

内容

1. 情報理論.....	2
1. 1 情報量.....	2
1. 2 エントロピー.....	2
1. 3 交差エントロピー（クロスエントロピー）	2
1. 4 KL ダイバージェンス.....	3
1. 5 連続型確率分布.....	4
1. 6 マルチヌーイ分布.....	5

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

$$\begin{aligned} I(A \cap B) &= -\log_2 P(A \cap B) \\ &= -\log_2 P(A)P(B) \\ &= -\log_2 P(A) - \log_2 P(B) \\ &= I(A) + I(B) \end{aligned}$$

確率で起きる事象の情報量は $-\log_2 1 = 0$

■ 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 $-\log_2 P(A)$

■

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 $A \in \Omega$ 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

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

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

$$D(p||q) = \sum_x p(x) \log_2 (p(x)/q(x))$$

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 と q の KL ダイバージェンスを $D_{KL}(p||q)$

$$\begin{aligned} H(p) + D_{KL}(p||q) &= -\sum p(x) \log_2 p(x) + \sum p(x) \log_2 (p(x)/q(x)) \\ &= -\sum p(x) \log_2 p(x) + \sum p(x) \log_2 p(x) - \sum p(x) \log_2 q(x) \\ &= -\sum p(x) \log_2 q(x) \\ &= H(p,q) \end{aligned}$$

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

$$\begin{aligned} D_{KL}(p||p) &= \sum p(x) \log_2 (p(x)/p(x)) \\ &= \sum p(x) \log_2 1 \\ &= 0 \end{aligned}$$

■ 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

$$D_{JS}(p||p) = 1/2 (\sum p(x) \log_2 (p(x)/r(x))) + \sum q(x) \log_2 (q(x)/r(x))$$

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 $D=\{x_1, x_2, \dots, x_n\}$ 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

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

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

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

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 $L_D(\theta) = \prod q(x, \theta)$ であるから、等式が成り立つ

$$\begin{aligned} \tilde{H}(p, q) &= -1/n \sum \log q(x_i; \theta) \\ &= -1/n \log \prod q(x_i; \theta) \\ &= -1/n \log L_D(\theta) \end{aligned}$$

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

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

■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

$$E[f(X)] = \int_a^b p(x) f(x) dx$$

■確率変数 X の観測として独立なデータ $D=\{x_1, x_2, \dots, x_n\}$ 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

$$E[f(X)] \doteq 1/n \sum_{i=1}^n f(x_i)$$

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 $p(\mathbf{x})$ を、マルチヌーイ分布

$$q(\mathbf{x}; \mu) = \prod_{j=1}^k \mu_j^{x_j}$$

によって推定することを考える。

ただし、 \mathbf{x}, μ の j 成分を x_j, μ_j とかく。

この時 $p(\mathbf{x})$ と $q(\mathbf{x}; \mu)$ の交差エントロピーは

$$H(p, q) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x}; \mu)$$

$$= - \sum p(\mathbf{x}) \log \prod_{j=1}^k \mu_j^{x_j}$$

$$= - \sum p(\mathbf{x}) \sum_{j=1}^k \log \mu_j$$

$$= - \sum p(\mathbf{x}) \sum x_j \log \mu_j$$

モンテカルロ積分を用いれば

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

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、 x_i 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。