

目 次

内容

1. 確率・統計.....	2
1. 1 ベルヌーイ分布 1	2
1. 2 ベルヌーイ分布 2	3
1. 3 マルチヌーイ分布(カテゴリカル分布)	4
1. 4 一変量正規分布.....	6
1. 5 条件付き確率 (ベイズの定理)	8
1. 6 母集団	9

1. 確率・統計

1. 1 ベルヌーイ分布 1

確率 p で表が出て、確率 $1 - p$ で裏が出るコイン投げを考える。
表が出た時に $X=1$ をとり、裏が出た時に $X=0$ を取るとする。

1) $X=x$ (ただし、 $x=0,1$)となる確率

$p^x (1-p)^{1-x}$ となる。

これをベルヌーイ分布という。

2) ベルヌーイ分布に従う確率変数の期待値

$$E[X] = \sum_{x=0}^1 x p^x (1-p)^{1-x} = p$$

3) 分散

$$\text{Var}[X] = E[X^2] - E[X]^2 = \sum_{x=0}^1 x^2 p^x (1-p)^{1-x} - p^2 = p \cdot p^2 = p(1-p)$$

■ベルヌーイ分布の確率関数は、 $P(X=x) = p^x (1-p)^{1-x}$

で与えられ、その期待値は p 、分散は $p(1-p)$ である。

■離散確率変数 X に対して、 $P(X=x)=f(x)$ である時、 $f(x)$ を確率関数（又は確率分布）と呼ぶ。

この時 X の期待値 $E[X]$ は

$$E[X] = \sum f(x) \text{ で定義され}$$

分散 $\text{Var}[X]$ は

$$\text{Var}[X] = E[(X - E[X])^2] \text{ で定義される}$$

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

1. 2 ベルヌーイ分布 2

{0,1}を取りうる 2 値データ $D=\{x_1, \dots, x_n\}$ がベルヌーイ分布 $f(x;p)=p^x \cdot (1-p)^{1-x}$ に独立に従うと仮定する。
このとき最尤法によりパラメータ p を決定する。

1) 尤度関数

$$LD(p) = \prod_{i=1}^n f(x_i; p) = \prod_{i=1}^n p^{x_i} \cdot (1-p)^{1-x_i}$$

2) 負の対数尤度

$$-\log LD(p) = -\log \prod_{i=1}^n f(x_i; p) = -\sum_{i=1}^n \log f(x_i; p) = -\sum_{i=1}^n \log p^{x_i} \cdot (1-p)^{1-x_i} = -\sum_{i=1}^n (x_i \log p + (1-x_i) \log(1-p))$$

この式は、2 クラス分類のニューラルネットワークの学習に適応されることの多い損失関数

3) p の最尤推定量

負の対数尤度の最小値が満たすべき必要条件は、

$$\frac{d}{dp}(-\log LD(p)) = 0$$

の解である。

$-\log p$ と $-\log(1-p)$ は凸関数なので、 $-\log LD(p)$ も凸関数であり、解は唯一で最小解である。

$$\begin{aligned} \frac{d}{dp}(-\log LD(p)) &= -\sum_{i=1}^n \left(x_i \frac{d}{dp} \log p + (1-x_i) \frac{d}{dp} \log(1-p) \right) = -\sum_{i=1}^n \left(x_i \cdot \frac{1}{p} - (1-x_i) \cdot \frac{1}{1-p} \right) = -\sum_{i=1}^n \frac{x_i - p}{p(1-p)} \\ &= -\frac{1}{p(1-p)} \left(\sum_{i=1}^n x_i - p \sum_{i=1}^n 1 \right) = \frac{1}{p(1-p)} \left(np - \sum_{i=1}^n x_i \right) \\ &= \frac{1}{p(1-p)} \left(np - \sum_{i=1}^n x_i \right) = 0 \end{aligned}$$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{p} \text{ はデータによって決定されるパラメータ } p \text{ の推定量}$$

・最尤推定量

「尤度関数が最大になる（負の対数尤度関数が最小になる）」ように決められる「確率分布がデータに最もよく当てはまる」ようなパラメータの推定量

■パラメータ θ によって定まる確率（密度）関数 $f(x; \theta)$ に対して、それに独立に従うと仮定されるデータ $D=\{x_1, x_2, \dots, x_n\}$ が与えられたとき、尤度関数は

$$LD(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

で定義される。

最尤推定量とは尤度関数を最大化、すなわち負の対数尤度関数 $-\log LD(\theta)$ を最小化するパラメータの推定量

$$\hat{\theta} = \theta$$

■ベルヌーイ分布 $f(x;p)=p \cdot (1-p)$ のパラメータ p の最尤推定量は、データの平均、すなわち「1 が出現する確率」で与えられる。

1. 3 マルチヌーイ分布(カテゴリカル分布)

1) ワンホットベクトル

ベクトル x の成分のうち、ただ 1 つが 1 であり、その他が全て 0 であるようなもの

2) マルチヌーイ分布

k 次元のワンホットベクトルで構成されるデータ

$$D=\{x_1, x_2, \dots, x_n\} \text{ がマルチヌーイ分布}$$

$$F(x;p) = \prod_{j=1}^k p_j^{x_j}$$

に従っていると仮定する。 x_i の j 成分を x_{ij} と書く。

ただし、

$$p=(p_1, p_2, \dots, p_k)^T, \quad \sum_{j=1}^k p_j = 1, \quad 0 \leq p_j \leq 1 (j=1, \dots, k)$$

出る目の確率が

$$p=(p_1, p_2, p_3, p_4, p_5, p_6)^T \quad (\text{ただし } \sum p_j = 1, \quad 0 \leq p_1, \dots, p_6 \leq 1)$$

で与えられる歪んだ 6 面サイコロについて、1 が出る事象を $(1, 0, 0, 0, 0, 0)^T$ 、2 が出る事象を $(0, 1, 0, 0, 0, 0)^T$...

6 が出る事象を $(0, 0, 0, 0, 0, 1)^T$ と、ワンホットベクトルの確率変数によって表現するすると、これに従う確率分布は

$$f(x;p) = \prod_{j=1}^6 p_j^{x_j}$$

と書ける。

マルチヌーイ分布は 3 値以上のカテゴリ変数に関するモデリングに用いるのに適した分布である。

ここで、尤度関数は確立関数にデータの各値を代入したものの積で書けるので、

この時尤度関数

$$LD(p) = \prod_{i=1}^n f(x_i; p) = \prod_{i=1}^n \prod_{j=1}^k p_j^{x_{ij}}$$

負の対数尤度

$$-\log LD(p) = -\log \prod_{i=1}^n f(x_i; p)$$

$$= -\sum \log \prod_{j=1}^k p_j^{x_{ij}}$$

$$= -\sum \sum \log p_j$$

$$= -\sum \sum x_{ij} \log p_j$$

この式は、分類問題のためのニューラルネットワークの学習に適応されることの多い損失関数である公差エントロピーである。

■ベクトル x の成分のうち、ただ 1 つが 1 であり、その他が全て 0 であるようなものをワンホットベクトルと呼ぶ。

■ワンホットベクトルの従う確率分布としてマルチヌーイ分布があり、その確率関数はパラメータベクトル p を用いて

$$F(x;p) = \prod p_j \quad (\text{ただし、} \sum p_j = 1, 0 \leq p_j \leq 1, j=1, \dots, k)$$

により与えられる。

■マルチヌーイ分布の最尤推定は、交差エントロピーの最小化に対応している。

■ マルチヌーイ分布の負の対数尤度

$-\log LD(\mathbf{p}) = -\sum_i \sum_j k_{ij} \log p_j$ を

制約 ($\sum_j p_j = 1, 0 \leq p_j \leq 1, j=1, \dots, k$) のもとで最小化する問題を解き、 \mathbf{p} の最尤推定量を $\hat{\mathbf{p}}$ と求めることができる。

求解にはラグランジェの未定乗数法を用いる。

つまり、ベルヌーイ分布と同様パラメータ \mathbf{p} の最尤推定量はデータの平均、すなわち「各次元において 1 が出現する頻度」で与えられる。

1. 4 一変量正規分布

平均が μ 、分散が 1 の 1 変量正規分布の確率密度関数は、

$$f(x;\mu)=1/\sqrt{2\pi}\exp(-1/2(x-\mu)^2)$$

よって与えられる。

この時、実数データ $D=\{x_1, x_2, \dots, x_n\}$ によって与えられ、データ D が平均 μ 、分散 1 の 1 変量正規分布に独立に従うとして、最尤推定によってパラメータ μ を求める。

1) 尤度関数

1 変量正規分布は、実数値確率変数が従う代表的な確率分布の 1 つで、左右対称なベルカーブを描く。

①. 確率密度関数

$$f(x;\mu, \sigma^2)=1/\sqrt{2\pi\sigma^2}\exp(-1/(2\sigma^2)*(x-\mu)^2)$$

で与えられる。

②. 分散 $\sigma^2=1$ で固定されているため、密度関数は

$$f(x;\mu)=1/\sqrt{2\pi}\exp(-1/2(x-\mu)^2)$$

③. 尤度関数

$$L(\mu)=\prod_{i=1}^n f(x_i;\mu)$$

$$=\prod_{i=1}^n (1/\sqrt{2\pi})\exp(-1/2(x_i-\mu)^2)$$

最尤推定では $L(\mu)$ の最大化問題を考えるが、ベルヌーイ分布と同じように、「負の対数尤度の最小化問題」に書き換える。

④. 負の対数尤度

$$-\log L(\mu)=-\log(\prod_{i=1}^n (1/\sqrt{2\pi})\exp(-1/2(x_i-\mu)^2))$$

$$=-\sum \log((1/\sqrt{2\pi})\exp(-1/2(x_i-\mu)^2))$$

$$=-\sum (\log(1/\sqrt{2\pi})-1/2(x_i-\mu)^2)$$

$$=-n\log(1/\sqrt{2\pi})+1/2\sum (x_i-\mu)^2$$

実際に最小化する関数は

$$g(\mu)=1/2\sum (x_i-\mu)^2$$

これは、回帰問題のモデルの学習に用いられる損失関数である、二乗和誤差に対応している。

ここで $g(\mu)$ は凸な二次関数の和なので、

$$\frac{d}{d\mu}g(\mu)=0$$

を解けば最小解が求められる。

$$\begin{aligned}\frac{d}{d\mu}g(\mu) &= 1/2\sum \frac{d}{d\mu}(x_i-\mu)^2 \\ &= 1/2\sum (-2(x_i-\mu)) \\ &= \sum \mu - \sum x_i\end{aligned}$$

$$=n\mu - \sum x_i$$

μ の最尤推定量、 $\hat{\mu}=1/n \sum x_i$ が得られる。

- ・「正規分布の最尤推定量は二乗和誤差の最小化」
- ・「正規分布の平均の最尤推定量はデータの平均」

■平均が μ 、分散が σ^2 である一変量正規分布は、次のようになる

$$f(x_i; \mu, \sigma^2) = 1/\sqrt{2\pi\sigma^2} \exp(-(1/2\sigma^2)(x_i - \mu)^2)$$

■正規分布をモデルとした最尤推定は、二乗和誤差の最小化に対応している。

■正規分布の平均の最尤推定量は、データの平均によって与えられる。

■平均ベクトルが μ 、分散共分散行列が Σ である多変量正規化分布の確率密度関数は、

$$f(\mathbf{x}; \mu, \Sigma) = 1/((\sqrt{2\pi})^n \sqrt{\det(\Sigma)}) \exp(-1/2(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu))$$

いま、分散が単位行列で与えられているとする（つまり、各変数間が無相関で、各変数の分散が1であることを意味する）。

このとき確率密度関数は、

$$f(\mathbf{x}; \mu, \Sigma) = 1/((\sqrt{2\pi})^n) \exp(-1/2(\mathbf{x} - \mu)^T (\mathbf{x} - \mu))$$

これについて、負の対数尤度を導出すると、一変数の場合と同様に多変量の場合の二乗和誤差を導出することができる。

1. 5 条件付き確率（ベイズの定理）

事象 A の起こる確率を $P(A)$ 、事象 B の起こる確率を $P(B)$ 、これらの同時確立を $P(A,B)$ と書く。

1) 条件付き確率の定義から $P(A|B)$ は

$$=P(A,B)/P(B)$$

2) $P(B|A)$

$$=P(A,B)/P(A)$$

3) ベイズの定理 $P(A|B)$

$$1) \text{ より } P(A,B)=P(A|B)P(B)$$

$$2) \text{ より } P(A,B)=P(B|A)P(A)$$

$$P(A|B)P(B)=P(B|A)P(A)$$

$$P(A|B)=(P(B|A)P(A))/P(B)$$

$$P(B|A)=(P(A|B)P(B))/P(A)$$

1. 6 母集団

母集団に属する人が疾患 X に罹患している確率を 0.010 とする。簡易検査薬 Y は、疾患 X に感染している人に適用した場合に確率 0.90 で陽性を示し、疾患 X に感染していない人に適用した場合に確率 0.10 に陽性を示すことが知られている。母集団に属する人のうち、ある 1 名 Z に対して簡易検査薬 Y を適用したところ、陽性を示した。この時、Z が疾患 X に罹患している確率は？

1) 母集団に属する人が疾患 X に罹患している確率 $P(\text{罹患})=0.010$

2) 罹患していない確率 $P(\text{非罹患})=0.99$

3) 疾患 X に罹患している人に対して簡易検査薬 Y が陽性を示す確率 $P(\text{陽性} | \text{罹患})=0.90$

4) 疾患 X に罹患していない人に対して簡易検査薬 Y が陽性を示す確率 $P(\text{陽性} | \text{非罹患})=0.10$

5) 疾患 X に「罹患していること」と「罹患していないこと」は排反事象

$$P(\text{陽性})=P(\text{陽性}, \text{罹患})+P(\text{陽性}, \text{非罹患})$$

$$=P(\text{陽性} | \text{罹患})P(\text{罹患})+P(\text{陽性} | \text{非罹患})P(\text{非罹患})$$

$$=0.90 \times 0.010 + 0.10 \times 0.99$$

$$P(\text{罹患} | \text{陽性})=P(\text{陽性} | \text{罹患})P(\text{罹患})/P(\text{陽性})$$

$$=0.90 \times 0.010 / (0.90 \times 0.010 + 0.10 \times 0.99)$$