

実データで学ぶ人工知能講座

演習課題：hoeffding_coins*

Matthew J. Holland[†]
大阪大学 産業科学研究所

実験条件の説明

この簡易な実験では、スライド中に見た「壺」が複数ある状況を考える。基本的な要素は以下の通りである。

- コインが 1000 枚ある
- 各コインを 10 回ずつ独立に投げて、「表」か「裏」か記録しておく

コイン自体はスライド中の「壺」に相当し、「コインを投げること」は、講義中の「壺から玉を取ること」に相当する。コインを投げた結果、「表」か「裏」かの 2 通りしかない。

- 「表」ならば $c = 1$
- 「裏」ならば $c = 0$

上記の具合に一つのコインの表裏の成否を表記する。また、コインは 1000 枚もあるので、コインごとの表裏については、 $c_1, c_2, \dots, c_{1000}$ という形式で表わすことにする。実験条件としては、各コインを何度も投げて、標本データを集めることになっている。標本数を n と表わすと、

- 1 番目のコイン： $c_1(1), \dots, c_1(n)$
- 2 番目のコイン： $c_2(1), \dots, c_2(n)$
- \vdots
- 1000 番目のコイン： $c_{1000}(1), \dots, c_{1000}(n)$

10 回投げるのであれば、 $n = 10$ である。最後に、コインごとに、投げた回数に占める「表」の割合を以下のように書くことにする。

$$\hat{c}_j := \frac{1}{n} \sum_{i=1}^n c_j(i) \quad (1)$$

上記のように、すべてのコインをひと通り 10 回ずつ投げるまでの一連の作業を、この実験の「1 試行」とする。次の節では演習課題の内容を説明する。

*この実験条件は Abu-Mostafa et al. [1, 2] の事例に基づくものである。

[†]作者の連絡先：matthew-h@ar.sanken.osaka-u.ac.jp.

演習課題の内容

問題 1. 乱数を発生させて、先述の実験を 10 万試行分、実施すること (numpy の random.binomial を使うと便利).

問題 2. 試行ごとに、次の 3 種類のコインの表の割合 (\hat{c}) を記録しておくこと.

- c_1 : 1 番目のコイン
- c_{rand} : 無作為に選んだコイン
- c_{min} : 全コインのなかで、表の数がもっとも少なかったコイン

問題 3. 上記の結果として、 \hat{c}_1 と \hat{c}_{rand} と \hat{c}_{min} それぞれ、10 万点からなるデータセットが得られる. その分布をヒストグラムで可視化すること (matplotlib.pyplot の hist を使うと便利).

問題 4. この 3 種類のコインそれぞれに対して:

- A. $\mathbf{P}\{|\hat{c} - \mathbf{E}c| > \varepsilon\}$ を近似し、 ε の関数としてそのグラフを描くこと.
- B. 同じプロットにおいて、Hoeffding の上界である $2e^{-2\varepsilon^2 n}$ のグラフも併せて表示すること.
- C. Hoeffding の不等式に従うコインと従わないコインはどれか. 従わないコインはなぜ従わないか自分の言葉で説明すること.

問題 5. 標本数 n (コインを投げる回数) を大幅に増やしてみること ($n = 500$ など). \hat{c}_{min} の分布が $n = 10$ のときと比べて、どのように変わるか. それはなぜか.

おまけ: 学習との関係

たとえばモデル $\mathcal{H} = \{h_1, \dots, h_k\}$ を検討しているとする. 各々の候補に対して、その標本内誤差を求めることができる. すると

$$\hat{R}(h_1), \hat{R}(h_2), \dots, \hat{R}(h_k)$$

を観測することができるので、当然その値がもっとも小さくなるような $h \in \mathcal{H}$ を突き止めることもできる. この候補を \hat{h}_{min} と書く.

対応関係は以下ようになる:

- コインの 1 試行は学習用のデータ $(x_1, y_1), \dots, (x_n, y_n)$ 一回分を得ることに相当する.
- \hat{c}_j は $\hat{R}(h_j)$ に相当する.
- \hat{c}_{min} は $\hat{R}(\hat{h}_{\text{min}})$ に相当する.
- \hat{c}_{min} が本当の表が出る確率よりも小さく偏るように、 $\hat{R}(\hat{h}_{\text{min}})$ が標本外誤差よりも小さく偏る.

したがって、 $|\hat{c}_{\text{min}} - \mathbf{E}c|$ が Hoeffding の不等式に従わず、誤差がより大きいのに同様に、汎化誤差 $|\hat{R}(\hat{h}_{\text{min}}) - R(\hat{h}_{\text{min}})|$ も Hoeffding の不等式ほどタイトな上界には従わない.

参考文献

- [1] Abu-Mostafa, Y., Song, X., Nicholson, A., and Magdon-Ismail, M. (2004). The bin model. Technical report, California Institute of Technology.
- [2] Abu-Mostafa, Y. S., Magdon-Ismail, M., and Lin, H.-T. (2012). *Learning From Data*.