

**维数灾难**（英语：curse of dimensionality，又名**维度的詛咒**）是一个最早由理查德·贝尔曼（Richard E. Bellman）在考虑优化问题时首次提出来的术语<sup>[1][2]</sup>，用来描述当（数学）空间维度增加时，分析和组织高维空间（通常有成百上千维），因体积指数增加而遇到各种问题场景。这样的难题在低维空间中不会遇到，如物理间通常只用三维来建模。

在很多领域中，如采样、组合数学、机器学习和数据挖掘都有提及到这个名字的现象。这些问题的共同特色是当维数提高时，空间的体积提高太快，因而可用数据变得很稀疏。稀疏性对于任何要求有统计学意义的方法而言都是一个**问题**，为了获得在统计学上正确并且有可靠的结果，用来支撑这一结果所需要的数据量通常随着维数的提高而呈指数级增长。而且，在组织和搜索数据时也有赖于检测对象区域，这些区域中的对象通过相似度属性而形成分组。然而在高维空间中，所有的数据都很稀疏，从很多角度看都不相似，因而平常使用的数据组织策略变得极其低效。

# 目录

- # 组合学

在一些问题中，每个变量都可取一系列离散值中的一个，或者可能值的范围被划分为有限个可能性。把这些变量放在一起，则必须考虑很多种值的组合方式，这后果就是常说的组合爆炸。即使在最简单的二元变量例子中，可能产生的组合总数就已经是在维数上呈现指数级的 $\mathcal{O}(2^d)$ 。一般而言，每个额外的维度都需要成倍地增加尝试所有组合方式的影响。

# 采样

当在数学空间上额外增加一个维度时，其体积会呈指数级的增长。如，点间距离不超过 $10^{-2}=0.01$ ， $10^2=100$ 个均匀间距的样本点足够采样到一个单位区间（“一个维度的立方体”）；一个10维单元超立方体的等价采样，其相邻两点间的距离为 $10^{-2}=0.01$ 则需要 $10^{20}$ 个样本点。一般而言，点距为 $10^{-n}$ 的10维超立方体所需要的样本点数量，是1维超立方体这样的单元区间的 $10^{n(10-1)}$ 倍。在上面的 $n=2$ 的例子中：当样本距离为0.01时，10维超立方体所需要的样本点数量会比单元区间多 $10^{18}$ 倍。这一影响就是上面所述组合学问题中的组合结果，距离函数问题将在下面介绍。

# 优化

当用数值逆向归纳法解决动态优化问题时，目标函数针对每个可能的组合都必须计算一遍，当状态变量的维度很大时，这是极其困难的。

# 机器学习

在机器学习问题中，需要在高维特征空间（每个特征都能够取一系列可能值）的有限数据样本中学习一种“自然状态”（可能是无穷分布），要求有相当数量的训练数据含有一些样本组合。给定固定数量的训练样本，其预测能力随着维度的增加而减小，这就是所谓的Hughes影响<sup>[4]</sup>或Hughes现象（以Gordon F. Hughes命名）。<sup>[5][6]</sup>

# 贝叶斯统计

在贝叶斯统计中维数灾难通常是一个难点，因为其后验分布通常都包含着许多参数。

然而，这一问题在基于模拟的贝叶斯推理（尤其是适应于很多实践问题的马尔科夫蒙特卡洛方法）出现后得到极大地克服，当然，基于模拟的方法收敛很慢，因此这也并不是解决高维问题的灵丹妙药。

# 距离函数

当一个度量，如欧几里德距离使用很多坐标来定义时，不同的样本对之间的距离已经基本上没有差别。

一种用来描述高维欧几里德空间的巨型性的方法是将超球体中半径 $r$ 和维数 $d$ 的比例，和超立方体中边长 $2r$ 和等值维数的比例相比较。这样一个球体的体积计算如下：
$$\frac{2r^d \pi^{d/2}}{d\Gamma(d/2)}$$

立方体的体积计算如下： $(2r)^d$

随着空间维度 $d$ 的增加，相对于超立方体的体积来说，超球体的体积就变得微不足道了。这一点可以从当 $d$ 趋于无穷时比较前面的比例清楚地看出：
$$\frac{\pi^{d/2}}{d2^{d-1}\Gamma(d/2)} \rightarrow 0$$

当 $d \rightarrow \infty$ 。因此，在某种意义上，几乎所有的高维空间都远离其中心，或者从另一个角度来看，高维单元空间可以说是几乎完全由超立方体的“边角”所组成的，没有“中部”，这对于理解卡方分布是很重要的直觉理解。给定一个单一分布，由于其最小值和最大值与最小值相比收敛于0，因此，其最小值和最大值的距离变得不可辨别。
$$\lim_{d \rightarrow \infty} \frac{\text{dist}_{\max} - \text{dist}_{\min}}{\text{dist}_{\min}} \rightarrow 0.$$

这通常被引证为距离函数在高维环境下失去其意义的例子。

# 最近邻搜索

最近邻搜索在高维空间中影响很大，因为其不可能使用其中一个坐标上的距离下界来快速地去掉一个候选项，因为该距离计算需要基于所有维度。<sup>[7][8]</sup>

然而，最近的研究表明仅仅一些数量的维度不一定会必然导致该问题，<sup>[3]</sup>因为相关的附加维度也能增加其相反项。另外，结果排序的方法仍然有助于辨别近处和远处的邻居。然而，不相关（“噪声”）维度也如期望一样会减少相反项，在时间序列分析中，数据一般都是高维的，只要信噪比足够高的话，其距离函数也同样能够可靠地工作。<sup>[9]</sup>

## **k近邻分类**

高维度在距离函数的另一个影响例子就是k近邻（k-NN）图，该图使用一些距离函数从数据集构造。当维度增加时，k-NN有向图的入度分页将会向右倾斜，从而导致中心的出现，很多的数据实例出现在其他许多实例（比预期多得多）的k-NN列表中。这一现象对很多技术，如分类（包括最近邻居法、半监督学习，和聚类分析）都有很大的影响。<sup>[10]</sup>，同时它也对信息检索问题有影响。<sup>[11]</sup>

## **延伸閱讀**

- |                  |                |                |
|------------------|----------------|----------------|
| ▪ <u>组合爆炸</u>    | ▪ <u>动态规划</u>  | ▪ <u>准随机</u>   |
| ▪ <u>相似度集中</u>   | ▪ <u>贝尔曼方程</u> | ▪ <u>聚类分析</u>  |
| ▪ <u>降维</u>      | ▪ <u>逆向归纳法</u> | ▪ <u>小波分析</u>  |
| ▪ <u>傅立葉變換列表</u> | ▪ <u>主成分分析</u> | ▪ <u>时间序列</u>  |
| ▪ <u>高维数据聚类</u>  | ▪ <u>最小二乘法</u> | ▪ <u>奇异值分解</u> |

## **參考資料**

- Richard Ernest Bellman; Rand CorporationDynamic programming Princeton University Press. 1957.ISBN 978-0-691-07951-6.  
Republished: Richard Ernest Bellman.Dynamic Programming Courier Dover Publications. 2003.ISBN 978-0-486-42809-3.
  - Richard Ernest Bellman.Adaptive control processes: a guided tour Princeton University Press. 1961.
  - doi: 10.1007/978-3-642-13818-8\_34([https://dx.doi.org/10.1007%2F978-3-642-13818-8\\_34](https://dx.doi.org/10.1007%2F978-3-642-13818-8_34))  
本引用來源將會在數十分鐘後自動完成。您可以[检查英文对应模板](#)或[手動擴充](#)
  - doi:10.1007/s11004-008-9156-6(<https://dx.doi.org/10.1007%2Fs11004-008-9156-6>)  
本引用來源將會在數十分鐘後自動完成。您可以[检查英文对应模板](#)或[手動擴充](#)
  - Hughes, G.F., 1968. "On the mean accuracy of statistical pattern recognizers", IEEE Transactions on Information Theory, IT-14:55-63.
  - Not to be confused with the unrelated, but similarly named Hughes effect in electromagnetism (named after Declan C. Hughes (<http://spiedl.aip.org/vsearch/servlet/VerityServlet?KEY=SPIEDL&possible1=Hughes%2C+Declan+C.&possible1zone=author&maxdisp=25&smode=strresults&pjournals=OPEGAR%2CJBOPFO%2CPSISDG%2CJEIME5%2CJMMMGF%2CJARSC4%2CJNOACQ&deliverytype=spiedl&aqs=true>)) which refers to an asymmetry in the hysteresis curves of laminated cores made of certain magnetic materials such as permalloy or mu-metal, in alternating magnetic fields.
  - R. B. Marimont and M. B. Shapiro, "Nearest Neighbour Searches and the Curse of Dimensionality"*Journal of the Institute of Mathematics and its Applications* 24, 1979, 59-70.
  - E. Chavez et al., "Searching in Metric Spaces"*ACM Computing Surveys* 33, 2001, 273-321.
  - doi: 10.1007/978-3-642-22922-0\_25([https://dx.doi.org/10.1007%2F978-3-642-22922-0\\_25](https://dx.doi.org/10.1007%2F978-3-642-22922-0_25))  
本引用來源將會在數十分鐘後自動完成。您可以[检查英文对应模板](#)或[手動擴充](#)
  - Radovanovi?, Milo?; Nanopoulos, Alexandros; Ivanovi?, MirjanaHubs in space: Popular nearest neighbors in high-dimensional data(PDF). *Journal of Machine Learning Research*. 2010**11**: 2487–2531
  - doi: 10.1145/1835449.1835482(<https://dx.doi.org/10.1145%2F1835449.1835482>)  
本引用來源將會在數十分鐘後自動完成。您可以[检查英文对应模板](#)或[手動擴充](#)
- Bellman, R.E. 1957.*Dynamic Programming* Princeton University Press, Princeton, NJ.
    - Republished 2003: Dover ISBN 0486428095
  - Bellman, R.E. 1961.*Adaptive Control Processes* Princeton University Press, Princeton, NJ.
  - Powell, Warren B. 2007.*Approximate Dynamic Programming: Solving the Curses of Dimensionality*Wiley, ISBN 0470171553.

---

取自“<https://zh.wikipedia.org/w/index.php?title=维数灾难&oldid=46355783>”

---

本页面最后修订于2017年9月27日 (星期三) 05:34。

本站的全部文字在[知识共享 署名-相同方式共享 3.0协议](#)之条款下提供，附加条款亦可能应用。（请参阅[使用条款](#)）  
Wikipedia®和维基百科标志是[维基媒体基金会](#)的注册商标；维基™是维基媒体基金会的商标。  
维基媒体基金会是在美国佛罗里达州登记的501(c)(3)[免税](#)、非营利、慈善机构。