



Gunosy | データ分析ブログ | LayerX | gunosy.fm | 採用情報

2017-11-30

Gunosyのパーソナライズを支える技術 -計算モデルとアーキテクチャ編-

この記事は Gunosy Advent Calendar 2017 の1日目の記事です(フライング)

Gunosy Advent Calendar 2017 - Qiita

Gunosyは情報を増やすことを目指しています。そんなGunosyで働くエンジニアが日々語っている[Gunosyデータ分析ブログ][http://data.gunosy...]

qiita.com 2 users

qiita.com

§1.はじめに

こんにちは。データ分析部ロジックチームの @mathetake です。いつもはデータ分析ブログにいるのでテックブログは初めてです。怖いです。Twitterとかやったことないですね。

最近は仕事でニュースバスというプロダクトの記事配信ロジックの改善を行っており、その一環としてパーソナライズロジックの開発プロジェクトに従事しています。

パーソナライズとはユーザーひとりひとりに対して別々の記事配信を行うことです。下記の記事でパーソナライズプロジェクト発足に至るまでの背景が語られているので、興味のある方はぜひご覧ください。

Q Gunosyのパーソナライズは私たちを知り、Gunosyを知る。

「ユーザーの課題をコードで解決する」Gunosyが考へるテックリード(Tech Lead)の役割...

こんな感じで広範なお悩みがあります。Gunosyではさまざまな難題を絶えず抱いていますが、今回はじめた「テックリード(Tech Lead)」の実験をするようになりました! テックリードがどういった仕事をするかなど...

2017-11-02 09:00

gunosy.gunosy.co.jp

この記事ではニュースバスの記事配信アルゴリズムのパーソナライズプロジェクトについて解説します。

パーソナライズの背景

・アーキテクチャ概説

・記事スコア計算モデルの概要

・パーソナライズされた記事リストを返すAPIサーバーの処理

について書きます。ほぼ毎月行われているGunosy Tech Night*2での発表内容を元にしています:

§2.背景とマイルストーン

パーソナライズpjの背景は、上記の記事の内容(a)をまとめてみると

・バッチ処理でユーザー属性の数だけリストを事前に生成することの限界

◦興味関心の取りこぼしが生まれる

◦組み合わせの数だけ処理が必要になりスケールしない

・技術的負債

◦最重要である記事配信ロジックの改善のサイクルが回りづらい

◦カジュアルにロジックのテストが行えない

◦複雑になり組んだコードにならざるため新しいメンバーのキャッチアップが大変

こんな感じです。

このような背景があり、パーソナライズプロジェクトの最初のマイルストーンとして例えば

・ユーザーのリクエストリガーで記事リストをリアルタイムに生成する

・リアルタイムにユーザーの興味嗜好の変化を反映する仕組みの導入

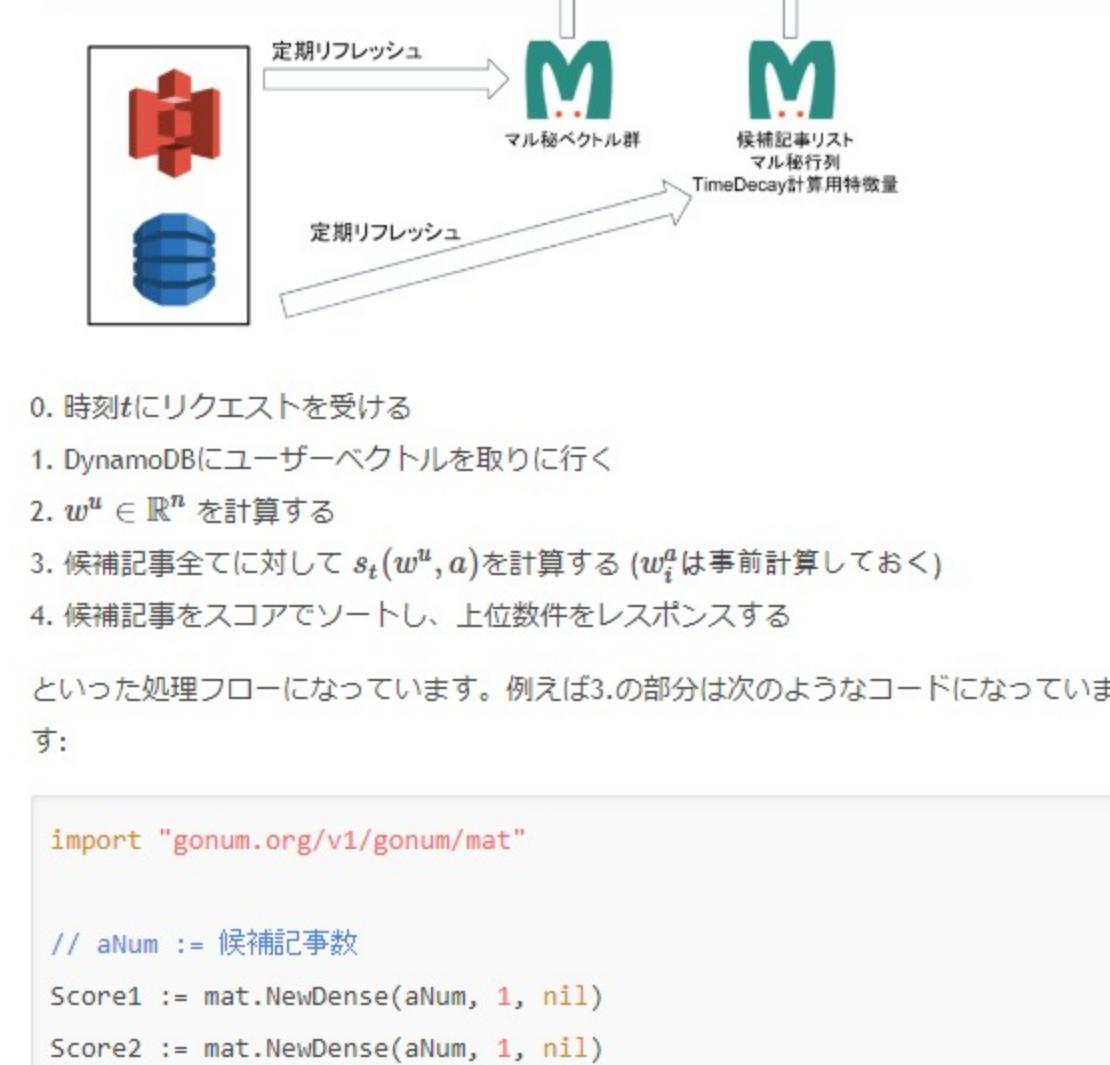
・既存の記事生成仕組みをわざわざゼロから作り直す

を設定しました。

§3.アーキテクチャ

まずパーソナライズシステム全体のアーキテクチャについて簡単に説明します。

大人の事情により勝手にマル秘と表示されている箇所があります。困った困った。



ざっくりまとめると

・後述の計算モデルに必要なものは全てDynamoDBとs3におく

・関連するほぼすべてのタスクのワークフローはdigdagにより管理

・例外はリアルタイム性を求めるユーザーべクトルを生成するAWS lambda

こんな感じです。digdagの運用とリアルタイムベクトル生成AWS lambdaについては後日別記事で解説します。

開発言語ですが

・Spark on EMRクラスタを用いてマル秘機械学習を行っている部分はScala

・パーソナライズリスト返す記事APIはgolang

・その他タスクは基本的にPython

と言った感じです。

クライアントからのリクエストは必ず一番右にある **Gateway API** を介し、ABテストのために複数存在する記事APIのうち、どのエンドポイントを叩くかを制御しています。これにより無駄なリクエストを減らすことができます。つまりソースコードを汚すことなくABテストを行うことができます。

以降では、記事APIがどのような計算モデルによって記事をスコアリングし、リストを生成するかについて述べてきます。

§4.スコア計算モデルとハイパーパラメータ

まず最初に、スコア、スコア計算モデルという言葉の定義をはっきりさせておきます。

スコアとは記事を出すべきかを表す数値のことです。候補となる記事全てに対してスコアを計算し、それを元にソートしてからユーザーに表示するというのが基本的なアルゴリズムになります。

スコア計算モデルとは、ユーザー u と記事 a を引数にしてスコアを算出する時間によって変化する関数 s_t

$$s_t : \{ \text{ユーザー} \} \times \{ \text{ニュース記事} \} \rightarrow R, (u, a) \mapsto s_t(u, a), t : \text{時間}$$

のことです。我々の目標はより良い s_t を構築する事ですが、その制約条件として

リアルタイムに数千からなる候補記事全てに対するスコア計算を行なう

があります。

そこで今回ver.1として、少しふわっとしてますが

・記事APIに行なう計算部分は行列演算のみで済む線形モデル

・線形モデルの各項に非線形効果が入るよう非同期なタスクを裏側で走らせておく

このようなスコア計算モデルを設計しました。

これにより、記事APIが行なう処理は行列演算のみなので現実的な時間内にレスポンスする

ことが可能になり、かつ非線形効果も入るのでそこそこの良いモデルになっています。

より詳細に(と言ってもマル秘を含みます)が解説するために、まず線形モデルに組み込まれるユーザーの特徴量が生成されるまでの流れを説明します:



1. ユーザーはクリックするたびに社内ではファインマンベクトルと呼ばれてる呼ばれているベクトルがリアルタイムで生成が更新される

2. 1.とは非同期的にある時点でのファインマンベクトルを大量に用いて、Sparkで実装されたマル秘機械学習アルゴリズムのおかげで、記事に随れて n 個のマル秘ベクトルを生成する

3. 記事APIがリクエストしたタイミングでのユーザーのファインマンベクトルと、2.で生成された n 個のマル秘ベクトルとの距離(みたいなもの)を計算

と言った感じです。ここで出来上がるユーザーの特徴量を $w^u \in \mathbb{R}^n$ としておきます。

次に線形モデルに入力される記事の特徴量についてです。ファインマンベクトルの生成方法とともに秘機械学習アルゴリズムのおかげで、記事に随れて n 個のマル秘ベクトルとその大きさを複数の尺度(m)で測る結果が出来ます。

つまり、複数の n 次元ベクトル $\{w_i^u\}_{i=1}^m \subset \mathbb{R}^n$ を生成することができます。

これらを用いてスコア計算モデル s_t を

$$s_t(u, a) := \text{TimeDecay}(t, a) \times \left(\sum_{i=1}^m \alpha_i \times \langle w_i^u, w_i^a \rangle \right), t : \text{リクエスト時間}$$

として設計します。このモデルは全てベクトルの内積と和の操作で完結するため、非常に高速に計算が出来ます。

ここで TimeDecay(t, a)は時間減衰関数(time decay function)と呼ばれるユーザーからリクエスト時間とニュース記事の情報を引数とした関数で、「基本的に時間が経つほど値が小さくなる関数」となっています。

例えば時間減衰関数に関しては次の記事が参考になります:

Gentle Intro to Function Scoring | Elastic

We'll cover the basics of scoring using functions while taking a look at some use cases where functional scoring techniques are highly useful and effective.

www.elastic.co 2 users

hrmos.co

また、 $\alpha_1, \alpha_2, \dots, \alpha_n$ はハイパーパラメータと呼ばれ、各独立したスコア(w^u, w^a)を最終的なスコアにどの程度寄与するかの重みです。

このハイパーはユーザー属性を上手く調整することで全く気色の異なる記事リストをレスポンスしたりすることができます。

もちろん各種計算アルゴリズムの重みを調整することで最も良いモデルになります。

そこで私はこのハイパーを実装する分野で最も優秀なエンジニアを絶賛募集中です! ご応募ください!

www.elastic.co

そこで記事とユーザーをえた時、ユーザーに表示すべきかどうかを表す数値のことです。候補となる記事全てに対してスコアを表示する時間によって

スコア計算モデルとは、ユーザー u と記事 a を引数にしてスコアを算出する時間によって

変化する関数 s_t

と言った感じです。

これにより、記事APIが行なう処理は行列演算のみなので現実的な時間内にレスポンスする

ことが可能になり、かつ非線形効果も入るのでそこそこの良いモデルになっています。

より詳細に(と言ってもマル秘を含みます)が解説するために、まず線形モデルに組み込まれるユーザーの特徴量が生成されるまでの流れを説明します:

1. ユーザーはクリックするたびに社内ではファインマンベクトルと呼ばれてる呼ばれているベクトルがリアルタイムで生成が更新される

2. 1.とは非同期的にある時点でのファインマンベクトルを大量に用いて、Sparkで実装されたマル秘機械学習アルゴリズムのおかげで、記事に随れて n 個のマル秘ベクトルを生成する

3. 記事APIがリクエストしたタイミングでのユーザーのファインマンベクトルと、2.で生成された n 個のマル秘ベクトルとの距離(みたいなもの)を計算

と言った感じです。ここで出来上がるユーザーの特徴量を $w^u \in \mathbb{R}^n$ としておきます。

次に線形モデルに入力される記事の特徴量についてです。ファインマンベクトルの生成方法とともに秘機械学習アルゴリズムのおかげで、記事に随れて n 個のマル秘ベクトルとその大きさを複数の尺度(m)で測る結果が出来ます。

つまり、複数の n 次元ベクトル $\{w_i^u\}_{i=1}^m \subset \mathbb{R}^n$ を生成することができます。

これらを用いてスコア計算モデル s_t を

$$s_t(u, a) := \text{TimeDecay}(t, a) \times \left(\sum_{i=1}^m \alpha_i \times \langle w_i^u, w_i^a \rangle \right), t : \text{リクエスト時間}$$

として設計します。このモデルは全てベクトルの内積と和の操作で完結するため、非常に高速に計算が出来ます。

ここで TimeDecay(t, a)は時間減衰関数(time decay function)と呼ばれるユーザーからリクエスト時間とニュース記事の情報を引数とした関数で、「基本的に時間が経つほど値が小さくなる関数」となっています。

例えば時間減衰関数に関しては次の記事が参考になります:

Gentle Intro to Function Scoring | Elastic

We'll cover the basics of scoring using functions while taking a look at some use cases where functional scoring techniques are highly useful and effective.

www.elastic.co 2 users

hrmos.co

また、 $\alpha_1, \alpha_2, \dots, \alpha_n$ はハイパーはマル秘と呼ばれ、各独立したスコア(w^u, w^a)を最終的なスコアにどの程度寄与するかの重みです。

このハイパーはユーザー属性を上手く調整することで全く気色の異なる記事リストをレスポンスしたりすることができます。

もちろん各種計算アルゴリズムの重みを調整することで最も優秀なエンジニアを絶賛募集中です! ご応募ください!

www.elastic.co

そこで記事とユーザーをえた時、ユーザーに表示すべきかどうかを表す数値のことです。候補となる記事全てに対してスコアを表示する時間によって

スコア計算モデルとは、ユーザー u と記事 a を引数にしてスコアを算出する時間によって

変化する関数 s_t

と言った感じです。

これにより、記事APIが行なう処理は行列演算のみなので現実的な時間内にレスポンスする

ことが可能になりました。

より詳細に(と言ってもマル秘を含みます)が解説するために、まず線形モデルに組み込まれるユーザーの特徴量が生成されるまでの流れを説明します:

1. ユーザーはクリックするたびに社内ではファインマンベクトルと呼ばれてる呼ばれている