

AI 中的数学

邓小铁 李翰禹

2024 年 9 月 28 日

目录

第零章 引言	1
第一部分 AI 的逻辑	2
第一章 合情推理	3
§1.1 命题逻辑的演绎推理	4
§1.2 合情推理的数学模型	12
§1.2.1 合情推理的基本假设，似然	14
§1.2.2 似然与概率	19
§1.2.3 先验与基率谬误	21
§1.3 合情推理的归纳强论证	23
§1.3.1 归纳强论证	23
§1.3.2 有效论证和归纳强论证的比较	28
§1.4 先验模型的存在性	33
§1.5 章末注记	35
§1.6 习题	36

第二章 Markov 链与模型	37
§2.1 Markov 链	38
§2.2 Markov 奖励过程 (MRP)	49
§2.3 Markov 决策过程 (MDP)	55
§2.4 隐 Markov 模型 (HMM)	64
§2.4.1 评估问题	67
§2.4.2 解释问题	69
§2.5 扩散模型	72
§2.5.1 采样逆向过程	77
§2.5.2 训练逆向过程	78
§2.6 章末注记	81
§2.7 习题	81
 第二部分 信息与数据	 82
第三章 熵与 Kullback-Leibler 散度	83
§3.1 熵	84
§3.1.1 概念的导出	84
§3.1.2 概念与性质	89
§3.2 Kullback-Leibler 散度	98
§3.2.1 定义	98
§3.2.2 两个关于信息的不等式	101
§3.3 编码理论	102
§3.3.1 熵与编码	102
§3.3.2 K-L 散度、交叉熵与编码	106

§3.4 在机器学习中的应用：语言生成模型	108
§3.5 附录：Shannon 定理的证明	110
§3.6 习题	113
§3.7 章末注记	115
第四章 高维几何，Johnson-Lindenstrauss 引理	118
§4.1 高维几何	120
§4.1.1 高维球体	120
§4.1.2 Stein 悖论	124
§4.1.3 为什么我们要正则化？远有潜龙，勿用	131
§4.2 集中不等式	132
§4.3 J-L 引理的陈述与证明	139
§4.4 J-L 引理的应用	144
§4.5 习题	147
§4.6 章末注记	148
第五章 差分隐私	149
§5.1 数据隐私问题	150
§5.2 差分隐私的定义与性质	154
§5.3 差分隐私的应用	162
§5.3.1 随机反应算法	162
§5.3.2 全局灵敏度与 Laplace 机制	165
§5.3.3 DP 版本 Llyod 算法	169
§5.4 习题	172
§5.5 章末注记	172

第三部分 决策与优化 173

第六章 凸分析 174

§6.1 决策与优化的基本原理	175
§6.1.1 统计决策理论	175
§6.1.2 优化问题	178
§6.1.3 例子: 网格搜索算法	185
§6.2 凸函数	189
§6.3 凸集	195
§6.3.1 基本定义和性质	196
§6.3.2 分离超平面定理	200
§6.4 习题	203
§6.5 章末注记	203

第七章 对偶理论 204

§7.1 约束的几何意义	207
§7.2 条件极值与 Lagrange 乘子法	215
§7.3 Karush–Kuhn–Tucker 条件	219
§7.4 Lagrange 对偶	224
§7.4.1 原始规划与对偶规划	224
§7.4.2 对偶的几何意义	229
§7.4.3 弱对偶定理	231
§7.4.4 Slater 条件, 强对偶定理	232
§7.5 应用: 支持向量机 (SVM)	237
§7.6 习题	240
§7.7 章末注记	240

第八章 不动点理论	241
§8.1 Banach 不动点定理	242
§8.2 Brouwer 不动点定理	253
§8.3 习题	260
§8.4 章末注记	260
 第四部分 博弈与逻辑	 261
第九章 逻辑与博弈	262
§9.1 博弈的基本语言：以井字棋为例	264
§9.2 输赢博弈	266
§9.2.1 博弈的不同维度	266
§9.2.2 Zermelo 定理与 AlphaGo Zero	268
§9.3 正则形式博弈	275
§9.3.1 定义	276
§9.3.2 理性与均衡	278
§9.3.3 生成对抗网络	281
§9.3.4 混合策略	284
§9.4 随机博弈 (Markov 博弈)	290
§9.5 习题	301
§9.6 章末注记	301
 第五部分 认知与逻辑	 302
第十章 共同知识, Bayes 博弈, Aumann 知识算子	303
§10.1 “泥泞的孩童”谜题	306

§10.2 不完全信息博弈 (Bayes 博弈)	311
§10.3 电子邮件博弈	321
§10.4 Aumann 知识算子	326
§10.5 习题	335
§10.6 章末注记	335
第十一章 模态逻辑, 知识的逻辑	336
§11.1 知识逻辑的形式语言	338
§11.2 Kripke 语义	343
§11.3 模态可定义性	350
§11.4 知识逻辑的基本模型与性质	353
§11.4.1 知识逻辑的 Kripke 模型与公理	353
§11.4.2 Kripke 模型与 Aumann 结构	360
§11.4.3 “泥泞的孩童”再回顾: 形式化解法	362
§11.5 对不一致达成一致	364
§11.5.1 模型	365
§11.5.2 定理及其证明	368
§11.6 习题	370
§11.7 章末注记	370
第六部分 附录: 预备知识	371
附录 A 线性代数基础	372
§A.1 线性空间	372
§A.2 线性映射	379
§A.3 矩阵	385

§A.4 双线性型与二次型	394
§A.5 带内积的线性空间	400
§A.6 行列式	409
§A.7 算子范数与谱理论	413
附录 B 微分学基础	422
§B.1 点集拓扑	422
§B.1.1 度量空间, 范数	422
§B.1.2 开集与闭集	427
§B.1.3 紧致性, 收敛性, 完备性	431
§B.1.4 连续映射	435
§B.1.5 与实数序有关的性质	440
§B.2 一元函数的微分学	443
§B.2.1 导数与微分的定义	444
§B.2.2 微分学基本定理	449
§B.3 多元函数的微分学	451
§B.3.1 微分、偏导数与导数的定义	451
§B.3.2 微分学基本定理	461
§B.3.3 隐函数定理	464
附录 C 概率论基础	470
§C.1 从朴素概率论到公理化概率论	470
§C.1.1 Kolmogorov 概率论	470
§C.1.2 条件概率, 独立性	476
§C.2 随机变量, 分布函数	482
§C.2.1 基本定义	482

§C.2.2 离散型随机变量	488
§C.2.3 连续型随机变量	489
§C.2.4 随机向量, 条件分布, 独立性	494
§C.2.5 随机变量 (向量) 的函数	501
§C.3 随机变量的数字特征, 条件数学期望	505
§C.3.1 数学期望, Lebesgue 积分	505
§C.3.2 数学期望的性质	512
§C.3.3 随机变量的内积空间	516
§C.3.4 特征函数	519
§C.3.5 条件数学期望	521
§C.4 多元正态分布 (Gauss 向量)	528
§C.5 大数定律	530

第一部分

AI 的逻辑

第二部分

信息与数据

第三章 熵与

Kullback-Leibler 散度

人脑和机器的区别是什么？我们有可能模拟人脑的功能吗？这些问题根植于认知科学和人工智能领域。在 20 世纪 50 年代，计算机科学、认知科学和人工智能仍然处于萌芽状态。就是在这个时候，基于对人类如何解决问题和决策的研究，Herbert A. Simon 提出，人脑其实是一个“信息处理器”，也就是输入信息（视觉、听觉等），进行处理，然后输出信息（动作、语言等）。从这个观点上，人脑和机器并无区别。

基于这样的观点，Simon 和 Allen Newell、J. C. Shaw 一起合作，制造了逻辑理论家（Logic Theorist）、通用问题求解器（General Problem Solver，GPS）等计算机程序。逻辑理论家可以证明《数学原理》（作者是 Whitehead 和 Russell）第二章前 52 个定理中的 38 个，而 GPS 则可以解决汉诺塔问题。这展现出“信息处理器”观点的巨大潜力。

时至今日，“信息处理器”的观点已经深入认知科学和人工智能研究者的心中。然而，信息是一个特别抽象的概念。它不像重量，可以从沉甸甸的铅块中直观感受到。那么，信息到底是什么？本章将要讨论这一

问题，并给出它在人工智能领域的应用。

§3.1 熵

§3.1.1 概念的导出

我们常说“恐惧来源于未知”，信息似乎代表着某种确定的东西，某种知识，因而和不确定性有相反的关系。更精确地说，消除不确定性的东西被称为信息。当然，这句话本身似乎是一种循环解释，它既没有回答信息是什么也没有回答不确定性是什么。所以我们进一步的问题是，给定一个“对象”，如何定量衡量它不确定性（或信息量）？

我们先从一个例子看起。

例 3.1 (信息论读本) 假设我们有一个信息论的读本（例如本章就是），我们想要衡量它的信息量。我们面临的第一个困难是，同样的内容对于不同的人来说，信息量是完全不同的。已经学过信息论的读者再看这一部分内容，他获得的信息会比没有学过的读者要少得多。因此，我们很难直接给单个对象衡量它的信息量。

但是，信息论读本的读者背景是多样的、不确定的，可能学过信息论，可能只学过概率论，也可能什么都没学过。要衡量这本书的信息量，我们可以考虑所有可能的读者背景，然后给出一个信息的概率分析。例如，读这本书的读者大概率不是信息论专家，但有一定概率论的背景，他们可以获得很多信息；而还有很少部分读者精通信息论，因此这本书给他们的信息量就很少。但综合来看，这本书的信息量依然是不少的。□

以上例子说明了这样一种思想：将世界视为不确定的，有多种可能的结果，然后考虑这一堆结果所带来的平均信息量。

我们可以用数学来表述上面的考虑, 假如我们进行一次试验, 一共有 n 种可能的结果, 第 i 种发生的概率为 p_i . 我们预测试验的结果, 如果越能正确地预测, 那么就说明我们对这个试验中包含的信息知道的越多.

- 假如 $p_1 = 1$, 那么我们完全确定试验一定会产生结果 1;
- 如果 $p_i = 1/n$, 那么我们完全无法预计试验的结果.

我们对试验结果的预期与试验结果的概率分布有密切联系. 因此概率分布给我们带来了信息, 使得我们能够产生不同的判断. 另一方面, 概率分布带来了不确定性, 使我们不能总是确信预言会成真.

我们遵循“信息论之父”Shannon 的思路, 为信息提供一个严格的数学模型: 熵. 假设随机变量 X 表示了所有可能的结果 (编号为 1 到 n), $\Pr(X = i) = p_i$, $p = (p_1, \dots, p_n)$, 有时候也把 p_i 写作 $p(i)$. 我们把不确定性度量记为 $H(p)$. Shannon 假设 H 满足以下三个性质:

1. H 是一个连续函数.
2. 事件结局可能数变多则不确定性增大: $p_i = 1/n$ 时, $H(p)$ 随 n 单调递增, n 是正整数.
3. 如果一个试验被分成了两个相继的试验, 那么原来的 H 应该等于分开之后的 H 的加权和.

前两个假设都比较好理解, 我们现在具体解释第三个假设.

如图图 3.1 所示, 假设我们有一个试验, 有三种可能的结果, 1, 2, 3, 概率分别为 $1/2, 1/3, 1/6$. 该试验的不确定性是 $H(1/2, 1/3, 1/6)$.

我们把试验分成两步相继的试验 (右图). 第一步试验有两种可能的结果, 概率分别都是 $1/2$.

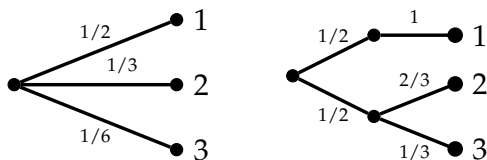


图 3.1: Shannon 的第三个假设

- 当第一步试验出现上面的结果时，第二步试验以概率 1 产生结果 1；
- 当第二步试验出现下面的结果时，第二步试验以概率 2/3 产生结果 2，以概率 1/3 产生结果 3。

我们可以看到，分成两步之后，第一步试验的不确定性是 $H(1/2, 1/2)$ ，第二步试验的不确定性有一半概率是 $H(1)$ （上面的分支），有一半概率是 $H(2/3, 1/3)$ （下面的分支），因而加权的确定性是 $1/2 \cdot 0 + 1/2 \cdot H(2/3, 1/3)$ 。因此第三个假设可以具体表述为

$$H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) = H\left(\frac{1}{2}, \frac{1}{2}\right) + \left[\frac{1}{2} \cdot H(1) + \frac{1}{2} \cdot H\left(\frac{2}{3}, \frac{1}{3}\right)\right].$$

这里，我们可以看出 Shannon 的哲学思想：不确定性只来自于概率分布而不是具体对象。他的考虑具有浓厚的工程意味，正如他自己针对通信的数学理论所说：“消息是具有含义的……然而，通信的语义层面并不是工程问题所关心的。”正是因为抽象掉了具体考虑的对象，信息论的应用才变得如此广泛。

基于上面三个假设，Shannon 证明了如下定理，这一定理直接给出了熵的概念。

定理 3.1 (Shannon 定理) H 满足三个假设当且仅当

$$H(p) = -C \sum_i p_i \log p_i,$$

其中 C 是正常数, $0 \log 0 = 0$.

这一定理的证明较长并且和后面的讨论关联较小, 所以我们在第 3.5 节中给出证明.

根据对数的换底公式, 可以将 $C \log p_i$ 写为 $\log_b p_i$, 这里 $C = 1/\log b$. 于是, Shannon 定理直接给出了熵的如下定义:

定义 3.1 (熵) 分布列 $p = (p_1, \dots, p_n)$ 的熵定义为

$$H(p) = - \sum_{i=1}^n p_i \log_b p_i.$$

其中 $b = e$ (自然对数底数), $0 \log 0 = 0$. 当 $b = 2$ 时, 我们记熵为 $H_2(p)$. □

通常来说, 使用 e 作为底数会使得数学推导简洁, 而用 2 为底数则常常是讨论信息量时的习惯. 在第 3.3 节中, 我们将讨论熵在通信中的含义, 以 2 为底的时候熵的实际意义会更清楚些. 如果没有特别强调, 我们在讨论时总是假设 $b = e$.

熵的定义还可以用数学期望的形式写出. 假设 X 的分布列是 p , $p(i) = \Pr(X = i)$, 那么我们也可以把熵写成期望的形式:

$$H(p) = -\mathbb{E}[\log p(X)].$$

每一个 (离散) 随机变量 X 会确定一个分布列 p_X , 因此我们也可以定

义随机变量的熵:

定义 3.2 (随机变量的熵) 随机变量 X 的熵定义为

$$H(X) = -\mathbb{E}[\log p_X(X)].$$

其中 p_X 是 X 的分布列, $0 \log 0 = 0$. □

尽管从信息论的角度我们可以唯一确定熵的定义, 但是熵的概念在物理学上早就已经存在. 下面我们给出统计力学中熵的推导过程.

在经典力学中, 物理系统的状态由粒子的位置和动量 (也就是速度) 完全确定, 将粒子位置和动量可能的值集合称为相空间, 于是物理系统的演化就是相空间中的粒子状态的变化.

将相空间等分成 m 个单元, 编号 1 到 m . 假设相空间中有 N 个可区分的粒子, 相互独立, 没有相互作用, 每个粒子等可能出现在每一个单元中. 如果单元 i 中有 N_i 个粒子, 那么按照粒子在单元中的分布来看, 系统处于某个特定状态的概率为

$$P = \frac{N!}{N_1! \dots N_m!} \left(\frac{1}{m} \right)^N.$$

这是一个多项分布. 两边取对数, 得

$$\log P = \log(N!) - \sum_i \log(N_i!) - N \log m.$$

考虑充分大的 N_i , 由 Stirling 公式, 有

$$\log(N_i!) \sim \log \left(\sqrt{2\pi N_i} \left(\frac{N_i}{e} \right)^{N_i} \right) \sim N_i \log N_i.$$

因此,

$$\log P \sim N \log N - \sum_i N_i \log N_i - N \log m \sim N \log N - \sum_i N_i \log N_i. \quad (3.1)$$

假设 N_i 充分大的时候, N_i/N 呈现固定的比例 p_i , 那么

$$\begin{aligned} N \log N - \sum_i N_i \log N_i &\sim N \log N - \sum_i N p_i \log(N p_i) \\ &= -N \sum_i p_i \log p_i. \end{aligned}$$

$\log P \sim -N \sum_i p_i \log p_i$. 于是我们证明了:

$$\frac{1}{N} \log P \rightarrow H(p_1, \dots, p_m), \quad N \rightarrow \infty.$$

因此, 熵刻画了充分多粒子的物理系统某种特定状态出现概率! 熵越大的系统越有可能达到. 更进一步, 在统计力学中有 Boltzmann H -定理: 孤立的粒子系统会向着熵 (H) 增加的方向演化, 并最终达到熵最大的状态. H -定理是热力学第二定律的微观解释, 熵越大的系统出现概率越大、越混乱、越接近均衡.

§3.1.2 概念与性质

现在, 我们将进一步探讨熵的若干拓展定义, 并讨论他们的性质.

首先, 我们考虑最简单的情形, 即分布列为 (p_1, p_2) , 此时, 我们不妨设 $p_1 = p$, $p_2 = 1 - p$, 那么熵就是

$$H(p_1, p_2) = H(p, 1 - p) = -p \log p - (1 - p) \log(1 - p).$$

H 是关于 p 的函数, 作图如图 3.2 所示.

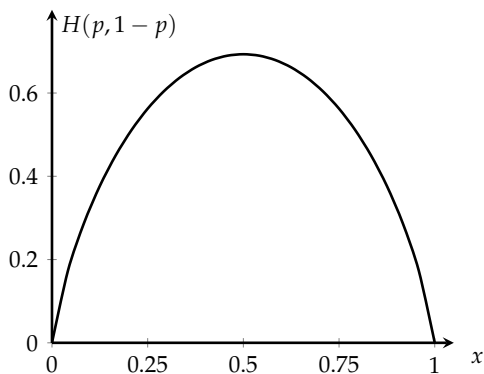


图 3.2: 熵 $H(p)$ 的图像.

利用导数的方法, 很容易证明:

命题 3.1 $H(p)$ 在 $p \in (0, 1/2)$ 严格单调递增, 在 $p \in (1/2, 1)$ 严格单调递减. 它的最小值是 0, 在 $p \in \{0, 1\}$ 取得; 它的最大值是 $\log 2$, 在 $p = 1/2$ 取得.

证明. 我们有

$$H'(p) = -\log p - 1 - \log(1-p) + 1 = \log(p^{-1} - 1),$$

$H'(p)$ 是一个关于 p 单调递减的函数. 我们有如下分类讨论:

- 当 $p \in (0, 1/2)$ 时, $p^{-1} - 1 > 1$, 所以 $H'(p) > 0$, $H(p)$ 是单调递增的;
- 当 $p \in (1/2, 1)$ 时, $p^{-1} - 1 < 1$, 所以 $H'(p) < 0$, $H(p)$ 是单调递减的;

• 当 $p = 1/2$ 时, $H'(p) = 0$, 结合前两点, $H(p)$ 取得最大值.

$H(0) = H(1) = 0$, $H(1/2) = \log 2$, 因此命题得证. \square

这与我们对于“不确定性”的直觉是相一致的: 当 p 接近 0 或 1 时, 我们对于 X 的取值几乎是确定的, 因此熵接近 0; 当 p 接近 $1/2$ 时, 我们对于 X 的取值几乎是完全不确定的, 因此熵接近最大值 $\log 2$.

实际上, 这样的性质对于一般的分布也是成立的, 我们分别将他们写在命题 3.2 和命题 3.4 中.

考虑一般分布的熵 $H(p) = H(p_1, \dots, p_n)$. 我们有如下性质:

命题 3.2 $H(p) \geq 0$, 等号成立当且仅当某个 $p_i = 1$.

证明. 这是一个典型的证明, 主要的技巧是使用熵的期望形式. 考虑随机变量 X , 其分布列为 p . 回忆 Jensen 不等式 (定理 C.17): 如果 f 是一个严格凸函数, 那么

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]).$$

等号成立当且仅当 X 是常数.

因为 $-\log(\cdot)$ 是严格凸函数, 所以根据 Jensen 不等式

$$H(X) = \mathbb{E}[-\log p(X)] \geq -\log \mathbb{E}[p(X)] \geq -\log 1 = 0.$$

等号成立当且仅当 X 是常数, 即对某个 i , $p(i) = 1$. \square

命题 3.3 p_i 朝着相等方向改变的时候 H 增加. 也就是说, 假设

$$p_i < p'_i \leq p'_j < p_j, \quad p_i + p_j = p'_i + p'_j,$$

那么, 用 p'_i 和 p'_j 代替原来的 p_i 和 p_j , H 将会变大.

证明. 为简化符号, 考虑 $i = 1$ 和 $j = 2$, 一般情况是一样的证明. 利用假设三, 第一步试验中, 将试验的结果 1 和结果 2 合并, 第二步试验再按照 $p_1/(p_1 + p_2)$ 和 $p_2/(p_1 + p_2)$ 的概率产生结果 1 和结果 2. 于是,

$$\begin{aligned}
 & H(p_1, p_2, \dots) \\
 &= H(p_1 + p_2, p_3, \dots) + (p_1 + p_2)H\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right) \quad (\text{假设三}) \\
 &\leq H(p_1 + p_2, p_3, \dots) + (p_1 + p_2) < H\left(\frac{p'_1}{p'_1 + p'_2}, \frac{p'_2}{p'_1 + p'_2}\right) \quad (\text{命题 3.1}) \\
 &= H(p'_1, p'_2, p_3, \dots). \quad (\text{假设三}) \quad \square
 \end{aligned}$$

命题 3.4 当且仅当 $p_1 = \dots = p_n = 1/n$ 时 H 取得最大值 $\log n$.

证明. 用反证法, 假设达到最大值的时候存在 $p_k \neq 1/n$, 那么, 因为 $\sum_i p_i/n = 1/n$, 根据鸽巢原理, 则必有 i, j 满足

$$p_i < 1/n < p_j.$$

根据命题 3.3, 我们可以将 p_i 和 p_j 替换为 $1/n$ 和 $p_i + p_j - 1/n$, 而 H 增大, 这与假设矛盾. 因此, $p_1 = \dots = p_n = 1/n$ 是 H 取得最大值的必要条件。

另一方面, 因为 H 连续, 所以根据 Weierstrass 最值定理 (??), H 一定有最大值, 所以 $p_1 = \dots = p_n = 1/n$ 也是 H 取得最大值的充分条件. □

至此, 命题 3.2 和命题 3.4 证明了一般情形的命题 3.1. 在等可能的

时候不确定性最大, 熵最大; 在确定事件的时候不确定性最小, 熵最小. 所以熵是符合直观的定义.

接下来, 我们讨论熵的拓展形式.

在一次试验中, 我们可以观察多个变量, 比如说 X 和 Y . 等价地, 我们其实只观察到了一个结果 (X, Y) , 只是这个结果是一个向量的形式, 服从分布 $p(i, j)$. 因此, 这一向量也有对应的熵, 这就是联合分布的熵:

$$H(X, Y) = -\mathbb{E}[\log p(X, Y)].$$

对应地, 我们也可以写成和的形式:

$$H(p) = -\sum_{i,j} p(i, j) \log p(i, j).$$

自然, 联合分布也可以引出边缘分布的熵:

$$H(X) = -\mathbb{E}[\log p_X(X)] = -\sum_i \sum_j p(i, j) \log \sum_j p(i, j).$$

$$H(Y) = -\mathbb{E}[\log p_Y(Y)] = -\sum_j \sum_i p(i, j) \log \sum_i p(i, j).$$

有了两个随机变量, 我们就可以讨论“条件”的概念. 具体来说, 我们可以把试验分为两步, 第一步观测 X , 第二步观测 Y , 那么, 第二步所产生的熵就是已经知道第一步结果之后的熵, 即:

$$H(Y|X=x) = -\mathbb{E}[\log p_{Y|X=x}(Y)|X=x] = -\sum_j p_{Y|X=x}(j) \log p_{Y|X=x}(j),$$

其中 $p_{Y|X=x}(j) = p(x, j)/p_X(x)$. 当我们知道了 $X=x$ 之后, 对 Y 的观

测就消除了部分的不确定性，因此根据我们对于不确定性和信息关系的讨论，从 $X = x$ 中获得的关于 Y 的信息是

$$I(X = x : Y) = H(Y) - H(Y|X = x).$$

考虑一个特殊情况， $Y = X$ ，那么刚刚的讨论就变成了自己从自己身上获得的信息，或者说知道 $X = x$ 带来的信息量。首先有

$$p_{X|X=x}(i) = \begin{cases} 1, & i = x \\ 0, & i \neq x. \end{cases}$$

因此，

$$H(X|X = x) = - \sum_j p_{X|X=x}(j) \log p_{X|X=x}(j) = -1 \log 1 = 0.$$

于是，

$$I(X = x : X) = H(X) - H(X|X = x) = H(X).$$

这正是定量版本的“消除不确定性的东西被称之为信息”！此外，我们之前说过，熵刻画的是一族可能对象的信息，这一点也反映在了这一公式中：只要知道了 X 的值，无论它具体是多少，我们得到的信息量是一样的！

再回到一般情况，还是同样的两步试验，我们定义给定 X 时 Y 的条件熵为

$$H(Y|X) = \mathbb{E}[H(Y|X = x)]$$

$$\begin{aligned}
&= -\mathbb{E}[\log p_{Y|X}(Y)] \\
&= -\sum_x p_X(x) \sum_j p_{Y|X=x}(j) \log p_{Y|X=x}(j) \\
&= -\sum_{x,j} p(x,j) \log p_{Y|X=x}(j).
\end{aligned}$$

换言之，我们现在进一步假定 X 也是不知道的，于是 $H(Y|X)$ 就是平均上来说第二步中 Y 的不确定性。条件熵和熵有着类似的性质：

命题 3.5 $H(Y|X) \geq 0$ ，等号成立当且仅当 Y 是退化的，即 Y 概率 1 只取一个值。

证明。 仿照命题 3.2 的证明即可。 □

类似地，我们可以考虑平均上 Y 中包含的关于 X 的信息量：

$$\mathbb{E}[I(X = x : Y)] = H(Y) - H(Y|X).$$

与之相对应地，平均上 X 中包含的关于 Y 的信息量为

$$\mathbb{E}[I(Y = y : X)] = H(X) - H(X|Y).$$

一个自然的问题是，二者相互包含的信息量是什么关系？根据概率的链式法则， $p(x,y) = p_{X|Y}(x|y)p_Y(y)$ ，带入 $H(X,Y)$ 的定义得熵的链式法则：

命题 3.6 对任意离散随机变量 X, Y ， $H(X,Y) = H(Y) + H(X|Y)$ 。

利用链式法则，我们注意到，

$$H(X) - H(X|Y) = H(X) - (H(X,Y) - H(Y))$$

$$\begin{aligned}
&= H(X) + H(Y) - H(X, Y) \\
&= H(Y) - (H(X, Y) - H(X)) \\
&= H(Y) - H(Y|X).
\end{aligned}$$

所以， X 中包含的 Y 的信息和 Y 中包含的 X 的信息是一样多的！

此外，直观上我们还应该觉得，信息量不能是负的，实际上的确如此：

命题 3.7 $H(X) - H(X|Y) \geq 0$ ，等号成立当且仅当 X 和 Y 相互独立。

我们将在第 3.2 节看到，命题 3.7 就是 K-L 散度信息不等式的一个特例，所以我们就不在这里给出证明了。命题 3.7 表明知道任何信息都不会增加不确定性，这个原理被称为“Information doesn’t hurt.”

根据以上讨论，我们可以自然地定义 X 和 Y 的互信息为

$$I(X; Y) = I(Y; X) = \mathbb{E}[I(X = x : Y)] = \mathbb{E}[I(Y = y : X)].$$

类似联合分布的熵，条件熵和互信息的概念也可以推广到多元情形。对于三个随机变量 X, Y, Z ，我们可以定义条件熵为

$$H(X, Y|Z) = H(X, Y, Z) - H(Z).$$

类似地，我们可以定义互信息为

$$I(X, Y; Z) = H(X, Y) - H(X, Y|Z).$$

他们的含义以及性质和二元情形类似。

同样，我们可以定义条件互信息为

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z),$$

表明 Z 已知时候 Y 给 X 带来的平均信息增益. 类似互信息，我们如下性质：

命题 3.8 条件互信息满足以下性质：

1. 非负性： $I(X; Y|Z) \geq 0$ ，等号成立当且仅当 X 和 Y 在给定 Z 的条件下相互独立.
2. 对称性： $I(X; Y|Z) = I(Y; X|Z)$.
3. 链式法则： $I(X, Y; Z) = I(X; Z|Y) + I(Y; Z)$.
4. 条件信息量： $I(X : X|Y) = H(X|Y) - H(X|X, Y) = H(X|Y)$.

最后一条性质说的其实是，在平均的意义下，给定 Y 的时候，知道 X 所能够得到的额外信息量就是 $H(X|Y)$. 这一命题的证明和前面都非常相似，见习题[\[lhy: 出一下\]](#).

最后，我们将各种熵以及信息量的关系总结为图 3.3. 在集合论中，这样的图被称为 Venn 图，所以我们可以用集合论来理解信息与熵. 对应关系可以总结为表 3.1.

信息论	集合论
$H(X)$	A
$H(Y)$	B
$H(X Y)$	$A \setminus B$
$H(X,Y)$	$A \cup B$
$I(X;Y)$	$A \cap B$

表 3.1: 信息论和集合论的对应关系.

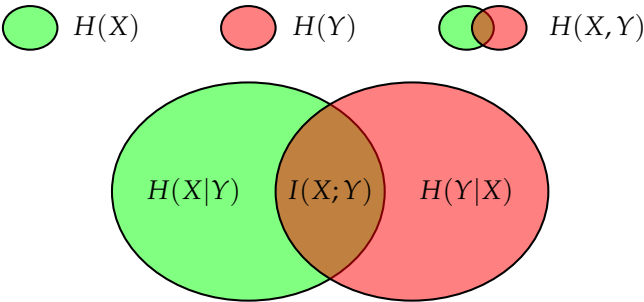


图 3.3: 熵和信息量的关系

§3.2 *Kullback-Leibler* 散度

§3.2.1 定义

为了引入 *K-L* 散度，我们从互信息出发. 它的定义是：

$$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) \\ &= -\sum_x p_X(x) \log p_X(x) + \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_Y(y)} \\ &= -\sum_{x,y} p(x,y) \log p_X(x) + \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_Y(y)} \end{aligned}$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_X(x)p_Y(y)}.$$

根据命题 3.7, $I(X;Y) \geq 0$, 等号成立当且仅当 X, Y 相互独立, 即 $p(x,y) = p_X(x)p_Y(y)$. X, Y 之间的互信息越大, 说明他们之间的关联越强, 分布越不独立, $p(x,y)$ 越不接近 $p_X(x)p_Y(y)$.

上面的推导说明, 互信息其实在用分布列的比值比较两个分布的接近程度, 这样的想法可以被推广到一般分布上. 考虑两个概率分布的似然函数 p_1 和 p_2 (也就是他们的分布列). 抽取一个样本 X , 考虑如下两个假设:

H_1 : 样本 X 来自 p_1 的分布 vs. H_2 : 样本 X 来自 p_2 的分布

为了判断哪个假设是更有可能的, 我们考虑两个假设分布的似然比 p_1/p_2 . 如果这个比值越大, 就越说明 p_1 的值更大, 因而更有可能, 倾向于接受 H_1 , 反之则越倾向于接受 H_2 . 这种方法被称之为似然比检验法.

从上述讨论出发, 我们定义区分 H_1 和 H_2 的检验量为对数似然比:

$$\log(p_1(x)/p_2(x)).$$

假设 H_1 是真的, 那么在 H_1 成立的世界里, 这个检验量的期望为

$$\mathbb{E}_{X \sim p_1}(\log(p_1(X)/p_2(X))) = \sum_i p_1(i) \log \frac{p_1(i)}{p_2(i)}.$$

期望越大, 说明 H_1 越有可能成立. 实际上, 上面的期望就是 K-L 散度的定义.

定义 3.3 (Kullback-Leibler 散度, 相对熵) 对于两个概率分布 p_1, p_2 , 他们的 **Kullback-Leibler 散度 (相对熵)** 定义为

$$D_{\text{KL}}(p_1 \| p_2) = \mathbb{E}_{X \sim p_1}(\log(p_1(X)/p_2(X))) = \sum_i p_1(i) \log \frac{p_1(i)}{p_2(i)}.$$

其中规定 $0 \log(0/0) = 0$, $0 \log(0/a) = 0$, $a \log(a/0) = +\infty$. □

我们马上知道, 互信息是 K-L 散度的一种特殊情况:

命题 3.9 对于两个随机变量 X, Y , 成立 $I(X; Y) = D_{\text{KL}}(p_{X,Y} \| p_X p_Y)$, 其中 $p_{X,Y}$ 是 X, Y 的联合分布列, p_X, p_Y 分别是 X, Y 的边缘分布列.

K-L 散度可以看成两个分布之间的区分衡量标准, 但他不是度量. 一般来说, 甚至连对称性都不成立. 例如, 设 p_1 和 p_2 都是定义在 $0, 1$ 上的 Bernoulli 分布, 参数分别为 $1/2$ 和 $1/4$. 于是

$$D_{\text{KL}}(p_1 \| p_2) = \frac{1}{2} \log \frac{1/2}{3/4} + \frac{1}{2} \log \frac{1/2}{1/4} = \frac{1}{2} \log \frac{4}{3}.$$

$$D_{\text{KL}}(p_2 \| p_1) = \frac{3}{4} \log \frac{3/4}{1/2} + \frac{1}{4} \log \frac{1/4}{1/2} = \frac{1}{2} \log \frac{3\sqrt{3}}{4}.$$

这两个值是不相等的. 进一步, 这个差值甚至可以到任意大 (见习题[\[lhy: 出一下\]](#)).

上面推导 K-L 散度的过程看起来有些捏造, 我们将在第 3.3 节中给出一个非常直观的理解方式. 现在, 我们先接受这个定义, 然后看一下 K-L 散度的一些性质.

§3.2.2 两个关于信息的不等式

利用 K-L 散度，我们可以给出两个关于信息的不等式，它们分别是信息不等式和数据处理不等式。

定理 3.2 (信息不等式) 对于两个概率分布列 p, q ，成立 $D_{\text{KL}}(p\|q) \geq 0$ ，当且仅当 $p = q$ 时取等号。

证明. 由于 $\log x$ 是凸函数，所以由 Jensen 不等式，我们有

$$D_{\text{KL}}(p\|q) = -\mathbb{E}_{X \sim p} \left[\log \frac{q(X)}{p(X)} \right] \geq -\log \mathbb{E}_{X \sim p} \left[\frac{q(X)}{p(X)} \right] = -\log \sum_i p(i) \cdot \frac{q(i)}{p(i)} = 0.$$

因此， $D_{\text{KL}}(p\|q) \geq 0$ ，当且仅当 $p = q$ 时取等号。 \square

信息不等式表明，K-L 散度虽然不是度量，但却是非负的，因而确实可以被作为熵，用来衡量“额外的不确定性”。此外，命题 3.7 是信息不等式的直接推论。利用类似的方法，我们可以证明条件互信息的非负性（即命题 3.8 中的第一条）。

接下来我们叙述并证明数据处理不等式。

定理 3.3 (数据处理不等式) 假设随机变量 X, Y, Z 形成了 Markov 链，那么 $I(X; Y) \geq I(X; Z)$ 。特别地，对任意函数 f ，成立 $I(X; Y) \geq I(X; f(Y))$ 。

证明. 根据互信息链式法则，

$$\begin{aligned} I(X; Y, Z) &= I(X; Z) + I(X; Y|Z) \\ &= I(X; Y) + I(X; Z|Y). \end{aligned}$$

根据 Markov 性，条件在 Y 上， X 和 Z 相互独立。因此， $I(X; Z|Y) = 0$ ，根据条件互信息的非负性， $I(X; Y|Z) \geq 0$ ，所以 $I(X; Y) \geq I(X; Z)$ 。

显然, $X, Y, f(Y)$ 也形成了 Markov 链, 所以 $I(X; Y) \geq I(X; f(Y))$. \square

数据处理不等式表明, 无论我们对随机变量 Y 进行了何种处理, 甚至是允许带随机的处理, 它的信息量都不会增加.

§3.3 编码理论

最早的时候, Shannon 建立信息论, 就是为了给通信和编码理论一个数学基础. 从编码的角度出发, 我们可以更本质地理解信息和熵.

§3.3.1 熵与编码

通信就是一个发射端和一个接收端, 中间有信道传递消息. 将所有可能要传递的消息集合记为 Ω (一个有限集), 我们现在考虑 Ω 所蕴含的信息量是多少. 注意, 根据 Shannon 的思想, Ω 里面具体是什么并不重要, 重要的是有多少个. 我们可以用自然数 $1, 2, \dots$ 表示集合 Ω 里的元素.

使用二进制编码, 我们至少需要 $\log_2 |\Omega|$ 个比特来表示 Ω 里的元素. 于是, 假如说随机变量 X 表示收到的消息, 那么 X 的熵就定义为 $H(X) = \log_2 |\Omega|$, 它衡量了接收端收到的消息的不确定性.

当我们选定了具体的消息 $m \in \Omega$, X 的不确定性被消除了, 于是 $X = a$ 的过程产生了 (或者说传递了) $\log_2 |\Omega|$ 比特的信息. 比如说, 我们发送一个长为 n 的二进制序列, 消息的集合大小就是 2^n , 发送任何一条具体的消息, 我们就传递了 n 比特的信息.

有时候, 我们会把消息看成一个序列. 具体来说, 我们可以发送独立的 k 条消息, 其中第 i 条 X_i 来自消息集合 Ω_i , $|\Omega_i| = n_i$, 那么

(X_1, \dots, X_k) 的熵就是

$$H(X_1, \dots, X_k) = \log_2 n_1 + \dots + \log_2 n_k,$$

它衡量了 k 条消息的不确定性.

在更常见的情况下, 每次发送的其实不是一条消息, 而是一个字母, 所有的字母组成了一个字母表, 我们用 $\Sigma = \{x_1, \dots, x_s\}$ 来表示. 于是, X_i 就是消息的第 i 个字母, 于是, 一条消息可以写作 $X_1 \dots X_k$, 其中每一个 X_i 都来自 Σ .

我们现在考虑更简单的情形, 即每个字母 X_i 其实是同一个随机变量 X 的独立采样. 如果具体知道某一个 x_i 出现的次数, 我们其实可以有更高效的传递信息的方式. 比如说, 在极端情况下, 如果只有 x_1 和 x_2 会出现, 那么我们其实只需要 $\log_2 2 = 1$ 比特就足够传递所有消息了: 0 表示 x_1 , 1 表示 x_2 .

在一般情况下, 考虑 Ω 中只包含长为 k 的消息, 并且 x_i 在消息中出现 k_i 次, 那么所有可能的消息数量为

$$|\Omega| = N(k) = \frac{k!}{k_1! \dots k_s!}.$$

假定我们需要 $h(\omega)$ 比特来具体确定发的消息是 ω , 我们来推导 $h(\omega)$ 的上下界.

- 无论如何编码, 我们需要能区分 Ω 中的不同元素, 这本身需要 $\log_2 \Omega$ 比特来表示.
- 在一些编码中, 我们还需要确定 (k_1, \dots, k_s) . 确定它的一种方式是按照顺序给出每一个 k_i . 因为 $k_i \leq k$, 每个 k_i 最多需要 $\log_2 k$

比特来表示, 所以按顺序表示所有的 k_i 至多需要 $s \log_2 k$ 比特.

于是, 我们需要的比特数为

$$\log_2 \frac{k!}{k_1! \cdots k_s!} \leq h(m) \leq s \log_2 k + \log \frac{k!}{k_1! \cdots k_s!}.$$

这刚好和我们在统计力学中推导熵的过程是一致的! 假设消息足够的长, x_i 出现的频率逐渐接近 p_i , 那么同样的推理我们可以知道,

$$h(m) \sim -k \sum_i p_i \log_2 p_i = k H_2(p_1, \dots, p_s).$$

因此, 如果知道字母的出现频率, 我们传递单位长度的消息至少需要 $H(p_1, \dots, p_s)$ 比特, 这完全给出了熵的具体含义, 而且, 我们现在也不难理解熵的形式为何会出现 \log 了: 熵就是期望上编码一个字母需要的比特数 (即 $\log(1/p(X))$).

那么, 是否有一种编码确实达到了这个理论上的编码长度下界呢? 答案是肯定的, 它被称为 *Huffman* 编码. 它的核心思想在于把出现频率高的字母用更短的编码表示. 类似的思想被用在了机器学习的决策树中, 作为选择节点非常常用的一种依据.

注. 决策树是一种常用的机器学习分类模型. 假设数据有很多属性 P_1, \dots, P_k , 这些属性共同决定了某一条数据的类别. 比如, 在银行的信用系统中, 给定了一个人的性别、是否已婚、是否负债等信息, 我们希望给他评估一个信用评级.

决策树的做法是, 将决策过程写成一棵树, 然后叶节点是决策类别的结果. 比如说, 我们会先看这个人是否负债, 如果不负债, 那么看是否已婚, 如果已婚, 那么我们信用评级就给 A. 但如果负债, 那么我们信用评

级就给 B.

那么, 每个节点应该判断什么属性呢? 树本身其实就是一种广义的消息, 沿着树, 从根节点走到叶节点得到的就是一条消息. 于是, 在这一观点下, 我们可以用熵与编码的关系来选择属性.

如果我们选择编码最短的属性, 这样我们的决策树就会更加简单. 一种近似的做法是, 对于每个节点, 都优先选择信息增益最高的属性. 这样的选择方式叫做 *ID3* 策略.

我们进一步的问题是, 为什么我们知道了每个字母的频次就可以压缩编码? 我们接下来将要说明, 其实长为 k 的消息中的“典型消息”的数量远远少于所有 k 长消息的数目, 因此我们实际上相当于只是针对一个子集进行编码. 注意到, 当 k 充分大的时候,

$$\log_2 N(k) \sim h(m) \sim kH_2(p_1, \dots, p_s).$$

因此,

$$N(k) \approx 2^{kH_2(p_1, \dots, p_s)} = e^{kH(p_1, \dots, p_s)}.$$

然而, 长为 k 的所有消息数目为

$$s^k = e^{k \log s}.$$

根据命题 3.4, 只有当所有 p_i 相等的时候 $N(k)$ 才会达到这一量级. 从这个意义上说, 熵所刻画的信息量定量刻画了数据压缩可能的极限.

以上关于信息编码下界以及数据压缩的讨论, 再更一般的情况下也成立, 此时这样的性质被称为渐近等分性. 而这一性质成立对应的结果被称为 *Shannon-McMillan-Breiman* 定理, 它的陈述以及证明都需要用到更多随机过程的知识, 这里就不再给出了.

§3.3.2 K-L 散度、交叉熵与编码

我们在 K-L 散度的定义中提到了它的另一个名字——相对熵. 实际上, 这可以从编码中看出来. 假设事实上消息中字母的分布是 p_1 , 那么期望上编码单位长度消息需要的比特数是

$$H(p_1) = \mathbb{E}_{X \sim p_1} [\log p_1(X)].$$

如果我们错误地认为消息中字母的分布是 p_2 并使用最优编码, 那么实际上期望编码单位长度消息需要的比特数是

$$\mathbb{E}_{X \sim p_1} [\log p_2(X)].$$

由于错误的认识所产生的额外编码长度是

$$\mathbb{E}_{X \sim p_1} [\log p_1(X) - \log p_2(X)] = D_{\text{KL}}(p_1 \| p_2).$$

根据本节中的讨论, 我们知道, 额外的编码长度代表的是额外的不确定性, 因而这一概念是某种“熵”的概念. 这正是“相对熵”的由来, $D_{\text{KL}}(p_1 \| p_2)$ 表示了当我们错误地把 p_1 当成 p_2 时带来的额外的不确定性, 或者说额外的信息损失.

在机器学习中, 比起讨论 K-L 散度, 更加常用的是直接讨论量 $\mathbb{E}_{X \sim p_1} [\log p_2(X)]$. 从机器学习的观点来说, p_1 是真实的分布, 而 p_2 是我们所学习到的分布. 根据刚刚的讨论, 这个量越小越说明 p_2 接近真实的 p_1 , 因此这又是一种衡量两个分布之间关系的量, 我们称之为交叉熵:

定义 3.4 (交叉熵) 给两个随机变量 X, Y , X 的分布为 p_X , Y 的分布为

p_Y , 则 X 的分布 p_X 和 Y 的分布 p_Y 的交叉熵¹为

$$CH(p_X, p_Y) = -\mathbb{E}_{X \sim p_X}[\log p_Y(X)] = -\sum_i p_X(i) \log p_Y(i). \quad \square$$

在机器学习的分类问题中, 我们希望学习到的分布 p_Y 尽可能地接近真实的分布 p_X , 所以我们训练的目标经常是最小化交叉熵 $CH(p_X, p_Y)$. 有趣的是, 从数理统计的角度来看, 最小化交叉熵等价于进行最大似然估计 (见习题[[lhy: 出一下](#)]), 因此这为最大似然估计提供了一种信息论意义下的理解.

注. 现代的主流信息论都是从 Shannon 发展起来的. 然而, 这一信息论也有很多问题.

- 信息论使用了概率论进行建模. 但我们已经看到, 概率要么是作为频率的近似理论 (频率学派), 要么反映了人们对未知的信念 (Bayes 学派). 无论哪种解释, 都将问题简化了. 正如 Kolmogorov 所说: “如果事情没有按照我们的预期发展, 那么问题一定出在我们对于概率和真实世界的随机之间关系不清晰的认识上.”
- 这一信息论考虑的是一族对象的信息. 我们是否能够用这样的方式来衡量单个对象的信息量呢? 比如, 我们要考虑这本书中包含的信息量, 是它放在所有可能的书的集合中去考虑呢, 还是把它的每一个章节分开考虑成一个随机序列呢? 因此, 信息论并不能很好地回答“单个对象”的信息量的问题.

现代概率论的奠基人 Kolmogorov 也非常严肃地考虑了这一问题. 他提出了被后世称为 **Kolmogorov 复杂度** 的概念, 旨在刻画一个随机字符串

¹文献中, 经常会直接写为 $H(p_X, p_Y)$, 但是在本书中为了区分熵, 我们使用了符号 CH .

的随机程度. 简单来说, 一个字符串的 Kolmogorov 复杂度就是描述输出它所需要的最短代码长度. 越随机的字符串就越需要更复杂的程序去描述它的输出方式.

例如, 尽管字符串 $x = 0101010101$ 看起来非常长, 但是我们可以用一个很短的程序来描述它: 输出 5 次“01”. 然而, 尽管字符串 $y = 011001$ 比 x 短得多, 我们却很难找到一个简短的程序来描述它. 因此, y 的 Kolmogorov 复杂度要比 x 大, 因而看起来更像是随机的.

利用这一概念, 我们可以将信息的概念变成一个对象自己的属性, 而不再需要把对象放在可能的一堆对象中去考虑. 这是信息论的另一种构建思路.

§3.4 在机器学习中的应用：语言生成模型

现如今, 机器学习中最为瞩目的成果之一就是大语言模型 (LLM), 它通过学习人类海量的高质量语料库来形成一个生成式的模型, 其中最为典型的例子是 ChatGPT.

从思路上来说, 大语言模型的核心思想非常简单: 给一段话, 将其中一些词掩盖掉, 让模型填出这些词来. 例如, 给出

“我在[mask]面条, 它真好吃.”

模型应该能够填出

“我在吃面条, 它真好吃.”

对于 GPT 模型来说, 这一思想更加简单: 永远只预测下一个词. 它的哲学是“通过预测下一个词, 可以理解世界.”

这样的思想，对于更一般的数据也是成立的：用（修改过的）数据本身作为输入，训练一个编码器，然后将编码器的输出送入解码器，而解码器的输出具有原始数据的格式，我们希望这一输出能够尽量匹配原始的输入。这正是本章开头 Simon 所说的“信息处理器”的具体实现。

在语言模型中，一个生成模型往往同时有编码器和解码器。比如说，图 3.4 展示的就是 BART [LLG⁺19] 的结构。

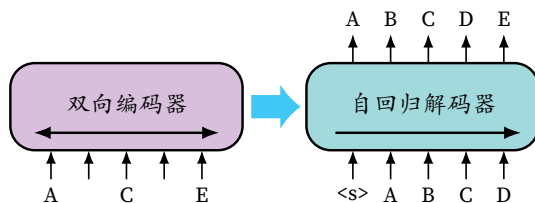


图 3.4: 生成式语言模型 BART 的示意图。

有的时候，编码器和解码器并不是显式给出的。例如，GPT 模型只有自回归解码器，没有编码器。然而，我们可以认为 GPT 的解码器实际上是一个编码器-解码器的结构：它总要处理输入的数据，因而需要编码器，同时也要输出数据，因而需要解码器。

在第 3.3 节我们指出，熵和编码有着密切的联系。从这个角度出发，我们很容易理解生成模型背后的思想：我们希望通过训练的方式得到一个由神经网络所表示的编码和解码规则，他要尽可能符合真实数据的分布。

我们可以用一种非常简单的模型去理解这一过程。假设所有的单词的集合为 Σ ，单词数为 k 的文本集合为 Ω 。我们希望训练一个生成模型 M ，给它输入 $k-1$ 个单词，它可以给出第 k 个单词的概率分布，我们选择出现概率最大的那个词作为预测。

在训练的时候，对于一个句子 ω ，我们只保留前 $k-1$ 个词，得到

$\omega[1:k-1]$, 然后将它输入到生成模型 M 中, 让它去预测第 k 个词.

对于这一个具体的句子 ω 来说, 理想的分布应该是一个 Dirac 分布² $\delta_{\omega[k]}$, 即以概率 1 取到 $\omega[k]$. 假如说生成模型的输出是一个概率分布 $M(\omega[1:k-1]) = p$, 那么, 我们可以用 K-L 散度去衡量这两个分布的差异.

因为 $H(\delta_{\omega[k]})$ 是固定的, 所以我们只用考虑交叉熵 $CH(\delta_{\omega[k]}, p)$. 一次训练会给多个样本, 所以我们的目标是同时最小化这些交叉熵的和. 假如训练集是 T , 我们的目标就是

$$\min_M \sum_{\omega \in T} CH(\delta_{\omega[k]}, M(\omega[1:k-1])).$$

实际上, 这个例子是有普适性的, 所有的监督训练的分类问题都可以用这种方式来建模. 而在第 6.1 节我们也会看到, 此时交叉熵实际上被作为了一种损失函数.

§3.5 附录: Shannon 定理的证明

我们在这一部分给出 Shannon 定理 (定理 3.1) 的证明. 整体上的思路是:

1. 证明如果 f 是单调函数, 对正整数 m, n 成立 $f(mn) = f(m) + f(n)$, 那么 $f(n) = C \log n$.
2. 求出 $H(1/n, \dots, 1/n)$ 的表达式.

²Dirac 分布是一个数学物理中更加常用的名字. 在概率论中, 这也被称为退化分布; 而在机器学习中, 分布经常会表示为一个概率向量, 文献中称为独热向量.

3. 假设 p_i 是有理数, 设 $p_i = n_i / \sum_j n_j$, 考虑 $\sum_j n_j$ 个等可能试验结果, 利用假设 3 推出 H 的表达式.
4. 利用有理数的稠密性和 H 的连续性推出一般情形.

最后一步是显然的, 我们只需要证明前三步即可.

对第一步, 我们需要证明的是, 如果 f 是单调函数, 对正整数 m, n 成立 $f(mn) = f(m) + f(n)$, 那么 $f(n) = C \log n$. 首先, 利用数学归纳法容易看出, 对正整数 n, k , 成立

$$f(n^k) = kf(n). \quad (3.2)$$

设 m, n 是任意两个大于 1 的整数, 再选任意的正整数 k , 从 m 进制数的性质可以看出, 总存在正整数 l 使得

$$m^l \leq n^k < m^{l+1}. \quad (3.3)$$

根据 f 的单调性, 我们有

$$f(m^l) \leq f(n^k) < f(m^{l+1}).$$

利用式 (3.2), 我们有

$$lf(m) \leq kf(n) < (l+1)f(m) \iff \frac{l}{k} \leq \frac{f(n)}{f(m)} < \frac{l+1}{k}.$$

将式 (3.3) 取对数, 得到

$$l \log m \leq k \log n < (l+1) \log m \iff \frac{l}{k} \leq \frac{\log n}{\log m} < \frac{l+1}{k}.$$

所以

$$\left| \frac{\log n}{\log m} - \frac{f(n)}{f(m)} \right| \leq \frac{1}{k}.$$

因为 k 可以是任意大的正整数, 取 $k \rightarrow \infty$, 我们就得到了

$$\frac{\log n}{\log m} = \frac{f(n)}{f(m)}.$$

由 m, n 的任意性, 取 $m = 2$, 我们就得到了 $f(n) = (f(2)/\log 2) \cdot \log n = C \log n$. 容易检验, $f(1) = 0 = C \log 1$, 因此这一等式对所有正整数 n 都成立.

对第二步, 我们需要求出 $f(n) = H(1/n, \dots, 1/n)$ 的表达式. 我们要利用第一步的结果, 首先, 根据假设二, $f(n)$ 是单调递增的函数. 然后, 考虑 mn 个等可能试验, 我们可以将它分成两步试验, 第一步有 m 中等可能的结果, 而在每一种结果之下, 第二步有 n 种等可能结果. 根据假设三,

$$f(mn) = f(m) + \frac{1}{n} \cdot n f(n) = f(m) + f(n).$$

所以 $f(n)$ 符合第一步的假设. 第二步就可以直接从第一步推出.

最后, 我们证明第三步. 设 p_1, \dots, p_n 都是有理数, 那么, 他们可以被写为

$$p_i = \frac{n_i}{\sum_{j=1}^n n_j}.$$

其中 n_i 是非负整数. 我们考虑 $\sum_j n_j$ 个等可能试验, 这个试验可以被看成两步的试验, 第一步有 n 种可能的结果, 第 i 种结果出现的概率是 p_i , 而在第 i 种结果之下, 第二步有 n_i 种等可能的结果. 根据假设三, 和证

明的第三步，我们有

$$C \log \sum_{j=1}^n n_j = H(p_1 + \cdots + p_n) + \sum_{i=1}^n p_i \cdot C \log n_i.$$

因此，

$$\begin{aligned} H(p_1, \dots, p_n) &= C \left(\log \sum_{j=1}^n n_j - \sum_{i=1}^n p_i \log n_i \right) \\ &= C \left(\log \sum_{j=1}^n n_j - \sum_{i=1}^n p_i \log \left(p_i \sum_{j=1}^n n_j \right) \right) \\ &= C \left(\log \sum_{j=1}^n n_j - \sum_{i=1}^n p_i \log p_i - \sum_{i=1}^n p_i \log \sum_{j=1}^n n_j \right) \\ &= -C \sum_{i=1}^n p_i \log p_i. \end{aligned}$$

这正是我们要证明的。于是，我们证明了 Shannon 定理。

§3.6 习题

1. 我们在熵以及 K-L 散度的定义中，都规定了一些无定义的量的值，这些值并不是随便规定的，他们实际上反映了熵或者 K-L 散度定义中的连续性。

- (1) 证明：对给定的 $a > 0$ ， $\lim_{x \rightarrow 0+} x \log(x/a) = 0$ ，因此我们规定了 $0 \log 0 = 0$ 以及 $0 \log(0/a) = 0$ 。
- (2) 证明：对给定的 $a > 0$ ， $\lim_{x \rightarrow 0+} x \log(a/x) = +\infty$ ，因此我们规定了 $0 \log(a/0) = +\infty$ 。

2. 考虑关于 n 的正实数序列 $a_1(n), \dots, a_k(n)$ 以及 $b_1(n), \dots, b_k(n)$, 假设对所有 i , 都成立 $\lim_{n \rightarrow \infty} a_i(n)/b_i(n) = 1$, 证明:

$$\lim_{n \rightarrow \infty} \frac{a_1(n) + \dots + a_k(n)}{b_1(n) + \dots + b_k(n)} = 1.$$

由此证明式 (3.1).

3. 证明命题 3.1.

4. 用 Lagrange 乘子法重新证明命题 3.4.

提示: 如果你不知道 Lagrange 乘子法, 可以参考 ??.

5. 证明命题 3.8.

6. [?] 仿照集合论的思路, 我们可以定义三个随机变量的互信息为:

$$I(X; Y; Z) = I(X; Y) - I(X; Y|Z).$$

(1) 证明对称性: $I(X; Y; Z) = I(Y; X; Z) = I(X; Z; Y)$.

(2) 举一个例子说明, 可能会有 $I(X; Y; Z) < 0$, 所以这样定义的互信息并不一定真的代表“信息量”.

7. 举一个例子说明, 即便 $D_{\text{KL}}(p_1 \| p_2)$ 很接近 0, $D_{\text{KL}}(p_2 \| p_1)$ 也可能很大.

8. (单变量数据处理不等式) 对任意离散随机变量 X 和函数 f , 证明: $H(X) \geq H(f(X))$.

9. 考虑二分类的学习问题, 此时对单个样本我们观察到的结果要么是 0 或 1, 假设在真实世界中样本总体服从参数为 θ 的 Bernoulli 分

布, 即 $\Pr(X = 1) = 1 - \Pr(X = 0) = \theta$. 假设我们的数据集是 $(x_1, y_1), \dots, (x_N, y_N)$, 他们是从总体中独立采样得到的.

- (1) 将问题考虑成一个数理统计问题, 估计 θ . 写出似然函数 $L(\theta; y_1, \dots, y_N)$.
- (2) 再将问题考虑为一个信息论问题, 写出每个样本的真实分布与估计分布之间的交叉熵之和 $CH(\theta; y_1, \dots, y_N)$.
- (3) 证明: $\max_{\theta} L(\theta; y_1, \dots, y_N) = \min_{\theta} CH(\theta; y_1, \dots, y_N)$, 也就是说, 最大似然估计等价于最小化交叉熵.

10. 请查找文献回答以下问题:

- (1) Fisher 信息量是什么? 它与 K-L 散度有什么样的关系?
- (2) 列举其他概率分布之间散度的概念, 他们是否是度量?
- (3) 列举概率分布之间的度量, 他们之间是否有关联?

§3.7 章末注记

信息一词的英文是“information”, 从动词“inform”来, 意思是告知、通知. 早在 15 世纪中叶, “information”一词的出现了义项“在通信中针对特定主题的知识”.[\[Inf\]](#) 这说明在那个时候人类就已经意识到, 通信会产生新的东西, 被称为知识或信息. 然而, 人类对信息的严谨探索起步晚得多. 关于信息的物理学讨论源自统计力学, Boltzmann 提出了著名的熵, 证明了 H 定理, 以此给出了热力学第二定律的微观解释. 关于 Boltzmann 的工作, 参见 [\[Uff22\]](#).

一般认为, 现代信息论的起源是 Shannon 的论文 [\[Sha48\]](#), 他在论文中提出了信息的数学定义, 以及信息的基本性质. 但是, Shannon 的工

作并不是孤立的，他的工作是在统计力学的基础上发展起来的。事实上，Shannon 在论文中也提到了 Boltzmann 的熵，这篇工作也被视为通信理论以及编码理论的奠基性工作。Shannon 在这篇论文中还给出了渐近意义下达到理论下界的最优编码，并且独立地被 Fano [Rob49] 以不同的形式发现，因此后世称为 Shannon-Fano 编码。但是 Shannon-Fano 编码并不是精确地达到下界，实际上，最优编码是 Huffman [Huf52] 给出的。Shannon 在这篇论文中还讨论了渐近等分性，后来 McMillan 的工作 [McM53] 和 Breiman 的工作 [Bre57] 拓展了这一结果，因此后世称为 Shannon-McMillan-Breiman 定理。

关于信息论与集合论的关系工作，可以参见 Hu Kuo Ting 的工作 [?]. 他的工作还给出了多个随机变量互信息的定义，在这一章习题中有涉及。

相对熵的概念依然是从 Shannon 的奠基性论文 [Sha48] 中提出的，但他只局限于通信的问题，更加一般的讨论是由 Kullback 和 Leibler 在 [KL51] 给出，他们的是一种数理统计的思路，但是他们也具体地讨论了这一概念与信息的关系。他们的论文中也讨论了交叉熵这一概念。

机器学习中编码器和解码器的思路，最早是由 Rumelhart, Hinton 和 Williams 在 [RHW86] 中提出，他们将编码器和解码器的整体称作自编码器。这篇工作几乎可以被视为深度学习的开山之作，它还提出了训练神经网络最常用的反向传播算法。

关于信息论的经典教科书，可以参见 [CT12]，此外，概率论的教材中也有很多很好的讨论，比如 [Jay02]，[Shi96] 以及 [李 10]。

关于 Kolmogorov 复杂度的讨论，可以参见专著 [?], 这本书对于随机、信息、编码、复杂度，乃至归纳推理等概念都有非常独到的见解，值得一读。

第三部分

决策与优化

第四部分

博弈与逻辑

第五部分

认知与逻辑

第六部分

附录：预备知识

参考文献

- [Bre57] Leo Breiman. The Individual Ergodic Theorem of Information Theory. *The Annals of Mathematical Statistics*, 28(3):809–811, 1957.
- [CT12] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- [Huf52] David A. Huffman. A Method for the Construction of Minimum-Redundancy Codes. *Proceedings of the IRE*, 40(9):1098–1101, September 1952.
- [Inf] Information | Etymology, origin and meaning of information by etymonline. <https://www.etymonline.com/word/information>. (accessed 2023-07-10).
- [Jay02] Edwin T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, 2002.

- [KL51] S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86, 1951.
- [LLG⁺19] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension, October 2019.
- [McM53] Brockway McMillan. The Basic Theorems of Information Theory. *The Annals of Mathematical Statistics*, 24(2):196–219, June 1953.
- [RHW86] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1: Foundations*, pages 318–362. MIT Press, Cambridge, MA, USA, January 1986.
- [Rob49] Robert M. Fano. *The Transmission of Information*. March 1949.
- [Sha48] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, July 1948.
- [Shi96] A. N. Shiryaev. *Probability*, volume 95 of *Graduate Texts in Mathematics*. Springer, New York, NY, 1996.

- [Uff22] Jos Uffink. Boltzmann's Work in Statistical Physics. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2022 edition, 2022.
- [李 10] 李贤平. 概率论基础. 高等教育出版社, 2010.