

# seis-ml-api 中間レポート

江畑 拓哉 (201611350)、栗本真太郎 (201511366)、畑中 智之 (201611402)

## Contents

### 1 情報特別演習概要

本演習は江畑、栗本、畑中の3名により実施する。それぞれ演習を進め、最終的に大規模データベースを用いた機械学習 API (seis-ml-api) を作成することが目標である。

#### 1.1 演習範囲に関して

以下に示す範囲を演習し、その成果を組み合わせることで seis-ml-api の作成を目指す。

Table 1: 演習範囲に関して

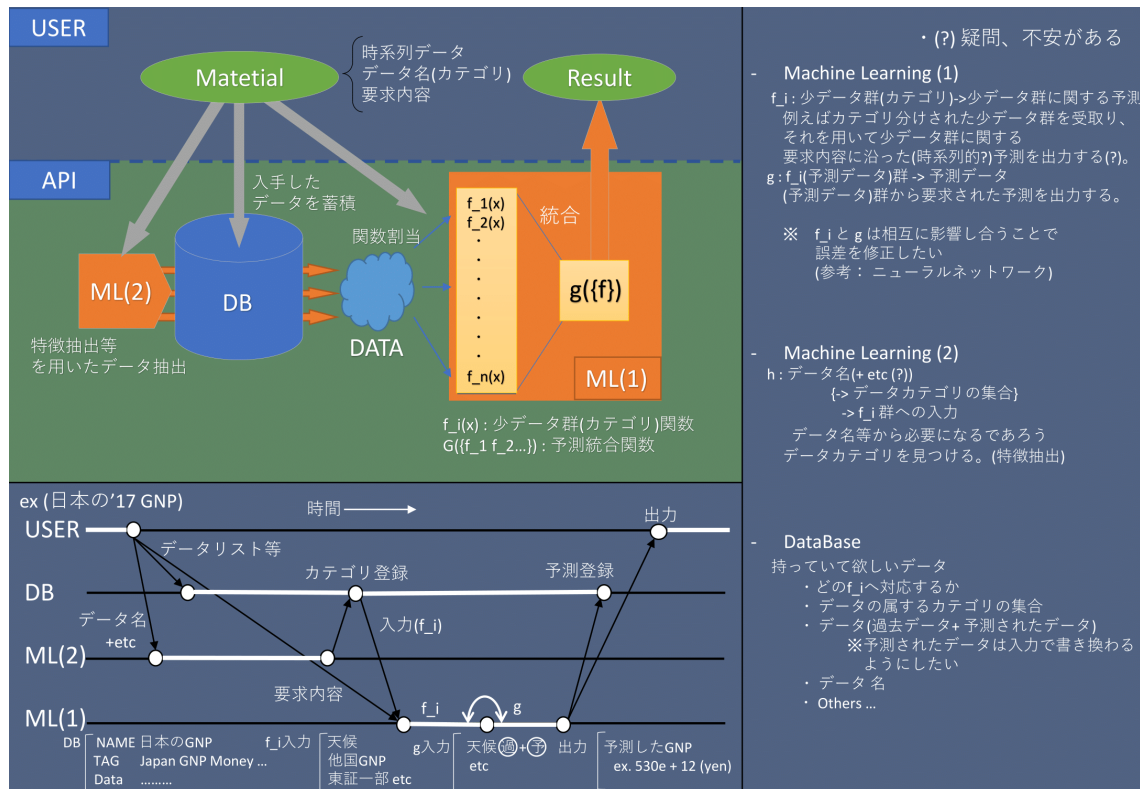
氏名	分野	内容
江畑	データベースと機械学習の統合	大規模データを利用した機械学習の作成
栗本	機械学習	機械学習モデルの調整
畑中	データベース	大規模データベースの作成

#### 1.2 seis-ml-api 概要

seis-ml-api は、機械学習を提供する Web API である。大規模データベースを採用する。

seis-ml-api のフローを Figure.1 に示す。

ユーザは API に自分の持っているデータを登録する。API は、API に登録されているデータを用いて機械学習処理を行い、その結果をユーザに返す。



・ (?) 疑問、不安がある

#### - Machine Learning (1)

$f_i$ : 少データ群(カテゴリ)→少データ群に関する予測  
例えばカテゴリ分けされた少データ群を受取り、  
それを用いて少データ群に関する  
要求内容に沿った(時系列的?)予測を出力する(?)。  
 $g: f_i(\text{予測データ}) \rightarrow \text{予測データ}$   
(予測データ)群から要求された予測を出力する。

※  $f_i$  と  $g$  は相互に影響し合うことで  
誤差を修正したい  
(参考: ニューラルネットワーク)

#### - Machine Learning (2)

$h$ : データ名(+ etc ?)  
{→ データカテゴリの集合}  
→  $f_i$  群への入力  
データ名等から必要になるであろう  
データカテゴリを見つける。(特徴抽出)

#### - DataBase

持っていて欲しいデータ

- ・ どの  $f_i$  へ対応するか
- ・ データの属するカテゴリの集合
- ・ データ(過去データ + 予測されたデータ)

※ 予測されたデータは入力で書き換わるようにしたい

- ・ データ名
- ・ Others ...

Figure 1: Figure.1

### 1.3 実験に用いるデータ

以下のデータをデータベースに登録し、実験段階においても用いたいと考えている。データの相関を求める関係上、これらの入手元から日本国内の経済についてのデータを集めたいと考えている。

- Quandl  
ほとんどすべてのデータはここで手に入る。ただし、長期間のデータは乏しいようである。主には、ここから得た株価データを用いて分析を行う予定である。
- google finance  
日本のデータを csv 形式で入手することは困難だが、海外のデータは容易に手に入る。
- 総務省統計データ  
めばしいデータは少ないが、ゼロではないため活用していきたい。

## 2 それぞれの進捗について

- 江畑  
機械学習アルゴリズムの策定を行った。また Java、Python、Clojure での HBase の利用方法について学習した。
- 栗本  
「機械学習」の履修及び関連書籍の学習を行った。また、主専攻実験「ヒューマンセンシング」の自主実験において、サポートベクターマシンを用いた簡易画像分類器を作成した。
- 畑中  
HBase、Hadoop を用いて疑似分散環境を構築した。また、Java から HBase にアクセスする方法を、実際にプログラムを動作させて確認した。

## 3 今後の演習について

- 江畑

策定した機械学習のモデルを実際のコードに実現する作業と HBase に入力されたデータを送る API を作成する。

- 栗本

機械学習モデルの理解を深め、Python によるデータ分析に慣れ、より良いモデル調整が可能のように学習していきたいと考えている。

- 畑中

完全分散環境を構築して、HBase の性能テストを行いたいと考えている。