



## Overview of Interestingness Measures for Data Mining

**Dr. Swati Joshi**

*Pune Vidyarthi Griha's College of Science and Commerce, Pune*

*Corresponding Author – Dr. Swati Joshi*

**DOI - 10.5281/zenodo.15533028**

### **Abstract:**

*In data mining, the primary goal is to discover hidden patterns and knowledge from large datasets. However, not all patterns discovered during the mining process are useful or relevant. To enhance the efficacy of data mining, interestingness measures are developed to evaluate the quality and significance of these patterns. This paper explores the concept of interestingness measures in data mining, categorizes existing approaches, and proposes new metrics for improving the relevance and usefulness of discovered patterns. We also discuss the implications of these measures for real-world applications in various domains, such as business, healthcare, and social sciences.*

### **Introduction:**

Data mining has become an essential tool for discovering patterns and knowledge from large datasets. However, with the vast amount of data generated in modern applications, not all discovered patterns are valuable or interesting. The concept of "interestingness" is central to this challenge. Interestingness measures serve as evaluation criteria for filtering patterns that are not only statistically significant but also relevant and actionable. This paper defines various types of interestingness measures, examines the existing research on the topic, and proposes new frameworks for evaluating pattern quality. The vast number of patterns generated during mining can overwhelm the data analyst. The ability to prioritize interesting patterns helps in focusing on insights that have practical implications, saving both time and resources.

### **Objectives:**

1. Define the concept of interestingness in data mining.

2. Classify existing interestingness measures.
3. Propose new measures that address current limitations.
4. Discuss the application of these measures in various domains.

### **Background and Related Work:**

#### **Data Mining and Knowledge Discovery:**

Data mining refers to the process of discovering patterns from large datasets, often through techniques like classification, clustering, association rule mining, and anomaly detection. Knowledge discovery, a broader field, includes all aspects of extracting useful information from data.

#### **Definition of Interestingness:**

Interestingness in the context of data mining refers to the quality of a pattern, rule, or model, which indicates its value for the end-user. Various definitions exist, and these definitions depend on the type of mining technique being used and the domain in which the data is applied. Various interestingness measures, such as lift,

correlation and all-confidence have been proposed for discovering useful association rules. Each measure has its own selection bias that justifies the rationale for preferring a set of association rules over another. As a result, selecting a right interestingness measure for mining association rules is a tricky problem.

#### Existing Approaches:

- **Statistical Measures:** The classic measures of support, confidence, and lift are commonly used in association rule mining. These measures evaluate the frequency and strength of relationships between items.
- **Actionability:** Some interestingness measures aim to prioritize patterns that lead to actionable insights. These measures focus on patterns that have a direct impact on decision-making.
- **Novelty:** Novelty measures evaluate the degree to which a pattern differs from known patterns, identifying new insights.
- **Utility-based Measures:** These measures focus on the utility or profitability of the discovered pattern. For instance, in market basket analysis, a pattern that predicts a profitable product combination is considered interesting.

#### Categories of Interestingness Measures:

Existing measures can be broadly classified into two types:

**1. Objective measures:** These are purely statistical and often include factors like

support, confidence, lift, and correlation. They focus on the mathematical properties of patterns. [9] Objective measures are based on the probability theory, statistics, or information theory. The objective measures do not require any prior knowledge about the user or domain. Objective measure measures the interestingness of an association rule in terms of the structure and the underlying data used in the discovery process. An objective measure is usually computed on the frequency counts tabulated in a contingency table.

**2. Subjective measures:** These depend on the application context and may include user preferences, domain-specific knowledge, and action ability of patterns. Subjective measures considers both the data and the user. A pattern is said to be subjectively interesting if it discovers unexpected information about the data or such knowledge which could lead to profitable results. To define a subjective measure, access to the user's domain or background knowledge about the data is required. Subjective measures recognize that a pattern of interest to one user may or may not be of interest to another user [2, 3]. In [3], proposed unexpectedness and action ability as the two measures of subjective interestingness. Negative encoding length and temporal description length have been used as subjective measures in [4] and [1], respectively.

An objective measure can be either symmetric or asymmetric. These measures are defined as below.

Measure	Formula
Correlation ( $\phi$ )	$\frac{Nf_{11}-f_{1+}f_{+1}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$
Odds ratio ( $\alpha$ )	$\frac{f_{11}f_{00}}{f_{10}f_{01}}$
Kappa ( $\kappa$ )	$\frac{Nf_{11}+Nf_{00}-f_{1+}f_{+1}-f_{0+}f_{+0}}{N^2-f_{1+}f_{+1}-f_{0+}f_{+0}}$
Lift ( $I$ )	$\frac{Nf_{11}}{f_{1+}f_{+1}}$
Cosine ( $IS$ )	$\frac{f_{11}}{\sqrt{f_{1+}f_{+1}}}$
Piatetsky-Shapiro ( $PS$ )	$\frac{f_{11}}{N} - \frac{f_{1+}f_{+1}}{N^2}$
Collective strength ( $S$ )	$\left(\frac{f_{11}+f_{00}}{f_{1+}f_{+1}+f_{0+}f_{+0}}\right) \times \left(\frac{N^2-f_{1+}f_{+1}-f_{0+}f_{+0}}{N-f_{11}-f_{00}}\right)$
All-confidence ( $h$ )	$\min \left[ \frac{f_{11}}{f_{1+}}, \frac{f_{11}}{f_{+1}} \right]$
Imbalance Ratio ( $IR$ )	$\frac{ f_{10}-f_{01} }{f_{11}+f_{10}+f_{01}}$
Jaccard ( $\zeta$ )	$\frac{f_{11}}{f_{1+}+f_{+1}-f_{11}}$

Symmetric interestingness measures

Measure	Formula
Confidence ( $conf$ )	$\frac{f_{11}}{f_{1+}}$
Goodman-Kruskal ( $\lambda$ )	$\frac{(\sum_j \max_k  f_{jk} - \max_k f_{jk} )}{N - \max_k f_{+k}}$
Mutual Information ( $M$ )	$\frac{(\sum_j \sum_i \frac{f_{ij}}{N} \log \frac{Nf_{ij}}{f_{i+}f_{+j}})}{(-\sum_i \frac{f_{i+}}{N} \log \frac{f_{i+}}{N})}$
J-Measure ( $J$ )	$\frac{f_{11}}{N} \log \frac{Nf_{11}}{f_{1+}f_{+1}} + \frac{f_{10}}{N} \log \frac{Nf_{10}}{f_{1+}f_{+0}}$
Gini index ( $G$ )	$\frac{f_{1+}}{N} \times \left[ \frac{(f_{11})^2 + (f_{10})^2}{(f_{1+})^2} - \left(\frac{f_{1+}}{N}\right)^2 \right] + \frac{f_{0+}}{N} \times \left[ \frac{(f_{01})^2 + (f_{00})^2}{(f_{0+})^2} - \left(\frac{f_{0+}}{N}\right)^2 \right]$
Laplace ( $L$ )	$(f_{11}+1)/(f_{1+}+2)$
Conviction ( $V$ )	$(f_{1+}f_{+0})/(Nf_{10})$
Certainty factor ( $F$ )	$(\frac{f_{11}}{f_{1+}} - \frac{f_{+1}}{N}) / (1 - \frac{f_{+1}}{N})$
Added Value ( $AV$ )	$\frac{f_{11}}{f_{1+}} - \frac{f_{+1}}{N}$

Asymmetric Interestingness Measures

### Challenges in Defining Interestingness:

Traditional interestingness measures are domain-independent and not reflecting real-world relevance hence often criticized. Defining interestingness in data mining is challenging because of subjective and context-dependent nature of "interesting." The concept of interestingness in data mining is crucial as it determines what patterns, insights, or knowledge should be considered valuable or noteworthy. Hence, several challenges arise in its definition and application.

Defining interestingness in data mining requires balancing various factors, such as user goals, domain-specific needs, computational constraints, and the subjective nature of interest. A flexible and context-sensitive approach is required to ensure that patterns deemed interesting are valuable and relevant to the end-users. Ethical implications and long-term consequences is vital in the process of defining interestingness.

### Subjectivity of Interestingness:

- **User Dependence:** Things for one person or business considers interesting may not be that much interesting to someone else. For example, for a retailer customer purchasing behavior may be interesting, while for a data

scientist might be patterns based on statistical significance.

- **Context Sensitivity:** Interestingness can change depending on the context. A pattern may seem interesting in one scenario (e.g., predicting customer churn) but may not be relevance in another context (e.g., predicting fraud).
- **Cultural and Domain Variability:** Different industries, domains, and cultures may interpret and define interestingness in multiple ways. What works for one dataset or application might not be applicable to another.

### Measuring Interestingness:

- **Objective Measures vs. Subjective Measures:** Defining a quantitative, objective measure of interestingness that can be universally accepted across all types of data is very difficult. Some of the Common measures like statistical significance, support, or confidence may not always align with what users consider "interesting." In contrast, subjective measures often rely on human judgment, which is inconsistent and difficult to formalize.
- **Lack of Universality:** Different mining algorithms might generate results with varying degrees of usefulness or relevance, making it hard to compare the "interestingness" of one pattern over another.

- **Complexity of Multi-dimensional Data:** In datasets with multiple attributes or features, the definition of interestingness can become more complicated. What is interesting in a high-dimensional space might not be obvious in lower-dimensional spaces, and vice versa.

#### Discovery of Unexpected Patterns:

- **Novelty and Surprise:** Patterns which are predictable or obvious may not be considered interesting. Unexpected, surprising findings could be interesting. However, determining whether a surprising pattern is truly valuable or just a statistical anomaly is difficult.
- **Balance between Novelty and Relevance:** It's challenging to balance novelty and relevance while defining interestingness. Focus on novelty could lead to patterns that are not practically useful, while focusing on relevance might overlook innovative findings that could be game-changing.

#### Dynamic and Evolving Data:

- **Changing Preferences over Time:** Interestingness can change over time as the data evolves. For instance, business goals might shift, or consumer behavior might change, making previous patterns less interesting or useful.
- **Scalability of Interestingness Metrics:** As the volume of data increases, the computational cost of determining interesting patterns becomes much higher. A pattern that is interesting in a small dataset may become irrelevant or hard to discern in a large dataset.

#### Evaluation and Validation:

- **Lack of Ground Truth:** It's often difficult to establish a benchmark for interestingness in data mining. Unlike classification or regression tasks, where accuracy or error metrics provide an objective measure, interestingness lacks a definitive validation process.

- **Evaluation Metrics:** Many algorithms depend on predefined metrics (e.g., support, confidence) to identify interesting patterns, but these metrics may not fully capture the value of the patterns. Evaluation frameworks for interestingness are often incomplete and context-dependent.

#### Trade-offs Between Simplicity and Complexity:

- **Overfitting vs. Underfitting:** More complex patterns might be interesting in terms of their novelty or complexity, but they could also be overfitting the data. On the other hand, overly simple patterns may fail to capture the true essence of the data.
- **Pattern Comprehensibility:** Complex Patterns are difficult to interpret and might not be considered interesting, even if they are statistically significant. Users often prefer simple and actionable insights, leading to challenges in balancing model complexity with human interpretability.
- **Novelty and Actionability:** One promising approach is to define interestingness in terms of novelty and actionability. A pattern that is both novel (not previously known) and actionable (capable of leading to a decision) is likely to be of greater interest to stakeholders.
- **Complexity of the Pattern:** Another interestingness measure could be based on the complexity or simplicity of the discovered pattern. Simpler patterns that are easier to interpret and explain tend to be more useful in real-world applications, while more complex patterns might be harder to translate into actionable insights.

#### Diversity in Mining Tasks:

- **Task-Specific Interestingness:** Different data mining tasks (e.g., clustering, classification, association

rule mining) may have different criteria of interestingness. For example, association rule mining is concerned with finding frequent patterns, whereas clustering is about discovering underlying structures. The notion of interestingness thus varies based on tasks to be performed.

- **Evaluation across Multiple Tasks:** Defining a uniform standard for interestingness across different data mining tasks is difficult and complex. What is interesting for one task (e.g., outlier detection) may not be interesting for another task and so on (e.g., feature selection).

#### Ethical and Societal Concerns:

- **Bias in Definition:** The definition of interestingness can be defined on the biases of data and algorithms. If the data used to mine patterns is biased, the patterns discovered may reinforce societal biases (e.g., in predictive policing or hiring practices).
- **Impact on Decision Making:** Deciding which patterns are considered interesting and which not, can have significant implications for real-world decision-making. For example, biased or ethically questionable patterns deemed interesting in data mining could lead to harmful outcomes when used in applications like healthcare or criminal justice.

#### Domain-Specific Metrics:

The value of a pattern often depends on the context in which it is discovered.

- In **healthcare**, a pattern that links certain patient characteristics to disease outbreaks might be highly valuable, while other patterns may not have the same impact.
- **Risk Factor Association:** It identifies associations between patient behaviors, conditions, or treatments and their risk of certain outcomes (e.g., disease, complications). For example in medical

data mining, interestingness may be defined by how strongly a certain behavior (e.g., smoking) is associated with a specific condition (e.g., lung cancer).

- **Predictive Accuracy (Outcome Prediction):** Measures how well a model predicts clinical outcomes, such as patient recovery, risk of complications, or survival rates. For example a model that predicts patient risk within 30 days would use predictive accuracy as an interestingness metric.
- **Medical Knowledge Enhancement:** Measures how well a discovered pattern can improve medical understanding or uncover previously unknown links. For example discovering an unknown relationship between two drugs' interactions could be considered highly interesting in medical research.
- In e-commerce, patterns predicting future customer purchases or identifying cross-selling opportunities can drive significant business value.
- **Lift (in association rule mining):** Measures the strength of an association rule relative to its frequency in the data. A higher lift means the items in the rule occur more often together than expected by chance. For example, In retail, if customers buying "bread" often buy "butter" together, this lift metric helps find combinations of products that increase sales.
  - **Sales Increase (Revenue Impact):** Measures the potential increase in revenue driven by discovered patterns or rules. For example discovering that "customers who buy running shoes are likely to buy athletic wear" could be used to design promotions that increase sales.



- **Customer Lifetime Value (CLV):** Measures the long-term value a customer brings, which can be linked to patterns in purchasing behavior. For example A pattern of frequent purchases by a specific customer group can inform strategies to retain high-value customers
- **Finance and Banking:**
  - **Fraud Detection Rate:** Measures how effectively a pattern or model identifies fraudulent activities. Example: A rule indicating that "customers who frequently change addresses tend to commit fraud" may be highly interesting if it can flag fraudulent activities.
  - **Risk Prediction (Default Likelihood):** In credit scoring, interesting patterns may be those that effectively predict loan defaults or financial risk. Example: Patterns showing that "individuals with a high ratio of credit utilization to income" are more likely to default could be critical for assessing creditworthiness.
  - **Anomaly Detection (Unusual Transactions):** Identifying transactions or behaviors that are rare but indicative of significant events, such as fraud. Example: A sudden large transfer of funds from a previously low-activity account may be flagged as an interesting anomaly.
- **Telecommunications:**
  - **Churn Prediction:** Measures how well patterns can predict customer churn (i.e., the likelihood of a customer leaving the service). For example identifying "customers who frequently contact customer support are more likely to cancel their subscription" is highly relevant in telecommunications.
  - **Service Usage Patterns:** Identifies patterns of how customers use telecom services (e.g., calling behavior, data usage) and their impact on service adoption. For example discovering that users who use mobile data heavily but not voice calling are more likely to switch to mobile-first carriers could inform service offerings.
- **Social Media and Online Platforms:**
  - **Engagement Metrics (Likes, Shares, Comments):** Measures how well content-related patterns predict user engagement (e.g., virality). For example a pattern where posts with certain keywords (e.g., sustainability") tend to get more shares and likes might be considered interesting in the context of social media marketing.
  - **Sentiment Analysis:** Measures the sentiment of text or social media posts and how it correlates with engagement or behavior. For example Patterns where posts with positive sentiment tend to increase user retention could be deemed interesting for platform optimization.
  - **Influencer Impact:** Identifies patterns related to influencer activity, such as correlations between influencers' posts and user actions. For example: If influencers in certain niches (e.g., tech gadgets) drive higher click-through rates or product purchases, such patterns may be highly valuable.
- **Manufacturing and Supply Chain:**
  - **Predictive Maintenance Metrics:** Measures the ability of patterns to predict equipment failure or the need for maintenance. For example: Discovering a pattern where machines that show specific vibration frequencies tend to fail after a certain number of operations could lead to predictive maintenance strategies.

- **Inventory Turnover:** Measures how efficiently patterns lead to the optimization of inventory, reducing stockouts or overstocking. For example: Patterns indicating that certain products have higher sales during specific months could lead to better forecasting of inventory needs.
- **Education and Learning Analytics:**
  - **Student Performance Prediction:** Measures how well patterns predict student success, dropout rates, or learning outcomes. For example: Identifying that students who engage in online quizzes are more likely to perform well on final exams could be an interesting finding.
  - **Learning Style Optimization:** Measures how patterns related to individual learning preferences improve educational outcomes. For example discovering students who prefer video content perform better in certain subjects might be considered interesting in developing personalized learning plans.
  - **Course Completion Rates:** Measures patterns related to course engagement and completion, helping to predict and improve retention. For example patterns showing that students attending live classes are more likely to complete the course could inform strategies for improving online education programs.
- **Government and Public Sector:**
  - **Crime Hotspot Detection:** Measures how well discovered patterns predict locations or times with higher crime rates. For example identifying that certain neighborhoods have a higher probability of certain crimes based on historical data could help optimize policing efforts.
  - **Public Health Trends:** Measures patterns of disease outbreaks, health interventions, or social behaviors related to public health. For example Identifying patterns in flu outbreaks based on certain weather conditions or travel behaviors could inform public health strategies.
  - **Resource Allocation:** Measures patterns that help with the optimal allocation of resources, such as emergency services or social welfare programs. For example discovering that certain demographics are more likely to request social services could help prioritize resource distribution.
- **Energy and Utilities:**
  - **Energy Consumption Patterns:** Measures how patterns in energy usage can inform efficiency or cost-saving strategies. For example discovering that certain regions or customer types tend to use more energy during specific months could inform demand-response strategies.
  - **Fault Detection in Power Grids:** Measures how effectively patterns can predict equipment failure or power outages in the grid. For example identifying patterns in voltage fluctuation that precede a grid failure could be critical for proactive grid management.

### Conclusion:

This paper has discussed the importance of interestingness measures in data mining, highlighting the challenges in defining relevant and useful patterns. We reviewed existing measures, proposed new hybrid metrics that combine statistical and domain-specific approaches, and explored their applicability in various domains. Future work should focus on refining these measures, incorporating more context and feedback from end-users, and developing frameworks for more robust evaluation.

**References:**

1. S. Chakrabarti, S. Sarawagi, and B. Dom. Mining surprising patterns using temporal description length. In VLDB, pages 606–617, 1998.
2. M. Kamber and R. Shinghal. Evaluating the interestingness of characteristic rules. In KDD, pages 263–266, 1996.
3. A. Silberschatz and A. Tuzhilin. On subjective measures of interestingness in knowledge discovery. In KDD, pages 275–281, 1995.
4. E. Suzuki. Negative encoding length as a subjective interestingness measure for groups of rules. In PAKDD, pages 220–231, 2009.
5. R. Agrawal, T. Imielinski, and A. N. Swami. Mining association rules between sets of items in large databases. In SIGMOD Conference, pages 207–216, 1993.
6. P.-N. Tan, V. Kumar, and J. Srivastava. Selecting the right interestingness measure for association patterns. In KDD, pages 32–41, 2002.
7. T. Wu, Y. Chen, and J. Han. Re-examination of interestingness measures in pattern mining: a unified framework. Data Mining and Knowledge Discovery, 2010.
8. X. Wu, C. Zhang, and S. Zhang. Efficient mining of both positive and negative association rules. ACM Trans. Inf. Syst., 22(3):381–405, 2004
9. Swati R. Ramdasi; Shailaja C. Shirwaikar; Vilas Kharat , Interestingness measures for quantified and ordered categorical attributes using fuzzy approach International Journal of Fuzzy Computation and Modelling (IJFCM), Vol. 2, No. 4, 2019
10. Swati R. Ramdasi; Shailaja C. Shirwaikar AICTC '16: Proceedings of the International Conference on Advances in Information Communication Technology & Computing, Article No.: 88, Pages 1 – 5
11. Hegland, M., Algorithms for Association Rules, Lecture Notes in Computer Science, Volume 2600, Jan 2003, Pages 226—234
12. Mojdeh Jalali-Heravi, Osmar R. Zaïane, A Study on Interestingness Measures for Associative Classifiers. Master's thesis, University of Alberta, 2009.