_i} \Big|_{b=b^k}$$

3.6 Мета-модель

Для ускорения подбора коэффициентов агрегирующей функции рассмотрим следующую нейронную сеть:



Процесс обучения модели осуществляется по траектории (1) с использованием алгоритма обратного распространения ошибки (метод вычисления градиента через оператор $A_{\mathcal{D}}(b)$ был представлен ранее). После завершения обучения модель функционирует по траектории (2) для генерации коэффициентов b , которые в дальнейшем используются для обучения градиентного бустинга.

3.7 Внешний критерий качества

В качестве функции потерь используется среднеквадратичная ошибка:

$$\mathcal{L}(y, F(x)) = \frac{1}{2}(y - F(x))^2$$

В качестве внешнего критерия качества выступает следующий функционал:

$$Q = \frac{\sum_{(x,y) \in \mathcal{D}} \mathcal{L}(y, F_N(x))}{N}$$

Целью является уменьшение числа базовых моделей N при сохранении точности, то есть минимизация Q .

Список литературы

- Justin Domke. Generic methods for optimization-based modeling. *Journal of Machine Learning Research*, 2012.
- Paolo Frasconi Massimiliano Pontil Luca Franceschi, Michele Donini. Forward and reverse gradient-based hyperparameter optimization. *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017.
- Anonymous. Differentiable hyper-parameter optimization. *ICLR 2022*, 2022.
- Ryan P Adams Dougal Maclaurin, David Duvenaud. Gradient-based hyperparameter optimization through reversible learning. *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- David Duvenaud Jonathan Lorraine, Paul Vicol. Optimizing millions of hyperparameters by implicit differentiation. *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- Nolan Hatch Byron Boots Amirreza Shaban, Ching-An Cheng. Truncated back-propagation for bilevel optimization. *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- Leo Breiman. Bagging predictors. *Machine Learning*, 1996.
- Robert E. Schapire Yoav Freund. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 1997.
- Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 2001.
- Carlos Guestrin Tianqi Chen. Xgboost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- Thomas Finley Taifeng Wang Wei Chen Weidong Ma Qiwei Ye Tie-Yan Liu Guolin Ke, Qi Meng. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems 30 (NIP 2017)*, 2017.
- Bin Yu Tong Zhang. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 2005.
- Torsten Hothorn Peter Bühlmann. Boosting algorithms: Regularization, prediction and model fitting. *Statistical Science*, 2007.
- Peter Bartlett Marcus Frean Llew Mason, Jonathan Baxter. Boosting algorithms as gradient descent. *Advances in Neural Information Processing Systems 12*, 2000.
- Jerome H Friedman. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 2002.
- Thomas G Dietterich Dragos D Margineantu. Pruning adaptive boosting. *Proceedings of the 14th International Conference on Machine Learning (ICML)*, 1997.
- Geoff Crew Alex Ksikes Rich Caruana, Alexandru Niculescu-Mizil. Ensemble selection from libraries of models. *Proceedings of the 21st International Conference (ICML 2004)*, 2004.
- Jerome Friedman Trevor Hastie, Robert Tibshirani. *The Elements of Statistical Learning: Data mining, Inference, and Prediction*. Springer, 2009.

-
- David H Wolpert. Stacked generalization. Neural Networks, 1992.
- I. H. Witten K. M. Ting. Issues in stacked generalization. Journal of Artificial Intelligence Research, 1999.
- D. Opitz R. Maclin. Popular ensemble methods: An empirical study. Journal of Artificial Intelligence Research, 1999.
- Vitaly Kuznetsov Mehryar Mohri Scott Yang Corinna Cortes, Javier Gonzalvo. Adanet: Adaptive structural learning of artificial neural networks. Proceedings of the 34th International Conference on Machine Learning (ICML), 2017.
- Fabian Pedregosa. Hyperparameter optimization with approximate gradient. Proceedings of the 33rd International Conference on Machine Learning (ICML), 2016.
- Alexey Andreyevich Radul Jeffrey Mark Siskind Atilim Gunes Baydin, Barak A Pearlmutter. Automatic differentiation in machine learning: A survey. Journal of Machine Learning Research, 2018.
- Ryan P Adams Jasper Snoek, Hugo Larochelle. Practical bayesian optimization of machine learning algorithms. Advances in Neural Information Processing Systems (NeurIPS), 2012.
- Balázs Kégl Y. Bengio James Bergstra, R Bardenet. Algorithms for hyper-parameter optimization. Advances in Neural Information Processing Systems (NeurIPS), 2011.