

# ◆ matplotlib の練習 ◆

と思ったら、正規分布の勉強もしてた。

## 1 緒言

matplotlib は Python における基本的な可視化ライブラリである。seaborn というライブラリもあり<sup>\*1</sup>、これは見た目が美しいことから、いずれは使ってみたいと考えているのだが、まずは基本中の基本である matplotlib を押さえることにした。

また、実行環境として Google Colaboratory を使用することにした。一応ローカルにも Anaconda を入れてあるので、Jupyter notebook を使えるのだが、今後みんな Python を勉強していくにあたって、環境構築が不要な Google Colab を使用したほうが何かと都合がいいだろう。

## 2 1 変数の基本的な関数

まずはデータではなく単純な関数をプロットしてみたい。変数が  $x$  の一つしかない関数を matplotlib を用いてプロットする。

とりあえず使うライブラリを import。

---

```
1 import numpy as np
2 from matplotlib import pyplot as plt
```

---

L<sup>A</sup>T<sub>E</sub>X だから無駄に数式を書いてみる。次の一次関数と二次関数を matplotlib にプロットさせよう。

$$y = 2x - 1 \tag{1}$$

$$y = x^2 - 2x - 1 \tag{2}$$

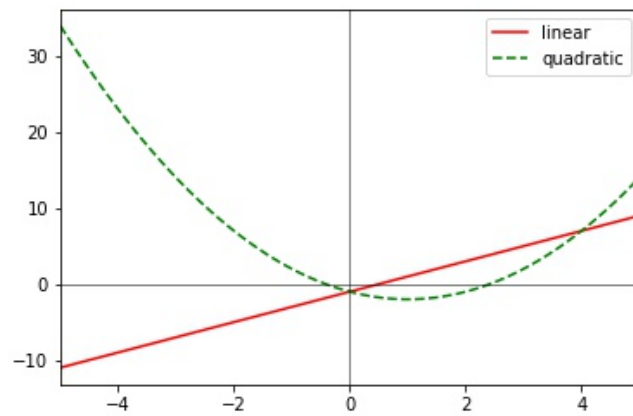
---

```
1 def linear(x):
2     return 2*x - 1
3
4 def quadratic(x):
5     return x**2 - 2*x - 1
6
7 x = np.linspace(-5, 5, 100) # x を numpy の linspace で list で与える。
8 y1, y2 = linear(x), quadratic(x)
9
10 plt.plot(x, y1, 'r', label='linear') # plot(x,y)を使う。'r'はマーカで赤の実線。
    label は凡例に使う。
11 plt.plot(x, y2, 'g--', label='quadratic') # 'g--'は緑の破線
12 plt.xlim(-5, 5) # x の範囲を絞る (デフォルトでは x の範囲外まで描画される)。
13 plt.axhline(y=0, linewidth=0.5, color='k') # x 軸
14 plt.axvline(x=0, linewidth=0.5, color='k') # y 軸
15 plt.legend() # 凡例を追加
16 plt.savefig('linear_and_quadratic.jpg') # 画像を保存。
    eps 形式だとなんか TEX にうまく貼り付けられなかった。
```

---

---

<sup>\*1</sup> ただし、これは matplotlib のラッパーである。



いい感じだ。三角関数  $y = \sin x$ ,  $y = \cos x$  や対数関数  $y = \log x$  もプロットしてみよう。

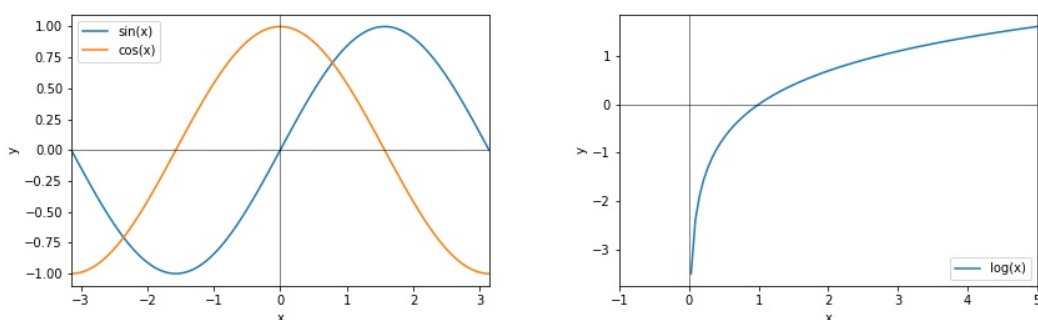
---

```

1 import math
2 pi = math.pi
3 x = np.linspace(-pi, pi, 100)
4
5 plt.plot(x, np.sin(x), label='sin(x)') # 三角関数は
    numpy のを使わないと怒られた。x が list なので。
6 plt.plot(x, np.cos(x), label='cos(x)')
7 plt.xlim(-pi, pi)
8 plt.xlabel('x') # ちゃんとラベルも追加した
9 plt.ylabel('y')
10 plt.axhline(y=0, linewidth=0.5, color='k')
11 plt.axvline(x=0, linewidth=0.5, color='k')
12 plt.legend()
13 plt.savefig('trigonometric.jpg')
14
15 x = np.linspace(-1, 5, 100)
16 plt.plot(x, np.log(x), label='log(x)')
17 plt.xlim(-1, 5)
18 plt.xlabel('x') # ちゃんとラベルも追加した
19 plt.ylabel('y')
20 plt.axhline(y=0, linewidth=0.5, color='k')
21 plt.axvline(x=0, linewidth=0.5, color='k')
22 plt.legend()
23 plt.savefig('logarithmic.jpg')

```

---



だいたい流れは分かった。 $x$  として、描きたい範囲の list を numpy の `linspace(start, end, slice)` で作り、matplotlib の `plot` に  $x$  と  $y$  の list を渡せばいいらしい。ただ、凡例やラベルや軸を書く書き方を忘れてしまいそう。

### 3 正規分布

2 変数の関数のプロットをしたいと思い、二変数の正規分布を描いてみたいと思ったのだが、その前にちょっと休憩。そもそも僕は、正規分布をよく知らない。

このページ<sup>\*2</sup>を見て、正規分布の式を 1 次元から多次元に拡張するところを勉強しようと思う。

#### 3.1 平均

確率変数  $x$  の平均とは  $x$  の期待値である。確率密度関数を  $p(x)$  として、 $x$  の平均  $\mu_x$  は、

$$\mu_x = \int_x xp(x)dx \quad (3)$$

で表される。うん、これはまあ言われてみれば分かる。

#### 3.2 分散

確率変数  $x$  の分散とは、平均  $\mu_x$  からの  $x$  の二乗誤差の期待値である。 $x$  の分散  $\sigma_x^2$  は、

$$\sigma_x^2 = \int_x (x - \mu_x)^2 p(x) dx \quad (4)$$

式 3 の  $x$  に  $x - \mu_x$  を代入した感じ？

#### 3.3 共分散

確率変数  $x$  と確率変数  $y$  の共分散は、同時確率密度関数<sup>\*3</sup>を  $p(x, y)$  として次の式で表される。

$$\sigma_{xy} = \int_x \int_y (x - \mu_x)(y - \mu_y)p(x, y) dy dx \quad (5)$$

……分からない。なんとなく式 4 に似ているのは分かる。 $x, y$  は、例えば数学のテストの点数と国語のテストの点数、だとか、身長と体重だとか、まあそういうのを思い浮かべるといいのだと思う。

<sup>\*2</sup> 一次元の正規分布から多次元正規分布へ - HELLO CYBERNETICS <https://www.hellocybernetics.tech/entry/2016/10/06/111153>

<sup>\*3</sup> 何ですかそれは。

この式のポイントは、 $x$ が大きくなると、 $y$ も大きくなるような関係がある場合、共分散は正の値、 $x$ が大きくなると  $y$  の値が小さくなるような関係がある場合、共分散は負の値を取ることである。いやあ、式を見るだけじゃよく分かんない……\*4。

### 3.4 相関係数

確率変数  $x$  と、確率変数  $y$  の相関係数  $\rho_{xy}$  は、次の式で表される。

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad (6)$$

共分散を正規化したものである。共分散は元の値の大小に引きずられてしまうので分かりにくい。この相関係数は  $[-1,1]$  の範囲に収まるので分かりやすい。相関係数はよく見るやつである。

### 3.5 独立性

確率変数  $x$  と確率変数  $y$  が独立であるとは、

$$p(x, y) = p(x)p(y) \quad (7)$$

が成り立つことである。一方の試行が他方の試行に影響を与えないことで、サイコロをふることなどが例として挙げられる。これが成り立てば、相関係数は 0 になる。

### 3.6 一次元の正規分布

一次元の正規分布は次の式で表される。

$$f(x) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x - \mu_x)^2}{2\sigma_x^2}\right) \quad (8)$$

へえ。 $\sigma_x$ (分散のルート、つまり標準偏差)と  $\mu_x^2$ (平均) が分かれば書けるのか\*5。matplotlib で書いてみるか。

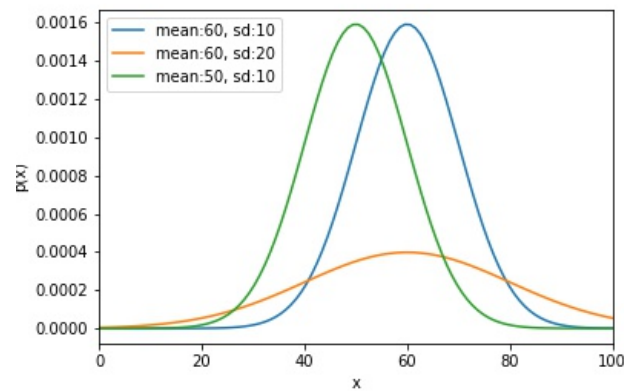
---

```
1 x = np.linspace(0, 100, 100)
2
3 def gaussian(x, mu, sigma):
4     e = np.exp(-(x-mu)**2/ (2*sigma**2))
5     return 1/(2*pi*sigma**2) * e # ルートを付けるの忘れてた
6
7 plt.plot(x, gaussian(x,60,10), label='mean:60, sd:10')
8 plt.plot(x, gaussian(x,60,20), label='mean:60, sd:20')
9 plt.plot(x, gaussian(x,50,10), label='mean:50, sd:10')
10 plt.xlim(0, 100)
11 plt.xlabel('x')
12 plt.ylabel('p(x)')
13 plt.legend()
```

---

\*4 今回  $x, y$  が確率変数だから余計に分らないのかも。

\*5 そらそうやな、という感じだが、改めて見ると「へえ」と思ってしまう。



標準偏差の値を 10 から 20 に変えるだけで、えらくグラフが平べったくなった\*6。  
 ん？ これは確率密度関数ということは全区間で積分すると 1 になる？

---

```

1 # gaussian にちゃんとルートを付け直した。
2 from scipy import integrate
3 integrate.quad(lambda x:gaussian(x, 60, 10), 0, 100)
4
5 # (0.9999683277715794, 3.6566249334013534e-09)
```

---

だいたい 1 になった。

$\mu_x = 0, \sigma_x^2 = 1$  であるような正規分布を標準正規分布という。

$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \quad (9)$$

### 3.7 特別な多次元標準正規分布

独立な確率変数  $x_i (i = 1, \dots, N)$  を考える。これをベクトルとして表す\*7。

$$\mathbf{x} = (x_1, \dots, x_N)^T \quad (10)$$

同時分布  $p(x_1, \dots, x_N) = p(\mathbf{x})$  は、次のように個々の確率密度変数の積になる。

$$p(\mathbf{x}) = p(x_1) \dots p(x_N) \quad (11)$$

なので、多変数の標準正規分布は、

$$f(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^N} \exp\left(-\sum_{i=1}^N \frac{x_i^2}{2}\right) \quad (12)$$

となる。指数の積の肩は和になるんだよね。

ここで！！ ベクトルの内積を使うとシグマを使わなくて済むのだ。すごい。

$$f(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^N} \exp\left(-\frac{\mathbf{x}^T \mathbf{x}}{2}\right) \quad (13)$$

---

\*6 上のグラフは正規分布の式の分母に  $\sqrt{\phantom{x}}$  をつけるのを忘れていたので、縦軸の値が間違っていますね……。

\*7 複数の変数をベクトル化して扱うのは Coursera の機械学習でだいぶ慣れた。

### 3.8 matplotlib で二変数正規分布をプロット

このあと、一般の正規分布について議論していたが、まずこの標準正規分布で2変数の場合について、matplotlib で実装してみよう。

まず、2変数の正規分布を返す関数と、 $x, y$ を表す ndarray<sup>\*8</sup>を作った。

---

```
1 import math
2 import numpy as np
3 pi = math.pi
4
5 def gaussian_3d(x,y):
6     return 1/np.sqrt(2 * pi)**2 * np.exp(-(x**2 + y**2) / 2)
7
8 x = np.linspace(-1,1,100)
9 y = np.linspace(-1,1,100)
```

---

3次元プロットをするためには、Axes3D とかいうライブラリを使うらしい。散布図の描画には scatter というメソッドを呼ぶ。

とりあえず、 $x, y$ を gaussian\_3d に渡して  $z$  を作り、scatter に突っ込んでみた。

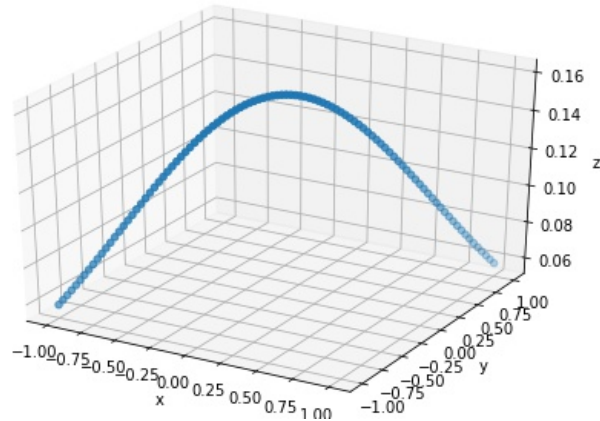
---

```
1 from mpl_toolkits.mplot3d import Axes3D
2 import matplotlib.pyplot as plt
3
4 z = gaussian_3d(x,y)
5
6 fig = plt.figure()
7 ax = Axes3D(fig)
8
9 ax.set_xlabel("x")
10 ax.set_ylabel("y")
11 ax.set_zlabel("z")
12
13 ax.scatter(x, y, z)
14
15 plt.show()
```

---

---

<sup>\*8</sup> 上で list とか言ってたやつも本当はこれ。



んんん……、線じゃなくて連続的な面として描いてほしいんだけどな。。。

いろいろググっているうちに、numpy の meshgrid というメソッドがあることが分かった。これは格子を作ってくれるメソッドらしい。もっばら三次元空間に何かをプロットしたいときに使うメソッドのようだ。

---

```

1 p = np.linspace(0,5,6)
2 q = np.linspace(0,5,6)
3
4 P, Q = np.meshgrid(p, q)
5
6 '''
7 P
8 [[0. 1. 2. 3. 4. 5.]
9  [0. 1. 2. 3. 4. 5.]
10 [0. 1. 2. 3. 4. 5.]
11 [0. 1. 2. 3. 4. 5.]
12 [0. 1. 2. 3. 4. 5.]
13 [0. 1. 2. 3. 4. 5.]]
14
15 Q
16 [[0. 0. 0. 0. 0. 0.]
17  [1. 1. 1. 1. 1. 1.]
18  [2. 2. 2. 2. 2. 2.]
19  [3. 3. 3. 3. 3. 3.]
20  [4. 4. 4. 4. 4. 4.]
21  [5. 5. 5. 5. 5. 5.]]
22 '''

```

---

なんか分かんけど、P と Q が行列になった。こいつを gaussian\_3d に渡すと、それぞれの要素同士の計算をしてくれるようである。なので当然戻り値も行列になる。要するに、P も Q も gaussian\_3d(P,Q) も面になったのだ。

---

```

1 R = gaussian_3d(P,Q)
2 np.round(R, 3) # 見やすいように桁数を制限
3
4 '''
5 array([[0.159, 0.097, 0.022, 0.002, 0. , 0. ],

```

---

```

6      [0.097, 0.059, 0.013, 0.001, 0. , 0. ],
7      [0.022, 0.013, 0.003, 0. , 0. , 0. ],
8      [0.002, 0.001, 0. , 0. , 0. , 0. ],
9      [0. , 0. , 0. , 0. , 0. , 0. ],
10     [0. , 0. , 0. , 0. , 0. , 0. ]])
11 '''

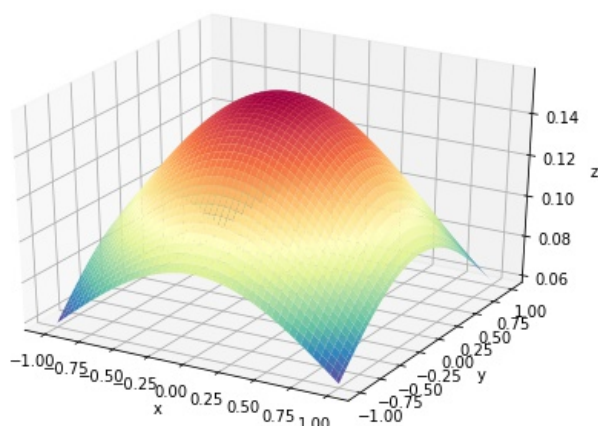
```

$x, y$  を meshgrid にして、それを gaussian\_3d に渡して、できた行列を Axes3D に渡してみた。plot\_surface というメソッドがあったので、そっちを使ってみることにした。

```

1 X, Y = np.meshgrid(x,y)
2 Z = gaussian_3d(X,Y)
3 ax.plot_surface(X, Y, Z, cmap='Spectral_r')

```

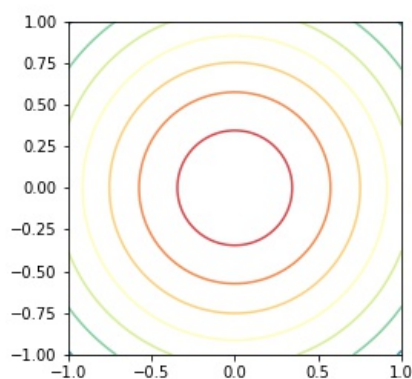


やったー。が、正直言うと 3 次元プロットは見た目はカッコいいが、向こう側が見えないという致命的欠陥がある。そこでよく使われるのが等高線図である。contour<sup>\*9</sup>というメソッドを使う。

```

1 plt.contour(X, Y, z, cmap='Spectral_r')
2 plt.gca().set_aspect('equal') # アス比を正方形に

```



\*9 輪郭という意味の英語



### 3.9 一般の多次元正規分布

閑話休題。上の二変数標準正規分布を一般の多次元正規分布に拡張する方法を学ぶ。ベクトルや行列を使って表記するといいらしい。

$\mathbf{x}$  は  $N$  個の確率変数を並べたベクトルである。よって確率変数  $x_i$  の平均  $\mu_i$  を同様にして並べたものは、次のように定義できる。

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)^T \quad (14)$$

次に、分散を表現したいところだが、その前に  $x_i$  と  $x_j$  の共分散の式は次のように表される。ここで、 $E$  は期待値を表す。分散/共分散とは、要は平均からの偏差の期待値なのである。

$$\sigma_{x_i x_j} = E((x_i - \mu_{x_i})(x_j - \mu_{x_j})) \quad (15)$$

分散とは  $i = j$  の場合の特別な共分散と考えられる。つまり、分散と共分散は同じ土俵にしまったほうが扱いが楽である。

$N$  個の確率変数があるとき、2つの  $i, j$  の組み合わせを考えると、全部で  $N^2$  個の組み合わせがあると考えられる。なので、分散は  $N \times N$  の行列で表したほうが楽なのだ。

次のような分散を表す行列  $\Sigma$  を考えてしまう\*10。

$$\Sigma = \begin{pmatrix} \sigma_{x_1}^2 & \sigma_{x_1 x_2} & \dots & \sigma_{x_1 x_N} \\ \sigma_{x_2 x_1} & \sigma_{x_2}^2 & \dots & \sigma_{x_2 x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_N x_1} & \sigma_{x_N x_2} & \dots & \sigma_{x_N}^2 \end{pmatrix} \quad (16)$$

行列の対角成分が分散、そのほかは共分散になっている感じ。これを分散共分散行列と呼ぶ。

$\Sigma$  の  $i, j$  成分は式 15 で表されるので、 $\Sigma$  をベクトルを使って書き直すと、

$$\Sigma = E((\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T) \quad (17)$$

になる。これ、縦ベクトル×横ベクトル=行列をうまく利用しているようなのだが、うーん、 $E$  が入っているからいまいちパツと分からない。式 15 を  $\Sigma$  に代入しちゃえ。

$$\mathbf{x} - \boldsymbol{\mu}_x = (x_1 - \mu_{x_1}, x_2 - \mu_{x_2}, \dots, x_N - \mu_{x_N})^T \quad (18)$$

$$\Sigma = E \begin{pmatrix} (x_1 - \mu_{x_1})(x_1 - \mu_{x_1}) & (x_1 - \mu_{x_1})(x_2 - \mu_{x_2}) & \dots & (x_1 - \mu_{x_1})(x_N - \mu_{x_N}) \\ (x_2 - \mu_{x_2})(x_1 - \mu_{x_1}) & (x_2 - \mu_{x_2})(x_2 - \mu_{x_2}) & \dots & (x_2 - \mu_{x_2})(x_N - \mu_{x_N}) \\ \vdots & \vdots & \ddots & \vdots \\ (x_N - \mu_{x_N})(x_1 - \mu_{x_1}) & (x_N - \mu_{x_N})(x_2 - \mu_{x_2}) & \dots & (x_N - \mu_{x_N})(x_N - \mu_{x_N}) \end{pmatrix} \quad (19)$$

確かに式 17 書くことができそう。

これらを使って、一般の多次元正規分布は次のように書けるそう。

$$f(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^m \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \quad (20)$$

ただし、 $|\Sigma|$  は  $\Sigma$  の行列式。

---

\*10 行列は TeX で書くのがとてもめんどくさい。

この式の導出については、本当にただひたすら式変形の話になってしまうので、天下りの受け入れることにしよう。

### 3.10 正規分布を返すメソッド

と、ここまで人力で実装してきたあれなのだが、`scipy` には正規分布を返すメソッドが用意されている。`norm` は正規分布。`pdf` は確率密度関数 (Probability Density Function)。出力されるグラフは当然同じなので省略する。

---

```
1 from scipy import stats
2
3 x = np.linspace(0, 100, 1000)
4 p = scipy.stats.norm.pdf(x, 60, 5)
5 plt.plot(x, p, label='mu=60, sd=5')
6 plt.legend()
```

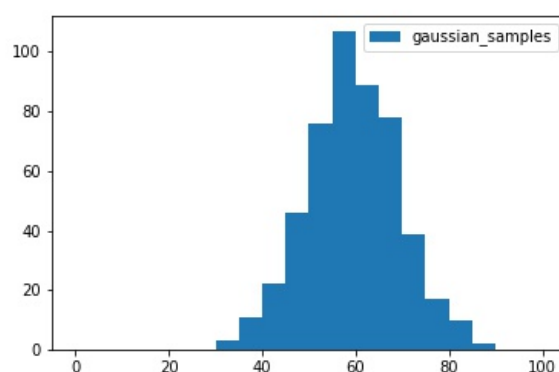
---

さらになんと次のように、正規分布に従う乱数の `ndarray` を返してくれるメソッドもある。ヒストグラムにしてみた。

---

```
1 g = np.random.normal(60, 10, 500)
2 plt.hist(g, bins=20, range=(0,100))
```

---



`np.random.normal` の引数は、(平均, 標準偏差, サンプル数) の意味。

## 4 結言

特になにか難しいことをしているわけではなかったが、真面目に `matplotlib` を触ったのがこれが初めてな気がする。結構呪文的なところが多く、ちゃんと覚えられないので、当分はコピペでやることになりそう。また、`matplotlib` や `numpy` にどんなメソッドがあるのかというのはきちんと覚えたい。しばらくはこれらの可視化系ライブラリの練習になりそう。また、今回は乱数や連続関数しかプロットしなかったが、今度は意味のあるデータ<sup>\*11</sup>をプロットしたり、`pandas` を用いて集計したりしてみたい。

また、正規分布の定義式を初めてじっくり眺めてみた。そこから多次元に拡張するのは、正直難しくてあまり

---

<sup>\*11</sup> 今考えているのは気象庁のアメダス観測データ。

よく分からなかったが、多変数を扱うのにベクトルや行列で記述すると効率的であるということが分かったので、今後線形代数やこれらを使った量の表記についても慣れていく必要があるなと思った。