

パターン認識と機械学習

1.6.1 相対エントロピーと 相互情報量

吉田拓夢

相対エントロピー

- 未知の分布 $p(x)$ ← 近似的に $q(x)$ でモデル化
- 真の分布 $p(x)$ の代わりに $q(x)$ で x を符号化
- x の値を特定するために必要な追加情報量の平均は

$$\begin{aligned} \text{KL}(p \parallel q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \end{aligned}$$

↑ 分布 $p(x)$ と $q(x)$ の間の相対エントロピー

相対エントロピー

$$\begin{aligned}\text{KL}(p \parallel q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x}\end{aligned}$$

↑相対エントロピー(relative entropy)

別称: KL(カルバック-ライブラー)ダイバージェンス
(Kullback-Leibler divergence)

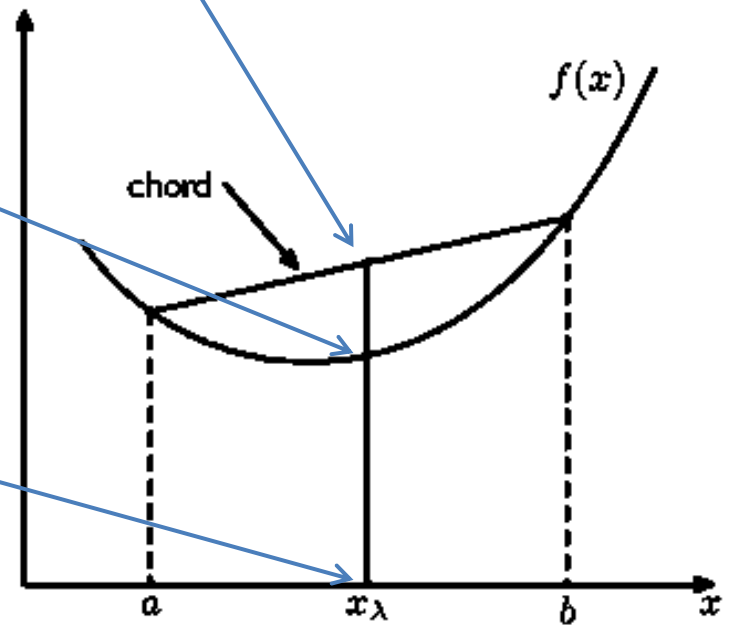
- 対称な量ではない $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$
- 右を満たす $\text{KL}(p \parallel q) \geq 0$

凸関数

- 凸性は次式で表現される

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$$

$$\lambda a + (1 - \lambda)b \quad 0 \leq \lambda \leq 1$$



- 凸関数(convex function)

$f(x)$ は全ての弦(chord)が関数に乗っているか
それより上にあるとき凸である

凸関数

- 凸関数の例: $x \ln x (x > 0)$, x^2
- $\lambda=0$ と $\lambda=1$ だけで成立 \rightarrow 真に凸(strictly convex)
- 凸性と逆の性質 \rightarrow 凹関数(concave function)
- $f(x)$ が凸関数 $\rightarrow -f(x)$ は凹関数

イェンセンの不等式

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad \lambda_i \geq 0 \quad \sum_i \lambda_i = 1$$

↑凸関数 $f(x)$ は任意の点集合 $\{x_i\}$ に対して上を満たす

λ_i を値 $\{x_i\}$ を取る離散確率変数 x 上の確率分布とすると

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)]$$

連続変数に対しては

$$f\left(\int \mathbf{x}p(\mathbf{x})d\mathbf{x}\right) \leq \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

イエンセンの不等式

- KLダイバージェンスに適用すると

($-\ln x$ が凸関数であること、規格化条件 $\int q(\mathbf{x})d\mathbf{x} = 1$ を使って)

$$\text{KL}(p \parallel q) = - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \geq -\ln \int q(\mathbf{x})d\mathbf{x} = 0$$

- 等式は全ての x について $p(x)=q(x)$ のときのみ成り立つ
- KLダイバージェンスを2つの分布 $p(x)$ と $q(x)$ の隔たりの尺度として扱える

データ圧縮と密度推定

- 最も効率的な圧縮 ← 真の分布が分かる
- 真の分布でない分布 → 非効率な符号化
追加情報量の平均 ≥ 2 分布間のKLダイバージェンス

データ圧縮と密度推定

- データが未知の分布 $p(x)$ から生成されるとき
それをモデル化

→可変パラメータ θ をもつパラメトリックな分布 $q(x|\theta)$ で近似

→ θ の決定: $p(x)$ と $q(x|\theta)$ 間のKLダイバージェンスを θ で最小化
 $p(x)$ から得られた有限個の訓練点の集合 $X_n(n=1, \dots, N)$ で近似

$$\text{KL}(p \parallel q) \simeq \frac{1}{N} \sum_{n=1}^N \{-\ln q(\mathbf{x}_n | \theta) + \ln p(\mathbf{x}_n)\}$$



訓練集合で評価した分布 $q(x|\theta)$ 下での θ の負の対数尤度

→KLダイバージェンスの最小化=尤度の最大化

相互情報量

- 2つの変数集合 x, y の同時分布 $p(x, y)$
→独立なら周辺分布の積に分解 $p(x, y) = p(x)p(y)$
- 独立でない→どれぐらい独立に「近い」か？
→同時分布と周辺分布の積の間のKLダイバージェンスを考える

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \parallel p(\mathbf{x})p(\mathbf{y})) \\ &= - \int \int p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x}d\mathbf{y} \end{aligned}$$

KLダイバージェンスの性質から $I(\mathbf{x}, \mathbf{y}) \geq 0$

相互情報量

- 相互情報量は条件付きエントロピーと関係

(確率の加法・乗法定理より)

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

- y の値を知ることで x に関する不確実性がどれだけ減少するかを表す
- $p(x)$: x の事前分布
 $p(x|y)$: 新たなデータ y を観測した後の事後分布
→ 新たな y を観測した結果として x に関する不確実性が減少した度合い