

卒研ゼミ「深層学習」

3 確率と情報理論

3.2 確率変数

確率変数 (random variable) とは起こりうる事象に対応した、ランダムな値を取れる変数のことである*1。確率変数を x のようにローマン体の小文字で書き*2、 x がとりうる値は x のようにイタリック体で書く。確率変数がベクトル値のときはそれらは太字で \mathbf{x} , \mathbf{x} のように書く。確率変数は離散値でも連続値でもよい。

3.3 確率分布

3.3.1 離散変数と確率質量関数

確率質量関数 (probability mass function, PMF) とは、離散型確率変数に対してその値をとるときの確率に対応させた関数のことであり、大文字 P で書き表す。たとえば確率 p で表 ($x = 1$)、確率 $1 - p$ で裏 ($x = 0$) が出るようなコインでコイントスをするとき、確率質量関数は

$$P(x = x) = \begin{cases} p, & (x = 1) \\ 1 - p & (x = 0) \end{cases}$$

とかける (これはベルヌーイ分布とよばれる)。このコイントスを n 回行ったうち k 回だけ表が出る確率は

$$P(y = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

となる (これは二項分布とよばれる)。このとき「確率変数 y は母数 (n, p) の二項分布 $B(n, p)$ に従う」とい、これを $y \sim B(n, p)$ と表記する。

$x = x$ かつ $y = y$ であるときの確率 $P(x = x, y = y) = P(x, y)$ は同時確率分布とよぶ。

関数 P が離散型確率変数 x の確率質量関数であるためには以下の性質を満たさなければならない。

- 定義域は確率変数のとりうる値すべての集合である。
- その集合を A とすると、 $\forall x \in A (0 \leq P(x) \leq 1)$ 。
- 正規化されている (normalized) : $\sum_{x \in A} P(x) = 1$ 。

たとえば、確率変数 x が k 個の異なる離散値 x_1, x_2, \dots, x_k をとるとき、確率質量関数を

$$P(x = x_i) = \frac{1}{k}$$

*1 たとえばサイコロを振ったときの「6の目が出る」という事象は $x = 6$ に対応させる、コイントスをしたときの「裏が出る」「表が出る」という事象をそれぞれ $x = 0$, $x = 1$ に対応させるなど。このように確率変数は、事象の集合 (標本空間) から数の集合への関数と考えることができる。

*2 多くの本では X のようにイタリック体の大文字で書かれることが多いようです。

とすれば一様分布を定義することができる。すべての i で和をとれば

$$\sum_{i=1}^k P(x = x_i) = \frac{1}{k} \sum_{i=1}^k 1 = \frac{1}{k} \cdot k = 1$$

となって、正規化されていることがわかる。

3.3.2 連続変数と確率密度関数

離散型確率変数に対して確率質量関数とよんでいたものを、連続型確率変数に対しては確率密度関数 (probability density function, PDF) とよぶ。確率密度関数 p は以下の性質を満たさなければならない。

- 定義域は確率変数のとりうる値すべての集合である。
- その集合を I とすると、 $\forall x \in I (0 \leq p(x))$ 。ただし $p(x) \leq 1$ である必要はない。
- $\int_I p(x) dx = 1$ 。

$x = x$ であるときの確率は $p(x)$ と直接得られるわけではなく、代わりに x が微小区間 $[x, x + \delta x]$ の値をとるときの確率が $p(x)\delta x$ で与えられる。区間 $[a, b]$ に x が存在する確率は $\int_{[a,b]} p(x) dx$ で求められる。

確率密度関数の例として、一様分布がある。

$$u(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b], \\ 0, & x \notin [a, b]. \end{cases}$$

積分すれば 1 になる。 x が区間 $[a, b]$ において一様分布にしたがうことを、 $x \sim U(a, b)$ と表す。