

La Maldición de la Dimensionalidad

Angel Luis Valdés Sánchez
Universidad Tecnológica de la Mixteca
Huajuapán de León, Oaxaca, México
angelluis2605@gs.utm.mx

I. ORIGEN DEL TÉRMINO

El concepto de “*maldición de la dimensionalidad*” se refiere a que muchos métodos algorítmicos en R^d se vuelven exponencialmente más difíciles conforme crece la dimensión d [5]. Esta expresión fue acuñada por el matemático Richard E. Bellman para describir el problema causado por el aumento exponencial del volumen al añadir dimensiones adicionales al espacio euclidiano [1]. Bellman introdujo el término a fines de la década de 1950, en el contexto de la optimización de múltiples variables y la programación dinámica. En particular, en su libro *Dynamic Programming* (1957) Bellman advirtió que al incrementar el número de variables de estado en un problema, el espacio de posibles estados crece de forma explosiva, dificultando los cálculos necesarios para encontrar soluciones óptimas [1]. Este fenómeno, descrito originalmente en problemas de control y decisiones dinámicas, fue posteriormente generalizado en su obra *Adaptive Control Processes* (1961) [2] y por otros investigadores, reconociéndose como un obstáculo fundamental en el análisis de datos de alta dimensionalidad [6].

II. CONCEPTO Y FUNDAMENTOS TEÓRICOS

En términos generales, la maldición de la dimensionalidad engloba una serie de fenómenos adversos que aparecen al analizar datos en espacios de alta dimensión y que no se presentan en entornos de baja dimensión [5]. Un tema común en estos problemas es la **escasez de datos**: a medida que la dimensionalidad aumenta, el volumen del espacio crece tan rápidamente que los datos disponibles se vuelven extremadamente dispersos o escasos en relación con el espacio total. Consecuentemente, para obtener resultados fiables, la cantidad de muestras requeridas crece de manera exponencial con el número de dimensiones del dato [6]. En otras palabras, la densidad efectiva de los puntos disminuye drásticamente al añadir nuevas variables, lo que implica que se necesitan muchas más observaciones para cubrir o muestrear adecuadamente el espacio en altas dimensiones.

Otro aspecto clave es la **geometría contraintuitiva** de las altas dimensiones. Ciertas propiedades geométricas cambian radicalmente con d : por ejemplo, el volumen de una hipersfera de radio fijo crece inicialmente con d pero después decrece y tiende a cero en dimensiones muy elevadas, volviéndose insignificante en comparación con el volumen de un hipercubo circunscrito [5]. De hecho, en dimensiones altas casi todo el volumen de un hipercubo se concentra en las “esquinas” lejanas del centro, lo que significa que puntos distribuidos aleatoriamente tienden a estar cerca de los bordes del espacio. Asimismo, las distancias euclídeas tienden a **concentrarse**: la diferencia entre la distancia al vecino más cercano y al más lejano se vuelve despreciable en altas dimensiones [4]. En la práctica, esto implica que la noción de cercanía relativa pierde significado cuando d es muy grande, ya que prácticamente todos los puntos están casi igual de lejos unos de otros.

Estas características conllevan **consecuencias negativas** para muchos algoritmos de minería de datos, aprendizaje automático y

análisis numérico [6]. Técnicas que funcionan bien en bajas dimensiones suelen volverse ineficientes o ineficaces en dimensiones altas. Por ejemplo, muchos métodos de búsqueda y organización de datos se basan en detectar regiones donde los objetos forman grupos con propiedades similares; bajo la maldición de la dimensionalidad, sin embargo, todos los objetos aparecen disímiles entre sí, dificultando la agrupación (*clustering*) y la indexación eficiente de los datos. En el campo del reconocimiento de patrones, un efecto conocido es el *fenómeno de Hughes*: al aumentar el número de características utilizadas en un clasificador, su desempeño inicialmente mejora, pero más allá de cierto punto comienza a degradarse si el número de muestras de entrenamiento es limitado [3]. Este comportamiento, analizado por Hughes en 1968, ilustra otra faceta de la maldición de la dimensionalidad: añadir dimensiones irrelevantes o poco informativas puede empeorar la generalización de un modelo debido a la escasez relativa de datos por dimensión adicional. De manera similar, Beyer *et al.* (1999) demostraron formalmente que en espacios de muy alta dimensión las consultas de vecino más cercano se vuelven inestables: con alta probabilidad, la distancia al punto más cercano es casi indistinguible de la distancia a cualquier otro punto tomado al azar [4]. Resultados como estos subrayan que, en dimensiones elevadas, muchas suposiciones habituales (por ejemplo, la existencia de vecinos “cercaños” significativos o de cúmulos compactos) dejan de ser válidas, a menos que se cuente con un volumen de datos exponencialmente mayor o se apliquen técnicas adecuadas de reducción de dimensionalidad.

III. CONCLUSIÓN

La maldición de la dimensionalidad representa un reto inevitable en la era del análisis de grandes volúmenes de datos. Este fenómeno muestra cómo al aumentar el número de dimensiones los datos se vuelven escasos, las distancias pierden significado y las propiedades geométricas se vuelven contraintuitivas. En otras palabras, lo que en espacios pequeños parece lógico y manejable, en espacios de alta dimensión se convierte en un obstáculo para algoritmos y métodos tradicionales.

REFERENCES

- [1] R. E. Bellman, *Dynamic Programming*. Princeton, NJ: Princeton University Press, 1957.
- [2] R. E. Bellman, *Adaptive Control Processes: A Guided Tour*. Princeton, NJ: Princeton University Press, 1961.
- [3] G. F. Hughes, “On the mean accuracy of statistical pattern recognizers,” *IEEE Trans. Inform. Theory*, vol. 14, no. 1, pp. 55–63, 1968.
- [4] K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft, “When is ‘nearest neighbor’ meaningful?,” in *Proc. 7th Int. Conf. on Database Theory (ICDT)*, 1999, pp. 217–235.
- [5] M. Verleysen and D. François, “The curse of dimensionality in data mining and time series prediction,” in *Proc. Int. Work-Conf. on Artificial Neural Networks (IWANN)*, LNCS 3512, 2005, pp. 758–770.
- [6] A. Zimek, E. Schubert, and H.-P. Kriegel, “A survey on unsupervised outlier detection in high-dimensional numerical data,” *Statistical Analysis and Data Mining*, vol. 5, no. 5, pp. 363–387, 2012.