

15. データサイエンス

プログラミング・データサイエンス I

2022/7/28

1 今日の目的

今日の目的

- データからの推計
- 機械学習

この講義は、「プログラミング・データサイエンス」というインターフェースプログラムの端緒となる科目です。この講義の最後として、データサイエンスの香りをちょっとだけ吸ってみましょう。

データサイエンスとは、データの分析に基づき、様々な課題解決を行うことを目指すものです。その基本には、統計学とプログラミング等の計算機科学があります。データサイエンスを活用する入り口として、具体的な例を見ることにします。今日のサンプルプログラムは、用意していますが、新しいパッケージが必要であり、動作しない人がいる懸念があります。また、大量の計算が必要であり、コンピュータによっては、負荷が大きくなりすぎます。話を聞くだけで結構です。

2 推計

初めに、データからその統計的特徴量を推計することを考えましょう。データは、母集団から得るサンプルに過ぎないことに注意します。ここでは、簡単なシミュレーションを通じた例を示します。この例は以下にあります。

<https://github.com/first-programming-saga/dataScience>

平均 μ 、標準偏差 σ の母集団を考えます。ここから n 個のサンプル $\{x_i\}$ を無作為抽出

しましょう。標本平均と標本標準偏差は

$$\bar{x} = \frac{1}{n} \sum_{i=0}^{n-1} x_i \quad (2.1)$$

$$s^2 = \frac{1}{n-1} (x_i - \bar{x})^2 \quad (2.2)$$

となります。

中心極限定理によると、標本平均は、平均 μ 、標本標準偏差 σ/\sqrt{n} の正規分布に従います。このことは、母集団の分布に依存しません。このような理由により、正規分布が偏在しているのです。

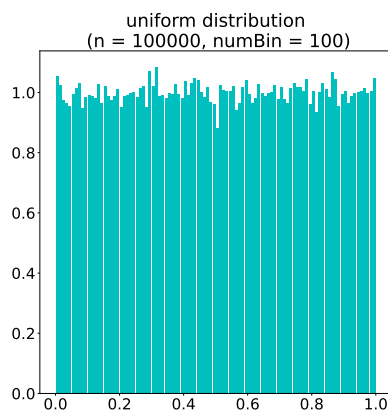


図1 [0, 1) の一様分布。データ数 10,000。bin の数 100。

実際に、[0, 1) の一様分布に対して、大きさサイズ 100 の標本を 10,000 個作って、標本平均の分布を求めてみましょう。一様分布とは、ある区間に標本が一様に分布するものです。[0, 1) の一様分布を図 1 に示します。

ソースコード 2.1 を見てください。こちらは、[0, 1) の一様分布に従う乱数から、標本平均を求めるものです。

母集団母集団の平均は $\mu = 1/2$ 、分散は $\sigma^2 = 1/12$ です。図 2 に示すように、平均 μ 、標準偏差 σ/\sqrt{n} の正規分布となっています。

本来、母集団の平均と標準偏差を直接に知ることはできません。しかし、前述のように標本平均が正規分布に従うことを使うと、母集団の平均を推計することができます。標本平均が、平均 μ 、標準偏差 s/\sqrt{n} の正規分布に従うことを出発点にしましょう。正規分布は、平均を中心としてある幅に入る確率を数値的に知ることができます。

ソースコード 2.1 標本平均を求める

```
1 def evalSampleMean(n: int) -> float:
2     """
3     大きさ  $n$  の標本に対して、平均を返す
4     """
5     x = 0
6     for i in range(n):
7         x += np.random.rand()#[0, 1) の乱数
8     return x / n
9
10 #標本平均を作る
11 sample_size = 100 #標本サイズ
12 num_samples = 10000 #標本数
13 sample_means = np.zeros(num_samples)
14 for i in range(num_samples):
15     sample_means[i] = evalSampleMean(sample_size)
```

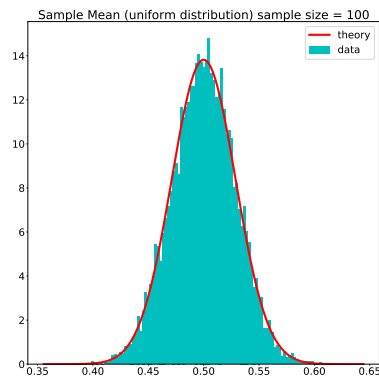


図 2 $[0, 1)$ の一様分布に対する標本平均の分布。標本サイズは 100

標本平均 \bar{x} と母集団平均 μ との差が、標本標準偏差のある範囲に入っているとしましょう。 $f > 0$ はその範囲を定めるパラメタです。

$$-f \frac{s}{\sqrt{n}} < \bar{x} - \mu < f \frac{s}{\sqrt{n}} \quad (2.3)$$

95% の確率 (これを信頼係数という) でこの範囲に入るとすると、 $f \simeq 1.96$ が該当します。これを変形することで、標本平均から信頼係数に依存して、母集団平均が入っている

区間を推定することができます。

$$\bar{x} - \frac{fs}{\sqrt{n}} < \mu < \bar{x} + \frac{fs}{\sqrt{n}} \quad (2.4)$$

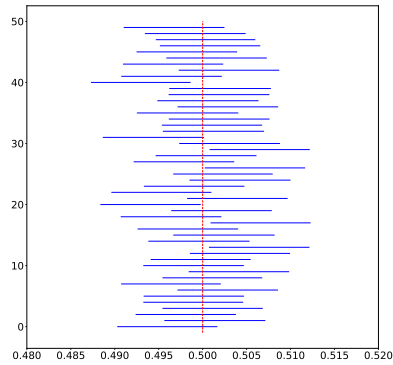


図3 [0, 1) の一様分布に対する標本平均から、母集団の平均の範囲を 95% で推計

実際に、[0, 1) の一様分布から 1000 個のサンプルを取り出し、95% の信頼度で、母集団平均を推計するイメージを示します。プログラムをソースコード 2.2 に示します。図 3 の横棒は、各サンプルで推定した母集団平均の範囲を表しています。縦の破線は、正しい母集団平均です。

ソースコード 2.2 標本平均から母集団平均を推定

```
1 def estimateRange(mean: float, s: float, n: int, f: float) ->
  ↳ dict[str, float]:
2     """
3     区間を推定
4     """
5     minX = mean - f * s / np.sqrt(n)
6     maxX = mean + f * s / np.sqrt(n)
7     return {'lower' : minX, 'higher' : maxX}
8
9 def evalSample(n : int) -> tuple[float, float]:
10    """
11    大きさ  $n$  の標本に対して、平均と分散を返す
12    """
13    d = np.zeros(n)
14    for i in range(n):
15        d[i] = np.random.rand()#[0,1) の乱数
16    x = 0
17    for i in range(n):
18        x += d[i]
19    mean = x / n
20    s = 0
21    for i in range(n):
22        s += (d[i] - mean) * (d[i] - mean)
23    s = s / (n - 1)
24    return mean, np.sqrt(s)
25
26 alpha = 0.05 #信頼度
27 f = norm.isf(alpha / 2) #信頼度に対応した標準正規分布の変位
28 sampleSize = 10000 #標本サイズ
29 num_trial = 50 #標本数
30 result = list()
31 for i in range(num_trial):
32     m, s = evalSample(sampleSize)
33     r = estimateRange(m, s, sampleSize, f)
34     result.append(r)
```

3 機械学習

次は機械学習を考えましょう。コンピュータに、大量の例を見せて、学習をさせる方法です。

最初に、ヒトが学習する過程を想像しましょう。小さなこどもが文字を覚える過程には、まだまだ分からないことがたくさんありますが、大まかには以下のような過程と考えることができるでしょう。子どもは、周囲にある書かれた文字を見て、あるいは大人が文字を書く様子を見て、その形や筆順をマネして書くことを繰り返します。ひらがなであれば、形と音を組み合わせて覚えていきます。時には、周囲の大人が、形や音を修正をします。このように、文字の形とそのラベル(ひらがなの場合には音)を、あるいはどれとどれが同じか異なるかを試行錯誤しながら覚えます。

成長してから漢字、特に画数の多い漢字を覚えることを想像しましょう。小さなこどもが、ひらがなを覚えるのとは少し違う方法をとっていると思われる。自分が知っている漢字の部品と同じものは覚えやすいのですが、知らない部品で出来ていると覚えるのが大変ですね。このときの学習過程は、すでに、文字という記号、つまり抽象化された概念を道具として使っているのです。

人工知能の研究が始まったころには、後者のように、既知の知識を集め、それに基づいて、コンピュータが何かを判断することが目標とされていました。狭くて明確な知識領域への応用では一定の成果を挙げましたが、知識領域が広い場合や、曖昧さを含む情報に基づくものへの応用は困難でした。

現在の人工知能の中核的手法である機械学習が目指すのは前者の方法です。データそのものから、アルゴリズムを使って特徴を学習し、それに基づき判断をすることを目指しています。そのため、アルゴリズムが学習した特徴を、ヒトが分かるように表現することは、困難になります。

機械学習は、データ分析の一つの手法であり、データサイエンスの重要な手法の一つです。今回は、手書きの数字を判別する機械学習の例を見ていきましょう。

本日のプログラムは以下にあります。

<https://github.com/first-programming-saga/MNIST>

4 手書き数字

アメリカ国立標準技術研究所 (National Institute of Standards and Technology, NIST) は、計量や標準化を担う連邦機関です。NIST が保有していたデータから、機械学習向けに変更したデータが MNIST データです。この中の手書き数字を例にしましょう。

Python には、データ分析のパッケージ `scikit-learn` があります。インストールには `scikit-learn` という名前を使うのですが、python コードに `import` する際には `sklearn` という名前を使います。

ソースコード 4.1 MNIST データ読み込み

```
1 from sklearn.datasets import fetch_openml
2
3 #データの取得
4 mnist = fetch_openml('mnist_784', version=1)
```

`scikit-learn` には、MNIST のデータが含まれています。ソースコード 4.1 が、データの読み込み部分です。この中には、 $28 \times 28 = 784$ 個のビットからなる、手書き文字画像が 70,000 個入っています。

5 0 4 1 9 2 1 3 1 4
3 5 3 6 1 7 2 8 6 9
4 0 9 1 1 2 4 3 2 7
3 8 6 9 0 5 6 0 7 6
1 8 7 9 3 9 8 5 9 3
3 0 7 4 9 8 0 9 4 1
4 4 6 0 4 5 6 7 0 0
1 7 1 6 3 0 2 1 1 7
8 0 2 6 7 8 3 9 0 4
6 7 4 6 8 0 7 8 3 1

図 4 イメージサンプル

これから行うのは、教師付き学習 (supervised learning) という機械学習です。つまり、各イメージを表す数字、つまり正答が対応ついています。この正答をラベルと言います。

サンプルを 100 個、図 4 に示します。ヒトが見ても、判断に迷うようなものも含まれています。

5 二値分類

機械学習では、データを訓練データと試験データに分割します。訓練データを使ってシステムを学習させます。その後に、システムが試験データに対して正答を答えることができれば成功です。70,000 個のうち 60,000 を訓練データとして、残りを試験データとすることにします (ソースコード 5.1)。

ソースコード 5.1 データを訓練データと試験データに分割

```
1 x_train, x_test =
  ↳ mnist['data'].iloc[:60000],mnist['data'].iloc[60000:]
2 y_train, y_test =
  ↳ mnist['target'].iloc[:60000],mnist['target'].iloc[60000:]
3 y_train = y_train.astype(np.int8)
4 y_test = y_test.astype(np.int8)
```

初めに、手書き文字を”5”とそれ以外に分ける学習を考えましょう。データを二つに分類する課題であるため、二値分類 (binary classification) と言います。

元の画像データは、 $28 \times 8 = 784$ 個の 0 と 1 の列です。これを 784 次元空間の点と考えることにします。”5”を表現する点のグループと、それ以外を区別する超平面、今の場合は $784 - 1$ 次元の平面を求めるのが、簡単な手法です。訓練データから、試行錯誤をしてこの面を求めます。ここでは、SGDClassifier というモジュールを使います (ソースコード 5.2)。

学習の後に、試験用の 10,000 個のデータを”5”とそれ以外に分類します。”5”をそれ以外と判断したものが 158 個、”5”でないものを”5”と判断したものが 182 個ありました。正答率にすると 96 パーセントです。簡単な手法なのに、そこそこ良い結果です。間違え

ソースコード 5.2 5 とそれ以外に分割して、学習を実行

```
1 y_train5 = (y_train == 5)
2 sgd_clf = SGDClassifier(random_state=48)
3 sgd_clf.fit(x_train,y_train5)
```


のあった文字の例を示します (図 5)。

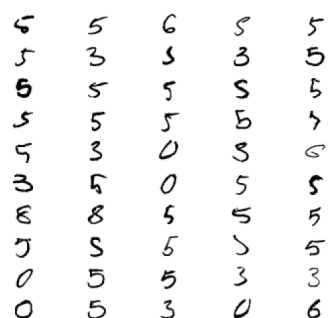


図 5 判断を誤った入力例

6 多値分類

一つの文字とそれ以外を区別することを、それぞれの文字を判別することに拡張することは、それほど困難ではありません。ある手書き文字をどの文字と判定するかの「確らしさ」を評価する方法があれば可能です。詳細は、省略します。

ラベル	データ数	正答率
0	980	0.98776
1	1135	0.95683
2	1032	0.63953
3	1010	0.96436
4	982	0.87271
5	892	0.55717
6	958	0.91127
7	1028	0.92023
8	974	0.81520
9	1009	0.85927

表 1 各試験データの正答率 (小数 6 桁切り捨て)

正答率を表 1 に示します。概ね良い結果を示していますが、“2”や“5”は誤りが非常に

多いことがわかります。結果を詳細に見ると、いずれの文字も”3”と誤って判定する
が多いことがわかります。

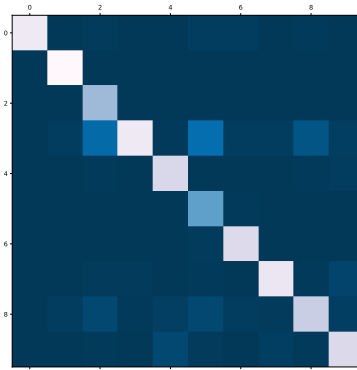


図 6 試験データの分析結果

結果を図にしたものを図 6 に示します。横軸は、試験データのラベル、つまり正答で
す。縦軸は、予想結果です。色が明るいほど件数が多いことを表しています。対角成分が
特に明るくなっていることは、多くの手書き文字を正しく判定していることを表していま
す。しかし、”2”や”5”は、判定が難しいことを表しています。縦軸の値が”3”であるところ
が対角成分以外で明るい色のところがあります。これは、”3”と誤って判断する場合は、
他の場合よりも多いことを示しています。

7 多値分類: その 2

学習手法を少し改善しましょう。784 次元の空間を曲面で区切ることを許しましょう。
また、訓練データからその局面をできるだけ離すような方法を使います。

正答率を表 2 に示します。各文字の正答率が大きく向上していることがわかります。誤
り率は 2% 程度になっています。

結果を図示しましょう (図 7)。先ほどの例 (図 6) と比べると、対角成分がさらに明るく
なっています。また、対角でない部分の明るい色の部分がなくなりました。予想精度が高
くなっていることがわかります。

ラベル	データ数	正答率
0	980	0.99286
1	1135	0.99207
2	1032	0.97481
3	1010	0.98515
4	982	0.97862
5	892	0.97646
6	958	0.98539
7	1028	0.96887
8	974	0.97536
9	1009	0.96135

表 2 各試験データの正答率 (小数 6 桁切り捨て)

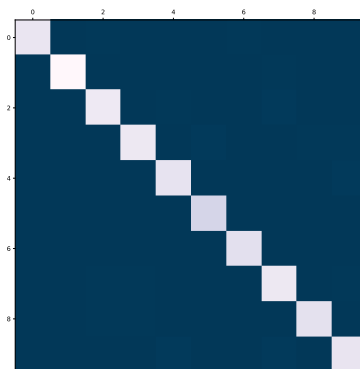


図 7 試験データの分析結果