

# co-ocurrencia – co-occurrence

Authored by  
**memjavad**

November 17, 2025

## RECOMMENDED CITATION

memjavad (2025). *co-ocurrencia – co-occurrence*. Spanish Psychological Databases.  
Retrieved from <https://spanish.arabpsychology.com/?p=4874>

## Co-ocurrencia

**Primary Disciplinary Field(s):** [Estadística](#), [Lingüística Computacional](#), [Minería de Datos](#), [Procesamiento del Lenguaje Natural \(PLN\)](#).

### 1. Definición Central

La **co-ocurrencia**, en su sentido más amplio y fundamental, se define como la manifestación conjunta, adyacente o correlacionada de dos o más elementos discretos dentro de un marco de referencia específico, ya sea espacial, temporal o contextual. Este concepto trasciende disciplinas, sirviendo como piedra angular en el análisis de patrones, la inferencia estadística y la modelización de sistemas complejos. Es crucial entender que la co-ocurrencia no implica inherentemente causalidad, sino meramente una **asociación observada**; la frecuencia con la que A y B aparecen juntos proporciona evidencia empírica para postular una relación subyacente que puede ser de dependencia funcional, similitud semántica o simple proximidad estructural. En el ámbito de la estadística descriptiva, la medición de la co-ocurrencia es el primer paso para determinar la intensidad y la naturaleza de la correlación, utilizando matrices de contingencia para cuantificar la distribución conjunta de las variables. Por ejemplo, en el análisis de texto, la co-ocurrencia de dos palabras dentro de una ventana de contexto limitada sugiere una relación semántica o sintáctica fuerte, mientras que en la epidemiología, la co-ocurrencia de ciertos síntomas puede indicar un síndrome o una patología específica. La definición precisa del "marco de referencia" o "ventana" es metodológicamente vital, ya que un marco demasiado amplio podría diluir las asociaciones significativas, mientras que uno demasiado estrecho podría ignorar dependencias importantes a largo alcance. Esta necesidad de contextualización hace que la co-ocurrencia sea un concepto dinámico cuya interpretación está intrínsecamente ligada al diseño experimental y al dominio de aplicación.

Desde una perspectiva formal, la co-ocurrencia se operacionaliza mediante la cuantificación de la probabilidad conjunta. Si  $P(A)$  es la probabilidad de observar el elemento A y  $P(B)$  es la probabilidad de observar el elemento B, la probabilidad de co-ocurrencia es  $P(A \cap B)$ , es decir, la probabilidad de que ambos A y B se manifiesten simultáneamente. Si la co-ocurrencia observada  $P(A \cap B)$  es significativamente mayor que el producto de sus probabilidades individuales  $P(A) * P(B)$  (asumiendo independencia), se infiere una asociación positiva. Esta desviación de la expectativa de independencia es lo que confiere utilidad analítica al concepto. En la práctica, especialmente en campos como la [minería de reglas de asociación](#), la co-ocurrencia se mide mediante métricas como el soporte (la frecuencia relativa de la aparición conjunta) y la confianza (la probabilidad condicional de B dado A), que permiten descubrir patrones ocultos en grandes conjuntos de datos, como los hábitos de compra de los consumidores o las interacciones genéticas. La robustez y la validez de las inferencias basadas en la co-ocurrencia dependen directamente de la calidad y el tamaño del corpus o conjunto de datos analizado, requiriendo

métodos estadísticos rigurosos para mitigar el riesgo de detectar asociaciones espurias o debidas únicamente al azar.

## 2. Origen y Desarrollo Histórico

Si bien el término "co-ocurrencia" como concepto estadístico y computacional es relativamente moderno, sus raíces filosóficas se encuentran en el empirismo y el asociacionismo, particularmente en las obras de pensadores como John Locke y David Hume, quienes exploraron cómo la mente humana establece conexiones entre ideas que se presentan repetidamente en proximidad temporal o espacial. El desarrollo formal del concepto de asociación y correlación en el siglo XIX, impulsado por figuras como [Francis Galton](#) y [Karl Pearson](#), sentó las bases matemáticas para medir la co-variación de variables, lo que es esencialmente una cuantificación de la co-ocurrencia en datos numéricos. Pearson, con su coeficiente de correlación, proporcionó una herramienta estandarizada para evaluar la fuerza y dirección de la relación lineal entre dos variables, formalizando así la idea de que la aparición conjunta frecuente o sistemática de valores específicos no es aleatoria.

El concepto experimentó una transformación significativa con el auge de la lingüística estructural y, posteriormente, la lingüística computacional en la segunda mitad del siglo XX. [Zellig Harris](#), en su trabajo sobre la distribución de elementos lingüísticos, destacó la importancia de los patrones de co-ocurrencia para definir el significado y la estructura. La hipótesis distribucional, clave en la semántica léxica, postula que las palabras que aparecen en contextos similares (es decir, co-ocurren con palabras similares) tienden a tener significados similares. Este principio fue fundamental para el desarrollo de modelos de vectores de palabras a partir de los años 80 y 90. Con la explosión de los datos digitales y la necesidad de procesar grandes volúmenes de texto (el **corpus**), la co-ocurrencia pasó de ser una herramienta teórica a una métrica práctica indispensable en el [PLN](#). La creación de matrices de co-ocurrencia de términos se convirtió en el método estándar para construir representaciones vectoriales densas del lenguaje, culminando en técnicas modernas de incrustación (embeddings) como Word2Vec y GloVe, que se basan intrínsecamente en la frecuencia y el contexto de la co-ocurrencia para capturar relaciones semánticas complejas.

## 3. Fundamentos Matemáticos y Estadísticos

Desde una óptica matemática, la co-ocurrencia se cuantifica a través de las **matrices de co-ocurrencia**. Estas matrices son estructuras bidimensionales donde las filas representan un conjunto de elementos (por ejemplo, términos, ítems, o eventos A) y las columnas representan otro conjunto de elementos (B), o el mismo conjunto si se analiza la co-ocurrencia interna. Cada celda (i, j) de la matriz almacena la frecuencia o el peso con el que el elemento i y el elemento j aparecen juntos dentro de la ventana de contexto definida. Si bien la matriz de co-ocurrencia más

simple registra el conteo bruto de apariciones conjuntas, su utilidad analítica a menudo requiere una normalización o ponderación. La normalización puede ser probabilística, transformando los conteos en probabilidades conjuntas  $P(i, j)$ , o puede involucrar técnicas de reducción de dimensionalidad como el Análisis de Componentes Principales (PCA) o la Factorización de Matrices No Negativas (NMF) para extraer los patrones de asociación más significativos y reducir el ruido inherente a los datos de alta dimensionalidad.

Un desafío estadístico fundamental es distinguir la co-ocurrencia significativa de la co-ocurrencia esperada por azar. Para ello, se emplean diversas medidas de asociación que contrastan la frecuencia observada con la frecuencia esperada bajo la hipótesis nula de independencia. El [Punto de Información Mutua \(PMI\)](#) es una de las métricas más utilizadas, especialmente en PLN. El PMI mide la dependencia entre dos variables, cuantificando cuánta información se obtiene sobre una variable al observar la otra. Un PMI alto indica que la co-ocurrencia observada es mucho más probable de lo que se esperaría si los elementos fueran independientes. Otras medidas incluyen el coeficiente de Dice, la razón de verosimilitud (likelihood ratio) y el estadístico chi-cuadrado, cada uno con sus propias fortalezas y debilidades para manejar la dispersión de datos y el problema de las frecuencias bajas (sparse data). La selección de la métrica adecuada es crucial, ya que el PMI, por ejemplo, tiende a sobreestimar la importancia de las co-ocurrencias raras, mientras que el soporte bruto puede verse dominado por elementos de alta frecuencia.

#### 4. Aplicaciones en Lingüística y Procesamiento del Lenguaje Natural

En la [Lingüística Computacional](#) y el PLN, la co-ocurrencia es posiblemente el concepto más importante para la representación del significado. La Semántica Distribucional se basa en la premisa de que el significado de una palabra se refleja en el conjunto de contextos en los que aparece. Al construir matrices de co-ocurrencia de palabras a partir de grandes corpus textuales, donde las filas son las palabras objetivo y las columnas son las palabras de contexto dentro de una ventana de tamaño fijo (por ejemplo, 4 palabras a la izquierda y 4 a la derecha), se genera un vector de características para cada palabra. Estos vectores, conocidos inicialmente como modelos de espacio vectorial o VSM (Vector Space Models), permiten calcular la similitud semántica entre palabras mediante métricas de distancia vectorial, como la similitud coseno. Palabras como "gato" y "felino" tendrán vectores de co-ocurrencia muy similares porque tienden a aparecer en los mismos contextos ("tiene bigotes", "caza ratones", "animal doméstico"), lo que permite a los sistemas informáticos inferir que sus significados están estrechamente relacionados.

Más allá de la simple similitud léxica, la co-ocurrencia es vital para la extracción de frases clave, la desambiguación del sentido de las palabras y el modelado de temas (topic modeling). En el caso del modelado de temas, algoritmos como la [Asignación Latente de Dirichlet \(LDA\)](#) identifican temas latentes examinando qué palabras co-ocurren consistentemente dentro de los mismos documentos. Un conjunto de palabras con alta co-ocurrencia (como "gen", "ADN", "célula") define

un tema (Biología). La evolución de estos métodos ha llevado a los modelos de incrustación neuronal (neural embeddings), como Word2Vec, GloVe (Global Vectors for Word Representation), y modelos basados en transformadores (BERT), que, aunque más sofisticados, siguen basándose fundamentalmente en la idea de que la información contextual de la co-ocurrencia es el proxy más efectivo para el significado. GloVe, por ejemplo, optimiza directamente la información global de las matrices de co-ocurrencia para crear representaciones vectoriales que no solo capturan la proximidad semántica, sino también relaciones analógicas complejas (por ejemplo, "Rey" es a "Hombre" como "Reina" es a "Mujer"), demostrando la profunda capacidad de la co-ocurrencia para codificar la estructura del lenguaje.

## 5. Aplicaciones en Minería de Datos y Análisis de Redes

En el campo de la **Minería de Datos**, la co-ocurrencia es la base operativa de la minería de reglas de asociación, una técnica fundamental para el análisis de cestas de mercado. El objetivo es descubrir patrones de co-ocurrencia frecuentes en grandes bases de datos transaccionales. Por ejemplo, si los ítems A (pan) y B (leche) co-ocurren frecuentemente en las mismas transacciones, se puede formular la regla de asociación "Si se compra A, es probable que se compre B". El algoritmo Apriori, pionero en este campo, utiliza el concepto de co-ocurrencia para identificar conjuntos de ítems frecuentes, empleando las métricas de soporte (frecuencia de co-ocurrencia) y confianza (probabilidad condicional) para filtrar y evaluar la utilidad de las reglas. Estas reglas de co-ocurrencia tienen un impacto directo en la toma de decisiones empresariales, influenciando la colocación de productos, las estrategias de ventas cruzadas y la personalización de recomendaciones en plataformas de comercio electrónico, demostrando cómo una simple estadística de aparición conjunta se traduce en valor económico y estratégico.

Además, la co-ocurrencia juega un papel crucial en el [Análisis de Redes](#). En este contexto, la co-ocurrencia define las conexiones o aristas entre los nodos. Por ejemplo, en las redes de colaboración científica, la co-ocurrencia de nombres de autores en el mismo artículo indica una colaboración. En las redes de conocimiento o redes semánticas, la co-ocurrencia de conceptos en el mismo texto o documento establece una relación de proximidad conceptual. La matriz de co-ocurrencia se convierte, de facto, en la matriz de adyacencia de la red, permitiendo la aplicación de algoritmos de teoría de grafos para medir la centralidad de los nodos, detectar comunidades (grupos de elementos con alta co-ocurrencia interna) y comprender la estructura general del sistema. Esta aplicación es vital en campos tan diversos como la sociología (análisis de redes sociales), la biología (interacciones proteína-proteína) y la bibliometría (mapeo de campos de investigación).

## 6. Tipos y Medidas de Co-ocurrencia

La co-ocurrencia se clasifica habitualmente en función de su alcance. La **co-ocurrencia local** se

refiere a la aparición conjunta de elementos dentro de una ventana de contexto muy restringida, como la proximidad inmediata de palabras en una oración o de genes en una región cromosómica. Este tipo es fundamental para capturar relaciones sintácticas o funcionales directas. Por otro lado, la **co-ocurrencia global** considera la aparición conjunta de elementos en unidades más grandes, como documentos, artículos o sesiones de compra completas. La co-ocurrencia global es más efectiva para detectar relaciones temáticas o asociaciones a largo plazo. Los modelos de incrustación modernos a menudo buscan un equilibrio entre estas dos perspectivas, utilizando técnicas que incorporan tanto la información de contexto local (ventanas deslizantes) como la estadística agregada global (matrices de conteo) para optimizar la calidad de las representaciones vectoriales.

En cuanto a las métricas específicas utilizadas para ponderar o medir la co-ocurrencia, la elección depende del objetivo analítico y de las propiedades estadísticas de los datos.

**Soporte (Support):** Es la métrica más básica, definida como la frecuencia relativa de la co-ocurrencia  $P(A \cap B)$ . Es esencial para identificar patrones frecuentes, pero no mide la fuerza real de la asociación, ya que es sensible a las frecuencias marginales de A y B.

**Confianza (Confidence):** Mide la probabilidad condicional  $P(B|A)$ . Es útil para reglas de asociación unidireccionales (Si A, entonces B), pero ignora la frecuencia de B, lo que puede llevar a reglas engañosas si B es extremadamente común.

**Elevación (Lift):** Corrige las deficiencias del soporte y la confianza al comparar la co-ocurrencia observada con la co-ocurrencia esperada bajo independencia:  $Lift(A, B) = P(A \cap B) / (P(A) * P(B))$ . Un valor de elevación mayor que 1.0 indica una asociación positiva fuerte, mientras que un valor menor que 1.0 sugiere una asociación negativa (exclusión mutua). El Lift es la métrica preferida para evaluar la verdadera fuerza de la asociación.

**Punto de Información Mutua (PMI):** Utilizado predominantemente en PLN, mide cuánto más probable es la co-ocurrencia de A y B de lo que se esperaría por azar, utilizando logaritmos. El PMI penaliza las co-ocurrencias muy comunes y favorece las asociaciones específicas y raras, siendo útil para identificar colocaciones o términos técnicos.

## 7. Limitaciones y Desafíos Metodológicos

Uno de los principales desafíos metodológicos en el análisis de co-ocurrencia es el problema de la **dispersión de datos** (sparsity). En conjuntos de datos grandes, como corpus de texto masivos, la mayoría de los pares de elementos posibles nunca co-ocurren, o lo hacen con una frecuencia extremadamente baja. Esto resulta en matrices de co-ocurrencia que son vastas pero casi vacías, lo que dificulta la estimación estadística precisa y puede llevar a que las métricas (como el PMI) sobrestimen la importancia de eventos raros. Las soluciones comunes incluyen la aplicación de técnicas de suavizado (smoothing), el uso de PMI positivo (PPMI), que reemplaza los valores negativos con cero, y la reducción drástica de la dimensionalidad, a menudo mediante

factorización matricial, para proyectar los datos en un espacio vectorial denso donde las relaciones latentes sean más evidentes y computacionalmente manejables.

Otro desafío crucial es la definición adecuada del **contexto**. La utilidad de la co-ocurrencia depende críticamente de cómo se define la "ventana" de aparición conjunta. En PLN, una ventana demasiado grande (por ejemplo, todo el documento) capturará relaciones temáticas generales, pero perderá las relaciones sintácticas finas. Una ventana demasiado pequeña (por ejemplo, una palabra adyacente) puede ignorar dependencias importantes a distancia. La elección del tamaño de la ventana es un hiperparámetro que debe ajustarse al objetivo específico del análisis. Además, la co-ocurrencia solo captura relaciones de primer orden; es decir, mide la relación directa entre A y B. Sin embargo, muchas relaciones interesantes son indirectas (A se relaciona con C a través de B). El análisis de redes y las incrustaciones vectoriales abordan esto al permitir inferir la similitud entre A y C si ambos co-ocurren frecuentemente con el mismo conjunto de elementos intermedios B, aprovechando así la información de co-ocurrencia de segundo orden.

Finalmente, la limitación conceptual más significativa es la distinción entre **co-ocurrencia y causalidad**. La co-ocurrencia es una medida de asociación; el hecho de que A y B aparezcan juntos frecuentemente no proporciona, por sí mismo, ninguna evidencia de que A cause B, o viceversa, o que ambos sean causados por una tercera variable latente C. Esta limitación requiere que los hallazgos de co-ocurrencia sean siempre interpretados a la luz de la teoría del dominio y complementados con métodos causales (como experimentos controlados o modelos causales bayesianos) para establecer la dirección y la naturaleza de la influencia. Ignorar esta distinción puede llevar a conclusiones erróneas, un problema conocido en estadística como la falacia de la correlación o la asociación espuria.

## 8. Conclusión e Impacto Interdisciplinario

La co-ocurrencia, partiendo de una simple observación empírica de proximidad, se ha consolidado como un principio organizador fundamental en la ciencia de datos moderna. Su sencillez conceptual, combinada con su poder predictivo cuando se aplica a grandes volúmenes de datos, la convierte en una herramienta indispensable. Desde la comprensión de cómo el cerebro humano asocia ideas hasta el desarrollo de la inteligencia artificial más avanzada, la capacidad de cuantificar la aparición conjunta de elementos ha permitido la transición de la mera recopilación de datos a la extracción de significado y la inferencia de estructuras latentes. El impacto de la co-ocurrencia es profundamente **interdisciplinario**, siendo un lenguaje común que conecta la estadística, la informática, la psicología, la sociología y la biología, proporcionando un marco unificado para el descubrimiento de patrones y la modelización de dependencias.

En el futuro, el concepto de co-ocurrencia seguirá evolucionando, especialmente a medida que los modelos de aprendizaje profundo busquen capturar no solo la frecuencia, sino también la

estructura jerárquica y las dependencias temporales y espaciales complejas. La tendencia actual se centra en desarrollar modelos que puedan distinguir diferentes tipos de co-ocurrencia (sintáctica vs. semántica, causal vs. incidental) de manera más explícita. Sin embargo, incluso en los sistemas más sofisticados, la matriz de co-ocurrencia, ya sea explícita o implícita, permanece como el cimiento sobre el cual se construyen la mayoría de las representaciones de conocimiento y los sistemas de recomendación, asegurando su relevancia continua como una de las herramientas analíticas más robustas y universales en la era del big data.

## Further Reading

[Estadística - Wikipedia](#)

[Lingüística Computacional - Wikipedia](#)

[Minería de Datos - Wikipedia](#)

[Procesamiento del Lenguaje Natural - Wikipedia](#)

[Minería de reglas de asociación - Wikipedia](#)

[Francis Galton - Wikipedia](#)

[Karl Pearson - Wikipedia](#)

[Zellig Harris - Wikipedia](#)

[Punto de Información Mutua - Wikipedia](#)

[Asignación Latente de Dirichlet - Wikipedia](#)

[Análisis de redes - Wikipedia](#)