

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ
ФЕДЕРАЦИИ

ГОСУДАРСТВЕННОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ
ВЫСШЕГО ПРОФЕССИОНАЛЬНОГО ОБРАЗОВАНИЯ
«МОСКОВСКИЙ ФИЗИКО-ТЕХНИЧЕСКИЙ ИНСТИТУТ
(ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ)»

ФАКУЛЬТЕТ ИННОВАЦИЙ И ВЫСОКИХ ТЕХНОЛОГИЙ
КАФЕДРА АНАЛИЗА ДАННЫХ

Выпускная квалификационная работа по направлению
01.03.02 "Прикладные математика и информатика"

На тему:

ГЛАДКАЯ МЕТРИКА ДЛЯ ЗАДАЧИ РАНЖИРОВАНИЯ

Студент _____ Левков Мирон Николаевич

Научный руководитель _____ Воронцов Александр Сергеевич

Зам. зав. кафедрой _____ Бунина Елена Игоревна

Москва, 2017

Гладкая метрика для задачи ранжирования

Рассматривается задача построения гладкой метрики ранжирования документов в поисковых запросах. Предложены различные способы сглаживания метрики DCG, также исследованы имеющиеся в данной области результаты. Проведен анализ зависимости качества сглаженных метрик от гиперпараметров и размера выборки поисковых запросов.

СОДЕРЖАНИЕ

1	Введение	4
1.1	Описание задачи ранжирования	4
1.2	Обозначения	4
2	Постановка задачи	6
3	Метрики качества результатов	7
4	Метрика SoftDCG	10
4.1	Описание метрики	10
4.2	Рассуждения о метрике	11
5	Метрика NoisedSoftDCG	12
5.1	Описание метрики	12
5.2	Рассуждения о метрике	13
6	Метрика FairSoftDCG	14
6.1	Описание метрики	14
6.2	Рассуждения о метрике	14
7	Применение результатов	16
8	Эксперименты	18
8.1	Зависимость гладкости от размера пула запросов	18

8.2	Зависимость качества аппроксимации от размера пула запросов	21
8.3	Изменение метрики при добавлении шума	22
8.4	Изменение метрики для комбинации двух формул	23
9	Проблема выбора гиперпараметра	25
10	Выводы	27
Список литературы		28

Введение

Описание задачи ранжирования

Пусть имеется некий поисковый запрос. Задача ранжирования состоит в том, чтобы из списка всех доступных документов выбрать наиболее хорошие (релевантные) и показать их пользователю в поисковой выдаче. Для того, чтобы измерять качество алгоритма, решающего данную задачу, используются оценки релевантности документов запросу, проставленные людьми. Метрика качества дает возможность получить более хороший ранжирующий алгоритм, улучшающий качество поиска.

Обозначения

Как правило метрику ранжирования измеряют не на одном запросе: берется большое количество запросов и вычисляется среднее значение метрики по ним.

Рассмотрим один поисковый запрос.

Обозначение 1.1. Обозначим $\{d_i\}_{i=1}^n$ - набор *документов*, которые нужно ранжировать по данному запросу; $\{r_i\}_{i=1}^n$ - *оценки* этих документов, проставленные ассесорами;
 $\{s_i\}_{i=1}^n$ - оценки релевантности (*скоры*), выданные ранжирующим алгоритмом

Определение 1.1. Метрика качества ранжирования **DCG** определяется по формуле:

$$DCG = \sum_{i=1}^k \frac{r_{p_i}}{\text{discount}(i)}$$

где

p_1, \dots, p_n - перестановка на множестве $\{1, \dots, n\}$, т.ч. $s_{p_1} > s_{p_2} > \dots > s_{p_n}$

k - количество документов, по которым считается метрика ($k \leq n$) - важен в том случае, когда нас интересуют лишь несколько первых документов в выдаче

$\text{discount}(i)$ - дисконтирующий фактор, как правило $\frac{1}{\log(i+1)}$ - однако в поисковых метриках Яндекса используется $\frac{1}{i}$, поэтому мы будем использовать его

Постановка задачи

Все DCG-подобные метрики ранжирования имеют такую конструкцию, в которой значение метрики напрямую не зависит от скоров алгоритма (при подсчете значения учитывается лишь перестановка на документах, заданная скорыми алгоритма).

Несложно заметить, что для такой конструкции при фиксированном наборе документов и их оценок данная метрика может принимать лишь конечное количество различных значений. При этом изменение значения метрики происходит лишь в случае перестановки местами двух документов в выдаче.

Как было замечено выше, основная проблема подобных метрик в том, что они не имеют гладкой зависимости от скоров. В то же самое время логично ожидать, что, если ранжирующий алгоритм выдал трем документам оценки $\{100, 1, 0.5\}$, а другим трем документам - $\{10, 1, 0.5\}$, то он считает первый документ из первой тройки сильно лучшим, чем первый документ из второй тройки. В добавок, гладкая относительно скоров метрика является более чувствительным инструментом для оценивания качества алгоритма и отслеживания случаев недообучения или переобучения.

В данной работе исследуются разные подходы к построению метрик ранжирования, с целью получения метрики, которая бы гладко зависела от скоров и хорошо коррелировала с DCG. Далее мы уточним эти критерии.

Метрики качества результатов

Хотелось бы, чтобы наша метрика обладала двумя свойствами. Во-первых, метрика должна быть "реалистичной". Т.е. ее локальные минимумы/максимумы должны соответствовать минимумам/максимумам DCG.

Определение 3.1. Пусть $\{y_i\}_{i=1}^k$ - значения метрики DCG на данном пуле запросов для k различных значений набора гиперпараметров. Пусть $\{\hat{y}_i\}_{i=1}^k$ - значения нашей метрики на том же пуле при тех же наборах гиперпараметров. Тогда **ошибкой аппроксимации** метрики DCG нашей метрикой назовем величину $\min_{\alpha, \beta} \frac{1}{k} \sum_{i=1}^k (\alpha \cdot \hat{y}_i + \beta - y_i)^2$

Обозначение 3.1. Далее будем обозначать *ошибку аппроксимации*, как $\text{err}_{\text{approx}}(\text{Metric Name})$

Также придуманная метрика должна быть гладкой. При этом необходим способ подсчета гладкости конкретной функции, если она задана не аналитически, а численно. Определим метрику гладкости.

Определение 3.2. Пусть $\{y_i\}_{i=1}^k$ - значения нашей метрики на данном пуле запросов для k различных значений набора гиперпараметров. Тогда **ошибкой гладкости** функции, посчитанной **через усреднение модуля разности** в соседних точках, будем называть величину

$$\text{err}_{\text{smooth}}^{\text{abs}}(y_1, \dots, y_k) = \sum_{i=2}^k |y_i - y_{i-1}| \cdot \frac{1}{|y_k - y_1|}$$

Определение 3.3. Пусть $\{y_i\}_{i=1}^k$ - значения нашей метрики на данном пуле запросов для k различных значений набора гиперпараметров. Тогда **ошибкой гладкости** функции, посчитанной **через нормировку**

дисперсии разности в соседних точках, будем называть величину

$$\text{err}_{\text{smooth}}^{\text{std}}(y_1, \dots, y_k) = \frac{\text{diff}_{\text{std}}}{|\text{diff}_{\text{mean}}|}$$

Здесь $\text{diff}_{\text{mean}} = \frac{1}{k-1} \sum_{i=2}^k (y_i - y_{i-1}), \quad \text{diff}_{\text{std}} = \frac{1}{k-1} \sum_{i=2}^k (y_i - y_{i-1} - \text{diff}_{\text{mean}})^2$

Определение 3.4. Пусть $\{y_i\}_{i=1}^k$ - значения нашей метрики на данном пуле запросов для k различных значений набора гиперпараметров. Пусть $\text{window} \in \mathbb{Z}, \text{deg} \in \mathbb{N}$. Тогда **ошибкой гладкости** функции, посчитанной **с помощью аппроксимации полиномами**, будем называть величину

$$\text{err}_{\text{smooth}}^{\text{poly}}(y_1, \dots, y_k) = \sum_{i=1}^{k-\text{window}+1} \frac{\left(\text{poly}_{\text{deg}}(y_i, \dots, y_{i+\text{window}-1})_{i+\lfloor \frac{\text{window}}{2} \rfloor} - y_{i+\lfloor \frac{\text{window}}{2} \rfloor} \right)^2}{k - \text{window}}$$

Здесь $\text{poly}_{\text{deg}}(y_i, \dots, y_{i+\text{window}-1})$ - аппроксимирующий полином степени deg , построенный по window точкам. Т.е. такой полином степени deg , что среднеквадратичное отклонение в точках $\{i, \dots, i+\text{window}-1\}$ минимально.

Возникает ряд проблем с тем, что выбор конкретного способа измерения гладкости для сравнения наших метрик между собой неочевиден. Первая метрика хорошо отображает гладкость в том случае, если функция монотонна (если не учитывать шумовые колебания) - однако же в иных случаях данная метрика может быть плоха из-за нормировки на разность значений в крайних точках.

Вторая метрика ведет себя лучше, однако тоже не является достаточно гибкой и интерпретируемой.

С последней метрикой встает проблема выбора параметров deg и window . Тем не менее эта метрика понятна и действительно способна достаточно хорошо отображать гладкость функции. После перебора разных вариантов выбор был сделан в пользу параметров $\text{deg} = 3, \text{window} = 11$. При этом

стоит учитывать, что гладкость метрики ранжирования измерялась на множествах размера порядка 100 точек.

Видно, что выбор метрики гладкости - отдельная задача, сложность которой заключается в том, что мы пытаемся мерять гладкость дискретно заданной функции. При этом метрика гладкости в работе выбиралась так, чтобы отразить наши представления о «хорошей» функции.

Метрика SoftDCG

Рассмотрим метрику SoftDCG. Данная метрика является попыткой сглаживания метрики DCG. Кратко разберем способ подсчета SoftDCG

Описание метрики

Пусть наш набор документов в поисковом запросе уже упорядочен по скорам. Рассмотрим конкретный документ d_j со скором s_j . Логично предположить, что скор не является точным определением того, на каком месте должен быть документ. Более точно: пусть $s_j - s_{j+1} < \varepsilon$, $s_{j-1} - s_j > \varepsilon$, тогда, если добавить к оценкам ранжирующего алгоритма случайный шум с дисперсией ε , документы (s_j, s_{j+1}) поменяются местами в выдаче более вероятно, чем (s_{j-1}, s_j) . Данный факт говорит о близости документов, однако DCG способен учитывать близость документов с точностью до 1 места.

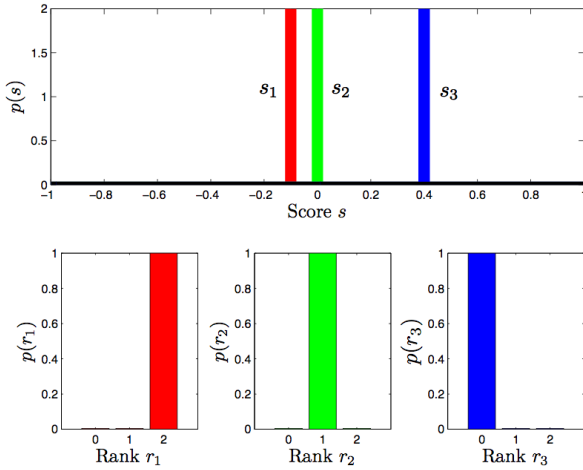


Figure 1: Deterministic scores and ranks: Three document scores as point (deterministic) values and their corresponding rank distributions. The lowest scoring document s_1 is certain to be ranked in the lowest position 2.

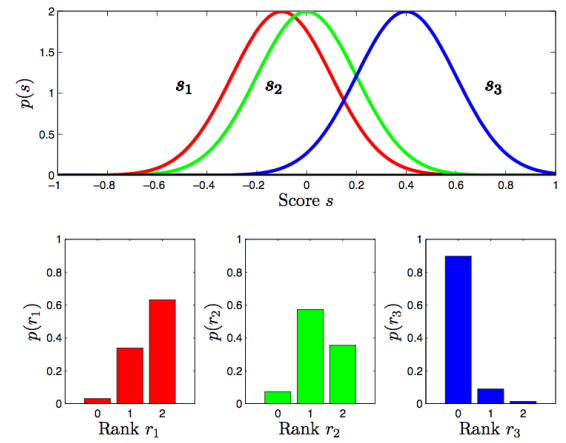


Figure 2: From score to rank distributions: Smoothed scores for 3 documents and the resulting 3 rank distributions.

Метрика SoftDCG предполагает следующее: вместо того, чтобы каждому документу сопоставлять конкретное место в поисковой выдаче, можно сопоставить ему нормальное распределение на скорях. Т.е. doc_i имеет скор t_i из распределения $t_i \sim \mathcal{N}(s_i, \sigma^2)$, где σ - гиперпараметр.

В такой модели

$$\pi_{ij} := P(doc_i \text{ выше } doc_j) = P(t_i - t_j > 0) = \int_{-\infty}^{\infty} \mathcal{N}(x|s_i - s_j, 2\sigma^2) dx$$

Из этого выводятся распределения на местах для каждого документа. При этом можно посчитать ожидаемую позицию документа:

$$\bar{pos}(j) = \sum_{i=1, i \neq j}^n \pi_{ij}$$

На этом строится подсчет среднего дискаунта каждого документа:

$$\hat{d}(j) = \sum_{i=1}^n p_j(i) \text{discount}(i) - \text{средний дискаунт } j\text{-го документа}$$

Далее средние дискаунты используются при подсчете SoftDCG:

$$\text{SoftDCG} = \sum_{i=1}^n \hat{d}_j \cdot r_j$$

Рассуждения о метрике

В то же самое время есть основания полагать, что данная метрика является не совсем «честной», т.к. средняя позиция документа вычисляется, как мат ожидание числа документов выше. Правильный способ подсчета средней позиции документа делается через введение вероятностного распределения на перестановках документов. Описанные в следующих разделах метрики по разному оценивают среднюю позицию документа.

Метрика NoisedSoftDCG

Рассмотрим другую метрику (назовем ее NoisedSoftDCG). Идея данной метрики берет свое начало в статье про SoftDCG, однако способ подсчета данной метрики несколько иной.

Описание метрики

В SoftDCG каждому документу сопоставлялось некоторое вероятностное распределение на множестве возможных занимаемых позиций. При этом следует отметить, что данное распределение вполне четко задавалось аналитически. Благодаря этому имелась возможность задать такой параметр, как $\hat{d}(j)$ (средний дискант), формулой.

В данном же случае при подсчете метрики принимается во внимание лишь тот факт, что при небольшом изменении скоров на случайные величины значение метрики может меняться в силу того, что близкие документы будут переставляться местами.

Алгоритм подсчета NoisedSoftDCG следующий:

- (1) Получим скоры ранжирующего алгоритма s_1, \dots, s_n
- (2) Сгенерируем случайный шум $\xi_1, \dots, \xi_n \sim \mathcal{N}(0, \sigma)$
- (3) $\hat{s}_1 = s_1 + \xi_1, \dots, \hat{s}_n = s_n + \xi_n$
- (4) Посчитаем значение метрики Value_{DCG}

- (5) Повторим шаги (2)-(4) достаточно большое количество раз (T) и вычислим $NoisedDCG = \sum_{i=1}^T Value_{DCG}^i$

Рассуждения о метрике

Можно дать некую интуицию по поводу того, почему метрика считается именно так. При достаточно большом T близкие документы будут часто меняться местами, а усреднение DCG при этом даст некий аналог среднего дискаунта. Таким образом ожидается, что при больших T данная метрика будет вести себя хорошо в плане гладкости

Метрика FairSoftDCG

Еще один вариант метрики, идея которой схожа с NoisedDCG. Заметим, что, основываясь на скорях, выданных ранжирующим алгоритмом, можно ввести вероятностное распределение на перестановках всех документов.

Описание метрики

При этом распределение вводится так, что $P(d_i \text{ выше } d_j) \sim \frac{e^{\sigma s_j}}{e^{\sigma s_i} + e^{\sigma s_j}}$ (данное распределение можно ввести единственным образом). Благодаря этому для каждой перестановки на документах p_1, \dots, p_n можно посчитать ее вероятность:

$$P(p_1, \dots, p_n) = \prod_{i=1}^{n-1} \frac{e^{\sigma s_{p_i}}}{\sum_{k=j}^n e^{\sigma s_{p_k}}}$$

Теперь перебрав все перестановки документов можно посчитать "честное" мат. ожидание DCG. Собственно это у будет значением FairSoft-DCG.

$$FairSoftDCG = \sum_{p_1, \dots, p_n} P(p_1, \dots, p_n) \cdot DCG(p_1, \dots, p_n)$$

Рассуждения о метрике

При этом возникает ряд проблем с данной метрикой, т.к. вычислительно данная задача достаточно сложная (всего имеется $n!$ возможных перестановок). В связи с этим был выбран способ подсчета близкого к

данной метрике значения - подсчета по топ-k документам:

$$FairSoftDCG = \sum_{p_1, \dots, p_k} P(p_1, \dots, p_k) \cdot DCG(p_1, \dots, p_k),$$

где сумма берется по всем возможным размещениям по k из n элементов. k берется таким, чтобы метрику можно было посчитать, т.к. n большое и перебор $n!$ перестановок - долгая операция.

Применение результатов

В данной части будут приведены предполагаемые способы использования искомой метрики. На данный момент идеальной - т.е. гладкой и хорошо аппроксимирующей DCG - метрики не найдено.

Итак. Первый и способ применения более гладкой метрики - это использование ее в качестве инструмента показывающего качество обучения модели. Благодаря большей чувствительности появляется возможность следить за тем, как малые изменения гиперпараметров влияют на качество модели.

Простой пример: благодаря этой метрики можно достаточно хорошо измерять влияние новых элементов ансамбля в моделях градиентного бустинга. Точно так же можно использовать данную метрику для active learning like задач. Благодаря высокой чувствительности можно будет немного дообучать алгоритм на небольших частях датасета и смотреть на изменение качества ранжирования.

Второй и, пожалуй, важный с практической точки зрения способ применения нашей метрики - смешивание нескольких моделей ранжирования. В наши дни в ранжирующих системах крупных компаний используются абсолютно разные по характеру формулы. Эти формулы имеют разный масштаб выдаваемых значений. Каждая из этих формул придумывалась для оптимизации конкретной метрики - причем метрики могут быть кардинально различными по смыслу. При этом хотелось бы уметь смешивать формулы, выданные несколькими различными моделями, в один ансамбль формул, получая таким образом более мощную формулу ранжирования.

Т.к. наша метрика является гладкой и хорошо аппроксимирует DCG, оптимальная комбинация формул согласно этой метрике должна быть в той же окрестности, что и оптимум с точки зрения DCG. В то же самое время естественно ожидать, что при взятии выпуклой комбинации двух формул должна получаться некая выпуклая гладкая кривая для значений метрики.

В качестве своеобразной sanity-check - добавления шума к скорам ранжирующего алгоритма. Здесь имеется ввиду следующее: если добавить к скорам модели шум с неким коэффициентом, это эквивалентно смешиванию двух моделей - нормальной и абсолютно бесполезной. Логично полагать, что оптимум качества в данном случае будет при коэффициенте, с которым берется шум, равном 0 (или близком к 0).

Эксперименты

Зависимость гладкости от размера пула запросов

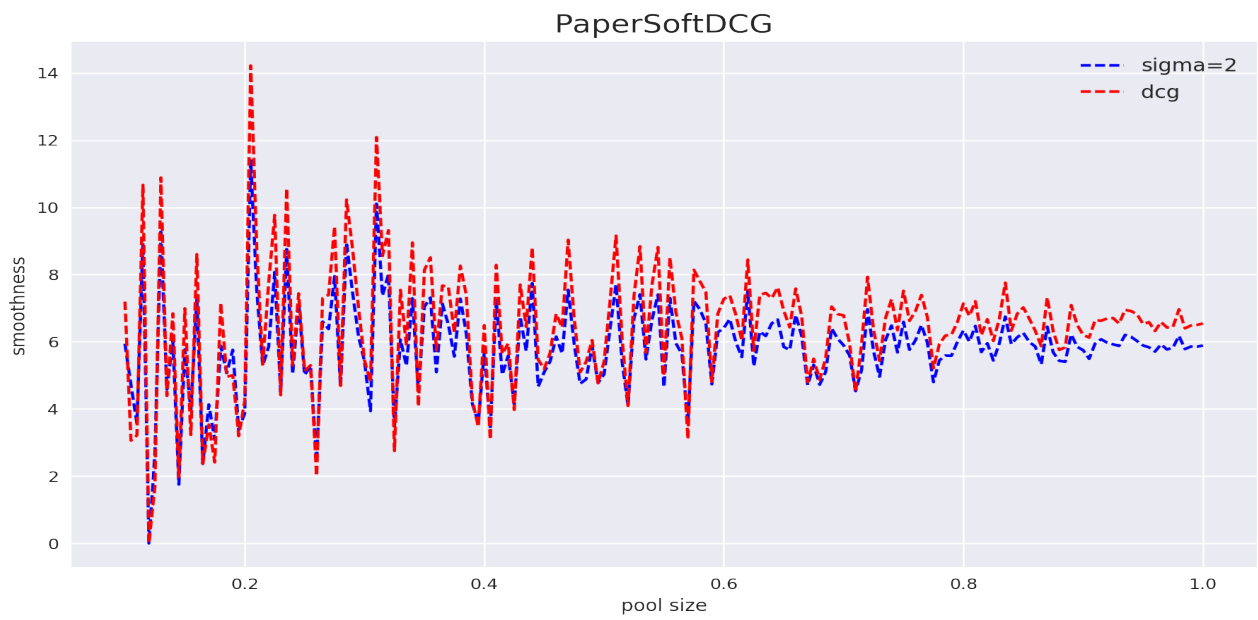
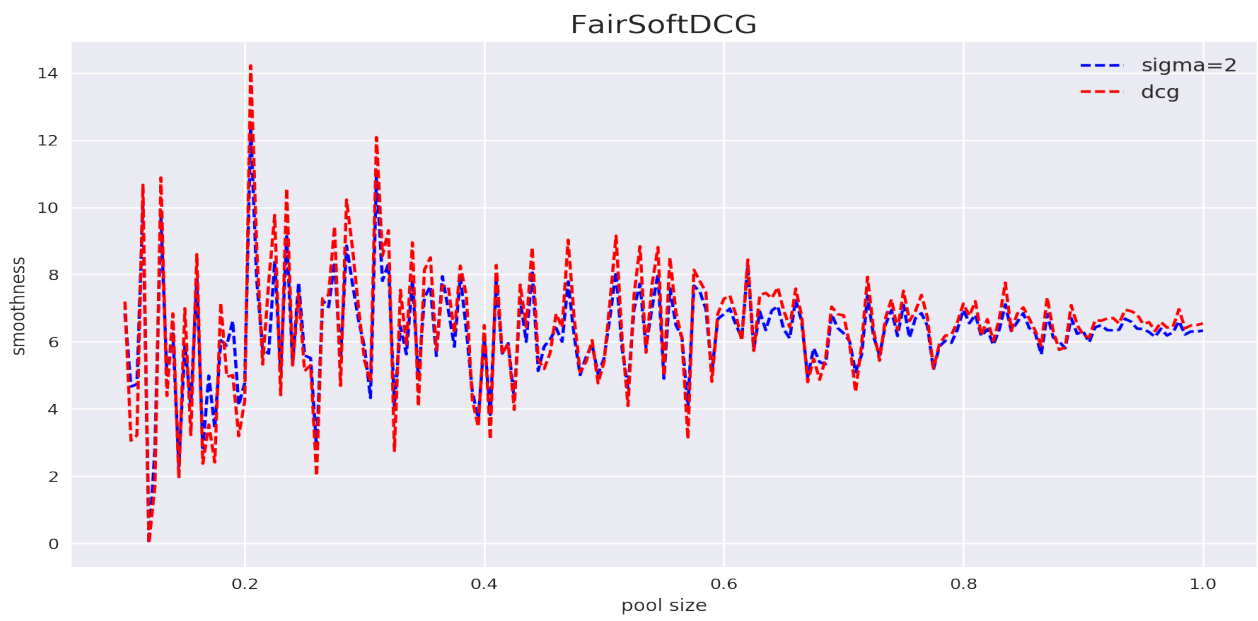
Первый эксперимент является одновременно способом проверить наш метод измерения гладкости на адекватность и способом сравнения между собой рассмотренных ранее метрик ранжирования. Здесь и далее рассуждения будут вестись для метрики DCG, но они точно так же верны для любой метрики из списка SoftDCG, NoisedSoftDCG, FairSoftDCG.

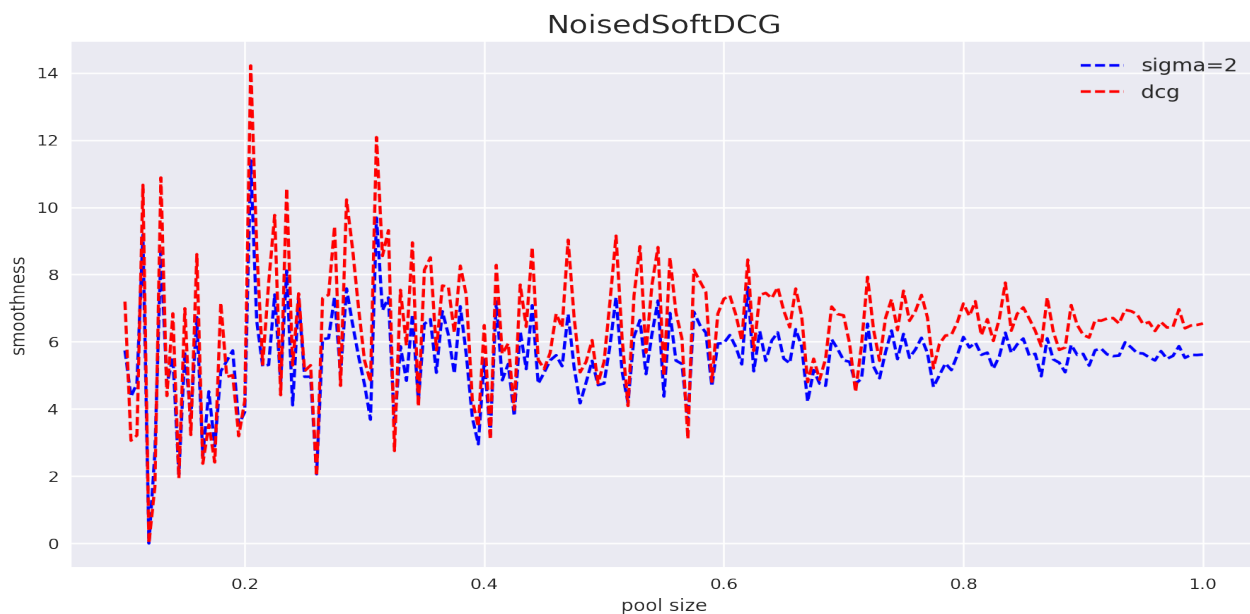
Для начала вспомним, что такое в нашей задаче "размер пула". Размер пула - это количество обрабатываемых поисковых запросов. Для каждого запроса мы считаем значение DCG, а потом производим усреднение данного значения по всем запросам. Все это происходит при каком-то зафиксированном наборе гиперпараметров.

111

Таким образом для каждого значения набора мы получаем какое-то значение метрики. Например, можно смотреть на значение метрики для ранжирующего алгоритма, который является выпуклой комбинацией двух других алгоритмов, в зависимости от коэффициента, с которым алгоритмы смешиваются. Если рассматривать данный график, на множестве коэффициентов $[0; 1]$, можно посчитать его гладкость. При этом, рост размера пула - т.е. увеличение количества запросов, по которому мы усредняем ответ, должно играть роль сглаживающего фактора.

Данный эксперимент состоял в том, чтобы посмотреть на изменение гладкости с увеличением роста пула. Хочется заметить, что т.к. увеличивать





пул мы возможности не имели, делалось эквивалентное действие - его уменьшение. Ожидалось, что для каждой метрики будет получен график, имеющий обратную пропорциональность - чем больше размер пула, тем меньше гладкость. Следует также напомнить, что в наших терминах гладкость кривой - тот параметр который хотелось бы уменьшить.

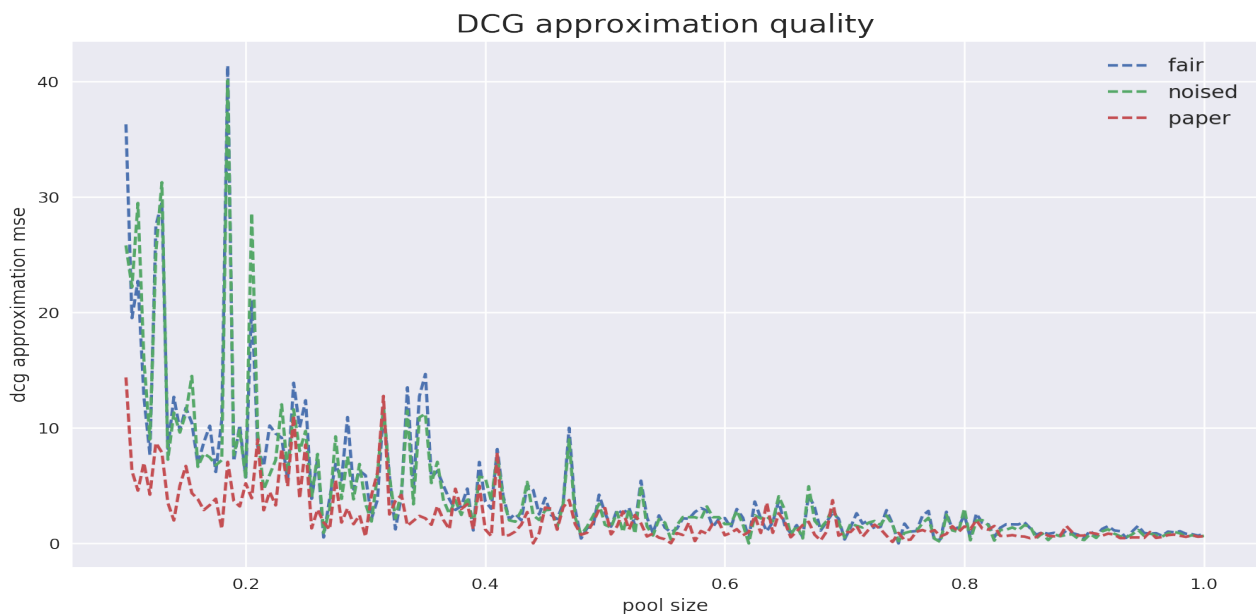
На графиках для сравнения приведена гладкость DCG и гладкость соответствующих метрик. Для наглядности значения метрик брались при одном и том же значении гиперпараметра $\sigma = 2$. Можно заметить, что самой гладкой в нашей терминологии являлась метрика NoisedSoftDCG. При совсем маленьких размерах пула (0.6 от изначального размера и меньше) графики, как и ожидалось, являются слишком шумными.

Зависимость качества аппроксимации от размера пула запросов

Аналогично предыдущему эксперименту логично было бы ожидать графики, на которых с увеличением размера пула аппроксимация DCG уменьшается.

Аппроксимация в нашей нотации имеет тот же смысл, что и гладкость: меньшие значения соответствуют более близкой к DCG метрике.

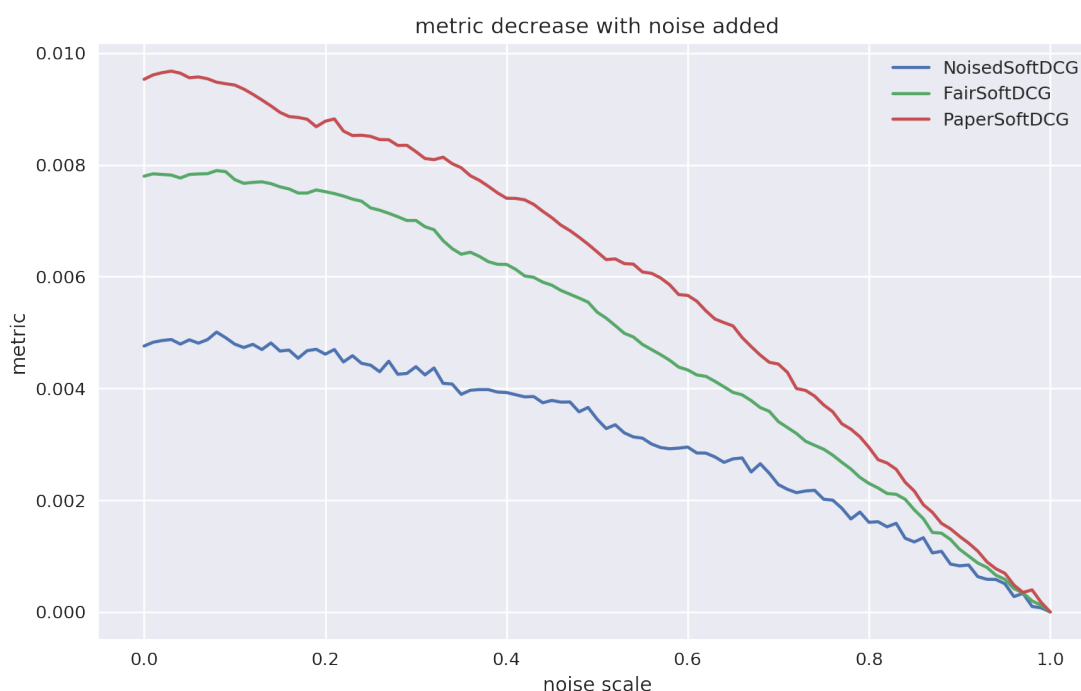
1



Что же касается аппроксимации DCG, можно заметить, что график для NoisedSoftDCG является наиболее шумным, а FairSoftDCG - наименее.

Изменение метрики при добавлении шума

Как упоминалось ранее, первоочередным действием в работе было желание убедиться в том, что метрика гладкости действительно несет ту смысловую нагрузку, которая на нее была возложена. Для этого проводился эксперимент по смешиванию двух ранжирующих формул, одна из которых была настоящей, а вторая - просто сгенерированным случайным шумом.

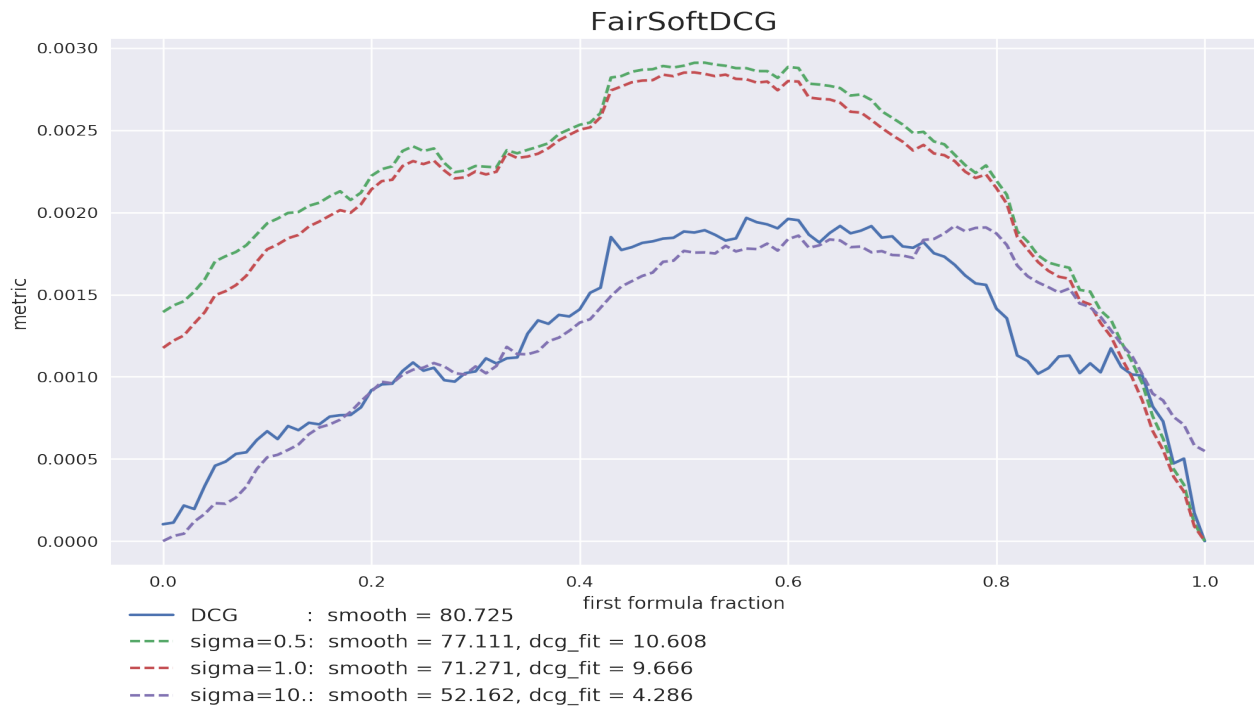


По оси x отложен коэффициент, с которым прибавляется случайный шум, по оси y - значения метрики. Можно заметить, что наши метрики убывают (т.е. они логичны) причем визуально видно, что они делают это почти монотонно.

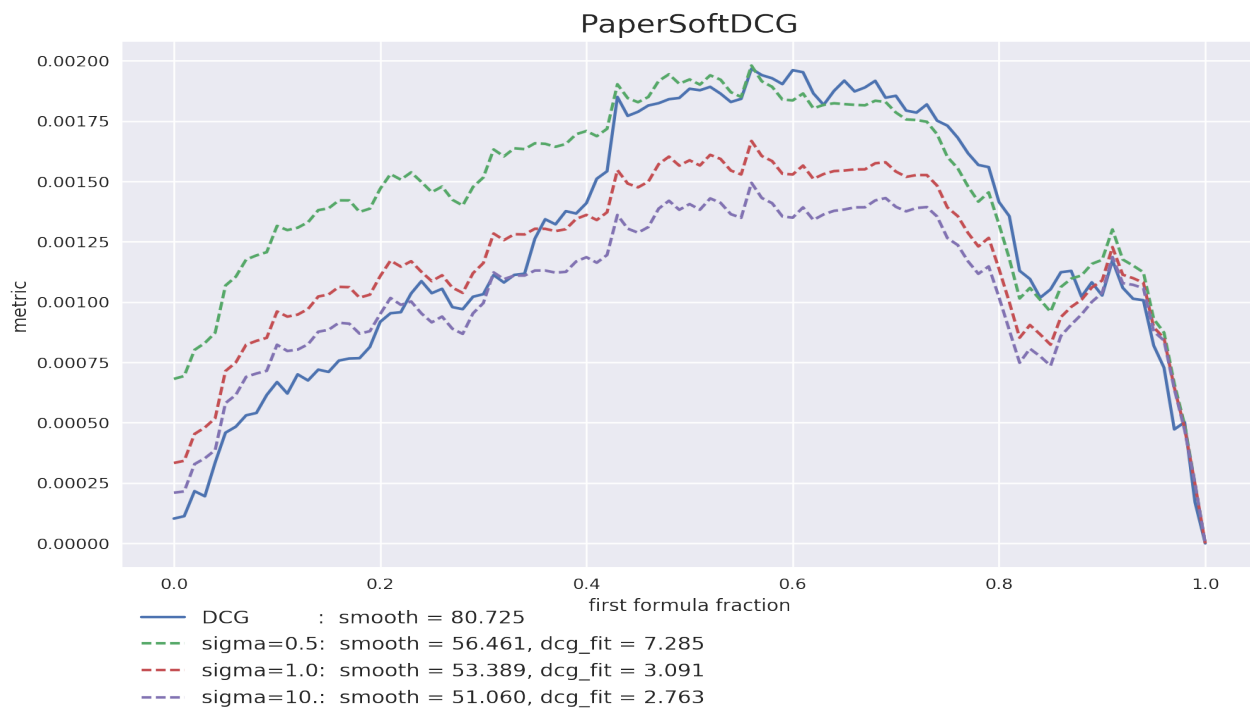
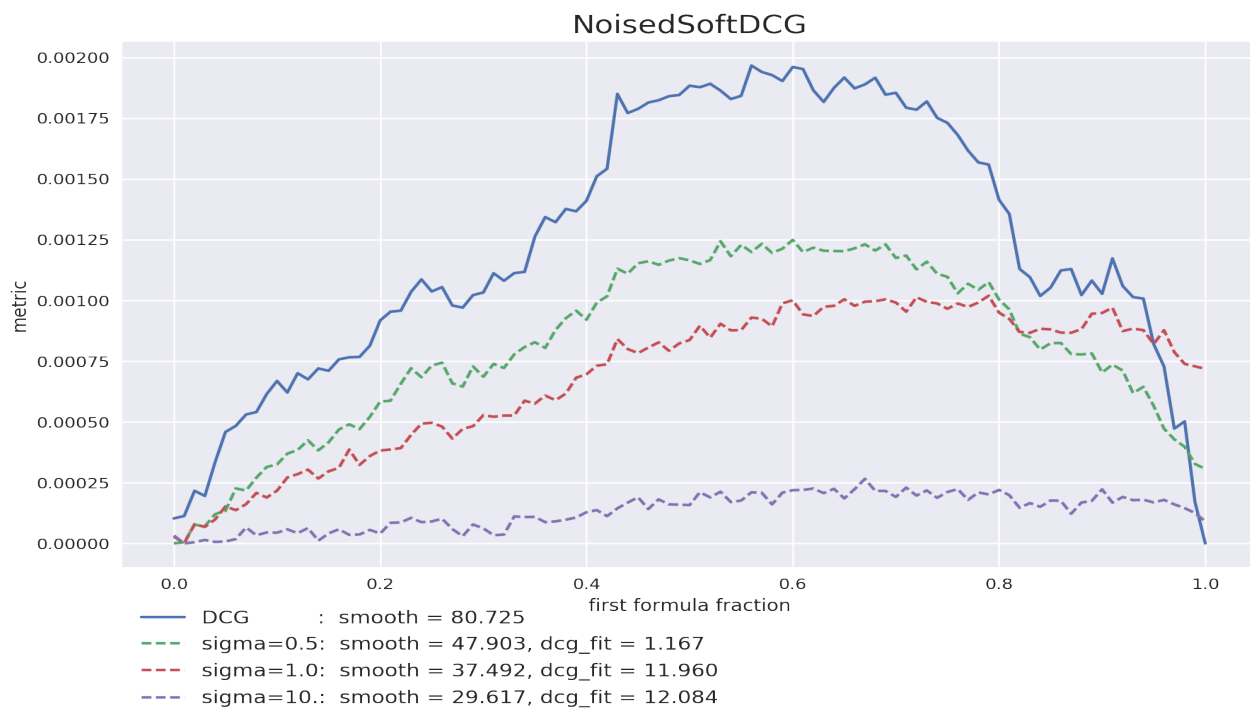
Подобное поведение в данном эксперименте является необходимым условием для класса искомых метрик.

Изменение метрики для комбинации двух формул

111



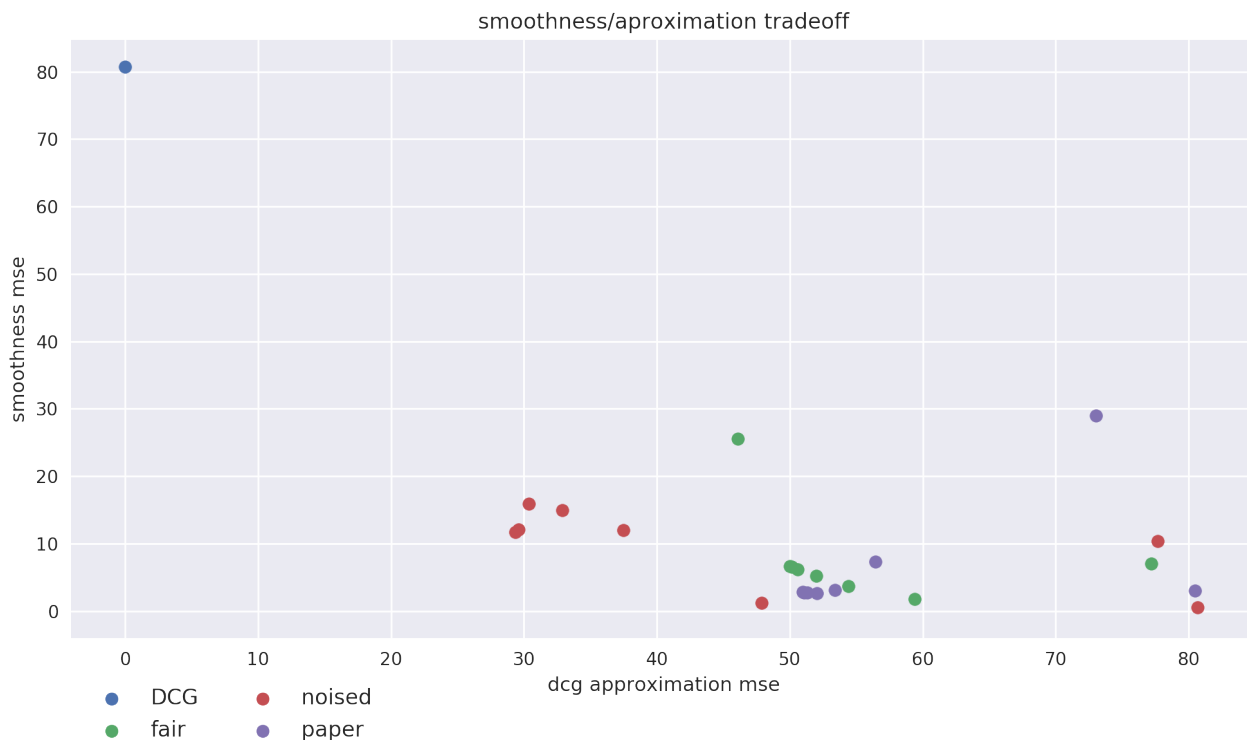
Самый важный с точки зрения практического применения эксперимент - эксперимент по смешиванию двух формул ранжирования. Цель эксперимента - проверить, какая из наших метрик лучше отображает реальный оптимум качества при использовании в качестве ранжирующей формулы выпуклую комбинацию двух формул.



Проблема выбора гиперпараметра

Так как главной целью работы было получение гладкой метрики, с помощью которой можно будет удобно и правильно подбирать коэффициенты в смеси ранжирующих формул, остановимся на этом эксперименте.

Как можно было заметить из графиков, приведенных выше, при разных значениях гиперпараметра σ получаются разные качества аппроксимации/гладкости. При этом, логично предположить следующее: чем лучше аппроксимация (т.е. чем ближе σ к 0), тем хуже гладкость кривой. Ниже приведен график в осях, соответствующих данным двум показателям, метрик с разными гиперпараметрами.



Как видно из графика, для каждой метрики σ можно варьировать этот гиперпараметр, получая тем самым желаемый компромисс между гладкостью метрики и хорошей аппроксимацией DCG.

Следует отметить, что выбор параметра σ - неочевидная задача. В итоге, возникает проблема определения σ . Она заключается в том, что необходимо выбирать данный параметр, опираясь, лишь на скоры, выданные каждой из формул. Однако скоры в каждом запросе свои. В добавок функция, которую мы при этом пытаемся минимизировать: $\alpha \cdot approx + \beta \cdot smoothness$, не является гладкой относительно σ .

Данную проблему можно описать следующим образом. Представим себе смесь двух формул ранжирования, где вторая формула совпадает в первой, но выдает значения в 10 раз больше. Очевидно вторая формула более устойчива к добавлению шума - из-за этого при том же самом ранжирующем алгоритме мы получаем другую σ .

Есть некоторые эвристики по подбору σ , но хорошего решения в ходе исследования данного вопроса найдено не было.

Выводы

В ходе исследования была рассмотрена реализация метрики SoftDCG, идея которой описана в [1]. Также были рассмотрены предложенные автором аналоги этой метрики: NoisedSoftDCG и FairSoftDCG.

При изучении этих метрик были введены такие параметры, как "гладкость" и "аппроксимация DCG". Было проведено исследование касательно изменения этих показателей при изменении различных гиперпараметров. Это дало возможность понять, что представленные метрики в целом являются неким продвижением в нужном направлении - но окончательного решения не дают.

Была исследована проблема выбора гиперпараметра σ , который соответствует тому, насколько сильно наша метрика стремится раздвинуть между собой документы (фактически - порядок разницы между скорями близких документов).

Список литературы

- [1] M. Taylor, J. Guiver, S. Robertson and T. Minka. SoftRank: Optimising Non-Smooth Rank Metrics. Microsoft Research Cambridge, 2016
- [2] Christopher J.C. Burges. From RankNet to LambdaRank to LambdaMART. Microsoft Research Technical Report, 2010
- [3] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, Hang Li. Learning to Rank: From Pairwise Approach to Listwise Approach. Microsoft Research Technical Report, 2007