第一章線形代数 関連記事

1. 線形代数

1. 1 固有值分解 1

$$A = \begin{bmatrix} 5 & 2 \\ 2 & 8 \end{bmatrix}$$
 対称行列の固有値分解公式 $A - \lambda$ $I = 0$ より固有値を求める

(1) 固有值

$$\begin{bmatrix} 5 - \lambda & 2 \\ 2 & 8 - \lambda \end{bmatrix} = (5 - \lambda)(8 - \lambda) - 2 + 2 = 40 - 13 + \lambda + \lambda + \lambda + \lambda + 4 = \lambda + \lambda - 13 + \lambda + 36 = 0$$

$$13 \pm \sqrt{(169 - 4 + 1 + 36)} \qquad 13 \pm \sqrt{25}$$

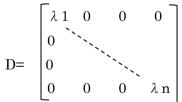
$$2 + 1 \qquad = 2 \qquad = 9,4$$

(2) 固有ベクトル

$$\begin{bmatrix} 5 & 2 \\ 2 & 8 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \end{bmatrix} = \begin{bmatrix} x1 \\ 2 & 8 \end{bmatrix} \begin{bmatrix}$$

$$D = \begin{bmatrix} 9 & 0 \\ 0 & 4 \end{bmatrix}$$

- ◎一般にn 次正方行列 A に対して $Av = \lambda v$ を満たすようなスカラー λ を固有値、非零ベクトルを固有ベクトル
- ◎固有値λは、Aの固有方程式 $det(\lambda A)=0$ の解となっている。
- ◎A の固有値 λ 1,...., λ n を対角成分に持つ対角行列



と、対応する固有ベクトルを列として並べた行列

に対いて、 $A=PDP^1$ が成立する。これを行列 A の「固有値分解」、又は「対角比」と呼ぶ。 さらに A が対称行列であるとき、対応する大きさ1の固有ベクトルを列として並べた行列 P は直行行列 $(P^{-1} = P^{T})$ になるため、 $A = PDP^{T}$ が成立する。

1. 2 固有值分解 2

$$A = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & -1 \\ 1 - 1 & 1 \end{bmatrix}$$
 対称行列の固有値分解

公式 $A - \lambda I = 0$ より固有値を求める

(1) 固有值

$$\begin{bmatrix} -\lambda & 0 & 1 \\ 0 - \lambda & -1 \\ 1 - 1 & 1 - \lambda \end{bmatrix} = -\lambda \begin{bmatrix} -\lambda & -1 \\ -1 & 1 - \lambda \end{bmatrix} = -\lambda \begin{bmatrix} 0 & 1 \\ -1 & 1 - \lambda \end{bmatrix} + 1 \begin{bmatrix} 0 & 1 \\ -\lambda & 1 \end{bmatrix} = -\lambda *((-\lambda + \lambda * \lambda) - (1)) + 1*(\lambda) = -\lambda *(\lambda * \lambda - \lambda - 1) + \lambda$$

$$= -\lambda * \lambda * \lambda + \lambda * \lambda + 2 * \lambda = 0$$

$$\lambda * \lambda * \lambda - \lambda * \lambda - 2 * \lambda = 0$$

$$\lambda (\lambda * \lambda - \lambda - 2) = \lambda ((\lambda + 1)(\lambda - 2)) = 0$$

$$\lambda = -1, 0, 2$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \\ x3 \end{bmatrix} = \begin{bmatrix} x1 \\ x2 \\ x3 \end{bmatrix}$$

(2) 固有値 $\lambda = 2$ に対する固有ベクトル $v=(1/\sqrt{6}, vy, 2/\sqrt{26})$

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & -1 \\ 1 - 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{6} \\ vy \\ 2/\sqrt{6} \end{bmatrix} = 2 \begin{bmatrix} 1/\sqrt{6} \\ vy \\ 2/\sqrt{6} \end{bmatrix}$$

$$\begin{bmatrix} 2/\sqrt{6} \\ -2/\sqrt{6} \\ 1/\sqrt{6} - vy + 2/\sqrt{6} \end{bmatrix} = \begin{bmatrix} 2/\sqrt{6} \\ 2vy \\ 4/\sqrt{6} \end{bmatrix}$$

$$2vy = -2/\sqrt{6} \quad vy = -1/\sqrt{6}$$

◎ 3 次正方行列の行列式は

これをサラスの公式を呼称する。

⑥高次方程式の解には因数分解を用いる。一般に高次方程式 $f(\mathbf{x})=0$ が $\mathbf{x}=\alpha$ を解に持つとき、 $f(\mathbf{x})$ は \mathbf{x} - α を因数に持つ。

2

1. 3 固有値分解3

$$\begin{bmatrix} 1 & 2 \\ A = \begin{bmatrix} 2 & 1 \end{bmatrix}$$
 対称行列の固有値分解

(1) 固有值

$$\begin{bmatrix} 1 - \lambda & 2 \\ 2 & 1 - \lambda \end{bmatrix} = (1 - \lambda)(1 - \lambda) - 4 = 1 - \lambda - \lambda + \lambda * \lambda - 4 = 0$$

 $\lambda * \lambda - 2* \lambda - 3 = (2 \mp \sqrt{4 - 4*1*(-3)})/2 = (2 \mp \sqrt{16})/2 = 3, -1$

固有値ベクトル

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \end{bmatrix} = 3 \begin{bmatrix} x1 \\ x2 \end{bmatrix}$$

$$X1+2x2=3x1 \quad 2x1=2x2 \qquad x1=x2 \qquad 1$$

$$2x1+x2=3x2 \quad 2x1=2x2 \qquad 1$$

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \end{bmatrix} = -1 \begin{bmatrix} x1 \\ x2 \end{bmatrix}$$

$$X1+2x2=-x1 \quad x1=-3x2 \qquad 1$$

$$2x1+x2=-x2 \quad x1=-x2 \qquad 1$$

1. 4 特異値分解 1

行列

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 1 \\ -1 & 1 & 0 \end{bmatrix}$$

を特異値分解し $A=U\Sigma V^T$ の形で表す

$$A A^{T} = \begin{bmatrix} 1 & 0 & 1 \\ -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 1+0+1 & -1+0+0 \\ -1+0+0 & 1+1+0 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix}$$

(1) 固有值分解

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 2 - \lambda & -1 \\ -1 & 2 - \lambda \end{bmatrix} = (2 - \lambda)(2 - \lambda) - 1 = 4 - 2\lambda - 2\lambda + \lambda * \lambda - 1 = \lambda * \lambda - 4\lambda + 3$$

 $\lambda = (4 \mp \sqrt{(16 - 4 * 3)})/2 = 3,1$

特異値は $\sqrt{3}$ 、 $\sqrt{1}$

⑥任意の零行列でない $\mathbf{m} \times \mathbf{n}$ 行列 \mathbf{A} に対して、 $\mathbf{A}\mathbf{v} = \sigma \mathbf{u}$ 、 $\mathbf{A}\mathbf{v} = \sigma \mathbf{v}$ を満たすような正の数 σ を特異値、 \mathbf{m} 次元ベクトル \mathbf{u} を左特異ベクトル、 \mathbf{n} 次元ベクトル \mathbf{v} を右特異ベクトルと呼ぶ。

また、定義式から $A^T_A u = {}^2_{\sigma} u$, $A^T_A Av = {}^2_{\sigma} v$ が得られるので、行列 $A^T_A e$ $A^T_A e$ の固有値ベクトルを求めることで、行列 A の特異値・(左右) 特異ベクトルを知ることができる。

このようにして求められた特異値を大きい順に(i,j)で成分に並べ、その他の成分を0で埋めた $m \times n$ 行列 Σ と、 左特異ベクトルを列として横に並べたm次正方行列U、右特異ベクトルを列として横に並べたn次正方行列Vを用いて $A=U^T$ Vと書ける。これをAの特異値分解と呼ぶ。なお特異値を降順に並べることにすると Σ は一意に定まるがU Σ びVは一意に定まらない。

1. 5 特異値分解 2

行列

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 2 \end{bmatrix}$$

を特異値分解し $A=U \Sigma V^T$ の形で表す

$$A A^{T} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1+0+0 & 0 \\ 0 & 0+1+4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix}$$

(1) 固有値分解

$$\begin{bmatrix} 1 - \lambda & 0 \\ 0 & 5 - \lambda \end{bmatrix} = (1 - \lambda)(5 - \lambda) = 0 \qquad \lambda = 1,5$$

特異値は $1,\sqrt{5}$

(2) 固有値ベクトル

$$\begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \end{bmatrix} = \begin{bmatrix} x1 \\ x2 \end{bmatrix}$$

X1=5*x1

5*x2=5*x2

$$\begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} x1 \\ x2 \end{bmatrix} = \begin{bmatrix} x1 \\ 1 & x2 \end{bmatrix}$$

X1=1*x1

5*x2=1*x2

(3)

A の特異値は $σ=1,\sqrt{5}$ なので、これを降順に並べれば

$$\Sigma = \begin{bmatrix} \sqrt{5} & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Uの第一列は AA^T の固有値5に対応する固有ベクトル $u=(ux,uy)^T$ のうち、大きさが1のもの

$$\begin{bmatrix} 1 & 0 & ux \\ 0 & 5 & uy \end{bmatrix} = \begin{bmatrix} ux \\ uy \end{bmatrix}$$

ux=5ux

ux=0,uy は任意の実数

5uy=5uy

V の各列は $\overset{T}{AA}$ の固有ベクトルだが、これらは正規直行性を満たす必要がある。 v1=(0,1/ $\sqrt{5}$,2/ $\sqrt{5}$),v3=(0,-2/ $\sqrt{5}$,vz)と置けばそれらの直交性から v^Tv3 =0 になる。 $0*0+1/\sqrt{5}*-2/\sqrt{5}+2/\sqrt{5}*vz$ =0

 $Vz=1/\sqrt{5}$

1. 6 L1ノルム、L2ノルム

(1) L1ノルム

n

$$x 1 = \Sigma |xi|$$

i=1

(2) L2ノルム

$$\begin{array}{c} x \ 2 = \sqrt{\sum_{i=1}^{n} \sum_{i=1}^{n}} \end{array}$$

(3) 最大値ノルム

$$x \infty = \max |xi|$$

Ι

1. 7 距離

(1) マンハッタン距離

n

$$d(x,y)=\Sigma |xi-yi|$$

i=1

(2) ユークリッド距離

$$d(x,y) = \sqrt{\frac{n}{\sum (xi-yi)2}}$$

$$i=1$$

(3) マハラノビス距離

$$D(x,y) = \sqrt{\sum (x-y)} T \sum -1 \, \not \equiv (x-y)$$

第二章 確率・統計 関連記事

- 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

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

(3) 分散

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

- ©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

 $E[X]=\sum f(x)$ で定義され

分散 Var[X]は

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

1. 2 ベルヌーイ分布2

 $\{0,1\}$ を取りうる 2 値データ $D=\{x1,...,xn\}$ がベルヌーイ分布 f(x;p)=px 乗(1-p)(1-x)乗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) =
$$\prod_{i=1}^{n} f(xi;p) = \prod_{i=1}^{n} pxi \# (1-p)(1-xi) \#$$

(2) 負の対数尤度

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、 $\frac{d}{dp}$ (-logLD(p))=0

の解である。

$$\begin{array}{l} \text{-logp} \ \& \text{-log}(1\text{-p}) \\ \text{は凸関数なので、-logLD(P)} \\ \text{-logLD(p)} \\ \text{-} \\ \text$$

$$= -1/p(1-p)(\sum_{i=1}^{n} x_i - \sum_{i=1}^{n} 1) = 1/p(1-p)(np - \sum_{i=1}^{n} x_i)$$

$$1/p(1-p)(np-\sum_{i=1}^{n}xi)=0$$

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

· 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」可能性がある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 $f(\mathbf{x}; \theta)$ に対して、それに独立に従うと仮定されるデータ $\mathbf{D}=\{\mathbf{x}1,\mathbf{x}2,...,\mathbf{x}n\}$ 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 $\theta=\stackrel{\wedge}{\theta}$

◎ベルヌーイ分布 f(x;p)=p (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

 $D={x1,x2,...,xn}$ がマルチヌーイ分布

$$F(x;p) = \prod_{i=1}^{k} p_i^{x_i^{i}}$$

j=1 に従っていると仮定する。Xi の j 成分を xij と書く。

ただし、

$$p=(p1,p2,...,pk)^{T}$$
 $\sum_{j=1}^{k} pj=1$ $0 \le pj \le 1(j=1,...,k)$

出る目の確率が

 $p=(p1,p2,p3,p4,p5,p6)^{T}$ (ただし $\Sigma pj=1$, $0 \le p1,....,p6 \le 1$)

 $f(x;p) = \prod_{i=1}^{6} p_i^{x_i^i}$

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

$$LD(p)=\prod_{\substack{i=1\\j=1}}^nf(xi;p)=\prod_{\substack{i=1\\j=1}}^n\prod_{j=1}^kxij$$

負の対数尤度

 $-\log LD(p) = -\log \prod f(xi;p)$

 $= - \sum \log \prod pj^{xij}$

=- $\Sigma \Sigma \log pj$

 $= - \sum \sum xijlogpj$

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 \mathbf{x} の成分のうち、ただ $\mathbf{1}$ つが $\mathbf{1}$ であり、その他が全て $\mathbf{0}$ 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)= Π pj (ただし、 Σ pj=1,0 \leq pj \leq 1、j=1,...,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

-logLD(p)=- Σ Σ kijlogpj $\stackrel{}{\mathcal{E}}$

制約 $(\Sigma pj=1,0 \le pj \le 1, j=1,...,k)$ のもとで最小化する問題を解き、p の最尤推定量を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 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 $D={x1,x2,...,xn}$ によって与えられ、データ D が平均 u、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;μ,
$$\sigma^2$$
)=1/ $\sqrt{2}$ π σ^2 exp(-1/(2 σ^2)*(x-μ)²) で与えられる。

②. 分散 σ^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

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

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

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

最尤推定では L(μ)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

$$-\log L(\mu) = -\log(\Pi(1/\sqrt{2}\pi) \exp(-1/2(xi-\mu)^{2}))$$

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

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

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

実際に最小化する関数は

$$g(\mu)=1/2 \Sigma (xi-\mu)^2$$

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

を解けば最小解が求められる。

$$\frac{d}{d\mu}g(\mu)=1/2 \sum_{\substack{d \\ d\mu}} (xi \cdot \mu)^2$$

$$=1/2 \sum (-2(xi \cdot \mu))$$

$$= \sum \mu \cdot \sum xi$$

$$=n\mu \cdot \sum xi$$

μ の最尤推定量、μ=1/nΣxi 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 \mathfrak{m}_{μ} 、分散 \mathfrak{m}_{σ}^2 である一変量正規分布は、次のようになる

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

- ◎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 ∑である多変量正規化分布の確率密度関数は、

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

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、

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

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) より P(A,B)=P(A|B)P(B)
 - (2) $\sharp V P(A,B)=P(B|A)P(A)$
 - $P(A \mid B)P(B)=P(B \mid A)P(A)$
 - $P(A \mid B) = (P(B \mid A)P(A))/P(B)$
 - $P(B \mid A) = (P(A \mid B)P(B))/P(A)$

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(罹患)=0.010
- (2) 罹患していない確率 P(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(陽性 |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(陽性 | 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(陽性)=P(陽性,罹患)+P(陽性,非罹患)

=P(陽性 | 罹患)P(罹患)+P(陽性 | 非罹患)P(非罹患)

 $=0.90\times0.010+0.10\times0.99$

P(罹患 | 陽性)=P(陽性 | 罹患)P(罹患)/P(陽性)

 $=0.90\times0.010/(0.90\times0.010+0.10\times0.99)$

第三章 情報理論 関連記事

1. 情報理論

1. 1 情報量

事象 A が起こる確率を P(A)とする。この時、事象 A が起こることの自己情報量は、

 $I(A) = -\log_2 P(A)$

- (1) 小さな確率の事象が起こることを、大きな情報量で表現したい 式が単調減少であること
- (2) 複数の事象が発生する確率は、積で表現されるが、情報量においては和で表現したい 事象 A,B が独立である時

$$\begin{split} \text{I}(\mathbf{A} \cap \mathbf{B}) &= \text{-log}_2 \ \text{P}(\mathbf{A} \cap \mathbf{B}) \\ &= \text{-log}_2 \ \text{P}(\mathbf{A}) \text{P}(\mathbf{B}) \\ &= \text{-log}_2 \ \text{P}(\mathbf{A}) \text{-log}_2 \ \text{P}(\mathbf{B}) \\ &= \text{I}(\mathbf{A}) \text{+I}(\mathbf{B}) \end{split}$$

確率で起きる事象の情報量は $-\log_2$ 1=0

②事象 A の起こる確率が P(A)の時、事象 A が起こることの情報量は $-\log_{9}$ P(A)

1. 2 エントロピー

離散確率変数 X において、X=x となる確率が p(x)で与えられたとする。

確率変数 X のエントロピーは?

$$H(X) = \sum_{x} p(x) \log_2 p(x)$$

エントロピー(平均情報量)は、情報量(すなわち、「事象の起こりにくさ・珍しさ」)の期待値で与えられる。 そのため、確率変数のランダム性の指標として用いられる。

◎事象の集合 Ω について A $∈ \Omega$ が起こる確率を P(A)とすると、P は確率分布であるとみなせる。確率分布 P のエントロピーは

$$H(A) = -\sum_{A} P(A) \log_2 P(A)$$

1. 3 交差エントロピー(クロスエントロピー)

2つの確率分布 p(x)と q(x)の交差エントロピーは

$$H(p,q) = -\sum_{X} p(x) \log_2 q(x)$$

分類問題を解くための損失関数として用いられる。

2つの確率分布が全く同じときに交差エントロピーが最小になる。 ここで p がデータによって近似される真の分布であり、 q がモデルの分布である。

1. 4 KL ダイバージェンス

2つの確率分布 p(x)と q(x)に対して

 $D(p \mid |q) = \sum_{x} p(x) \log_{2} (p(x)/q(x))$

KL(カルバック・ライブラー)ダイバージェンスと呼ぶ

KL ダイバージェンスは、2つの確率分布の近さを表現する最も基本的な量で、統計学・情報理論において非常に重要な役割を果たすものである。

- ・pのエントロピーを H(p)
- ・pとqの交差エントロピーを H(p,q)
- ・ $p \ge q O KL ダイバージェンスを D_{KT}(p||q)$

 $H(p) + D_{KL}(p \mid q) = -\sum p(x)log_{2} p(x) + \sum p(x)log_{2} (p(x)/q(x))$

- $= -\sum p(x)\log_2 p(x) + \sum p(x)\log_2 p(x) \sum p(x)\log_2 q(x)$
- $= -\sum p(x)\log_2 q(x)$
- =H(p,q)

p と q が全く同じ分布である時に交差エントロピーは最小になり、同時に KL ダイバージェンスは最小になる。 このとき

 $\begin{aligned} \mathbf{D}_{\mathrm{KL}}(\mathbf{p} \mid | \mathbf{p}) &= \sum \mathbf{p}(\mathbf{x}) \log_2 \ (\mathbf{p}(\mathbf{x})/\mathbf{p}(\mathbf{x})) \\ &= \sum \mathbf{p}(\mathbf{x}) \log_2 \ 1 \\ &= 0 \end{aligned}$

◎ 2 つの確率分布 p(x)と q(x)に対して JS ダイバージェンスは次のようになる

 D_{JS} (p | | p)=1/2(Σ p(x)log $_2$ (p(x)/r(x)))+ Σ q(x)log $_2$ (q(x)/r(x)) ただし

r(x) = (p(x) + q(x))/2

JS ダイバージェンスは KL ダイバージェンスとは異なり、 2 つの分布に対して対称な量である。 JS ダイバージェンスは敵対的ネットワークの損失関数に用いられる。

◎一般化 KL ダイバージェンスや IS ダイバージェンスは、正規化されていない(合計が1にならない)ような分布同士の近さを測る量である。

どちらも非負値行列因子分解の損失関数として、音響信号処理の領域でよく用いられる。

1.5 連続型確率分布

データ $D=\{x1,x2,...,xn\}$ が、p(x)を確率密度関数とする連続型確率分布に独立に従っている。 この時、モデル $q(x;\theta)$ によって、交差エントロピーが最小になるように p(x)を推定することを考える。

(1) P(x)と $q(x;\theta)$ の交差エントロピーは

連続型確率分布 p(x)と $q(x;\theta)$ の交差エントロピーは

$$H(p,q)=-\int_{X} p(x)\log q(x;\theta)dx$$

エントロピー、KLダイバージェンスについても同様の拡張が可能である。

(2) 真の分布 p(x)での期待値をデータ D による平均に置き換えた量は モンテカルロ積分

$$\widetilde{H}(p,q)=-1/n\sum_{i=1}^{n}\log q(ki;\theta)$$

これに置き換えることで、実際にはわからない真の分布 p(x)を含むような計算を回避している。

(3) 尤度関数は \mathbf{L} (θ)= $\Pi \mathbf{q}(\mathbf{x},\theta)$ であるから、等式が成り立つ

$$\widetilde{H}(p,q)=-1/n \sum \log q(xi; \theta)$$

$$=-1/n \log \Pi q(xi; \theta)$$

$$=-1/n \log L_D(\theta)$$

交差エントロピーが最小となる推定は、「負の対数尤度が最小となる尤度」すなわち、最尤推定と等価である。

◎連続型確率分布 p のエントロピーは

$$H(p) = \int p(x) \log p(x) dx$$

◎連続型確率分布 p と q の交差エントロピー

$$H(p,q)=-\int p(x)\log q(x)dx$$

◎連続型確率分布 p と q の KL ダイバージェンス

$$D_{KL}(p \mid q) = -\int p(x) \log(q(x)/p(x)) dx$$

- ②定義域[a,b]の実数値確率変数 X において、X=x となる確率密度関数を p(x)と書くとき、関数 f(X)の期待値は $E[f(X)]=\int\limits_{-b}^{b}p(x)f(x)dx$
- の確率変数 X の観測として独立なデータ $D=\{x1,x2,...,xn\}$ が与えられたとき、E[f(X)]はモンテカルロ積によって近似できる

$$E[f(X)] \stackrel{\text{n}}{=} 1/n \sum_{i=1}^{n} f(xi)$$

1. 6 マルチヌーイ分布

k 次元ワンホットベクトルが従う確率分布 p(x)を、マルチヌーイ分布 $q(x;\mu)=\prod_{j=1}^{n}\mu_j^{x_j}$ によって推定することを考える。

ただし、 \mathbf{x} , μ の \mathbf{j} 成分を \mathbf{x} \mathbf{j} , μ \mathbf{j} とかく。 この時 $\mathbf{p}(\mathbf{x})$ と $\mathbf{q}(\mathbf{x};\mu)$ の交差エントロピーは

$$\begin{split} H(p,q) &= \sum_{x} p(x) log q(x;\mu) \\ &= -\sum_{x} p(x) log \prod_{j=1}^{k} \mu j^{xj} \\ &= -\sum_{j=1}^{k} p(x) \sum_{j=1}^{k} log \mu j \\ &= -\sum_{j=1}^{k} p(x) \sum_{j=1}^{k} x j log \mu j \end{split}$$

モンテカルロ積分を用いれば

$$H(p,q) = -1/n \sum_{i=1}^{n} \sum_{j=1}^{k} x_{i} i_{j} \log \mu_{j}$$

ここでμをロジスティック回帰やニューラルネットワークの出力、xiを正解ラベルと置き換えれば、交差エントロピーは分類問題において典型的な損失関数となる。