

## Bernoulli 分布

Bernoulli (ベルヌーイ) 分布は、「成功か失敗か」、「表か裏か」、「勝ちか負けか」のように2種類のみ結果しか得られない試行 (Bernoulli 試行) の結果を 0 と 1 で表した分布です。式で表すと、次のようになります。

$$P(X = 1) = \phi$$

$$P(X = 0) = 1 - \phi$$

$$P(X = x) = \phi^x(1 - \phi)^{1-x}$$

以下では、確率分布の基本的なことから上の式を説明していきます。また、Bernoulli 分布の期待値や分散などの計算も行なっていきます。

### ♣ 確率分布とは

例として、「サイコロ 2 個振ったときの出た目の和  $X$ 」を考えましょう。この  $X$  は、どのような値となるかが確率法則によって決まるような変数で、一般に**確率変数**と呼ばれます。この確率変数  $X$  に対する確率分布 (確率変数をとる値とその値をとる確率の対応の様子) は次の表のようになります。

$X$	2	3	4	5	6	7	8	9	10	11	12
確率	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$

サイコロ 2 個の目の出かたの場合の数は、全部で  $6 \times 6 = 36$  通りです。これより、

$$X = 2: (1, 1) \quad \rightarrow \frac{1}{36}$$

$$X = 3: (1, 2), (2, 1) \quad \rightarrow \frac{2}{36}$$

$$X = 4: (1, 3), (2, 2), (3, 1) \quad \rightarrow \frac{3}{36}$$

⋮

といった要領で、 $X$  の取りうる値、それぞれに対する確率を求めることができます。「確率変数  $X$  が  $X = 2$  となる確率は  $\frac{1}{36}$  である」ということを、数式では次のように表現します。

$$P(X = 2) = \frac{1}{36}$$

$P$ は、確率 (probability) の頭文字となっています。確率分布で最も重要なことは、確率分布では全体の合計が1となるということです。

$$P(X = 2) + P(X = 3) + \cdots + P(X = 12) = \sum_{k=2}^{12} P(X = k) = 1$$

$P(X = 2)$  から  $P(X = 12)$  までの和を、総和記号  $\sum$  を用いて表現しました<sup>26</sup>。

### ♣ Bernoulli 分布とは

Bernoulli 分布は、上で説明した確率分布の1つの例です。先ほどは「サイコロ2個振ったときの出た目の和」という結果が複数あるような試行 (起こりうる結果がいくつかあり、どれが起こるかは偶然で決まる手順のこと) を考えました。ここでは、起こりうる結果が2つしかない、Bernoulli 試行と呼ばれる試行を考えます。例えば、コインを投げるという試行は、結果が表が出るか裏が出るかの2つなので、Bernoulli 試行となっています。コインを投げた結果、表が出ることを  $X = 1$ 、裏が出ることを  $X = 0$  と表しましょう。表が出る確率を  $\phi$  (ファイと読みます) とすると<sup>27</sup>、数式では

$$P(X = 1) = \phi$$

となります。 $\phi$  は確率を表すので、必ず0以上1以下の数です。このことを数学の記号で

$$\phi \in [0, 1]$$

と表します。 $\phi$  が0以上1以下の領域 ( $[0, 1]$ ) に含まれている ( $\in$ ) という意味です。確率分布では全体の合計 (今の場合は  $X = 0$  と  $X = 1$  の2つの合計) は1となるので、

$$P(X = 1) + P(X = 0) = 1$$

です。これより、

$$P(X = 0) = 1 - \phi$$

と分かります。この  $P(X = 1)$  と  $P(X = 0)$  の式を1つにまとめると、次のように表すことができます<sup>28</sup>。

$$P(X = x) = \phi^x (1 - \phi)^{1-x}$$

実際に  $x = 1$ 、 $x = 0$  を入れてみれば、上の2つの式が再現できることが確認できます。

$$P(X = 1) = \phi^1 (1 - \phi)^{1-1} = \phi$$

$$P(X = 0) = \phi^0 (1 - \phi)^{1-0} = 1 - \phi$$

<sup>26</sup>総和記号については、softmax 関数のところで詳しく説明しています。

<sup>27</sup>完璧なコインなら  $\phi = 0.5$  ですが、コインが歪んでいたりすると0.5からズレるでしょう。

<sup>28</sup>この  $P(X = x)$  のことを「確率質量関数 (probability mass function)」と言います。

0乗は1となることに注意して下さい<sup>29</sup>。以上で Bernoulli 分布の説明はおしまいです。最初式を見た時は難しそうに見えたかもしれませんが、実は非常にシンプルな分布なのです。

## ♣ 期待値と分散

確率分布を特徴づける量として、「期待値 (expected value)」と「分散 (variance)」があります。それぞれがどのような量か、確認しておきましょう。

期待値とは「確率変数を取る値を、確率によって重み付けした平均値」です。例えば、300 円の宝くじ 1 枚の期待値が 100 円であった場合、その宝くじには 100 円の価値が期待できるということです。当たったくじ、外れたくじの総合で見て、平均すると 1 枚あたり 100 円の価値であったということです。期待値という名前は、確率変数を取ると「期待」される値であることから名付けられています。数式では、期待値は次のように定義されます。

$$E[X] = \sum_k kP(X = k)$$

確率変数  $X$  の取り得る全ての値  $k$  について、その値とその値が出る確率を掛け合わせたものを足し合わせる、ということです<sup>30</sup>。例えば、最初に考えた「サイコロ 2 個振ったときの出た目の和  $X$ 」の期待値を計算しましょう。上の式をこの例の場合に適用すると、

$$\begin{aligned} E[X] &= \sum_{k=2}^{12} kP(X = k) \\ &= 2 \times P(X = 2) + 3 \times P(X = 3) + \dots + 12 \times P(X = 12) \\ &= 2 \times \frac{1}{36} + 3 \times \frac{2}{36} + \dots + 12 \times \frac{1}{36} = 7 \end{aligned}$$

となります。つまり、サイコロ 2 個振ったときの出た目の和の期待値は 7 と分かります。

次に分散を確認します。分散とは「確率変数のばらつき具合を表すための指標」です。分散が大きいほどばらつきが大きく、逆に分散が小さいほど確率変数の散らばりが少なく平均の近くにあることとなります。確率変数  $X$  の分散は、数式では次のように表現されます。

$$\text{Var}[X] = E[(X - E[X])^2]$$

<sup>29</sup>softmax 関数のところで説明しています。

<sup>30</sup>上の式では添え字  $k$  の動く範囲が一般には異なるので、特に範囲は明記していません。

つまり、確率変数  $X$  の期待値  $E[X]$  からのズレの 2 乗の期待値です。この式では分かりにくいので、総和記号を使って表すと、次のようになります。

$$\text{Var}[X] = \sum_k (k - E[X])^2 P(X = k)$$

これは次のように計算を進めることができます。

$$\begin{aligned} \text{Var}[X] &= \sum_k (k^2 - 2E[X]k + E[X]^2) P(X = k) \\ &= \sum_k k^2 P(X = k) - 2E[X] \sum_k k P(X = k) + E[X]^2 \sum_k P(X = k) \end{aligned}$$

ここで、1 つ目の項は期待値の定義より、

$$\sum_k k^2 P(X = k) = E[X^2]$$

つまり、確率変数  $X$  の 2 乗の期待値となります。2 つ目の項も期待値の定義より、

$$2E[X] \sum_k k P(X = k) = 2E[X] \times E[X] = 2E[X]^2$$

最後の項は、確率の合計が 1 であることを使うと、

$$E[X]^2 \sum_k P(X = k) = E[X]^2 \times 1 = E[X]^2$$

これより、分散は次のようになることが分かります。

$$\text{Var}[X] = E[X^2] - E[X]^2$$

つまり、2 乗の期待値から期待値の 2 乗を引いた量となっています。分散を計算する際には、最初の定義の式で計算しても良いですし、この 2 乗の期待値から期待値の 2 乗を引く方法で計算しても同じ結果となります。状況に応じて、計算が簡単な方で計算を行えば良いでしょう。

## ♣ Bernoulli 分布の期待値と分散

それでは、最後に Bernoulli 分布の期待値と分散を計算してみましょう。Bernoulli 分布は

$$P(X = 1) = \phi$$

$$P(X = 0) = 1 - \phi$$

でした。1か0しかないので、計算はとても簡単です。先ず期待値は、前のセクションで説明した定義に従って計算すると、

$$E[X] = \sum_{k=0}^1 kP(X = k) = 0 \times P(X = 0) + 1 \times P(X = 1) = \phi$$

となり、Bernoulli 分布の期待値は  $\phi$  であることが分かります。次に分散を計算します。ここでは、2乗の期待値から期待値の2乗を引く方法で分散を計算しましょう。2乗の期待値は、

$$E[X^2] = \sum_{k=0}^1 k^2P(X = k) = 0^2 \times P(X = 0) + 1^2 \times P(X = 1) = \phi$$

となり、2乗の期待値も  $\phi$  であることが分かります。これより、

$$\text{Var}[X] = E[X^2] - E[X]^2 = \phi - \phi^2 = \phi(1 - \phi)$$

となり、Bernoulli 分布の分散は  $\phi(1 - \phi)$  であることが分かります。

# Gauss 分布

一般的な Gauss (ガウス) 分布は、次の式で与えられます。

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\beta}(\mathbf{x} - \boldsymbol{\mu})\right)$$

この分布は「**多変量 Gauss 分布**」や「**多変量正規分布**」とも呼ばれます。複雑な上記の式の意味を理解することが、この解説プリントの目標です。

## ♣ 1 変数の Gauss 分布

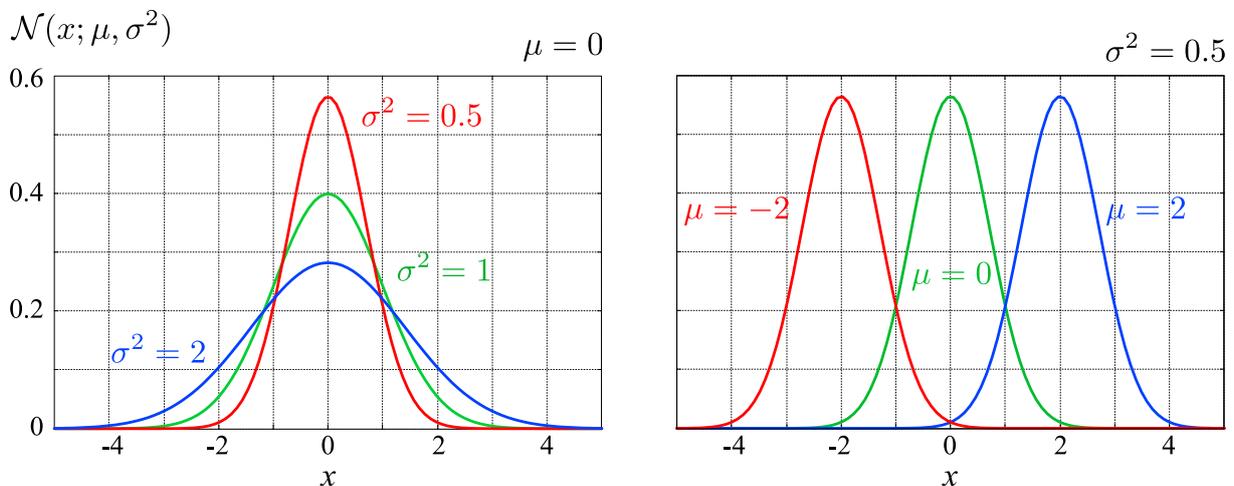
最初に書いた一般的な Gauss 分布は、複数の変数 (多変数) に対する確率分布を表すものです。ここでは先ず、変数が1つの場合を考えましょう。これを基本に、多変数の場合に拡張していくことで、最初の式の意味を理解することができます。Gauss 分布 (正規分布ともいいます) は、次の式で与えられる確率分布です。

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

$\exp$  というのは、底がネイピア数の指数関数のことです。つまり、

$$\exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) = e^{-\frac{1}{2\sigma^2}(x - \mu)^2}$$

です。特に指数が複雑な場合には、 $e^A$  よりも  $\exp(A)$  という書き方の方が見やすいため、この書き方がよく使われます。 $x$  は確率変数です。Bernoulli 分布の場合には、確率変数は0か1のみを取る「**離散型**」の変数でした<sup>31</sup>。この Gauss 分布では連続的にあらゆる値を取る「**連続型**」の変数を考えるため、Gauss 分布も次のように連続的なグラフとなります。



<sup>31</sup>離散というのは、とびとび、という意味です。

$\sigma^2$  と  $\mu$  (ミュー) の値を変化させると分布の様子が変化することが分かります。このように  $\sigma^2$  と  $\mu$  は分布を特徴づける量となっており<sup>32</sup>、 $\sigma^2$  は「分散」、 $\mu$  は「平均」と呼ばれています。分散  $\sigma^2$  を大きくすると、曲線の幅が広がり、確率変数の取り得る値のばらつきが大きくなることが分かります。つまり、 $\sigma^2$  は確率変数の取り得る値のばらつきをコントロールするパラメータです。一方で平均  $\mu$  を変化させると、曲線の中心の位置がずれることが分かります。Gauss 分布は、この平均  $\mu$  を中心に左右対称な曲線となります。

以上で1変数の Gauss 分布の基本的な事項はおしまいです。最後にいくつか細かなコメントをしておきます。まず、Gauss 分布の表し方

$$\mathcal{N}(x; \mu, \sigma^2)$$

の表記法についてです。 $\mathcal{N}$  は正規分布 (normal distribution) の頭文字です。そして連続変数  $x$  に対する確率分布なので、 $\mathcal{N}(x)$  と表します。ただ、Gauss 分布は2つのパラメータ  $\mu$  と  $\sigma^2$  によって分布が変化します。そのため、これらのパラメータの値が与えられた上での、 $x$  に対する確率分布という意味で  $\mathcal{N}(x; \mu, \sigma^2)$  と表しています。もう1つのコメントは、

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

についている  $\sqrt{\frac{1}{2\pi\sigma^2}}$  の部分についてです。この部分は定数なので、単に分布全体の大きさを変化させるだけで、分布の分散や平均を変化させることはない、いわばどうでもよい部分です。なのに何やら複雑な係数となっているのは、 $\mathcal{N}(x; \mu, \sigma^2)$  が確率分布の基本ルール「確率分布では全体の合計が1となる」を守るようにするためです。少しややこしいのは、Gauss 分布の場合には確率変数が連続なので、単純な足し算で「全体の合計」を計算することはできない点です。このような連続的な変数の場合の全体の合計の計算には、積分と呼ばれる特殊な“足し算”を使います。分からない場合には無視してもらって大丈夫ですが、数式で表すと、

$$\int_{-\infty}^{\infty} \mathcal{N}(x; \mu, \sigma^2) dx = 1$$

となるようになっています。とにかく、 $\mathcal{N}(x; \mu, \sigma^2)$  が確率分布の基本ルールを守るようにするために、なにやら複雑な係数がついているのだと思っておけば十分です<sup>33</sup>。

<sup>32</sup>このような変数をパラメータ (parameter) と呼びます。

<sup>33</sup>このような係数のことを、一般に「規格化係数」と言います。

## ♣ 1 変数から多変数へ

一般に何か測定したとすると、その測定結果は様々な要素に依存しているでしょう。こういった様々な要素をまとめて扱うためには、1変数の確率分布のみでは不十分で、多変数の確率分布を考える必要が出てきます。ここでは、前のセクションで説明した1変数の Gauss 分布を多変数の場合に拡張するための準備を行なっていきます。

### ベクトル

$n$  個の確率変数  $x_i$  ( $i = 1, \dots, n$ ) があつたとします。このように多くの変数を扱う際には、「**ベクトル (vector)**」と呼ばれる量を導入すると便利です<sup>34</sup>。

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$n$  個の変数  $x_1$  から  $x_n$  を縦に並べて、それをひとつの「グループ」とみなしています。このベクトルは、プログラミング言語でいうところの「配列」です。 $n$  個の変数を全てバラバラの変数名で扱うと大変なので、配列を導入して  $n$  個の変数を1つのグループとして扱えるようにしたのと同じことです。このベクトル  $\mathbf{x}$  に対応する量で、次のような量を導入しておきましょう<sup>35</sup>。

$$\mathbf{x}^T = (x_1 \ x_2 \ \cdots \ x_n)$$

$\mathbf{x}$  は縦に変数を並んでいるため「縦ベクトル」、 $\mathbf{x}^T$  では変数が横に並んでいるため「横ベクトル」と呼ぶことにしましょう<sup>36</sup>。

ベクトルは、複数の数を1つのグループにまとめた量なので、通常の数とは計算のルールが異なります。例として、2つのベクトルを考えます。

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

この2つのベクトルの足し算と引き算は、次のようになります。

<sup>34</sup>ここではベクトルは太字で表すことにしますが、 $\vec{x}$  と書く場合もあります。

<sup>35</sup>T は転置 (transpose) という操作を表す記号ですが、あまり気にしなくて大丈夫です

<sup>36</sup>本当は、列ベクトル、行ベクトルと呼ばれますが、行と列がややこしいので横と縦にしています。

$$\mathbf{x} + \mathbf{y} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{pmatrix}$$

$$\mathbf{x} - \mathbf{y} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 - y_1 \\ x_2 - y_2 \\ \vdots \\ x_n - y_n \end{pmatrix}$$

これがベクトルの足し算と引き算のルールです。同じ位置の変数を足せば(引けば)よいだけです。次に掛け算についてです。ベクトルの場合、実は普通の数とは違って掛け算といっても色々な種類の掛け算があるのですが<sup>37</sup>、ここでは多変量 Gauss 分布を理解するのに必要な掛け算のみを確認しましょう(単に掛け算といった場合には、通常この計算を指します)。

$$\mathbf{x}^T \mathbf{y} = (x_1 \ x_2 \ \cdots \ x_n) \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n$$

いくつか注意点を述べます。まず、このベクトルの計算は、横ベクトルと縦ベクトルの積のみに適用できる計算です。ですから、

$$\mathbf{x} \mathbf{y} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{x}^T \mathbf{y}^T = (x_1 \ x_2 \ \cdots \ x_n) (y_1 \ y_2 \ \cdots \ y_n)$$

といった掛け算はありません<sup>38</sup>。また、横ベクトルと縦ベクトルの順番にも気をつけて下さい。

$$\mathbf{x} \mathbf{y}^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} (y_1 \ y_2 \ \cdots \ y_n)$$

みたいな計算も今回は扱いません<sup>39</sup>。普通の数と違って、掛け算の順番も大事なのがベクトル

<sup>37</sup>内積、外積、テンソル積など色々あります。

<sup>38</sup>内積  $\mathbf{x} \cdot \mathbf{y}$  や外積  $\mathbf{x} \times \mathbf{y}$  みたいな計算はありますが、ここでは気にしなくて大丈夫です。

<sup>39</sup>テンソル積と呼ばれる掛け算なのですが、気にしなくてよいです。

の計算の大きな特徴です。もう1つ重要な注意点は、ベクトルの掛け算  $\mathbf{x}^T \mathbf{y}$  の結果は、ベクトルではないという点です。

$$\mathbf{x}^T \mathbf{y} = (x_1 \ x_2 \ \cdots \ x_n) \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n$$

において、 $x_1 y_1 + x_2 y_2 + \cdots + x_n y_n$  は  $n$  個の数を足した単なる数です。このような数を、ベクトルと区別するために「**スカラー (scalar)**」と呼びます。別にスカラーは何も新しい数ではなく、 $2+3=5$  のような普通の計算はスカラーの足し算に他なりません。ですから、ベクトル同士の掛け算の結果は、普通の数になるのだと思っておけば十分です。

## 行列

ベクトルをさらに発展させたものとして「**行列 (matrix)**」と呼ばれる量があります。これは例えば次のような量です<sup>40</sup>。

$$\mathbf{M} = \begin{pmatrix} M_{11} & M_{12} & \cdots & M_{1n} \\ M_{21} & M_{22} & \cdots & M_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ M_{n1} & M_{n2} & \cdots & M_{nn} \end{pmatrix}$$

上のように複数の数を四角状に並べたものです。なぜこんな量を考えるのでしょうか？ 行列のありがたみを感じるために、例として、3人の英語、数学、国語の点数が下の表のようになっている場合を考えましょう。

	英語	数学	国語
Aさん	80	90	70
Bさん	55	100	65
Cさん	85	60	95

行列を使うと、このデータを次のように1つの文字  $\mathbf{M}$  で表すことができるのです。

$$\mathbf{M} = \begin{pmatrix} M_{A, \text{英語}} & M_{A, \text{数学}} & M_{A, \text{国語}} \\ M_{B, \text{英語}} & M_{B, \text{数学}} & M_{B, \text{国語}} \\ M_{C, \text{英語}} & M_{C, \text{数学}} & M_{C, \text{国語}} \end{pmatrix} = \begin{pmatrix} 80 & 90 & 70 \\ 55 & 100 & 65 \\ 85 & 60 & 95 \end{pmatrix}$$

このように、複数のデータを1つにまとめにして扱えるため、多変数のデータを扱う場合にこの行列という量が大活躍するのです。

<sup>40</sup>下の行列は特に、 $n \times n$  行列と呼ばれます。大きさ (サイズ) は  $n \times n$  のように表されます。

行列に関しても様々な計算のルールがあります。ここでは多変量 Gauss 分布を理解するのに必要となる、ベクトルと行列の掛け算を確認しましょう。大きな行列、ベクトルを考えると大変なので、例として次の  $2 \times 2$  の行列とベクトルを考えましょう。

$$\mathbf{M} = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

ベクトルと行列の掛け算としては、次の 2 パターンがあります。

$$\mathbf{x}^T \mathbf{M} = (x_1 \ x_2) \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} = (x_1 M_{11} + x_2 M_{21} \quad x_1 M_{12} + x_2 M_{22})$$

$$\mathbf{M} \mathbf{x} = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x_1 M_{11} + x_2 M_{12} \\ x_1 M_{21} + x_2 M_{22} \end{pmatrix}$$

ややこしいですが、これが計算のルールです。よく見ると、どのような規則で計算が行われているか分かると思います。大事なのは、 $\mathbf{x}^T \mathbf{M}$  の結果は横ベクトルに、 $\mathbf{M} \mathbf{x}$  の結果は縦ベクトルになる点です。

上の 2 つのパターンのベクトルと行列の掛け算を組み合わせた、次のような式を考えましょう。

$$\mathbf{x}^T \mathbf{M} \mathbf{x} = (x_1 \ x_2) \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

このような量は、 $\mathbf{x}^T \mathbf{M}$ 、 $\mathbf{M} \mathbf{x}$ 、どちらの部分から計算しても構いません。

$$\mathbf{x}^T \mathbf{M} \mathbf{x} = (\mathbf{x}^T \mathbf{M}) \mathbf{x} = (x_1 M_{11} + x_2 M_{21} \quad x_1 M_{12} + x_2 M_{22}) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\mathbf{x}^T \mathbf{M} \mathbf{x} = \mathbf{x}^T (\mathbf{M} \mathbf{x}) = (x_1 \ x_2) \begin{pmatrix} x_1 M_{11} + x_2 M_{12} \\ x_1 M_{21} + x_2 M_{22} \end{pmatrix}$$

いずれにしても、先ほど確認した、横ベクトルと縦ベクトルの積となります。そのため、どちらで計算しても、

$$\mathbf{x}^T \mathbf{M} \mathbf{x} = M_{11} x_1^2 + (M_{12} + M_{21}) x_1 x_2 + M_{22} x_2^2$$

となります。大事なことは、 $\mathbf{x}^T \mathbf{M} \mathbf{x}$  の計算結果は、行列でもベクトルでもなく、普通の数 (スカラー) となるという点です。ここでは、 $2 \times 2$  の行列を例に説明しましたが、この計算はそのまま  $n \times n$  の行列での計算に応用できます。

## ♣ 多変量 Gauss 分布の式を理解していく

以上で基本的な準備は完了しました。もう1度、最初の式を見てみましょう。

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\beta}(\mathbf{x} - \boldsymbol{\mu})\right)$$

以下では1変数の場合の式

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

と比較しながら、最初の式を理解していきます。先ず、確率変数が1つの場合には単に  $x$  と書いていましたが、複数 ( $n$  個) の確率変数が含まれる場合には、 $n$  個の確率変数  $x_i$  ( $i = 1, \dots, n$ ) をまとめたベクトル  $\mathbf{x}$  を導入すると便利です。

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

確率変数が  $n$  個あると、それぞれのデータに対して平均値  $\mu$  が存在するため、 $\mu$  も次のようなベクトルとなっています。

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$

ですから、2つのベクトルの引き算は、

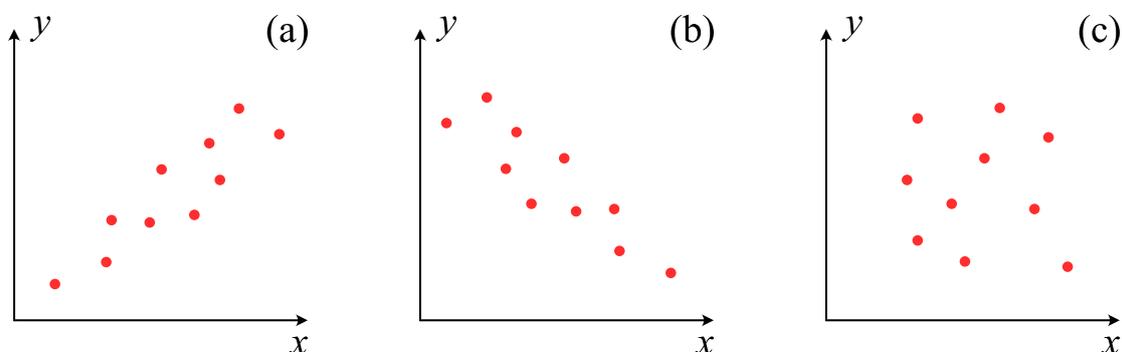
$$\mathbf{x} - \boldsymbol{\mu} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_n - \mu_n \end{pmatrix}$$

となっています。多変量 Gauss 分布に含まれる  $\boldsymbol{\beta}$  は「精度行列 (precision matrix)」と呼ばれる  $n \times n$  の行列です<sup>41</sup>。この行列は、確率変数の散らばり具合を表す量で、共分散と呼ばれる量を一般化したものだと考えられます。2つの確率変数  $x, y$  に対して、

$$\text{Cov}(x, y) = E[(x - E[x])(y - E[y])] = E[xy] - E[x]E[y]$$

<sup>41</sup>逆共分散行列 (inverse covariance matrix) とも呼ばれます

を  $x$  と  $y$  の共分散といいます。  $E[x]$  は Bernoulli 分布のところで説明した確率変数  $x$  の期待値です。なにやらややこしい量ですが、この量は2つの変数  $x, y$  の「相関」を表しています。相関というのは、2つの変数の関係性のことです。例えば、確率変数  $x, y$  のデータ点を10点ほどプロットしたら、次のようになったとしましょう。



(a) の場合には、  $x$  が増加すると  $y$  も増加する傾向にあることが分かります。この場合には、  $x$  と  $y$  の間には「正の相関がある」と言います。一方、(b) の場合には、  $x$  が増加すると  $y$  が減少する傾向にあることが分かります。この場合には、  $x$  と  $y$  の間には「負の相関がある」と言います。そして(c) の場合には、特に  $x$  と  $y$  の間には特に関係がなくバラバラにデータが広がっています。このような場合を「無相関」と言います。共分散  $\text{Cov}(x, y)$  の値は、正の相関がある場合には正、負の相関がある場合には負の数となります。そして無相関の場合には、共分散の値は0となります。このように共分散  $\text{Cov}(x, y)$  は、  $x$  と  $y$  の相関を反映した量となっています。

共分散について最後にいくつかコメントしておきます。共分散に関しては、定義から次の関係式が成り立ちます。

$$\text{Cov}(x_1, x_2) = \text{Cov}(x_2, x_1)$$

$x_1$  と  $x_2$  を入れ替えても値は変わりません。また、  $x$  と  $x$  に関する共分散は、

$$\text{Cov}(x, x) = E[x^2] - E[x]^2$$

となります。2乗の期待値から期待値の2乗を引いた量は、 Bernoulli 分布のところで出てきた「分散」に他なりません。つまり、

$$\text{Cov}(x, x) = \text{Var}(x)$$

ということになります。

この共分散を用いて、精度行列  $\beta$  (の逆行列) は次のように与えられます。

$$\beta^{-1} = \begin{pmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) & \cdots & \text{Cov}(x_1, x_n) \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) & \cdots & \text{Cov}(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_n, x_1) & \text{Cov}(x_n, x_2) & \cdots & \text{Cov}(x_n, x_n) \end{pmatrix}$$

つまり  $n$  個の確率変数  $x_i$  の関係を調べるためには、2つの変数を選び出してそれらの共分散を考えることとなりますが、2つの変数の選び方にはたくさんパターンがあります。それらの情報を1つにまとめてしまったのが精度行列  $\beta$  なのです。尚、上の式では  $\beta$  の表式ではなく、 $\beta$  の「逆行列」と呼ばれる行列  $\beta^{-1}$  の方が綺麗な形となるので、逆行列の形で書いています<sup>42</sup>。一般のサイズの行列の逆行列の計算の仕方は、大学レベルの数学となり、かなりややこしいのでここでは説明を省略します。ただ実際には、逆行列の計算はコンピュータで簡単にできてしまうので、使う分には計算アルゴリズムの詳細を知らなくても全然問題ありません。

以上で必要な行列  $\beta$  とベクトル  $\mathbf{x}$ ,  $\boldsymbol{\mu}$  が揃いました。これらの行列、ベクトルを用いて多変量 Gauss 分布は構成されています。

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \beta^{-1}) = \sqrt{\frac{\det(\beta)}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \beta (\mathbf{x} - \boldsymbol{\mu})\right)$$

$(\mathbf{x} - \boldsymbol{\mu})^T \beta (\mathbf{x} - \boldsymbol{\mu})$  の部分は、行列とベクトルの積となっています。先ほどのベクトルと行列のところで説明しましたが、このような横ベクトル、行列、縦ベクトルの積は、普通の数 (スカラー) となります。つまり、 $(\mathbf{x} - \boldsymbol{\mu})^T \beta (\mathbf{x} - \boldsymbol{\mu})$  もルールに従って計算すれば、スカラーの量となります。そして  $\exp$  の前についている

$$\sqrt{\frac{\det(\beta)}{(2\pi)^n}}$$

は、1変数の場合の  $\sqrt{\frac{1}{2\pi\sigma^2}}$  に対応する係数です。  $\det(\beta)$  というのは、 $\beta$  の「**行列式 (determinant)**」と呼ばれる量です。一般のサイズの行列の  $\det$  を求めるアルゴリズムも逆行列と同様、かなり複雑なので (大学レベルの数学となります)、ここでは説明を省略しますが、行列に対してスカラーを返すような計算です。つまり、 $\det(\beta)$  は普通の数であり、パッケージやライブラリを利用すればどのプログラミング言語でも簡単に計算できる量となっています。この係数も1変数の場合と同様に、分布  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \beta^{-1})$  が確率分布の基本ルール「全て合計したら1になる」を満たすようにするためにくっついています。

<sup>42</sup>精度行列  $\beta$  の逆行列  $\beta^{-1}$  は、分散共分散行列と呼ばれ、一般には  $\Sigma$  と表します。

以上で、多変量 Gauss 分布に関する説明はおしまいです。一般の  $n$  変数の場合だとこのようにややこしくて分かりにくいので、最後に 2 変数の場合の多変量 Gauss 分布を具体的に考えます。具体例をみれば、より理解が深まるかと思えます。

### ♣ 2 変数の多変量 Gauss 分布

では、具体的に  $n = 2$  の場合、つまり 2 変数に対する Gauss 分布を考えてみましょう。確率変数のベクトル  $\mathbf{x}$ 、平均値のベクトル  $\boldsymbol{\mu}$  は、この場合、

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

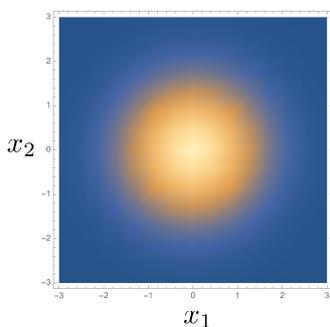
となります。平均は 1 変数の場合に説明したように、全体の位置をずらすだけなので、ここでは簡単のために、 $\mu_1 = \mu_2 = 0$  としておきましょう。精度行列  $\boldsymbol{\beta}$  の逆行列 (つまり分散共分散行列) は、この場合

$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) \end{pmatrix} = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_1, x_2) & \text{Var}(x_2) \end{pmatrix}$$

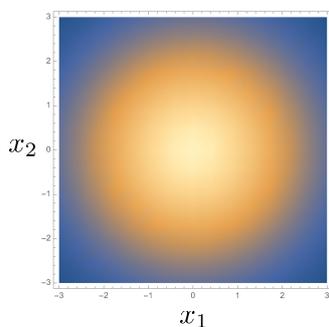
となっています。前のセクションで説明したように、 $\text{Cov}(x_1, x_2) = \text{Cov}(x_2, x_1)$  であること、 $\text{Cov}(x_1, x_1) = \text{Var}(x_1)$  であることを用いています。以下では先ず、 $x_1$  と  $x_2$  が無相関である場合を考えましょう。この場合には、共分散  $\text{Cov}(x_1, x_2)$  は 0 となるので、

$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} \text{Var}(x_1) & 0 \\ 0 & \text{Var}(x_2) \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$$

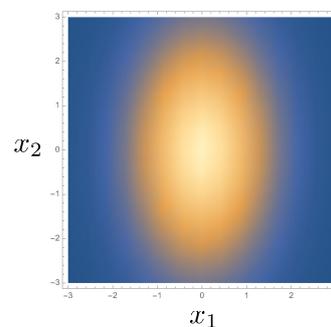
となります。ここでは、1 変数の場合の Gauss 分布で分散を  $\sigma^2$  と表現していたため、それを真似して  $\text{Var}(x_1) = \sigma_1^2$ 、 $\text{Var}(x_2) = \sigma_2^2$  としました。 $\sigma_1^2$ 、 $\sigma_2^2$  の値をいくつか変えて多変量 Gauss 分布のグラフを描いてみると、次のようになります。



$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} 3 & 0 \\ 0 & 3 \end{pmatrix}$$



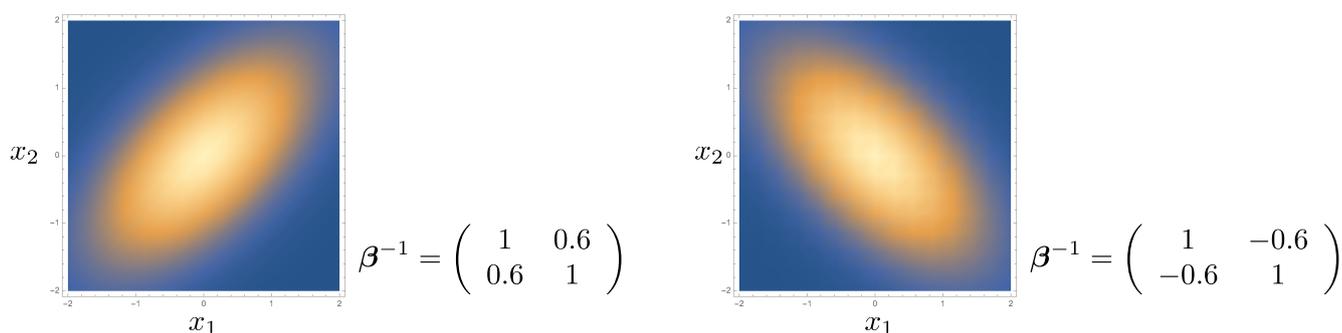
$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & 3 \end{pmatrix}$$

上の図では、色が薄い(白っぽい)ほど確率分布の値が大きく、濃い(青っぽい)ほど小さくなっています。1変数の場合と同様に、 $\sigma_1^2$ と $\sigma_2^2$ は、それぞれの変数方向の分布の広がり具合を変化させるパラメータとなっていることが分かります。 $\sigma_1^2$ を大きくすると、 $x_1$ 方向に分布が広がる様子が見て取れます。

次に、 $x_1$ と $x_2$ が相関を持っており、 $\text{Cov}(x_1, x_2)$ が0でない場合を考えましょう。例として、分散共分散行列が

$$\boldsymbol{\beta}^{-1} = \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix}, \quad \boldsymbol{\beta}^{-1} = \begin{pmatrix} 1 & -0.6 \\ -0.6 & 1 \end{pmatrix}$$

の2つの場合の多変量 Gauss 分布のグラフを描いてみると、次のようになります。



$\text{Cov}(x_1, x_2) = 0$ の場合とは異なり、分布が歪んで $\text{Cov}(x_1, x_2) > 0$ の場合には正の相関、 $\text{Cov}(x_1, x_2) < 0$ の場合には負の相関を反映した分布となっていることが分かります。このように、多変量 Gauss 分布では精度行列(分散共分散行列)の値を調整することで、分布の広がり具合や歪み具合を調整できるようになっています。