

信息论

信号传输与处理的理论基础

Markov序列和数据处理不等式
连续随机变量的信息熵

教程阅读：2.8和2.10节；8.1, 8.4, 8.5和8.6节



互信息量的凸性和凹性(1)

- * 将条件概率公式 $p(x,y)=p(y|x)p(x)$ 代入互信息量的定义式得
- * $I(X;Y) \equiv \sum_{x,y} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right) = \sum_{x,y} p(y|x)p(x) \log\left(\frac{p(y|x)}{p(y)}\right)$
- * 其中 $p(y) = \sum_x p(y|x)p(x)$ 。

* 针对 $I(X;Y)$ 的两种函数观点

- * (1) 将条件概率参数 $p(y|x)$ 视为常数、 $p(x)$ 视为变量，相应的函数记为 **C**;
- * (2) 将条件概率参数 $p(y|x)$ 视为变量、 $p(x)$ 视为常数，相应的函数记为 **D**;

*

C是凹函数；D是凸函数

- * 注：本讲义在下面所给的证明方法与教程中的不同，请同时阅读教程中的证明，
- * 并比较两者的特点。

互信息量的凸性和凹性(2)

若将条件概率参数 $p(y|x)$ 视为常数、 $p(x)$ 视为变量，则

*
$$I(X;Y) = \sum_{x,y} p(y|x)p(x) \log\left(\frac{p(y|x)}{p(y)}\right)$$

* 是 $p(x)$ 的凹函数，其中 $p(y) = \sum_x p(y|x)p(x)$ 。

(1) $I(X;Y) = \sum_{x,y} p_{y|x} \log\left(\frac{p_{y|x}}{p_y}\right)$, 其中 $p_y = \sum_x p_{y|x} p_x$, 故 $\partial p_y / \partial p_x = p_{y|x}, \forall x, y$.

$$= \sum_{x,y} p_{y|x} \cdot p_x \log(p_{y|x}/p_y)$$

(1) $I(X;Y)$ 是 (p_x) 的凹函数:

$$\frac{\partial I}{\partial p_x} = \sum_y p_{y|x} \log\left(\frac{p_{y|x}}{p_y}\right) - \sum_{x',y} p_{y|x'} \cdot p_{x'} \cdot \frac{\partial \log p_y}{\partial p_x}$$

$$= \sum_y p_{y|x} \log p_{y|x} - \sum_{y|x} p_{y|x} \log p_y - \sum_{x',y} p_{y|x'} \cdot p_{x'} \cdot \frac{1}{p_y} \frac{\partial p_y}{\partial p_{y|x}} = \frac{p_x}{p_y}$$

$$= \sum_y p_{y|x} \log p_{y|x} - \sum_{y|x} p_{y|x} \log p_y - \sum_{x',y} p_{y|x'} \cdot p_{x'} \cdot \frac{p_x}{p_y}$$

$$= \sum_y p_{y|x} \log p_{y|x} - \sum_{y|x} p_{y|x} \log p_y - p_x \sum_y \frac{1}{p_y} \cdot p_y = p_x$$

与 (p_x) 无关常数

$$\frac{\partial^2 I}{\partial p_x \partial p_{x'}} = - \sum_y p_{y|x} \cdot \frac{p_{y|x'}}{p_y} = -\delta_{xx'}$$

$$\sum_{x,x'} \delta_x \delta_{x'} \frac{\partial^2 I}{\partial p_x \partial p_{x'}} = - \sum_y \frac{1}{p_y} \left(\sum_x \delta_x \cdot p_{y|x} \right)^2 - \sum_x \delta_x^2 \leq 0 \text{ 且 } = 0 \text{ 当 } \delta_x = 0, \forall x$$

故 $I(X;Y)$ 是 (p_x) 的凹函数且严格凹. [注] 这里导出 $\frac{\partial I}{\partial p_x}$ 和 $\frac{\partial^2 I}{\partial p_x \partial p_{x'}}$ 表达式对计算 Shannon 容量数值算法有用.

(2) $I(X;Y)$ 是 $(p_{y|x})$ 的凹函数



互信息量的凸性和凹性(3)

若将条件概率参数 $p(y|x)$ 视为变量、 $p(x)$ 视为常量，则

* $I(X;Y) = \sum_{x,y} p(y|x)p(x) \log\left(\frac{p(y|x)}{p(y)}\right)$

* 是 $p(y|x)$ 的凸函数，其中 $p(y) = \sum_x p(y|x)p(x)$ 。

(2) $I(X;Y) = \sum_{x,y} p_{xy} \log\left(\frac{p_{xy}}{p_x p_y}\right) = \sum_{x,y} p_{y|x} p_x \log p_{y|x} - \sum_{x,y} p_{y|x} p_x \log p_y$
 是 $(p_{y|x})$ 的凹函数。
 $p_y = \sum_x p_{y|x} p_x$, $\partial p_y / \partial p_{y|x} = p_x$, $\partial p_y / \partial p_{y'|x'} = 0$, $\forall y' \neq y$.
 计算 $\partial I / \partial p_{y|x} = p_x (1 + \log p_{y|x}) - p_x \log p_y - \sum_{x'} p_x p_{x'} p_{y|x'} / p_y$.
 $\partial^2 I / \partial^2 p_{y|x} \partial p_{y'|x'} = \frac{p_x \delta_{xx'} \delta_{yy'}}{p_{y|x}} - \frac{p_x p_{x'} \delta_{yy'}}{p_y} = \frac{p_x}{p_y} \delta_{xx'} \delta_{yy'} - \frac{p_x p_{x'} \delta_{yy'}}{p_y}$
 $\sum_{x,y,x',y'} \xi_{y|x} \xi_{y'|x'} \frac{\partial^2 I}{\partial p_{y|x} \partial p_{y'|x'}} = \sum_y \left[\sum_x \frac{p_x \xi_{y|x}^2}{p_{y|x}} - \frac{1}{p_y} \left(\sum_x \xi_{y|x} p_x \right)^2 \right] \geq 0$.

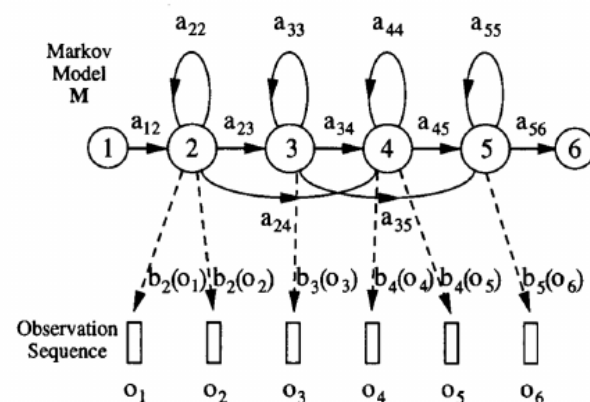
注：该凸性结论本课程今后用不到，仅用到关于凹性的结论。

习题（选做）：证明该不等式。





Markov序列(1)



* Markov序列

* 随机序列 $\dots, X_{n-1}, X_n, X_{n+1}, \dots$ 定义为Markov序列，是指

* 对任何 $n > n_1 > n_2 > \dots > n_p$ 恒有

$$P[X_n | X_{n_1}, X_{n_2}, \dots, X_{n_p}] = P[X_n | X_{n_1}]$$

* Markov性质的含义

* 基于任何历史数据 $(X_{n_1}, X_{n_2}, \dots, X_{n_p})$ 推断未来状态 X_n ，其效果总是等价于基于距离未来时刻最近的历史数据 X_{n_1} 推断 X_n 。

*

* 记号：在上述意义下的Markov序列，记为 $\dots \rightarrow X_{n-1} \rightarrow X_n \rightarrow X_{n+1} \rightarrow \dots$ 。



Markov序列(2)

* Markov序列的基本性质

* (1) 对称性

* 若随机序列 $\dots, X_{n-1}, X_n, X_{n+1}, \dots$ 是Markov序列，则沿相反的时间方向也是Markov序列。

* 证明：以三个时间点的情况为例，

$$\begin{aligned} & P[X_{n-1}|X_n, X_{n+1}] = P[X_{n-1}, X_n, X_{n+1}]/P[X_n, X_{n+1}] \quad (\text{条件概率的定义}) \\ & = P[X_{n+1}|X_n, X_{n-1}]P[X_n, X_{n-1}]/P[X_n, X_{n+1}] \quad (\text{根据条件概率公式重写} P[X_{n-1}, X_n, X_{n+1}]) \\ & = P[X_{n+1}|X_n]P[X_n, X_{n-1}]/P[X_n, X_{n+1}] \quad (\rightarrow X_{n-1} \rightarrow X_n \rightarrow X_{n+1} \rightarrow) \\ & = P[X_{n+1}|X_n]P[X_{n-1}|X_n]P[X_n]/P[X_{n+1}|X_n]P[X_n] \\ & \quad \text{根据条件概率公式重写} P[X_n, X_{n-1}] \text{ 和 } P[X_n, X_{n+1}] \\ & = P[X_{n-1}|X_n] \end{aligned}$$

* 习题：设 $n < n_1 < n_2 < \dots < n_p$ ，按以上思路证明 $P[X_n|X_{n_1}, X_{n_2}, \dots, X_{n_p}] = P[X_n|X_{n_1}]$ 。

* 以上分析表明基于任何未来的数据 $(X_{n_1}, X_{n_2}, \dots, X_{n_p})$ 推断过去的状态 X_n ，其效果总是等价于基于距离过去最近的未来数据 X_{n_1} 推断 X_n 。

* 采用前面的记号，就是 $\dots \leftarrow X_{n-1} \leftarrow X_n \leftarrow X_{n+1} \leftarrow \dots$ 。



Markov序列(3)

Markov序列的基本性质

* (2) 继承性

* 若随机序列 $\dots, X_{n-1}, X_n, X_{n+1}, \dots$ 是Markov序列，则其任何子序列也是Markov序列。

* 思考题：为什么？

* (3) 数据处理不等式

若 $X \rightarrow Y \rightarrow Z$ ，则 $I(X;Z) \leq I(X;Y)$

* 证明： $I(X;Z) - I(X;Y) = \sum_{xz} p(x,z) \log[p(x,z)/p(x)p(z)]$
* $\quad - \sum_{xy} p(x,y) \log[p(x,y)/p(x)p(y)]$
* $= \sum_{xyz} p(x,y,z) \log[p(x,z)/p(x)p(z)] - \sum_{xy} p(x,y,z) \log[p(x,y)/p(x)p(y)]$ (为什么?)
* $= \sum_{xyz} p(x,y,z) \log[p(x,z)p(y)/p(x,y)p(z)]$
* $= \sum_{xyz} p(x|y,z)p(y,z) \log[p(x|z)/p(x|y)] \quad (p(x,y,z)=p(x|y,z)p(y,z))$
* $= \sum_{xyz} p(x|y)p(y,z) \log[p(x|z)/p(x|y)] \quad (X \leftarrow Y \leftarrow Z)$
* $= - \sum_{xyz} p(y,z) p(x|y) \log[p(x|y)/p(x|z)] \leq 0$ (为什么?) 其他的证明参见2.8节



Markov序列(4)

* Markov序列的基本性质

* (4) 狭义数据处理不等式

* g 是任何函数, X 和 Y 是任意的随机变量, 则 $X \rightarrow Y \rightarrow g(Y)$,
因此 $I(X; g(Y)) \leq I(Y; g(Y))$

*

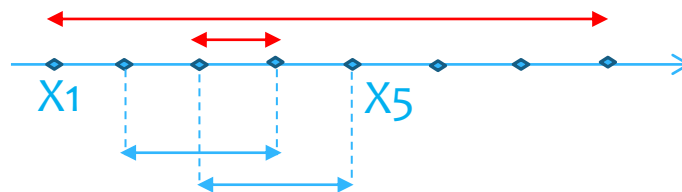
* 小结: 对Markov序列 $\dots \rightarrow X_{n-1} \rightarrow X_n \rightarrow X_{n+1} \rightarrow \dots$, “距离近”的状态变量的互信息量, 总是不低于“距离远”的状态量之间的互信息量。

* 习题 对以上Markov序列, 确定以下互信息量的相对大小, 或指出两者不能比较:

* $I(X_2; X_4)$ 和 $I(X_3, X_4)$;

* $I(X_1; X_{100})$ 和 $I(X_3, X_4)$;

* $I(X_2; X_4)$ 和 $I(X_3, X_5)$.



* 答案: 小于; 小于; 不能比较 (相对大小不确定)。



小结与习题（教程第一版第二章）

讲义未包含/阅读教程的内容：

- * 条件互信息量 $I(X:Y|Z)$ 及其基本性质，参阅 2.5 节定义式 (2.60-61)。该量在后续应用很少，属于一个有用但不是本质性的工具。凡是应用 $I(X:Y|Z)$ 分析和论证的问题，都可以等效地采用熵和条件熵的已经学过的关系来论证。
- * 2.9 节：充分统计量，可略去。
- * 2.10 节：核心内容是 Fano 不等式，该不等式重要但在后续应用很少。该节没有新概念，可作为前述已学知识的一个有趣的应用练习。
- * 习题： 2.2, 2.3, 2.4, 2.5, 2.10, 2.15, 2.21, 2.26, 2.27, 2.28, 2.29, 2.30（较难/选做）。
- * 部分参考答案或简要提示：
 - * 2.2 暂时将 X 作为取值离散的随机变量，下同：(a) $H(X) = H(2^X)$ ；(b) $H(X) \geq H(\cos X)$ ；
 - * 思考：如果 Y 是 X 的函数，什么情况下 $H(X)=H(Y)$ ？



部分习题选解（概要）

* 习题2.10

- * 因为 X_1 、 X_2 的取值范围不相交, 可以定义一个随机变量 X , 其取值
- * 范围是 X_1 、 X_2 值域的并集, 概率 $P[X=X_1]=a$, 概率 $P[X=X_2]=1-a$.
- * 再定义一个函数 $y=f(X)$, $f(X)=1$, 若 $X=X_1$; $f(X)=2$, 若 $X=X_2$ 。
- * 于是 $H(X) = H(X, f(X)) = H(y) + H(X|y)$
- * $= H(y) + p[y=1]H(X|y=1) + p[y=2]H(X|y=2)$
- * $= H(a, 1-a) + aH(X_1) + (1-a)H(X_2)$
- * 其中 $H(a, 1-a) = -a \log a - (1-a) \log(1-a)$. 【补充完整的细节】



部分习题选解（概要）

* 习题2.15

* 对Markov序列 $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \dots \rightarrow X_n$

* 有 $P[X_1, \dots, X_n] = P[X_n | X_{n-1}, \dots, X_1] P[X_{n-1}, \dots, X_1]$

* $= P[X_n | X_{n-1}] P[X_{n-1} | X_{n-2}, \dots, X_1] P[X_{n-2}, \dots, X_1]$

* $= P[X_n | X_{n-1}] P[X_{n-1} | X_{n-2}] P[X_{n-2} | X_{n-3}, \dots, X_1]$

* $= \dots = P[X_n | X_{n-1}] P[X_{n-1} | X_{n-2}] \dots P[X_2 | X_1] P[X_1]$

* 代入 $I(X_1; X_2, \dots, X_n)$

* $= \sum P[X_n, \dots, X_1] \log \{ P[X_n, \dots, X_1] / P[X_1] P[X_2, \dots, X_n] \}$

* 计算化简得

* $I(X_1; X_2, \dots, X_n) = I(X_1; X_2)$

* 【完成完整的计算，并思考：该结论的含义是什么？】



部分习题选解（概要）

* 习题2.21

* 证明概要

$$\begin{aligned} P(p(X) < d) \log \frac{1}{d} &= \sum_{x:p(x) < d} p(x) \log \frac{1}{d} \\ &\leq \sum_{x:p(x) < d} p(x) \log \frac{1}{p(x)} \\ &\leq \sum_x p(x) \log \frac{1}{p(x)} \\ &= H(X) \end{aligned}$$



部分习题选解（概要）

* 习题2.29 推导提示

- * (a) $H(X, Y | Z) = H(X|Z) + H(Y | X, Z) > H(X|Z)$
- * (b) $I(X, Y ; Z) = I(X; Z) + I(Y ; Z|X) > I(X; Z)$
- * (c) $H(X, Y, Z) - H(X, Y) = H(Z|X, Y) = H(Z|X) - I(Y ; Z|X)$
* $< H(Z|X) = H(X, Z) - H(X)$
- * (d) $I(X; Z|Y) + I(Z; Y) = I(X, Y ; Z) = I(Z; Y | X) + I(X; Z)$
* 和
- * $I(X; Z|Y) = I(Z; Y | X) - I(Z; Y) + I(X; Z)$
- * 注意该题的结论实际上是一个等式。



对连续随机变量的推广(1)

前面建立的全部概念和关系，都可以推广到连续型随机变量。

【在学习这部分时，注意连续随机变量情形特有的一些性质，例如信息熵的数值可以为负】

【阅读教程第8章 微分熵 8.1、8.2（在学习第三章后阅读）、8.4、8.5、8.6（略去定理8.6.6）】

- 1 连续随机变量 X ，取值范围记为 S ，概率密度为 $f(x)$ ，定义其

信息熵

$$h(X) = - \int_S f(x) \log f(x) dx,$$

- 2 连续随机变量 X 和 Y ，其联合概率密度是 $f(x,y)$ ，定义其联合熵

$$h(X,Y) = - \int dx dy f(x,y) \log f(x,y)$$

- 3 连续随机变量 X 和 Y ，定义其互信息量

$$I(X;Y) = - \int dx dy f(x,y) \log \{f(x,y)/f(x)f(y)\}$$



对连续随机变量的推广(2)

4 基本性质

$$h(X,Y) = h(X) + h(Y|X) = h(Y) + h(X|Y)$$

$$h(X|Y) \leq h(X)$$

$h(X,Y) \leq h(X) + h(Y)$ 并且上界在X和Y概率独立时达到

$$I(X;Y) \geq 0$$

$$I(X; Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$$

【习题】 仿照离散变量的情形检验上述关系。



对连续随机变量的推广(3)

* 5 基本实例

* 5.1 一维Gauss变量X

* 概率密度 $N(m, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp(-\frac{1}{2\sigma^2}(x-m)^2)$

* 信息熵 $H(X)$

*
$$= - \int_{-\infty}^{+\infty} dx (2\pi\sigma^2)^{-1/2} \exp(-\frac{1}{2\sigma^2}(x-m)^2) \log_e[(2\pi\sigma^2)^{-1/2} \exp(-\frac{1}{2\sigma^2}(x-m)^2)]$$

*
$$= - \int_{-\infty}^{+\infty} dx (2\pi\sigma^2)^{-1/2} \exp(-\frac{1}{2\sigma^2}(x-m)^2) \log_e[(2\pi\sigma^2)^{-1/2}]$$

*
$$+ \frac{1}{2\sigma^2} \int_{-\infty}^{+\infty} dx (x-m)^2 (2\pi\sigma^2)^{-1/2} \exp(-\frac{1}{2\sigma^2}(x-m)^2)$$

*
$$= \log_e[(2\pi\sigma^2)^{1/2}] + \frac{1}{2} = \log_e[(2\pi e\sigma^2)^{1/2}]$$



对连续随机变量的推广(4)

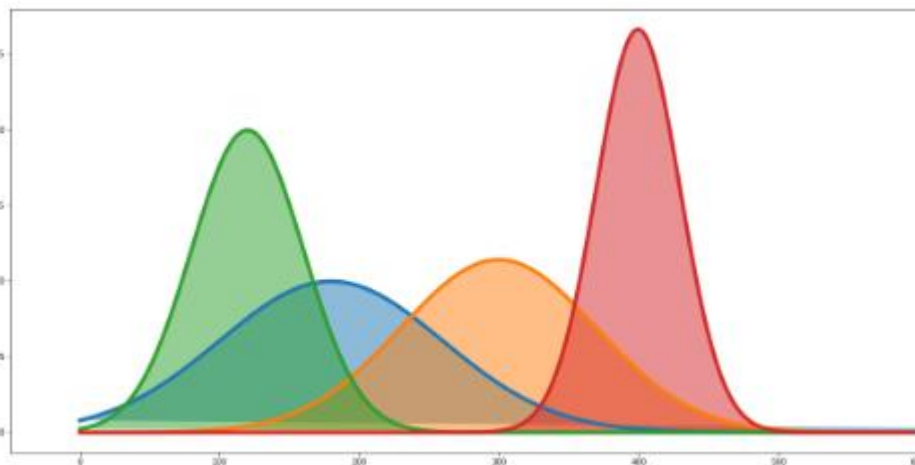
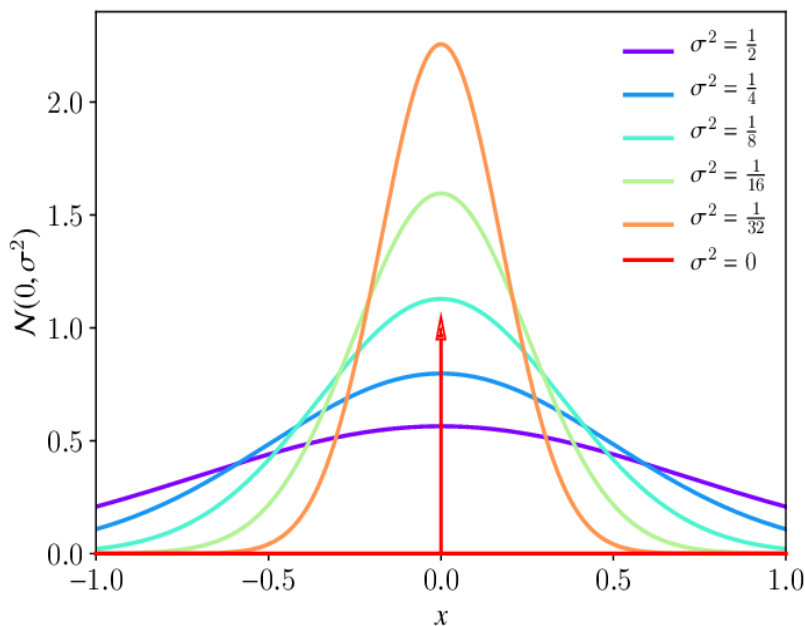
* 5 基本实例

* 5.1 (续) 一维Gauss变量 $X \sim N(m, \sigma^2)$ 的信息熵

*
$$H(X) = \frac{1}{2} \log_e(2\pi e \sigma^2)$$

* 注: Gauss变量的熵 $H(X)$ 同均值 m 无关, 仅与方差 σ^2 有关, 并是方差的增函数。

* 你能定性地解释为什么如此吗?



对连续随机变量的推广(5)

* 5 基本实例

* 5.2 n维Gauss随机向量 $\mathbf{X}=(X_1,\dots,X_n)\sim N(\mathbf{m},\mathbf{Q})$

$$* \quad N(\mathbf{m},\mathbf{Q}) = (2\pi|\mathbf{Q}|)^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \mathbf{Q}^{-1}(\mathbf{x}-\mathbf{m})\right)$$

* \mathbf{m} 是均值向量, $|\mathbf{Q}|$ 表示 \mathbf{Q} 的行列式, \mathbf{Q} 是正定对称的协方差矩阵:

$$* \quad Q_{ij} = E[(x_i - m_i)(x_j - m_j)]$$

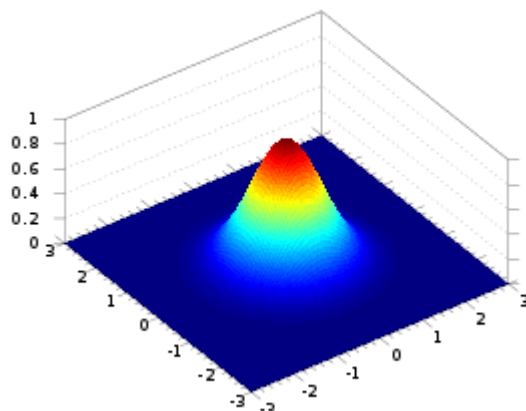
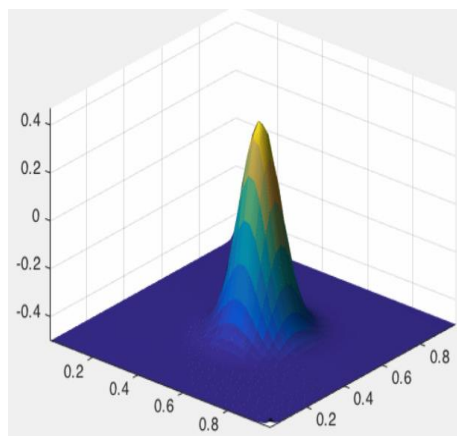
$$* \quad \text{信息熵} \quad H(\mathbf{X}) = \frac{1}{2} \log_e((2\pi e)^n |\mathbf{Q}|)$$

【习题】

(1)根据熵的定义推导该结果;

(2)当 X_1,\dots,X_n 是独立的Gauss变量, 方差分别为 $\sigma_1^2, \dots, \sigma_n^2$, 这时 $H(\mathbf{X})$ 的表达式是什么?

结合普遍的关系式 $h(X_1,\dots,X_n) \leq h(X_1) + \dots + h(X_n)$, 能由此导出著名的Hardamard不等式(8.64)。



对连续随机变量的推广(6)

- * $I(X;Y) = - \int dx dy f(x,y) \log \{f(x,y)/f(x)f(y)\}$

- * $= - \int dx dy f(y|x)f(x) \log \{f(y|x)/f(y)\}$

- * $f(y) = \int dx f(y|x)f(x)$

- * $I(X;Y)$ 是 $f(y|x)$ 的凸泛函、 $f(x)$ 的凹泛函

* 补充知识：泛函(functional)

- * (1) 所谓泛函，就是一个随不同函数而变化的实数或复数，即函数的函数。

- * (2) 典型的实例： $h(X)$, 和 $I(X;Y)$, 两者是概率密度的函数；

- * Dirac泛函： $\delta_a[f(x)] = f(a)$; $\delta_a^{(m)}[f(x)] = (-1)^m f^{(m)}(a)$

- * (3) 当函数变化时，可以考虑泛函的变分(variation)，即泛函微分：

- * $F[f(x)+\Delta f(x)] \approx F[f(x)] + L_f[\Delta f(x)]$

- * 其中 $\Delta f(x)$ 是幅度很小的函数、 L_f 是 $\Delta f(x)$ 的所谓线性泛函。则 L_f 定义做泛函 F 在 $f(x)$ 处的变分。

- * (4) 如果在 $f^*(x)$ 处 $L_{f^*}[\Delta f(x)] = 0$ 对任何 $\Delta f(x)$ 成立，则 $f^*(x)$ 是使泛函达到某种极值的函数。

- * 【以上概念参见下面的例子。在Gauss信道容量一章将继续出现泛函及其优化的例子】



对连续随机变量的推广(7)

* 7 关于Gauss概率分布的一个极值性质 (定理8.6.5)

* 在均值为给定的向量 \mathbf{m} 、协方差矩阵为给定的正定对称矩阵 \mathbf{Q} 的各种 n 维随机变量 \mathbf{X} 中, Gauss变量的熵 $h(\mathbf{X})$ 最大。

* 证明: 第一步: 考虑带约束的极值问题

*
$$\max - \int d\mathbf{x} f(\mathbf{x}) \log f(\mathbf{x})$$

*
$$\text{s.t. } \int d\mathbf{x} \mathbf{x} f(\mathbf{x}) = \mathbf{m}, \quad \int d\mathbf{x} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T f(\mathbf{x}) = \mathbf{Q}$$

* 第二步: 针对约束引进乘子向量 \mathbf{a} 和乘子矩阵 \mathbf{A} , 计算广义目标泛函

*
$$F[f(\mathbf{x})] = - \int d\mathbf{x} f(\mathbf{x}) \log f(\mathbf{x}) + \mathbf{a}^T (\int d\mathbf{x} \mathbf{x} f(\mathbf{x}) - \mathbf{m})$$

*
$$+ \text{tr}(\mathbf{A} (\int d\mathbf{x} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T f(\mathbf{x}) - \mathbf{Q}))$$

* 对 $f(\mathbf{x})$ 的变分。

* 【符号】 $\int d\mathbf{x}$ 是 n 重积分 $\int dx_1 \dots dx_n$, $\text{tr}(\mathbf{M})$ 表示矩阵的迹, 即对角线元素之和。

* 【注】 这里对定理8.6.5的推导与教程不同, 请阅读教程上的推导, 并比较两种方法的特点。



对连续随机变量的推广(8)

* 7 关于Gauss概率分布的一个极值性质 (定理8.6.5)

* 证明 (续) 计算泛函 $F[f(x)] = - \int dx f(x) \log f(x) + a^T (\int dx x f(x) - m)$
* $+ \text{tr}(\Lambda (\int dx (x-m)(x-m)^T f(x) - Q))$ 对 $f(x)$ 的变分:

* 仿照微分算法计算 $\Delta f(x)$ 无穷小时的差, 得到

*
$$F[f(x) + \Delta f(x)] - F[f(x)] \approx - \int dx \Delta f(x) (1 + \log f(x)) + a^T (\int dx x \Delta f(x))$$

*
$$+ \text{tr}(\Lambda (\int dx (x-m)(x-m)^T \Delta f(x)))$$

* 【习题】计算上述结果。注意将第一项对比普通的微分公式 $d(u \log u) = (1 + \log u) du$ 。

* 第三步: 令以上表达式为零, 得到一个关于最优分布的积分等式。由于该等式需
要对其中任意的函数 $\Delta f(x)$ 被满足, 因此被积函数中 $\Delta f(x)$ 的权因子必须为0:

*
$$1 + \log f^*(x) + a^T x + \text{tr}(\Lambda (x-m)(x-m)^T) = 0$$

* 这就是使得熵达到极大值的概率分布 $f^*(x)$ 应满足的方程。

* 第四步: 容易解出 $f^*(x) = \exp(x \text{ 的二次型})$, 这种函数形式的概率分布正是
* Gauss分布。证毕。

【习题】应用初等概率论的知识, 完成最后一步的论证。



对连续随机变量的推广(9)

* 7 小结 (定理8.6.5)

- * 在均值为给定的向量 \mathbf{m} 、协方差矩阵为给定的正定对称矩阵 \mathbf{Q} 的各种 n 维随机变量 \mathbf{X} 中，Gauss变量的熵 $h(\mathbf{X})$ 最大。
- * 特例：如果一个标量随机信号 x 的均值为零、功率 P 给定，则Gauss信号的熵最大。
- * 【习题】证明上述结论。注意一维信号若均值为0，则方差恰是其功率。
- * 【注】改组结论将在学习Gauss信道（第9章）和MIMO信道的容量性能时用到。
- * 【习题】用类似的方法证明：如果一维随机信号 X 的均值为已知的正数 m ，
- * （协方差任意）则指数分布的随机信号具有最大熵 $h(X)$ 。



习题（第8章 微分熵）

* 习题8.1、8.5

* 参考答案：

* 8.1(a) 指数随机变量的熵

$$\begin{aligned}h(f) &= -\int_0^{\infty} \lambda e^{-\lambda x} [\ln \lambda - \lambda x] dx \\&= -\ln \lambda + 1 \text{ nats.} \\&= \log \frac{e}{\lambda} \text{ bits.}\end{aligned}$$

* (b) Laplace随机变量的熵 $h(X) = \log(2e/\lambda)$

* (c) 独立Gauss变量之和的熵 $h(X_1+X_2) = \frac{1}{2} \log_e(2\pi e(\sigma_1^2 + \sigma_2^2))$

* 8.5 应用积分变量变换公式计算得：

$$\begin{aligned}h(AX) &= -\int g(y) \ln g(y) dy \\&= -\int \frac{1}{|A|} f(A^{-1}y) [\ln f(A^{-1}y) - \log |A|] dy \\&= -\int \frac{1}{|A|} f(x) [\ln f(x) - \log |A|] |A| dx \\&= h(X) + \log |A|.\end{aligned}$$



下一讲：典型集合与渐进均分性质

- * 预习：教程第三章 渐进均分性质
- * 基础：概率的大数定律

