

高维向量空间中的密度估计算法综述

1. 基础挑战：高维空间中的密度

问题导论

在高维数据分析领域，密度估计是一项基础且关键的任务。本文旨在探讨在处理大规模、高维度数据集（例如，包含数百万个3072维向量）时进行密度估计的有效算法。具体而言，我们关注两个核心问题：1) 针对一个新的数据点，如何计算其在现有数据库中的局部密度；2) 如何为整个数据集构建一个全局的概率密度函数 (Probability Density Function, PDF)，即一个输入为3072维向量、输出为该点概率密度的函数。

这两个问题都面临一个共同的、根本性的障碍，即由Richard Bellman提出的“维度灾难”(Curse of Dimensionality)。该理论指出，随着数据维度的增加，为了维持统计上的可靠性，所需的数据量会呈指数级增长。¹ 这一现象从根本上改变了我们对空间、距离和密度的直观理解，使得低维空间中行之有效的方法在高维空间中彻底失效。

高维空间的几何特性

为了理解高维密度估计的困难，必须首先摒弃基于二维或三维空间的直观经验，并审视高维空间独特的、有悖常理的几何特性。

指数级体积增长与数据稀疏性

当维度 d 增加时，空间的体积会以惊人的速度膨胀。例如，在一个单位超立方体中，为了保持数据点之间的平均距离不变，所需的数据点数量会随着维度的增加而指数级增长¹。这种现象也被称

为组合爆炸(combinatorial explosion)¹。对于一个包含数百万个样本的有限数据集,无论其规模多大,当维度高达3072时,数据点都会变得极其稀疏。这意味着任何一个数据点的“局部邻域”(local neighborhood)几乎都是空的。因此,依赖于在局部区域内进行数据点计数或平均的传统密度估计方法,从一开始就失去了其理论基础。

测度集中现象

高维空间的几何结构还表现出一种被称为“测度集中”(Concentration of Measure)的悖论现象。对于一个高维超球体,其绝大部分体积都集中在靠近其表面的一个薄壳区域内³。一个形象的比喻是“高维橙子”:几乎所有的果肉都在果皮里,内部几乎是空的³。更令人惊讶的是,一个单位半径的d维球体的体积会随着维度

d趋向于无穷大而趋向于零³。同样,对于一个高维超立方体,其大部分体积都集中在远离中心的“角落”里³。

这一现象对数据分布有着深远的影响:在高维空间中,数据点不仅彼此之间相距遥远,而且它们距离其自身分布的中心也非常遥远⁸。每个数据点都像是一个离群点,这使得区分真正异常的低密度点和由于空间稀疏性而显得孤立的正常点变得异常困难。

距离度量的失效

维度灾难最直接的后果之一是传统距离度量(如欧氏距离)的失效。随着维度d的增加,任意两个随机数据点之间的距离趋于相等⁸。一个查询点的最近邻和最远邻之间的距离差会趋近于零,这使得“最近邻”这一概念失去了区分能力¹⁰。

这种距离度量的失效,不仅仅是一个计算上的问题,它动摇了许多机器学习算法的根基。传统的非参数统计方法,如核密度估计(KDE)和k近邻(k-NN)密度估计,都建立在一个核心假设之上:空间上的“邻近性”等同于性质上的“相似性”。在高维空间中,由于所有点之间的距离都变得相似,这个假设被彻底打破。任何试图通过观察一个点的“局部邻域”来估计其密度的尝试都注定会失败,因为这个“局部”邻域要么是空的,要么必须扩展到包含整个数据集的很大一部分,从而失去了“局部”的意义¹¹。因此,任何希望在高维空间中成功进行密度估计的算法,都必须采用新的策略:要么找到一种方法在更低维的有效空间中恢复“局部性”的概念(例如,通过流形学习或结构性假设),要么采用不依赖于距离度量的方法,要么足够强大以至于能够直接学习整个分布的全局结构,而无需依赖局部假设。

2. 经典方法的必然失效

基于前一章节对高维空间几何特性的分析，本节将详细阐述为何两种最经典的非参数密度估计方法——核密度估计(KDE)和k近邻(k-NN)密度估计——在高维场景下会不可避免地失效。

核密度估计(KDE):固定带宽方法

KDE通过在每个数据点 x_i 上放置一个核函数 K (通常是高斯核)，并将它们叠加起来，来估计概率密度函数。其数学表达式为：

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N K_h(x - x_i)$$

其中， N 是样本总数， h 是至关重要的带宽参数，它控制着核函数的平滑程度¹²。

在低维空间中，KDE是一种强大而灵活的工具。然而，在高维空间中，它会遭遇灾难性的失败。其均方积分误差(Mean Integrated Squared Error, MISE)的收敛速度会随着维度 d 的增加而急剧下降，最优收敛率约为 $O(n^{-4/(d+4)})$ ¹²。当

$d=3072$ 时，这个收敛速度慢到几乎可以忽略不计，意味着需要天文数字般的样本量才能获得一个稍有意义的估计。

带宽 h 的选择变成了一个无法解决的难题。如果 h 太小，查询点 x 的周围几乎没有任何数据点会落入核函数的支撑域内，导致估计出的密度几乎处处为零。如果 h 太大，则会将整个分布过度平滑成一个毫无信息量的均匀分布¹²。近期的研究进一步证实，在高维极限下，KDE的行为会进入一个新的统计学范畴，其估计值由极值统计(extreme value statistics)所主导，即只有极少数几个数据点对最终的密度估计起决定性作用，这完全违背了KDE旨在利用所有数据进行平滑估计的初衷¹⁵。此外，高维KDE还容易产生“幻影模式”(phantom modes)，即在真实分布中不存在的地方凭空制造出虚假的密度峰值¹⁷。

k-近邻(k-NN)密度估计:自适应带宽方法

为了克服KDE固定带宽的缺陷，k-NN密度估计应运而生。它不固定邻域的体积，而是固定邻域内的样本数量 k 。其密度估计公式为：

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^k \frac{1}{V_k(x)}$$

其中, $V_k(x)$ 是以查询点 x 为中心、半径为到其第 k 个最近邻的距离 $R_k(x)$ 的 d 维超球体的体积¹⁸。 k -NN方法的核心优势在于其自适应性:在数据稀疏的区域, $V_k(x)$ 会自动增大, 从而得到一个较低的密度估计;在数据稠密的区域, $V_k(x)$ 会相应减小, 得到较高的密度估计²¹。这种特性使其在理论上比KDE更适合处理多尺度分布。

然而, 尽管具有自适应性, k -NN密度估计最终也无法逃脱维度灾难的魔咒。正如第一节所分析的, 在高维空间中, 第 k 个最近邻实际上并不“近”。为了包含 k 个邻居, 超球体的半径 $R_k(x)$ 必须增长到非常大, 以至于覆盖了整个数据空间的很大一部分¹¹。这直接违反了 k -NN方法的一个基本假设:在

$V_k(x)$ 这个体积内, 真实的概率密度是近似恒定的¹⁹。当体积变得巨大时, 这个假设显然不再成立。研究表明, 尽管 k -NN估计量的方差不随维度变化, 但其偏差会随维度增加而改变, 导致估计结果的系统性失真¹⁹。

3. 局部密度的实用解决方案:k-NN与近似最近邻搜索的结合

对于用户提出的第一个问题——计算一个新向量的局部密度——直接应用前述的经典密度估计算法既不准确也无法扩展。然而, 我们可以将问题重新定义为一个更具实践意义的任务:计算一个“典型性得分”(typicality score), 该得分能够反映一个点是位于稠密区域还是稀疏区域。

重新定义问题:从概率密度到典型性得分

k -NN密度估计的核心思想 $p(x) \propto 1/V_k(x)$ 为我们提供了一个理想的出发点。我们可以定义一个与局部密度成正比的得分, 它与到第 k 个最近邻的距离 $R_k(x)$ 的 d 次方成反比:

$$\text{Score}(x) \propto R_k(x)^{-d}$$

这个得分直观地捕捉了局部密度: $R_k(x)$ 越小, 意味着邻居越近, 局部密度越高, 得分也越高。这个定义回避了计算一个严格的、归一化的概率值, 转而提供一个可用于排序和比较的相对度量。然而, 主要的技术瓶颈在于计算 $R_k(x)$ 。在一个包含 N 个 d 维向量的数据集中, 对单个查询点进行精确的 k -NN搜索(即暴力搜索)需要计算 N 次距离, 其计算复杂度为 $O(N \cdot d)$ ²³。对于百万级规模的数据集和3072的维度, 这种方法的计算成本是无法接受的。

近似最近邻(ANN)搜索:实现可扩展性的关键

近似最近邻(Approximate Nearest Neighbor, ANN)搜索技术是解决这一瓶颈的关键。ANN的核心思想是通过牺牲一定的精度(即返回的邻居可能不是绝对精确的最近邻)来换取查询速度的巨大提升,通常能将查询复杂度降低到亚线性甚至对数级别²⁴。

基于图的方法(HNSW)

层级式可导航小世界图(Hierarchical Navigable Small World, HNSW)是目前性能最优、应用最广的ANN算法之一²⁷。

- 原理:HNSW为数据集构建一个多层的邻近图(proximity graph)²⁷。顶层图包含的节点最少,但连接的边是“长程”的,用于在搜索初期进行快速、粗略的定位,这类似于数据结构中的跳表(skip list)³⁰。底层图则包含所有数据点,连接的边是“短程”的,用于在搜索后期进行精确的查找³¹。
- 机制:搜索过程从顶层的某个入口点开始,在当前层贪婪地遍历图,每一步都移动到离查询向量最近的邻居节点。当在当前层找到一个局部最优解后,算法会以下降到下一层,以该最优解为新的入口点,继续进行更精细的搜索,直至到达最底层³³。
- 参数与权衡:HNSW的性能受到几个关键参数的显著影响:M(图中每个节点的最大连接数)、efConstruction(构建索引时的搜索深度)和efSearch(查询时的搜索深度)。增加这些参数通常会提高搜索的召回率(准确性),但同时也会增加索引的内存占用和构建/查询时间²⁸。

哈希与混合方法(LSH & ScaNN)

- 局部敏感哈希(LSH):LSH的基本思想是设计一组特殊的哈希函数,使得原始空间中相似的向量有很高的概率被哈希到同一个“桶”中,而距离较远的向量则有很高的概率被哈希到不同的桶中³⁶。查询时,只需在查询向量所在的桶内进行搜索,从而大大减少了需要比较的向量数量。基于随机投影的哈希是实现这一目标的一种常用技术³⁸。尽管LSH在理论上具有重要意义,但在极高维度下,其性能通常不如基于图的方法⁴⁰。
- ScaNN(可扩展最近邻):ScaNN是谷歌开发的一种先进的混合ANN方法,它结合了多种技术的优点,在性能上达到了业界领先水平⁴¹。其架构通常包括三个阶段:1)分区(Partitioning):首先将整个数据集划分为多个簇(cluster),查询时仅需搜索与查询向量最接近的少数几个簇。2)量化(Quantization):对每个分区内的向量进行有损压缩(如乘积量化),使得距离计算可以在压缩后的低维空间中快速进行。3)重排(Re-ranking):在通过量化搜索得到一批候选向量后,使用它们的全精度原始向量与查询向量计算精确距离,并进

行重新排序，以提高最终结果的准确性⁴¹。ScaNN针对现代CPU指令集(如AVX2)进行了深度优化，性能极高⁴⁵。

使用Faiss进行大规模实现

Faiss(Facebook AI Similarity Search)是一个由Meta AI开发的高性能向量相似性搜索库，已成为实现上述ANN算法的工业标准⁴⁶。它提供了C++核心实现、完整的Python接口以及强大的GPU支持，非常适合处理百万级规模的数据集⁴⁷。Faiss库中包含了多种索引类型，如

IndexFlatL2(暴力搜索)、IndexHNSWFlat(HNSW索引)以及IndexIVFPQ(结合了分区和乘积量化的复合索引)，为用户提供了在速度、内存和精度之间进行权衡的丰富选项⁴⁸。

结合k-NN的密度定义和ANN的快速搜索，我们得到的是一个计算上可行且非常有效的局部密度或“典型性”的启发式得分。这个过程并未产生一个严格意义上的概率密度函数。首先，ANN算法引入了近似误差，找到的邻居可能并非真正的最近邻。其次，即使找到了精确的k个最近邻， $k / (N * V_k(x))$ 这个值本身也只是一个点估计，它所属的函数在整个空间上的积分并不为1，因此不是一个有效的PDF。然而，这个得分对于诸如异常检测(低分表示异常)或根据典型性对项目进行排序等任务来说是极其有用的。明确这一区别至关重要：这个得分不能直接用于需要有效似然值的下游概率模型中。

下表对主流的ANN算法进行了比较，为在实际应用中做出工程决策提供了参考。

表1: 近似最近邻(ANN)算法对比分析

算法	类型	查询复杂度	构建复杂度	内存使用	关键参数	精度-速度权衡	3072维数据适用性
暴力搜索	精确	$O(N \cdot d)$	$O(1)$	$O(N \cdot d)$	无	100%精度, 速度最慢	仅适用于小规模数据
HNSW	基于图	近似 $O(\log N)$	$O(N \log N)$	高	M, efConstruction, efSearch	极佳, 高召回率下仍能保持高速	非常高, 业界领先

LSH	基于哈希	亚线性	亚线性	中到高	哈希函数数量, 表数量	依赖参数, 高维下性能可能下降	中等, 对参数敏感
ScaNN	混合	近似 $O(\log N)$	$O(N \log N)$	中等	分区数, 量化参数	极佳, 通过重排机制平衡	非常高, 尤其在CPU上

4. 使用深度生成模型对全局概率分布进行建模

对于用户提出的第二个、更具挑战性的问题——学习一个全局的概率密度函数 $p(x)$, 其中 $x \in \mathbb{R}^{3072}$ ——经典方法已无能为力。现代深度生成模型(Deep Generative Models)为此类任务提供了最前沿的解决方案, 它们能够学习并表示高维空间中极其复杂的概率分布¹⁷。

正则化流(Normalizing Flows): 精确密度估计

正则化流(Normalizing Flows, NFs)是一类能够提供精确似然估计的生成模型, 这使其在纯密度估计任务中尤为突出。

- 核心原理: NFs的核心思想是学习一个可逆的、可微的双射变换(bijective mapping) $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, 这个变换能够将一个简单的基础分布 $p_z(z)$ (例如, 标准正态分布)转化为复杂的数据分布 $p_x(x)$ 。根据变量替换公式, 我们可以精确地计算出任意数据点 x 的对数似然值:

$$\log p_x(x) = \log p_z(f^{-1}(x)) + \log |\det(J(f^{-1}(x)))|$$

其中 J 表示雅可比矩阵, \det 是其行列式⁵¹。

- 架构创新: NFs面临的主要挑战在于如何设计变换 f , 使其雅可比行列式对于高维数据(如 $d=3072$)易于计算。主流架构通过巧妙的设计解决了这个问题:
 - **RealNVP & Glow**: 这类模型采用“耦合层”(coupling layers)。它们将输入向量 x 分成两部分 x_a 和 x_b 。变换只作用于其中一部分, 且变换的参数由另一部分决定, 例如: $y_a = x_a$, $y_b = s(x_a) \odot x_b + t(x_a)$ 。这种设计的雅可比矩阵是三角矩阵, 其行列式就是对角线元素的乘积, 即缩放因子 s 的乘积, 计算非常高效⁵³。Glow模型在此基础上引入了

可逆的1x1卷积，用于在不同层之间置换和混合维度，增强了模型的表达能力⁵⁶。

- 掩码自回归流(MAF): MAF利用掩码自编码器(Masked Autoencoders, MADE)来强制实现一种自回归结构，即每个输出维度 x_i 只依赖于其之前的维度 $x_{1:i-1}$ 。这种结构同样能产生一个三角雅可比矩阵，从而简化了行列式的计算⁵⁸。MAF的密度评估速度很快，但采样过程是串行的，速度较慢。

变分自编码器(VAEs): 近似密度估计

变分自编码器(Variational Autoencoders, VAEs)是一类基于隐变量的生成模型，它通过一个低维的隐空间 z 来学习数据的分布。

- 核心原理: VAEs假设数据 x 是由隐变量 z 生成的，其边际似然为 $p(x) = \int p(x|z)p(z)dz$ 。由于这个积分通常是难以计算的(intractable)，VAEs引入了一个编码器网络 $q(z|x)$ 来近似真实的后验分布 $p(z|x)$ ，并通过最大化证据下界(Evidence Lower Bound, ELBO)来间接优化对数似然：
$$\log p(x) \geq \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x) \parallel p(z))$$

$p(z)$

这个目标函数由两部分组成：第一部分是“重构项”，鼓励解码器 $p(x|z)$ 能够从隐变量 z 中恢复原始数据 x ；第二部分是KL散度“正则项”，它促使编码器产生的后验分布 $q(z|x)$ 与预设的先验分布 $p(z)$ （通常是标准正态分布）保持一致⁶⁰。

- 重要性加权自编码器(IWAE): 标准的VAEs优化的是一个单样本的ELBO，这个下界可能比较松。IWAE通过从近似后验 $q(z|x)$ 中抽取 k 个样本，并计算一个重要性加权的似然估计，从而提供了一个比标准ELBO更紧的下界。更紧的下界通常能带来更好的密度模型和更丰富的隐变量表示⁶³。
- 关键挑战: 后验坍塌(Posterior Collapse)
 - 定义: 后验坍塌是指在训练过程中，模型找到一个平凡的局部最优解，其中近似后验 $q(z|x)$ 完全等于先验 $p(z)$ 。这意味着隐变量 z 中没有包含任何关于输入 x 的信息，编码器被完全忽略⁶⁵。
 - 成因: 这个问题在使用一个非常强大、表达能力强的解码器(例如，自回归模型)时尤为突出。强大的解码器可以直接对 $p(x)$ 进行建模，而不需要依赖 z 提供的信息，就能获得很高的重构似然 $p(x|z)$ 。这使得模型可以轻易地将KL散度项优化为零，从而陷入后验坍塌的困境⁶⁶。
 - 缓解策略: 为了对抗后验坍塌，研究者们提出了多种策略，其中最常用的是: 1) **KL退火(KL Annealing)**: 在训练初期，将KL散度项的权重设置为0，然后逐渐增加到1。这给了模型足够的时间先学习如何利用隐变量进行有效重构，之后再施加使其后验接近先验的正则化约束⁶⁵。 2) **自由位(Free Bits)**: 这种方法修改了目标函数，仅当每个隐变量维度的KL散度超过一个预设的阈值 λ 时才对其进行惩罚。这相当于给了每个隐变量维度一定的“信息容量预算”，鼓励模型至少利用这些容量来编码信息⁶⁵。

基于能量的模型(EBMs):非归一化密度

基于能量的模型(Energy-Based Models, EBMs)提供了最大的模型灵活性,但其训练过程也最具挑战性。

- 核心原理:EBMs通过一个能量函数 $E(x)$ 来隐式地定义概率分布。能量越低的点,其对应的概率密度越高。这种关系通过吉布斯-玻尔兹曼分布来形式化:

$$p(x)=Z\exp(-E(x))$$

其中 Z 是归一化常数(也称配分函数),它需要对整个数据空间进行积分,因此通常是难以计算的⁷⁰。

- 训练挑战与方法:由于 Z 的存在,直接最大化似然是不可行的。对数似然的梯度可以分解为两部分:一部分降低数据点的能量(“正相”),另一部分则提高由模型当前分布产生的样本的能量(“负相”)。负相需要从 $p(x)$ 中采样,这本身就是一个难题⁷²。通常,这是通过马尔可夫链蒙特卡洛(MCMC)方法(如朗之万动力学)来实现的,即通过迭代采样来近似从 $p(x)$ 中抽取样本⁷²。为了降低计算成本,一种广泛使用的近似方法是对比散度(Contrastive Divergence, CD),它从一个真实数据点开始,只运行少数几步MCMC,从而得到一个有偏但计算上可行的梯度估计⁷²。

这三类深度生成模型并非相互孤立,而是代表了在概率建模中一个根本性的权衡谱。正则化流(NFs)将可计算性放在首位,通过严格的架构约束(可逆性、易于计算的雅可比行列式)来保证能够得到精确的似然值,这使其成为纯密度估计任务的理想选择。变分自编码器(VAEs)则寻求一种平衡,它放宽了严格的可逆性约束,允许更灵活的编码器-解码器结构(例如,在隐空间中实现降维),但代价是只能得到似然的一个下界,使得密度估计是近似的。基于能量的模型(EBMs)则追求极致的表达能力,任何可微函数都可以作为能量函数,几乎不对模型形式施加任何限制。其代价是最大的不可计算性:似然是未归一化的,需要依赖昂贵且复杂的MCMC采样进行训练。因此,模型的选择应基于首要目标:如果精确的概率值至关重要,NFs是首选;如果学习丰富的低维表示更为重要且可接受近似密度,VAEs是合适的;如果目标是捕捉尽可能复杂的数据分布且计算成本是次要考虑,EBMs则最为强大。

5. 表征数据密度的替代范式

除了直接建模概率密度函数外,还存在一些替代方法,它们通过不同的视角来刻画数据点的密度或典型性,并且在某些情况下能更好地规避维度灾难。

孤立森林:一种无需距离的密度代理

孤立森林(Isolation Forest)是一种最初为异常检测设计的无监督算法,但其核心机制可以被巧妙地用于高效地估计密度的代理指标⁷³。

- 原理:该算法基于一个简单而强大的直觉:异常点(即低密度点)由于其“稀少且不同”的特性,比正常点更容易被“孤立”出来⁷⁵。
- 机制:孤立森林构建了一个由多个“孤立树”(iTrees)组成的集成模型。每棵树都是通过对数据的随机子样本进行递归分区来构建的。在树的每个节点,算法会随机选择一个特征,然后在这个特征的最大值和最小值之间随机选择一个分割点,将数据一分为二⁷⁷。这个过程持续进行,直到每个数据点都被单独划分到一个叶子节点。
- 异常得分作为密度代理:一个数据点 x 从树的根节点到其所在的叶子节点的路径长度 $h(x)$ 被记录下来。直观上,低密度点会很快被分割开,因此它们的平均路径长度 $E[h(x)]$ 会很短⁷⁹。算法将这个平均路径长度通过以下公式转换为一个介于0和1之间的异常得分

$s(x)$:

$$s(x)=2-c(\psi)E[h(x)]$$

其中 $c(\psi)$ 是与子样本大小 ψ 相关的归一化常数。得分越接近1,表示路径越短,该点越可能是异常点(密度越低);得分越接近0.5,表示路径长度接近平均水平,该点越可能是正常点(密度越高)⁸⁰。

- 复杂度与优势:孤立森林的训练复杂度为 $O(t \cdot \psi \cdot \log(\psi))$,而对新数据点的评分(推理)复杂度仅为 $O(t \cdot \log(\psi))$,其中 t 是树的数量, ψ 是子采样大小。这使得它在计算上极为高效⁷⁹。最关键的是,该算法完全不使用任何距离度量,因此它天然地免疫于第一节中描述的距离度量失效问题,在高维数据上表现出很强的鲁棒性⁷³。

利用结构性假设缓解维度灾难

如果关于数据内在结构存在先验知识,我们可以利用这些知识来显著缓解维度灾难的影响。

- 流形假设:如果3072维的数据实际上分布在一个维度远低于3072的低维流形(manifold)上或其附近,那么密度估计问题就可以在这个低维流形上进行,从而将问题的有效维度 d_{eff} 大大降低⁵⁰。流形学习流(Manifold learning flows)就是利用这一思想的生成模型之一⁸²。
- 条件独立性:如果数据的3072个维度之间并非完全相互依赖,而是存在某些条件独立关系,我们就可以使用能够利用这种稀疏依赖结构的模型。
 - 藤蔓杯葛(Vine Copulas):这种方法可以将一个 d 维的联合密度分解为 $d(d-1)/2$ 个二元(二维)的copula密度。这成功地将一个棘手的高维问题转化为了多个易于处理的低维

问题¹⁴。

- 图模型: 马尔可夫随机场 (Markov Random Fields) 等图模型可以显式地编码变量之间的条件独立性。在这种情况下, 密度估计的样本复杂度不再取决于环境维度 d , 而是取决于一个更小的、与图结构相关的参数, 称为“图弹性” (graph resilience)⁵⁰。

6. 综合分析 & 策略建议

执行摘要

本报告深入探讨了在高维空间 ($d=3072$) 和大规模数据集 ($N \sim$ 百万级) 下进行密度估计的挑战与可行方案。分析表明, 由于“维度灾难”, 传统的非参数方法 (如 KDE 和 k -NN) 会失效。针对用户的两个具体需求, 我们提出了不同的策略路径。对于计算新向量的局部密度, 推荐采用基于 k -NN 的密度得分, 并通过近似最近邻 (ANN) 搜索技术 (如 HNSW 或 ScaNN) 实现可扩展计算。对于学习全局概率密度函数 (**PDF**), 这是一个更复杂的研发挑战, 深度生成模型是当前最先进的解决方案, 但不同模型家族 (正则化流、VAEs、EBMs) 在精确性、可计算性和模型灵活性之间存在显著的权衡。

下表对本文讨论的关键方法进行了全面评估, 旨在为技术选型提供一个清晰的决策矩阵。

表2: 高维密度估计方法综合评估

方法论	任务 (局部/ 全局)	密度类 型 (代 理/近 似/精 确/非 归一 化)	核心原 理	训练可 扩展性	推理可 扩展性	高维鲁 棒性	主要优 势	主要劣 势/挑 战
k-NN (+ANN)	局部	代理得 分	邻域体 积的倒 数	N/A (惰性 学习)	极高 (近 $O(\log N)$)	中 (依 赖距 离)	实现简 单, 结 果直观 , 可扩 展性强	产生 的是得 分而非 概率, 对 距离度 量敏感

孤立森林	局部/全局	代理得分	隔离点的难易程度	极高 ($O(t\psi \log \psi)$)	极高 ($O(t\log \psi)$)	高 (不依赖距离)	速度极快, 内存占用低, 对维度灾难鲁棒	得分粒度可能较粗, 非严格概率密度
正则化流 (NFs)	全局	精确概率	可逆变换与变量替换公式	中到低	高	中	提供精确的对数似然, 理论完备	架构约束强, 可能难以拟合某些复杂分布
VAE / IWAE	全局	近似概率 (下界)	隐变量模型与变分推断	中	高	中	可学习有意义的低维表示, 模型灵活	仅提供似然下界, 易发生后验坍塌
EBMs	全局	非归一化概率	能量函数与玻尔兹曼分布	低	低 (需MCMC采样)	高	模型表达能力极强, 约束最少	训练不稳定且计算昂贵, 似然不可直接计算

局部密度计算建议(问题1)

对于计算新向量局部密度的任务，最实用和可扩展的方案是：

- 首选推荐: 使用基于**k-NN**的密度得分，该得分是关于到第k个最近邻距离 $R_k(x)$ 的函数。
- 实现方案: 必须使用近似最近邻(**ANN**)索引来加速**k-NN**搜索。我们推荐采用**HNSW**算法，并通过**Faiss**库进行实现，因为它在速度和召回率之间取得了当前最佳的平衡²⁷。谷歌的**ScaNN**也是一个顶级的替代方案，尤其是在CPU推理场景下表现出色⁴¹。
- 理由: 该方法在计算上可以轻松扩展到数百万向量，直接满足了用户对逐点密度度量的需求

，并且有成熟的库支持，实现相对直接。其结果是一个鲁棒的典型性得分，非常适用于排序或异常检测等应用。

- 备选方案：如果追求极致的速度，孤立森林是一个绝佳的选择。它对距离度量问题的免疫性以及线性时间复杂度使其非常鲁棒和高效⁷³。

全局PDF估计建议(问题2)

学习一个全局的概率密度函数是一项艰巨的研究与工程挑战。最终选择取决于对近似误差的容忍度、计算资源以及实现复杂性。

- 追求精确似然：如果任务的核心是纯粹的密度估计，且获得精确的概率值至关重要，那么正则化流(NFs)是理论上最合理的选择⁵²。像Glow或MAF这样的架构是强大的起点。
- 追求灵活建模与表示学习：如果目标不仅是密度估计，还包括学习数据有意义的低维表示，并且可以接受一个近似的密度(似然的下界)，那么推荐使用变分自编码器，特别是重要性加权自编码器(IWAE)⁶³。
重要提醒：必须准备好诊断并积极采用KL退火等技术来缓解后验坍塌问题，以确保模型学到有意义的表示。
- 追求极致表达能力：如果预期底层数据分布异常复杂，且计算资源不是主要限制因素，**基于能量的模型(EBMs)**提供了最大的模型灵活性。然而，这需要应对复杂且可能不稳定的基于MCMC的训练过程⁷²。

结论与展望

为高维数据选择合适的密度估计算法，需要深刻理解维度灾难带来的根本性挑战，并根据具体应用场景(局部得分 vs. 全局PDF)做出战略性权衡。对于局部密度查询，ANN加速的k-NN得分提供了一个兼具实用性和可扩展性的解决方案。对于全局PDF建模，深度生成模型开辟了新的可能性，但要求使用者在模型的可计算性、近似程度和表达能力之间做出明智的选择。该领域仍在飞速发展，诸如扩散模型(Diffusion Models)等新兴模型在生成任务上已展现出最先进的性能，其在显式密度估计方面的应用是一个值得关注的活跃研究方向⁵⁰。

引用的著作

1. Curse of dimensionality - Wikipedia, 访问时间为 九月 7, 2025, https://en.wikipedia.org/wiki/Curse_of_dimensionality
2. Curse of Dimensionality: Challenges of High-Dimensional Data | by Tiya Vaj - Medium, 访问时间为 九月 7, 2025,

- <https://vtiya.medium.com/curse-of-dimensionality-challenges-of-high-dimensional-data-115e67e1b5e6>
3. 1 Surprises in high dimensions - UC Davis Math, 访问时间为 九月 7, 2025, <https://www.math.ucdavis.edu/~strohmer/courses/180BigData/180lecture1.pdf>
 4. Why is volume of a high-D ball concentrated near its surface?, 访问时间为 九月 7, 2025, <https://math.stackexchange.com/questions/2118098/why-is-volume-of-a-high-d-ball-concentrated-near-its-surface>
 5. concentration of volume of hypersphere - Mathematics Stack Exchange, 访问时间为 九月 7, 2025, <https://math.stackexchange.com/questions/1648376/concentration-of-volume-of-hypersphere>
 6. Volume of an n-ball - Wikipedia, 访问时间为 九月 7, 2025, https://en.wikipedia.org/wiki/Volume_of_an_n-ball
 7. History of the high-dimensional volume paradox - MathOverflow, 访问时间为 九月 7, 2025, <https://mathoverflow.net/questions/128786/history-of-the-high-dimensional-volume-paradox>
 8. Curse of dimensionality: Challenges & impact in high-dimensional data - Domino Data Lab, 访问时间为 九月 7, 2025, <https://domino.ai/blog/the-curse-of-dimensionality>
 9. K-Nearest Neighbors and Curse of Dimensionality - GeeksforGeeks, 访问时间为 九月 7, 2025, <https://www.geeksforgeeks.org/machine-learning/k-nearest-neighbors-and-curse-of-dimensionality/>
 10. The Effects of Dimensionality Curse in High Dimensional kNN Search - ResearchGate, 访问时间为 九月 7, 2025, https://www.researchgate.net/publication/224265352_The_Effects_of_Dimensionality_Curse_in_High_Dimensional_kNN_Search
 11. Lecture 2: k-nearest neighbors / Curse of Dimensionality - Cornell: Computer Science, 访问时间为 九月 7, 2025, https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html
 12. Kernel density estimation - Wikipedia, 访问时间为 九月 7, 2025, https://en.wikipedia.org/wiki/Kernel_density_estimation
 13. 2.8. Density Estimation — scikit-learn 1.7.1 documentation, 访问时间为 九月 7, 2025, <https://scikit-learn.org/stable/modules/density.html>
 14. Evading the curse of dimensionality in nonparametric density estimation with simplified vine copulas - arXiv, 访问时间为 九月 7, 2025, <https://arxiv.org/pdf/1503.03305>
 15. Kernel Density Estimators in Large Dimensions - arXiv, 访问时间为 九月 7, 2025, <https://arxiv.org/html/2408.05807v1>
 16. [2408.05807] Kernel Density Estimators in Large Dimensions - arXiv, 访问时间为 九月 7, 2025, <https://arxiv.org/abs/2408.05807>
 17. Nonparametric Density Estimation for High-Dimensional Data ... - arXiv, 访问时间为

- 为 九月 7, 2025, <https://arxiv.org/pdf/1904.00176>
18. KnnDensityEstimation - Project Rhea, 访问时间为 九月 7, 2025, <https://www.projectrhea.org/rhea/index.php/KnnDensityEstimation>
 19. Lecture 7: Density Estimation: k-Nearest Neighbor and Basis Approach - faculty.washington.edu, 访问时间为 九月 7, 2025, http://faculty.washington.edu/yenchic/18W_425/Lec7_knn_basis.pdf
 20. Nearest Neighbors and Non-parametric Density Estimation, 访问时间为 九月 7, 2025, <https://www.csc.kth.se/utbildning/kth/kurser/DD2427/bik08/LectureNotes/Lecture6.pdf>
 21. 11 Nearest Neighbor Methods, 访问时间为 九月 7, 2025, <https://www.ssc.wisc.edu/~bhansen/718/NonParametrics10.pdf>
 22. L8: Nearest neighbors, 访问时间为 九月 7, 2025, https://people.engr.tamu.edu/rgutier/lectures/pr/pr_l8.pdf
 23. 1.6. Nearest Neighbors — scikit-learn 1.7.1 documentation, 访问时间为 九月 7, 2025, <https://scikit-learn.org/stable/modules/neighbors.html>
 24. Nearest neighbors in high-dimensional data? - Stack Overflow, 访问时间为 九月 7, 2025, <https://stackoverflow.com/questions/5751114/nearest-neighbors-in-high-dimensional-data>
 25. Approximate Nearest Neighbor Methods | by Abdullah Şamil Güser - Medium, 访问时间为 九月 7, 2025, <https://medium.com/@abdullahsamilguser/approximate-nearest-neighbor-methods-713dcfa8518f>
 26. Unifying Speed-Accuracy Trade-Off and Cost-Benefit Trade-Off in Human Reaching Movements - Frontiers, 访问时间为 九月 7, 2025, <https://www.frontiersin.org/journals/human-neuroscience/articles/10.3389/fnhum.2017.00615/full>
 27. Hierarchical navigable small world - Wikipedia, 访问时间为 九月 7, 2025, https://en.wikipedia.org/wiki/Hierarchical_navigable_small_world
 28. How hierarchical navigable small world (HNSW) algorithms can improve search - Redis, 访问时间为 九月 7, 2025, <https://redis.io/blog/how-hnsw-algorithms-can-improve-search/>
 29. What is a Hierarchical Navigable Small World - MongoDB, 访问时间为 九月 7, 2025, <https://www.mongodb.com/resources/basics/hierarchical-navigable-small-world>
 30. Hierarchical Navigable Small Worlds (HNSW) - Pinecone, 访问时间为 九月 7, 2025, <https://www.pinecone.io/learn/series/faiss/hnsw/>
 31. Understanding Hierarchical Navigable Small Worlds (HNSW) - Zilliz Learn, 访问时间为 九月 7, 2025, <https://zilliz.com/learn/hierarchical-navigable-small-worlds-HNSW>
 32. Introduction to HNSW: Hierarchical Navigable Small World - Analytics Vidhya, 访问时间为 九月 7, 2025, <https://www.analyticsvidhya.com/blog/2023/10/introduction-to-hnsw-hierarchical-navigable-small-world/>

33. HNSW indexing in Vector Databases: Simple explanation and code | by Will Tai - Medium, 访问时间为 九月 7, 2025, <https://medium.com/@wtaisen/hnsw-indexing-in-vector-databases-simple-explanation-and-code-3ef59d9c1920>
34. Similarity Search, Part 4: Hierarchical Navigable Small World (HNSW), 访问时间为 九月 7, 2025, <https://towardsdatascience.com/similarity-search-part-4-hierarchical-navigable-small-world-hnsw-2aad4fe87d37/>
35. The Hierarchical Navigable Small Worlds (HNSW) Algorithm | Lantern Blog, 访问时间为 九月 7, 2025, <https://lantern.dev/blog/hnsw>
36. Locality-sensitive hashing - Wikipedia, 访问时间为 九月 7, 2025, https://en.wikipedia.org/wiki/Locality-sensitive_hashing
37. Locality Sensitive Hashing (LSH): The Illustrated Guide - Pinecone, 访问时间为 九月 7, 2025, <https://www.pinecone.io/learn/series/faiss/locality-sensitive-hashing/>
38. Proximity in High Dimensions: The Power of Locality Sensitive Hashing | by Shubham Sangole | CodeX | Medium, 访问时间为 九月 7, 2025, <https://medium.com/codex/proximity-in-high-dimensions-the-power-of-locality-sensitive-hashing-3f4e37865fd7>
39. Locality Sensitive Hashing for Fast Search in High Dimension Data. | by sid dhuri - Medium, 访问时间为 九月 7, 2025, <https://medium.com/geekculture/locality-sensitive-hashing-for-fast-search-in-high-dimension-data-a2cdef1b6eff>
40. Nearest Neighbor Indexes for Similarity Search - Pinecone, 访问时间为 九月 7, 2025, <https://www.pinecone.io/learn/series/faiss/vector-indexes/>
41. What is ScaNN (Scalable Nearest Neighbors)? - Zilliz Learn, 访问时间为 九月 7, 2025, <https://zilliz.com/learn/what-is-scann-scalable-nearest-neighbors-google>
42. Faster retrieval with Scalable Nearest Neighbours (ScANN) - Keras, 访问时间为 九月 7, 2025, https://keras.io/keras_rs/examples/scann/
43. State-of-the-Art Approximate Nearest Neighbor Search with Google's ScaNN and Facebook's FAISS | by ANURAG DIXIT | Medium, 访问时间为 九月 7, 2025, <https://medium.com/@DataPlayer/scalable-approximate-nearest-neighbour-search-using-googles-scann-and-facebook-s-faiss-3e84df25ba>
44. Similarity Search: ScaNN and 4-bit PQ | by Takuma Yamaguchi (Kumon) | Medium, 访问时间为 九月 7, 2025, <https://medium.com/@kumon/similarity-search-scann-and-4-bit-pq-ab98766b32bd>
45. ScaNN - Python LangChain, 访问时间为 九月 7, 2025, <https://python.langchain.com/docs/integrations/vectorstores/scann/>
46. Welcome to Faiss Documentation — Faiss documentation, 访问时间为 九月 7, 2025, <https://faiss.ai/>
47. facebookresearch/faiss: A library for efficient similarity search and clustering of dense vectors. - GitHub, 访问时间为 九月 7, 2025, <https://github.com/facebookresearch/faiss>
48. Faiss indexes · facebookresearch/faiss Wiki - GitHub, 访问时间为 九月 7, 2025, <https://github.com/facebookresearch/faiss/wiki/Faiss-indexes>

49. HNSW : Semantic Search Using FAISS - BakingAI Blog, 访问时间为 九月 7, 2025,
<https://bakingai.com/blog/hnsw-semantic-search-faiss-integration/>
50. Breaking the curse of dimensionality in structured density estimation - OpenReview, 访问时间为 九月 7, 2025,
<https://openreview.net/pdf?id=dWwin2uGYE>
51. Understanding Normalizing Flows - TU Delft Repository, 访问时间为 九月 7, 2025,
https://repository.tudelft.nl/file/File_897ea75b-9851-45fe-a6c7-c5ce012e0ea1
52. Normalizing Flows for Probabilistic Modeling and Inference, 访问时间为 九月 7, 2025,
<https://jmlr.org/papers/volume22/19-1028/19-1028.pdf>
53. Flow-based Deep Generative Models | Lil'Log, 访问时间为 九月 7, 2025,
<https://lilianweng.github.io/posts/2018-10-13-flow-models/>
54. Density estimation using Real NVP - Keras, 访问时间为 九月 7, 2025,
https://keras.io/examples/generative/real_nvp/
55. GLOW: Generative flow - Amélie Royer, 访问时间为 九月 7, 2025,
<https://ameroyer.github.io/portfolio/2021-04-12-Glow/>
56. Flow-based generative model - Wikipedia, 访问时间为 九月 7, 2025,
https://en.wikipedia.org/wiki/Flow-based_generative_model
57. Glow: Generative Flow with Invertible 1x1 Convolutions - NIPS, 访问时间为 九月 7, 2025,
<https://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions>
58. Masked Autoregressive Flow for Density Estimation - The University of Edinburgh, 访问时间为 九月 7, 2025,
<https://homepages.inf.ed.ac.uk/imurray2/pub/17maf/maf.pdf>
59. Masked Autoregressive Flow for Density Estimation - NIPS, 访问时间为 九月 7, 2025,
<https://papers.nips.cc/paper/6828-masked-autoregressive-flow-for-density-estimation>
60. Student-t Variational Autoencoder for Robust Density Estimation - IJCAI, 访问时间为 九月 7, 2025,
<https://www.ijcai.org/proceedings/2018/0374.pdf>
61. Variational autoencoder - Wikipedia, 访问时间为 九月 7, 2025,
https://en.wikipedia.org/wiki/Variational_autoencoder
62. Evidence lower bound - Wikipedia, 访问时间为 九月 7, 2025,
https://en.wikipedia.org/wiki/Evidence_lower_bound
63. Importance Weighted Autoencoder (IWAE) - Emergent Mind, 访问时间为 九月 7, 2025,
<https://www.emergentmind.com/topics/importance-weighted-autoencoder-iwae>
64. Importance Weighted Autoencoders, 访问时间为 九月 7, 2025,
https://www.mlmi.eng.cam.ac.uk/files/d421a_poster_importance_weighted_autoencoders.pdf
65. Understanding Posterior Collapse in Generative Latent Variable Models - OpenReview, 访问时间为 九月 7, 2025,
<https://openreview.net/pdf/729562a11b8fe6b0af7244d73dea98ec6c5f8376.pdf>
66. Detecting Posterior Collapse in Conditional and Hierarchical Variational Autoencoders - arXiv, 访问时间为 九月 7, 2025,

- <https://arxiv.org/html/2306.05023v3>
67. What is "posterior collapse" phenomenon? - Data Science Stack Exchange, 访问时间为 九月 7, 2025,
<https://datascience.stackexchange.com/questions/48962/what-is-posterior-collapse-phenomenon>
 68. Variational Auto-Encoders - Deep learning courses at UC Berkeley, 访问时间为 九月 7, 2025,
<https://berkeley-deep-learning.github.io/cs294-131-s17/slides/VAE%20talk.compressed.pdf>
 69. Mastering KL Divergence in PyTorch | by Amit Yadav | We Talk Data - Medium, 访问时间为 九月 7, 2025,
<https://medium.com/we-talk-data/mastering-kl-divergence-in-pytorch-4d0be6d7b6e3>
 70. Energy-Based Models in Machine Learning - GeeksforGeeks, 访问时间为 九月 7, 2025,
<https://www.geeksforgeeks.org/machine-learning/energy-based-models-in-machine-learning/>
 71. LEARNING ENERGY-BASED MODELS BY SELF-NORMALISING THE LIKELIHOOD - OpenReview, 访问时间为 九月 7, 2025, <https://openreview.net/pdf?id=zrxISviRqC>
 72. How to Train Your Energy-Based Models, 访问时间为 九月 7, 2025,
<https://arxiv.org/abs/2101.03288>
 73. An Introduction to Isolation Forests, 访问时间为 九月 7, 2025,
https://cran.r-project.org/web/packages/isotree/vignettes/An_Introduction_to_Isolation_Forests.html
 74. Isolation Forest For Anomaly Detection Made Easy & How To Tutorial - Spot Intelligence, 访问时间为 九月 7, 2025,
<https://spotintelligence.com/2024/05/21/isolation-forest/>
 75. How to perform anomaly detection with the Isolation Forest algorithm, 访问时间为 九月 7, 2025,
<https://towardsdatascience.com/how-to-perform-anomaly-detection-with-the-isolation-forest-algorithm-e8c8372520bc/>
 76. Isolation Forest algorithm for anomaly detection | by Arpit - Medium, 访问时间为 九月 7, 2025,
<https://medium.com/@arpitbhayani/isolation-forest-algorithm-for-anomaly-detection-f88af2d5518d>
 77. What is Isolation Forest? - GeeksforGeeks, 访问时间为 九月 7, 2025,
<https://www.geeksforgeeks.org/machine-learning/what-is-isolation-forest/>
 78. Isolation Forest - Auto Anomaly Detection with Python - Towards Data Science, 访问时间为 九月 7, 2025,
<https://towardsdatascience.com/isolation-forest-auto-anomaly-detection-with-python-e7a8559d4562/>
 79. Isolation Forest - LAMDA, 访问时间为 九月 7, 2025,
<http://www.lamda.nju.edu.cn/publication/icdm08b.pdf>
 80. (PDF) Isolation Forest - ResearchGate, 访问时间为 九月 7, 2025,
https://www.researchgate.net/publication/224384174_Isolation_Forest

81. Isolation Forest - LAMDA, 访问时间为 九月 7, 2025,
<https://cs.nju.edu.cn/zhoush/zhoush.files/publication/icdm08b.pdf?q=isolation-forest>
82. NeurIPS Poster Canonical normalizing flows for manifold learning, 访问时间为 九月 7, 2025, <https://neurips.cc/virtual/2023/poster/69924>
83. Breaking the curse of dimensionality in structured density estimation - OpenReview, 访问时间为 九月 7, 2025,
[https://openreview.net/forum?id=dWwin2uGYE&referrer=%5Bthe%20profile%20of%20Bryon%20Aragam%5D\(%2Fprofile%3Fid%3D~Bryon_Aragam1\)](https://openreview.net/forum?id=dWwin2uGYE&referrer=%5Bthe%20profile%20of%20Bryon%20Aragam%5D(%2Fprofile%3Fid%3D~Bryon_Aragam1))
84. Breaking the curse of dimensionality in structured density estimation, 访问时间为 九月 7, 2025, <https://neurips.cc/media/neurips-2024/Slides/94337.pdf>