

計算機科学実験及演習 4 エージェント 課題 4

1029-28-2473 二見 颯

2018 年 11 月 16 日

1 課題 4-1

1.1 プログラム概要

リスティングデータによって与えられる類似の宿泊施設の提示価格を予測して、それによって自己の宿泊施設の価格推薦を行うことを考える。民泊市場にはリスティングデータに従う貸し手 1 人と価格推薦エージェント 1 体という条件で、価格推薦エージェントはリスティングデータの価格を予測して、それより少し価格を提示するとする。

1.2 外部仕様

リスティングデータとして、San Francisco-listings.csv を用いることとする。

課題 3 と同様に、96 の特徴量のうち 'latitude', 'longitude', 'accommodates', 'property_type', 'room_type', 'number_of_reviews', 'review_scores_rating' を説明変数として、'price' を目的変数とする。ただし、'price' が 500 以下のデータのみを利用する。

bayes_opt.py

pip install scikit-optimize が必要。

Grid Search の代わりにベイズ最適化の手法を用いて、SVR の最適なハイパーパラメータを探索する。

agent.py

予測価格の何 % とするかをコマンドライン引数として取り、SVR による価格予測を用いて価格推薦を行う 1 体のエージェントをシミュレーションする。プログラムの実行は

./agent.py [percentage]

とする。以下は agent.py の実行例である。

```
1 $ ./agent.py 90
2 load sanfrancisco: read 1000 / 6196 data
3 start SVR fit...
4 cvxopt.solvers.qp: optimization succeeded
5 start SVR predict...
6 revenue by SVR agent: 24257.254376011988
7 the agent wins 298 times
8 mean of y_train: 172.424
9 revenue of mean price offering: 17343.636000000005
```

1.3 内部仕様

bayes_opt.py

load_sanfrancisco 関数によって San Francisco-listing.csv を読み込み、データの前処理ののちに X_pro, y_pro を得る。

spaces には SVR の各パラメータ (p, C, eps) の定義域を与えて、skopt.gp-minimize により bayes_opt 関数を最小化するようなパラメータを探索する。

bayes_opt では与えられたパラメータの SVR を構成して、交差検証によって得た決定係数 (5 回の平均) に-1 を掛けたものを返す。(決定係数は 1 に近づくように最大化したいため)

agent.py

load_sanfrancisco 関数によって得た X, y について、前半を訓練データ、後半を評価データとする。

bayes_opt.py によって発見した最適なパラメータで SVR を構成して、訓練データによって SVR を学習させ、それによって評価データを予測する。

コマンドライン引数を m として、SVR による予測価格の $m\%$ が提示価格となる。

リストの価格の 30% は原価として考えて、各提示価格について提示価格がリストの価格を下回ったときに (提示価格 - リストの価格 * 0.3) が利益となる。これらの総利益をエージェントの評価基準とする。利益はマイナス (すなわち赤字) となることも考えられる。

1.4 評価結果

選択した説明変数のうち数値データ ('latitude', 'longitude', 'accommodates', 'number_of_reviews', 'review_scores_rating') と目的変数 'price' との相関を図 1 に示す。

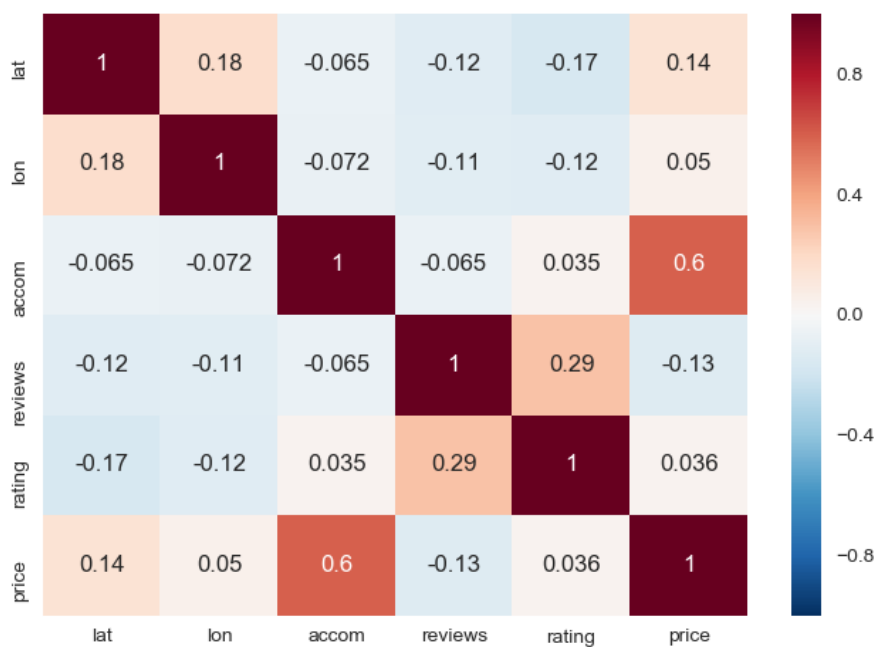


図 1 Correlation coefficient

SVR の予測精度 (決定係数) についてベイズ最適化により最適なハイパーパラメータ C, σ, ϵ を探索した。探索範囲は $2^{-1} \leq \sigma \leq 2^5$, $2^{-5} \leq C \leq 2^{15}$, $2^{-10} \leq \epsilon \leq 2^{-3}$ である。図 2 にベイズ最適化で探索したパラメータ C, σ とそのときの決定係数の値を示す。(データ数は 500 で 100 回探索した)

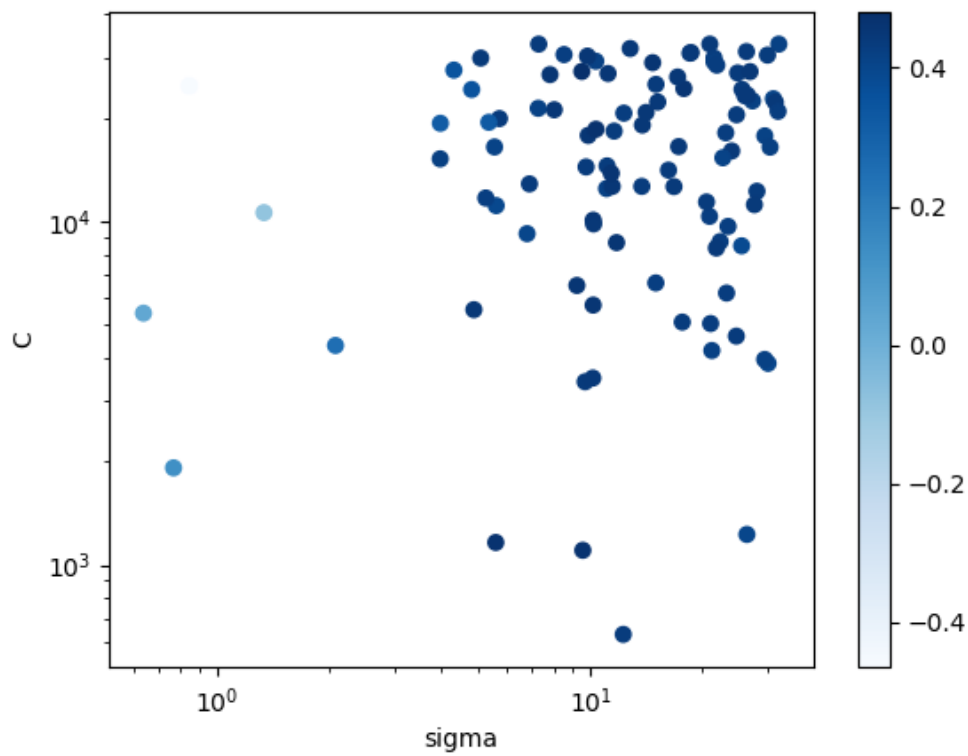


図2 Bayes Optimization

決定係数は $\sigma = 9.50, C = 2.73 \times 10^4, \epsilon = 9.27 \times 10^{-2}$ のときに 0.481 と最大になった。

データ数 2000 を SanFrancisco リスティングデータからサンプリングして、その前半 1000 個を訓練データ、後半 1000 個を評価データとした。提示価格を予測価格の何 % とするかを変化させた時の総利益の変化を表 1 に示す。

(%)	総利益	成功回数	効率 (総利益/成功回数)
70	44470	837	53.1
75	47421	793	59.7
80	48558	739	65.7
85	49322	689	71.5
90	48622	629	77.3
95	47511	568	83.6
100	46029	510	90.2
105	42665	449	95.0
110	39205	394	99.5
115	37179	351	105.9
120	35160	313	112.3
mean	33871	359	94.3

表 1 agent by percentage

2 課題 4-2

2.1 プログラム概要

リスティングデータに従う貸し手 1 人と価格推薦エージェント 1, エージェント 2 が存在するとする。エージェントが最安の値段をつけたときに借り手がついて、収入が得られるとする。ペイズ最適化で最も精度が良かったパラメータの SVR を利用したものをエージェント 1、あまり精度が良くない SVR を利用したものをエージェント 2 とする。

2.2 外部仕様

`./multiagent.py` で実行する。以下は `multiagent.py` の実行例である。エージェント 1, 2 の予測精度 (決定係数), 総利益, 成功回数, 効率 (総利益/成功回数) を出力する。

```

1 $ ./multi_agent.py
2 load sanfrancisco: read 2000 / 6196 data
3 === SVRAgent start ===
4 1 th cross validation...
5 cvxopt.solvers.qp: optimization succeeded
6 R2 score: 0.4797468020396928
7 2 th cross validation...
8 cvxopt.solvers.qp: optimization succeeded
9 R2 score: 0.33178028654882397
10 3 th cross validation...
11 cvxopt.solvers.qp: optimization succeeded
12 R2 score: 0.4714137117838899
13 4 th cross validation...
```

```
14 cvxopt.solvers.qp: optimization succeeded
15 R2 score: 0.49906353860044683
16 5 th cross validation...
17 cvxopt.solvers.qp: optimization succeeded
18 R2 score: 0.31606765685527194
19 5-fold average score: 0.41961439916562504
20 cross val score: 0.41961439916562504
21 start SVR fit...
22 cvxopt.solvers.qp: optimization succeeded
23 start SVR predict...
24 === SVRAgent start ===
25 1 th cross validation...
26 cvxopt.solvers.qp: optimization succeeded
27 R2 score: -0.05695700854191843
28 2 th cross validation...
29 cvxopt.solvers.qp: optimization succeeded
30 R2 score: 0.03379843202281596
31 3 th cross validation...
32 cvxopt.solvers.qp: optimization succeeded
33 R2 score: -0.00822695472188295
34 4 th cross validation...
35 cvxopt.solvers.qp: optimization succeeded
36 R2 score: 0.06128181455545334
37 5 th cross validation...
38 cvxopt.solvers.qp: optimization succeeded
39 R2 score: 0.017475693763498
40 5-fold average score: 0.009474395415593185
41 cross val score: 0.009474395415593185
42 start SVR fit...
43 cvxopt.solvers.qp: optimization succeeded
44 start SVR predict...
45 revenue by SVR agent1: 16242.43831357982
46 the agent1 wins 340/1000 times
47 revenue1 per 1: 47.77187739288182
48 revenue by SVR agent2: 23875.955321474925
49 the agent2 wins 436/1000 times
50 revenue2 per 1: 54.761365416226894
```

2.3 内部仕様

multi_agent.py

load_sanfrancisco 関数によって San Francisco-listing.csv を読み込み、データの前処理ののちに X, y を得る。

load_sanfrancisco 関数によって得た X, y について、前半を訓練データ、後半を評価データとする。

SVRAgent 関数はパラメータを引数にとって、交差検証で精度を示したのち訓練データで学習、評価データの予測を行う。提示価格は予測価格の 85% とする。

エージェント 1, エージェント 2, リスティングデータの価格のうち、最小のものに借り手がつくとする。エージェントに借り手がついた場合、(提示価格 - リストの価格 * 0.3) がエージェントの利益となる。

2.4 評価結果

精度の良い Agent1($\sigma = 9.50, C = 2.70 * 10^4, \epsilon = 9.00 * 10^{-2}$) と精度のあまり良くない Agent2($\sigma = 1.0, C = 1.0, \epsilon = 0.1$) の総利益, 成功回数, 効率は表 2 のようになった。サンプリングされたデータを変更して 5 回実行した。

(Agent1) 決定係数	総利益	成功回数	(Agent2) 決定係数	総利益	成功回数
0.419	16242	340	0.00947	23875	436
0.473	19088	386	0.0000350	23035	396
0.528	17584	384	0.0115	24607	422
0.486	17119	338	0.0175	24666	439
0.442	16432	330	-0.0013	24807	450

表 2 Agent1 vs Agent2

3 考察 (課題 4-1 + 課題 4-2)

課題 4-1 では、表 1 より提示価格を予測価格の 85% とした場合が総利益が最大になった。提示価格の予測価格に対する割合を下げることで、予測価格がリスティングデータの価格より高い場合にも提示価格が予測価格を下回る場合が多くなり、成功回数は上がるが、一回の提供あたりの利益が小さくなる。そのトレードオフを考慮して 85% あたりが望ましい。

課題 4-2 で 2 つのエージェントが存在する場合、どちらのエージェントが他より利益を上げることができるかは単純に決定することはできない。(予測精度が高い方が総利益が常に高くなるとはいえない) 価格の予測精度が低くても他のエージェントとリスティングデータの価格より低い価格を提示することができれば、利益を得ることができて、他のエージェントはそのとき利益を逃すからである。原価 (30%) を下回らない範囲で最低の価格を提示するようなエージェントがあれば、他のエージェントはどんなに正確に価格予測してリスティングデータに対し大きな利益が得られるような提示を行っても、それを得る機会を失い、最終的に原価に近い価格を提示し続けたエージェントが最大の総利益を得る。

しかし、このように価格競争が発生して提示価格は限りなく原価に近づいた場合、一方のエージェントが価格競争に勝利して全ての借り手に民泊を提供できたとしてもほとんど利益を上げられなくなる。それぞれのエージェントが自身の利益を追求すると、エージェント全体の利益は最小化されてしまう。(囚人のジレンマといえる)

一方、お互いのエージェントが交代で市場に参加すれば、それぞれのエージェントは借り手を独占することはできないが、結果的に総利益を最大にできることがある。

3.1 ベイズ最適化

ハイパーパラメータの探索において、一般的な手法としてグリッドサーチがある。グリッドサーチは離散的なパラメータの組の候補について全探索を行う。一方、ベイズ最適化では、ある連続的な定義域をもつ入力(パラメータ) x に対して、ブラックボックス関数 $y = f(x)$ を最小化 (あるいは最大化) するような x を逐次的に求める手法である。次の観測点 x の候補の求め方は全探索ではなく、これまでの関数の評価値に基づいて判断する。

次の観測点を決定する戦略として、ここでは EI 戦略を用いた。EI(期待改善量) 戦略では、評価値 y とこれまでで最高の評価値 y_{best} として、 $y - y_{best}$ の期待値 (獲得関数) が最も高くなる点を次に観測する。ブラックボックス関数はガウス過程に従うと仮定して、これまでの観測点から信頼区間を形成して、関数の形をある程度予想することを考える。(図 3, 参考文献 5 より)

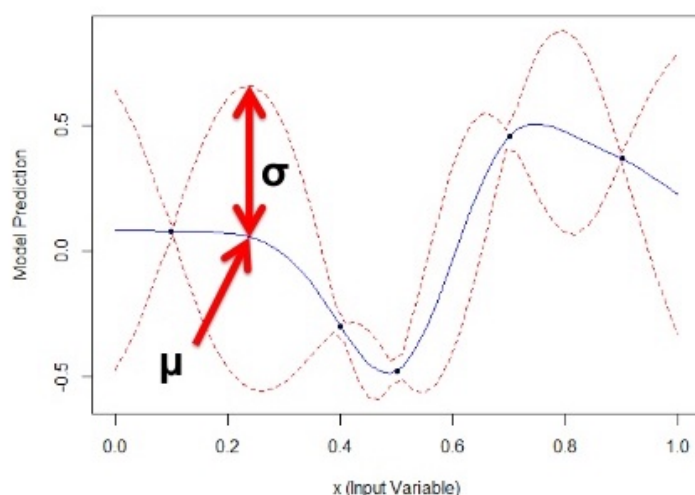


図 3 Gaussian process

未観測点の期待値 μ と分散 σ^2 を計算すると、 μ によって、未観測点の評価値の大きさ、 σ によって周囲の x (パラメータ) がどの程度観測されているかを知ることができる。 μ が最大の点を次の観測点に選ぶと、評価値が大きくなることが予想できるが、局所解に陥る可能性があるので適宜 σ が大きい点も観測する。

3.2 回帰モデルの評価

決定係数は

$$R^2 = 1 - \frac{\sum_{i=0} (y_i - \hat{y}_i)^2}{\sum_{i=0} (y_i - \bar{y})^2} \quad (1)$$

で表すことができる。 R^2 の分子は回帰モデルでは説明できない目的変数のばらつき (変動)、分母は目的変数の平均値のまわりにおけるばらつき (変動) を表していて、 R^2 は目的変数自体の変動のうち、回帰モデルで説明できるものの割合を示しているといえる (決定係数が 1 に近づくほど良いモデルであるといえる) 決定係数 (R^2) は分子に 2 乗誤差 (MSE) $\sum_{i=0} (y_i - \hat{y}_i)^2$ が存在していて、MSE に対して、 R^2 は単調減少する。

平均二乗誤差 (MSE) は

$$\frac{1}{n} \sum_{i=0} (y_i - \hat{y}_i)^2$$

平均絶対誤差 (MAE) は

$$\frac{1}{n} \sum_{i=0} \|y_i - \hat{y}_i\|$$

と表される。

MSE の最小化は、誤差に正規分布を仮定した場合の最尤推定といえる。
正規分布は、説明変数を x_i , 目的変数を y_i として

$$p(y_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - f(x_i))^2}{2\sigma^2}\right)$$

その対数尤度関数は

$$\begin{aligned} \log \prod_i p(y_i) &= \sum_i \log p(y_i) \\ &= C - \frac{1}{2\sigma^2} \sum_i (y_i - f(x_i))^2 \end{aligned}$$

となり、その最大化は $\sum_i (y_i - f(x_i))$ の最小化となる。

MAE の最小化は、誤差にラプラス分布を仮定した場合の最尤推定といえる。
ラプラス分布は

$$p(y_i) = \frac{1}{2\phi} \exp\left(-\frac{\|y_i - f(x_i)\|}{\phi}\right)$$

その対数尤度関数は

$$\log \prod_i p(y_i) = C - \frac{1}{\phi} \sum_i \|y_i - f(x_i)\|$$

となり、その最大化は $\sum_i \|y_i - f(x_i)\|$ の最小化となる。

一般的に、正規分布よりもラプラス分布の方が分布の裾野が広いことが知られているので、多くの外れ値が存在するデータの誤差を評価するときには、MSE より MAE のほうが指標として適切である。

誤差の分布からモデルの良し悪しを評価することを考える。

誤差の絶対値を $e_i = \|y_i - f(x_i)\|$ として、RMSE(Root Mean Squared Error) を \sqrt{MSE} と定義すると、

$$RMSE^2 = \frac{\sum_i e_i^2}{n}$$
$$MAE^2 = \frac{(\sum_i e_i)^2}{n^2}$$

ここで

$$RMSE^2 - MAE^2 = Var(e_i)$$

が成立する。さらに誤差の絶対値の平均 (\bar{e}_i) は MAE と一致するので、

$$\frac{RMSE}{MAE} = \sqrt{1 + \frac{Var(e_i)}{\bar{e}_i^2}}$$

誤差が平均 0、標準偏差 σ の正規分布に従う場合、誤差の絶対値の確率密度関数は $0 \leq e$ で

$$p(e) = \frac{2}{\sqrt{2\pi}\sigma} \exp\left(-\frac{e^2}{2\sigma^2}\right)$$

\bar{e} および $Var(e)$ は、

$$\bar{e} = \int_0^\infty ep(e)de$$
$$Var(e) = \int_0^\infty (e - \bar{e})^2 p(e)de$$

これらを用いて、

$$\frac{RMSE}{MAE} = \sqrt{\frac{\pi}{2}} = 1.253...$$

誤差が平均 0、分散 $2\phi^2$ のラプラス分布に従う場合、

$$\frac{RMSE}{MAE} = \sqrt{1 + \frac{\phi^2}{(\phi)^2}} = \sqrt{2} = 1.414...$$

SanFrancisco リスティングデータから 500 ドルより大きいデータを取り除いて 1000 個を利用するデータとして、最適化した SVR($\sigma = 9.50, C = 2.70 * 10^4, \epsilon = 9.00 * 10^{-2}$) を交差検証して評価した。このとき、 $RMSE = 70.67$, $MAE = 49.71$ となり、 $RMSE / MAE$ は 1.4216 となった。モデルの予測とデータとの誤差はラプラス分布に近いことがわかった。(外れ値が比較的多い)

4 参考文献

1. scikit-learn と Tensorflow による実践機械学習 Aurelien Geron 著
2. Python 機械学習プログラミング Sebastian Raschka 著
3. ベイズ最適化 (<https://www.kumilog.net/entry/bayesian-optimization>)
4. 精度評価指標と回帰モデルの評価 (<https://funatsu-lab.github.io/open-course-ware/basic-theory/accuracy-index/>)
5. 機械学習のためのベイズ最適化入門 (https://www.slideshare.net/hoxo_m/ss-77421091)