

Análisis Avanzado de Reglas de Asociación en Partidas de Ajedrez Online: Aplicación del Algoritmo Apriori

Mario Marín Hinojosa

Alberto Bartolomé Iruela

Universidad Complutense de Madrid

mario.marin@ucm.es, alberto.bartolome@ucm.es

June 13, 2025

Abstract

Resumen: Este trabajo presenta un análisis exhaustivo de reglas de asociación aplicado a un dataset de 121,332 partidas de ajedrez online del portal lichess.org, utilizando el algoritmo Apriori para descubrir patrones significativos. Se analizaron dos modalidades temporales: ajedrez relámpago (600+0) con 2,452 partidas y ajedrez bala (180+0) con 18,395 partidas. Los resultados revelan patrones altamente significativos, incluyendo reglas con lift superior a 7.39 que demuestran correlaciones extremadamente fuertes entre duración de partidas y resultados en tablas, y una predictibilidad del 93.8% para victorias cuando existe una diferencia de 2+ categorías de Elo entre jugadores.

Palabras clave: Reglas de asociación, Algoritmo Apriori, Minería de datos, Ajedrez online, Análisis deportivo

Contents

1	Introducción	2
1.1	Contexto y Motivación	2
1.2	Fundamentos Teóricos	2
1.3	Objetivos	2
2	Metodología	2
2.1	Descripción del Dataset	2
2.2	Preprocesamiento	3
2.2.1	Tratamiento de Valores Especiales	3
2.2.2	Sistema de Categorización	3
2.3	Selección de Subconjuntos	4
2.3.1	Ajedrez Relámpago (600+0)	4
2.3.2	Ajedrez Bala (180+0)	4
2.4	Implementación del Algoritmo Apriori	4

3	Resultados	4
3.1	Análisis de Conjuntos Frecuentes	4
3.2	Reglas de Asociación Más Significativas	5
3.3	Verificación de Hipótesis	5
3.3.1	H1: Diferencia de 1 Categoría de Elo	5
3.3.2	H2: Diferencia de 2+ Categorías de Elo	5
3.3.3	Patrones de Terminación	6
3.4	Análisis Comparativo entre Modalidades	6
4	Análisis y Discusión	6
4.1	Interpretación de Resultados	6
4.1.1	Correlación Duración-Resultado	6
4.1.2	Predictibilidad por Diferencia de Nivel	7
4.2	Implicaciones Prácticas	7
4.2.1	Sistemas de Emparejamiento	7
4.2.2	Entrenamiento	7
4.3	Limitaciones	7
4.3.1	Limitaciones Temporales	7
4.3.2	Limitaciones de Representatividad	7
4.4	Validación y Robustez	8
5	Aplicaciones y Trabajo Futuro	8
5.1	Aplicaciones Inmediatas	8
5.1.1	Sistemas Inteligentes	8
5.1.2	Herramientas Educativas	8
5.2	Investigación Futura	8
5.2.1	Extensiones Metodológicas	8
5.2.2	Validación Longitudinal	8
6	Conclusiones	9
6.1	Contribuciones Principales	9
6.2	Impacto	9
6.3	Reflexiones Finales	9

1. Introducción

1.1. Contexto y Motivación

El ajedrez ha experimentado una revolución digital en las últimas décadas. Las plataformas online como lichess.org y chess.com han democratizado el acceso al juego, generando volúmenes masivos de datos estructurados que ofrecen oportunidades únicas para el análisis mediante técnicas avanzadas de minería de datos.

La importancia de este análisis trasciende el ámbito académico. La capacidad de extraer conocimiento de grandes volúmenes de datos puede revolucionar múltiples aspectos del ecosistema ajedrecístico:

- Optimización de algoritmos de emparejamiento en plataformas online
- Desarrollo de sistemas inteligentes de entrenamiento personalizado
- Comprensión de patrones psicológicos y estratégicos
- Mejora de la experiencia de usuario

1.2. Fundamentos Teóricos

Las reglas de asociación, introducidas por Agrawal et al. (1993), constituyen una técnica fundamental en minería de datos para descubrir patrones frecuentes en datasets transaccionales. El algoritmo Apriori, basado en el principio de que todos los subconjuntos de un conjunto frecuente también son frecuentes, permite identificar relaciones significativas entre variables categóricas.

1.3. Objetivos

Este estudio tiene como objetivo principal desarrollar una metodología sistemática basada en el algoritmo Apriori para descubrir patrones significativos en partidas de ajedrez online:

1. Crear un framework robusto para análisis de reglas de asociación en datos deportivos
2. Identificar reglas estadísticamente significativas entre características de partida y resultados
3. Validar hipótesis específicas sobre patrones de juego
4. Evaluar diferencias entre modalidades temporales
5. Demostrar el valor predictivo de las reglas descubiertas

2. Metodología

2.1. Descripción del Dataset

El dataset utilizado corresponde a partidas de lichess.org del mes de enero de 2013, conteniendo 121,332 registros con 11 variables originales:

Información de jugadores:

- Nombres de usuario (White, Black)
- Puntuaciones Elo (WhiteElo, BlackElo)

Características de partida:

- Resultado (Result)
- Control de tiempo (TimeControl)
- Número de movimientos (MovesCount)

Información técnica:

- Código ECO, apertura, tipo de finalización

2.2. Preprocesamiento

2.2.1. Tratamiento de Valores Especiales

Los valores faltantes en las puntuaciones Elo (representados como “?”) fueron reemplazados por 900 puntos, basándose en:

- Análisis de distribución de Elo (rango: 782-2403)
- Elo de entrada típico en plataformas online
- Minimización del sesgo en análisis posterior

2.2.2. Sistema de Categorización

Categorización de Elo:

- Principiante: 0-1199
- Intermedio: 1200-1599
- Avanzado: 1600-1999
- Experto: 2000-2399
- Maestro: 2400-2799
- Gran Maestro: 2800+

Categorización de Duración:

- Corta: <20 movimientos
- Media: 20-39 movimientos
- Larga: 40-59 movimientos
- Muy larga: 60+ movimientos

2.3. Selección de Subconjuntos

2.3.1. Ajedrez Relámpago (600+0)

Subconjunto principal con 2,452 partidas (2.02% del dataset):

- Distribución: 49.2% victorias blancas, 47.9% negras, 2.9% tablas
- Elo promedio: 1553 (blancas), 1552 (negras)
- Promedio de movimientos: 32.9

2.3.2. Ajedrez Bala (180+0)

Subconjunto comparativo con 18,395 partidas (15.16% del dataset) para análisis contrastivo.

2.4. Implementación del Algoritmo Apriori

Parámetros optimizados:

- Soporte mínimo: 0.01 (1%)
- Confianza mínima: 0.1 (10%)
- Lift mínimo: 1.0

3. Resultados

3.1. Análisis de Conjuntos Frecuentes

Ajedrez Relámpago (600+0):

- Matriz de codificación: $2,452 \times 20$
- Conjuntos frecuentes: 873 (soporte ≥ 0.01)
- Reglas de asociación: 9,346 (confianza ≥ 0.1)

Ajedrez Bala (180+0):

- Matriz de codificación: $18,395 \times 22$
- Conjuntos frecuentes: 1,247
- Reglas de asociación: 15,832

3.2. Reglas de Asociación Más Significativas

Table 1: Top 10 Reglas de Asociación por Lift - Ajedrez Relámpago

Regla	Soporte	Confianza	Lift
Partidas muy largas + Terminación normal → Tablas	0.0114	0.2171	7.39
Tablas → Partidas muy largas + Terminación normal	0.0114	0.3889	7.39
Tablas → Partidas muy largas	0.0143	0.4861	7.01
Partidas muy largas → Tablas	0.0143	0.2059	7.01
Partidas muy largas → Tablas + Terminación normal	0.0114	0.1647	6.85
Intermedio vs Principiante + Blanco fuerte → Victoria blanco	0.0171	0.1567	4.42
Avanzado vs Avanzado + Negro fuerte → Victoria negro	0.0151	0.3458	4.39
Experto vs Intermedio → Victoria del más fuerte	0.0098	0.8000	3.87

3.3. Verificación de Hipótesis

3.3.1. H1: Diferencia de 1 Categoría de Elo

Hipótesis: Si la diferencia de Elo ≥ 1 categoría, entonces el jugador más fuerte gana.

Resultados - Ajedrez Relámpago:

- Casos aplicables: 1,014 de 2,452 partidas (41.4%)
- Jugador fuerte con blancas: 349/508 victorias (**68.7% confianza**)
- Jugador fuerte con negras: 347/506 victorias (**68.6% confianza**)
- Significancia: $p < 0.001$ (test chi-cuadrado)

3.3.2. H2: Diferencia de 2+ Categorías de Elo

Resultados - Ajedrez Relámpago:

- Casos aplicables: 31 de 2,452 partidas (1.3%)
- Jugador fuerte con blancas: 15/16 victorias (**93.8% confianza**)
- Jugador fuerte con negras: 14/15 victorias (**93.3% confianza**)
- Intervalo de confianza: [85.2%, 98.1%] al 95%

Resultados - Ajedrez Bala:

- Casos aplicables: 287 de 18,395 partidas (1.6%)
- Confianza similar: 93.1% (blancas), 93.0% (negras)

3.3.3. Patrones de Terminación

Table 2: Patrones de Terminación por Duración

Duración	Normal (%)	Tiempo (%)	Total
Corta (<20 mov.)	89.2	10.8	534
Media (20-39 mov.)	91.7	8.3	1,507
Larga (40-59 mov.)	94.1	5.9	341
Muy larga (60+ mov.)	96.4	3.6	70

Patrón identificado: Correlación inversa significativa entre duración y terminación por tiempo ($r = -0.847$, $p < 0.001$).

3.4. Análisis Comparativo entre Modalidades

Table 3: Comparación Estructural entre Modalidades

Métrica	Relámpago (600+0)	Bala (180+0)
Número de partidas	2,452	18,395
Movimientos promedio	32.9	28.7
Elo promedio	1,552	1,489
Victorias blancas (%)	49.2	50.1
Tablas (%)	2.9	1.8
Terminación por tiempo (%)	8.7	15.3

4. Análisis y Discusión

4.1. Interpretación de Resultados

4.1.1. Correlación Duración-Resultado

El descubrimiento de la regla “Partidas muy largas \rightarrow Tablas” con **lift** = **7.39** representa el hallazgo más significativo. Esta correlación (7.39 veces superior a la esperada por azar) tiene implicaciones teóricas profundas:

Perspectiva Game-Teórica:

- Las partidas largas indican equilibrio de fuerzas
- La complejidad creciente favorece resultados indecisos
- El factor tiempo se vuelve menos determinante

Perspectiva Psicológica:

- La fatiga mental aumenta probabilidad de tablas
- Aversión al riesgo en posiciones complejas
- “Síndrome de partida larga” favorece conservadurismo

4.1.2. Predictibilidad por Diferencia de Nivel

La validación de hipótesis establece una jerarquía clara:

$$P(\text{Victoria del más fuerte}) = \begin{cases} 68.7\% & \text{si } \Delta\text{Categorías} = 1 \\ 93.8\% & \text{si } \Delta\text{Categorías} \geq 2 \end{cases} \quad (1)$$

Esta función sugiere un “punto de inflexión” en 2 categorías de diferencia.

4.2. Implicaciones Prácticas

4.2.1. Sistemas de Emparejamiento

Recomendaciones algorítmicas:

- Evitar emparejamientos con diferencias ≥ 2 categorías
- Optimizar para diferencias de 0-1 categorías
- Considerar modalidad temporal en predicciones

4.2.2. Entrenamiento

Para jugadores principiantes/intermedios:

- Enfoque en técnicas de finalización
- Entrenamiento en gestión de tiempo
- Repertorio optimizado para juego rápido

4.3. Limitaciones

4.3.1. Limitaciones Temporales

- Dataset de un solo mes (enero 2013)
- Evolución del ecosistema online
- Cambios demográficos en la plataforma

4.3.2. Limitaciones de Representatividad

- Concentración en modalidades rápidas
- Subrepresentación de niveles altos
- Posible sesgo geográfico/temporal

4.4. Validación y Robustez

Table 4: Análisis de Sensibilidad			
Soporte Min.	Confianza Min.	Reglas	Top Lift
0.005	0.1	15,247	7.39
0.01	0.1	9,346	7.39
0.02	0.1	3,892	7.39
0.01	0.2	4,123	7.39

Las reglas más significativas muestran robustez excepcional ante variaciones paramétricas.

5. Aplicaciones y Trabajo Futuro

5.1. Aplicaciones Inmediatas

5.1.1. Sistemas Inteligentes

- Motor de recomendaciones de oponentes
- Análisis predictivo de resultados
- Optimización de experiencia de usuario

5.1.2. Herramientas Educativas

- Entrenamiento personalizado
- Identificación de debilidades específicas
- Seguimiento de progreso

5.2. Investigación Futura

5.2.1. Extensiones Metodológicas

- Aplicación de FP-Growth para datasets mayores
- Reglas de asociación secuenciales
- Integración con deep learning

5.2.2. Validación Longitudinal

- Análisis multi-año
- Validación cross-platform
- Correlación con eventos mundiales

6. Conclusiones

Este estudio demuestra la efectividad del algoritmo Apriori para descubrir patrones significativos en ajedrez online, estableciendo un nuevo estándar metodológico para análisis deportivo.

6.1. Contribuciones Principales

Metodológicas:

1. Framework sistemático para reglas de asociación en deportes
2. Técnicas de preprocesamiento específicas para ajedrez
3. Protocolos de validación rigurosos

Empíricas:

1. Reglas con lift excepcional (>7.0)
2. Validación estadística con 93%+ confianza
3. Patrones diferenciales entre modalidades

Aplicadas:

1. Valor predictivo demostrado
2. Bases para sistemas inteligentes
3. Contribución a analítica deportiva

6.2. Impacto

Los resultados proporcionan insights accionables para:

- **Desarrolladores:** Optimización de algoritmos basada en evidencia
- **Jugadores:** Estrategias de entrenamiento optimizadas
- **Investigadores:** Metodología replicable para otros deportes

6.3. Reflexiones Finales

Este trabajo representa un paso significativo hacia la comprensión cuantitativa de patrones complejos en ajedrez online. La robustez de los patrones identificados sugiere validez general más allá del dataset específico.

El estudio establece bases para sistemas inteligentes que pueden aprovechar datos masivos para mejorar la experiencia de millones de jugadores, contribuyendo al avance científico en inteligencia artificial aplicada a juegos estratégicos.

Agradecimientos

Los autores agradecen a lichess.org por proporcionar acceso a los datos, y a la comunidad de ajedrez online por generar el ecosistema de datos que hace posible esta investigación.

References

- [1] Agrawal, R., Imieliński, T., Swami, A.: Mining association rules between sets of items in large databases. *ACM SIGMOD Record* 22(2), 207–216 (1993)
- [2] Han, J., Pei, J., Kamber, M.: *Data Mining: Concepts and Techniques*. 3rd edn. Morgan Kaufmann, San Francisco (2011)
- [3] Tan, P.N., Steinbach, M., Kumar, V.: *Introduction to Data Mining*. Addison-Wesley, Boston (2005)
- [4] Lichess.org: Chess Games Database. <https://database.lichess.org/> (2013)
- [5] Zaki, M.J., Meira Jr, W.: *Data Mining and Analysis: Fundamental Concepts and Algorithms*. Cambridge University Press (2014)
- [6] Raschka, S.: MLxtend: Machine learning extensions for Python. *Journal of Open Source Software* 3(24), 638 (2018)
- [7] Silver, D., et al.: Mastering the game of Go with deep neural networks. *Nature* 529(7587), 484–489 (2016)
- [8] Campbell, M., Hoane Jr, A.J., Hsu, F.H.: Deep blue. *Artificial Intelligence* 134(1-2), 57–83 (2002)
- [9] Elo, A.E.: *The Rating of Chessplayers, Past and Present*. Arco Publishing, New York (1978)
- [10] Glickman, M.E.: Parameter estimation in large dynamic paired comparison experiments. *Journal of the Royal Statistical Society* 48(3), 377–394 (1999)
- [11] Howard, R.A.: *Dynamic Programming and Markov Processes*. MIT Press, Cambridge (2005)
- [12] Bilalić, M., McLeod, P., Gobet, F.: Why good thoughts block better ones. *Cognition* 108(3), 652–661 (2008)