





# 組織に創発現象を起こす クラウド型データ分析環境

西川大亮(新日鉄住金ソリューションズ(株))

本稿では, KDD Cup 2015 で 2 位に入賞したチーム (FEG&NSSOL@DataVeraci) の取り組み内容におい て、他チームと比べて特徴的な部分を中心に解説する.

### グループ企業内での合同チーム

筆者の所属チームは、システムインテグレータ の新日鉄住金ソリューションズ (株) (NSSOL <sup>☆ 1</sup>) 5 名と、NSSOL の 100% 子会社で金融分野を中心 にデータ分析やモデル化を主業とする(株)金融工 ンジニアリング・グループ (FEG) 9名の計 14名の 合同で結成された. チームの中心メンバは、NSSOL, FEG ともにデータマイニングの実務経験を持ち、一 部は KDD Cup や Kaggle などで開催されるオープン コンペティションへの参加経験があった.

例年の KDD Cupでは、技能向上を目的として両企 業それぞれで主に個人で参加していた。 今回 14 名の 合同チームを結成した目的は、従来の技能向上に加 え今後データマイニング分野の市場が拡大することを 想定し、課題を個人ではなくチームで解決するための 知見の獲得と、チーム戦に必要なクラウドサービスの 実証検証であった.

# データ分析統合環境

チーム名にある Data Veraci (ダータヴェラーチ) <sup>☆1</sup> は、NSSOL が 2014 年から開発/運用している、デー タ分析を共同で行うためのクラウドサービス型の統合 環境である. 特徴は単なるデータ分析アプリケーショ ンではなく、タスク管理や文書検索、テキストチャット などデータ分析業務に必要な IT 環境全体を用意して

☆1 新日鉄住金ソリューションズの登録商標

いる点である。これによりデータ分析プロジェクトの素 早い立ち上げと、遠隔地からの支援、IT 環境の統一に よる効率化、セキュリティの確保 <sup>☆ 2</sup> などを狙っている.

今回は地理的に離れたオフィス間での共同作業がで きた点や、検討時のソースコードがそのまま別のメン バで実行できることで他メンバの検討状況が容易に把 握できた点で、クラウド型の統合環境が有効であった. また情報共有により重複のない仮説検証や、課題解決 までの時間短縮にも効果があった、結果として個人で は思い浮かばないような仮説を立て、各メンバが自然 に役割を分担し実行できる体制が構築できた.

### 特徴抽出に注力

今年の KDD Cup のデータセットで主要なテーブル はユーザの受講ログと教材マスタのみであり、さらに 教材マスタにはその種類と階層関係しかなく、この限 られた項目からいかに特徴を抽出するかが重要なポ イントとなった  $^{4}$ 3.

人数は多いが作業場所と作業時間はまちまちであ る今回のチームでは、各人が Data Veraci 内で行われ ている仮説検証状況を確認しつつ, 新たな仮説を発 案し検証するという形で、特徴抽出を分担して行った.

最初にチームで共有した特徴量は、受講ログの項目 定義から機械的に集計パターンを網羅した約500次元 の "ベースライン"特徴量である(図-1). 各メンバは この特徴量をもとに、異なる観点からの仮説を立て特 徴量を追加していった. 特に以下の特徴抽出が精度向 上に大きく貢献した.

実務上、委託元からデータ格納場所を社内に限定される場合が多く、機 器を社内で保有管理するクラウドサービスとしている.

 $<sup>^{\</sup>diamond \, 3}$  今回の KDD Cup の課題については,本特集「編集にあたって」を参照.

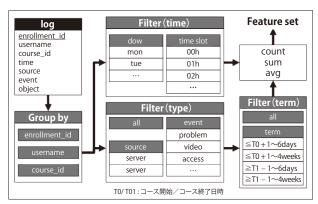


図-1 "ベースライン"特徴量

- 学習進捗:ユーザごとの教材のアクセス順序を集計 し、平均的なアクセス順序やタイミングを定義、進 捗状況や平均からのズレを求める
- 未来ログ<sup>☆4</sup>:今回は、講座に10日間アクセスがない状態を脱落と定義しているが、講座の開催期間に重複があるため、その10日間に別の講座を受けている場合があり、その回数を求める
- 周期性:受講口グの単純集計で失われるアクセスの 周期性を、フーリエ変換により求める
- アクセスパターン:テキストを見る,ビデオを見る, 問題を解くなどの受講ログの出現順序を状態遷移 と見なし、その遷移確率や時間を求める
- **講座別モデル**: 講座単位で各教材のアクセスなどを 特徴量化し講座単位で識別器を作成, その結果を 特徴量として用いる

上記の取り組みの結果, 最終的には約 2,500 次元 からなるデータセットが構築できた. この特徴量に対して xgboost ☆5 (勾配ブースティング木の一実装) 単体で作成したモデルのプライベートスコアは 0.9087 であり, ほぼ特徴抽出のみで入賞できたことになる.

### 自発的な分担

今回は、チームメンバに事前に役割を設定せず、メンバ間の調整で分担を決めていた。このような管理で十分であった理由としては、メンバそれぞれが経験

ある実務者として、重要な課題を発見し解決する能力があったため、Data Veraci 内での議論やドキュメント、ソースコードの確認により、メンバの対応余力や力量を理解し共有できたからである。

今回はメンバの能力で解決したが、より組織的にデータマイニングを行うためには、スキルやプロセスを定義しその中で役割を決めていく必要があると考えている。

#### チーム戦でのモチベーション

多人数チームで参加する場合、個人で参加する場合と 比べて結果に対する貢献度合いが不明確になる。加え て通常業務の合間に実行していることと、明確な役割 がないことが影響して、モチベーションの維持が難しい。

今回はチームマージ期限までは個人で参加し、メンバ間で競争させる対策をとった。それでもチームマージ後のモチベーション維持が課題になるが、今回はチームマージの途中で1位になったため、順位の維持が明確な目的となり、モチベーションが維持しやすかったと考える。

また,他メンバの活動を見て自分もやる気になる点や,グループ企業内とはいえ別組織との合同チームであることによる,取り組みへの責任感もモチベーション維持に寄与したと考える.

# 今後の課題と取り組み

今回は惜しくも2位であったが、参加によって共同作業の進め方など、業務にフィードバックできる多くの知見を得ることができた。今後も企業内のデータ分析実務者として、新手法のキャッチアップ、IT環境の整備、プロセスの洗練などを継続して行う中で、KDD Cup に代表されるオープンコンペティションに積極的にチャレンジしていきたいと考えている。

(2015年10月30日受付)

西川大亮 ■ nishikawa.daisuke.8tx@jp.nssol.nssmc.com

2001年北海道大学大学院工学研究科博士後期課程修了.博士(工学).現在,新日鉄住金ソリューションズ(株)システム研究開発センターデータ分析グループリーダー.

<sup>☆4</sup> 実務上は予測時点では利用できないデータだが、競技であるため使用している。

<sup>&</sup>lt;sup>☆ 5</sup> https://github.com/dmlc/xgboost