

A Probabilistic Comparison of Commonly Used Interest Measures for Association Rules

Michael Hahsler

Contents

About this Document	2
Introduction	2
Measures for Itemsets	3
Support	3
Support Count	4
All-Confidence	4
Cross-Support Ratio	4
Measures for Rules	4
Contingency Table	4
Confidence	5
Added Value	5
Casual Confidence	5
Casual Support	6
Centered Confidence	6
Certainty Factor	6
Chi-Squared	6
Collective Strength	7
Confidence Boost	7
Conviction	8
Cosine	8
Coverage	8
Descriptive Confirmed Confidence	8
Difference of Confidence	9
Example and Counter-Example Rate	9
Fisher's Exact Test	9
Gini Index	9
Hyper-Confidence	10
Hyper-Lift	10
Imbalance Ratio	10
Implication Index	11
Importance	11
Improvement	11
Jaccard Coefficient	11
J-Measure	12
Kappa	12
Kloggen	12
Kulczynski	12

Lambda	13
Laplace Corrected Confidence	13
Least Contradiction	13
Lerman Similarity	13
Leverage	14
Lift	14
MaxConfidence	14
Mutual Information	15
Odds Ratio	15
Phi Correlation Coefficient	15
Ralambondrainy	16
Relative Linkage Disequilibrium	16
Relative Risk	16
Rule Power Factor	16
Right-Hand-Side Support	17
Sebag-Schoenauer	17
Standardized Lift	17
Varying Rates Liaison	18
Yule's Q	18
Yule's Y	18

About this Document

This work is licensed under the Creative Commons Attribution Share Alike 4.0 International License.

Please cite this document as **Michael Hahsler, A Probabilistic Comparison of Commonly Used Interest Measures for Association Rules, 2015**, URL: <https://mhahsler.github.io/arules/docs/measures>

A PDF version of the document can be found here: <https://mhahsler.github.io/arules/docs/measures.pdf>

Introduction

This document contains a collection of commonly used measures of significance and interestingness (sometimes also called strength) for association rules and itemsets. Agrawal, Imielinski, and Swami (1993) introduced the problem of association rule mining in the following way:

Let $I = \{i_1, i_2, \dots, i_m\}$ be a set of m binary attributes called **items**. Let $D = \{t_1, t_2, \dots, t_n\}$ be a set of transactions called the **database**. Each transaction $t \in D$ has a unique transaction ID and contains a subset of the items in I , i. e., $t \subseteq I$. A **rule** is defined as an implication of the form $X \Rightarrow Y$ where $X, Y \subseteq I$ and $X \cap Y = \emptyset$. The sets of items (for short **itemsets**) X and Y are called antecedent (left-hand side or LHS) and consequent (right-hand side or RHS) of the rule, respectively. Measures of importance (interest) can be defined for itemsets and rules. The support-confidence framework defines the measures support and confidence. Rules that satisfy a user-specified minimum thresholds on support and confidence are called **association rules**.

Interest measures are usually defined in terms of itemset support, here we also present them using probabilities and, where appropriate, counts. The probability $P(E_X)$ of the event that all items in itemset X are contained in an arbitrarily chosen transaction can be estimated from a database D using maximum likelihood estimation (MLE) by

$$\hat{P}(E_X) = \frac{|\{t \in D; X \subseteq t\}|}{n}$$

where $n_X = |\{t \in D; X \subseteq t\}|$ is the count of the number of transactions that contain the itemset X and $n = |D|$ is the size (number of transactions) of the database. For conciseness of notation, we will drop the hat and the E from the notation for probabilities. We will use in the following $P(X)$ to mean $\hat{P}(E_X)$ and $P(X \cap Y)$ to mean $\hat{P}(E_X \cap E_Y) = \hat{P}(E_{X \cup Y})$, the probability of the intersection of the events E_X and E_Y representing the probability of the event that a transaction contains all items in the union of the itemsets X and Y . The event notation should not be confused with the set notation used in measures like support, where $\text{supp}(X \cup Y)$ means the support of the union of the itemsets X and Y .

Note on probability estimation: The used probability estimates will be very poor for itemsets with low observed frequencies. This needs to be always taken into account since it affects most measured discussed below.

Note on null-transactions: Transaction datasets typically contain a large number of transactions that do not contain either X or Y . These transactions are called null-transactions, and it is desirable that measures of rule strength are not influenced by a change in the number of null-transactions. However, most measures are affected by the number of null-transactions since the total number of transactions is used for probability estimation. Measures that are not influenced by a change in the number of null-transactions are called null-invariant (see Tan et al., 2004 and Wu et al., 2010).

Good overview articles about different association rule measures are provided by Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293-313, 2004 and Liqiang Geng and Howard J. Hamilton. Interestingness measures for data mining: A survey. ACM Computing Surveys, 38(3):9, 2006. Michael Hahsler maintains a comprehensive Annotated Bibliography on Association Rule Mining.

Implementation

All measures discussed on this page are implemented in the freely available R-extension package `arules` in function `interestMeasure()`.

Corrections and Feedback

For corrections and missing measures on this page or in `arules`, please open an issue here or contact me directly.

Measures for Itemsets

Support

Reference: R. Agrawal, T. Imielinski, and A. Swami. Mining associations between sets of items in large databases. In Proc. of the ACM SIGMOD Int'l Conference on Management of Data, pages 207-216, Washington D.C., May 1993.

$$\text{supp}(X) = \frac{n_X}{n} = P(X)$$

Support is defined on itemsets and gives the proportion of transactions that contain X . It is used as a measure of significance (importance) of an itemset. Since it uses the count of transactions, it is often called a **frequency constraint**. An itemset with support greater than a set minimum support threshold, $\text{supp}(X) > \sigma$, is called a **frequent or large itemset**.

For rules the support defined as the support of all items in the rule, i.e., $\text{supp}(X \Rightarrow Y) = \text{supp}(X \cup Y) = P(X \cap Y)$.

Support's main feature is that it possesses the **downward closure property (anti-monotonicity)**, which means that all subsets of a frequent set are also frequent. This property (actually, the fact that no superset

of an infrequent set can be frequent) is used to prune the search space (usually thought of as a lattice or tree of itemsets with increasing size) in level-wise algorithms (e.g., the Apriori algorithm).

The disadvantage of support is the **rare item problem**. Items that occur very infrequently in the data set are pruned, although they would still produce interesting and potentially valuable rules. The rare item problem is important for transaction data which usually have a very uneven distribution of support for the individual items (typical is a power-law distribution where few items are used all the time and most items are rarely used).

Range: $[0, 1]$

Support Count

Alias: Absolute Support Count

Range: $[0, n]$ where n is the number of transactions.

All-Confidence

Reference: Edward R. Omiecinski. Alternative interest measures for mining associations in databases. IEEE Transactions on Knowledge and Data Engineering, 15(1):57-69, Jan/Feb 2003.

All-confidence is defined on itemsets (not rules) as

$$\text{all-confidence}(X) = \frac{\text{supp}(X)}{\max_{x \in X}(\text{supp}(x))} = \frac{P(X)}{\max_{x \in X}(P(x))} = \min\{P(X|Y), P(Y|X)\}$$

where $\max_{x \in X}(\text{supp}(x \in X))$ is the support of the item with the highest support in X . All-confidence means that all rules which can be generated from itemset X have at least a confidence of $\text{all-confidence}(X)$. All-confidence possesses the downward-closed closure property and thus can be effectively used inside mining algorithms. All-confidence is null-invariant.

Range: $[0, 1]$

Cross-Support Ratio

Reference: Hui Xiong, Pang-Ning Tan, and Vipin Kumar. Mining strong affinity association patterns in data sets with skewed support distribution. In Bart Goethals and Mohammed J. Zaki, editors, Proceedings of the IEEE International Conference on Data Mining, November 19-22, 2003, Melbourne, Florida, pages 387-394, November 2003.

Defined on itemsets as the ratio of the support of the least frequent item to the support of the most frequent item, i.e.,

$$\text{cross-support}(X) = \frac{\min_{x \in X}(\text{supp}(x))}{\max_{x \in X}(\text{supp}(x))}$$

a ratio smaller than a set threshold. Normally many found patterns are cross-support patterns which contain frequent as well as rare items. Such patterns often tend to be spurious.

Range: $[0, 1]$

Measures for Rules

Contingency Table

A 2×2 contingency table with counts for rule $X \Rightarrow Y$ in the transaction dataset. The counts are:

	Y	\bar{Y}
X	n_{XY}	$n_{X\bar{Y}}$
\bar{X}	$n_{\bar{X}Y}$	$n_{\bar{X}\bar{Y}}$

n_{XY} is the number of transactions that contain all items in X and Y . All other measures for rules can be calculated using these counts.

Confidence

Alias: Strength

Reference: R. Agrawal, T. Imielinski, and A. Swami. Mining associations between sets of items in large databases. In Proc. of the ACM SIGMOD Int'l Conference on Management of Data, pages 207-216, Washington D.C., May 1993.

$$conf(X \Rightarrow Y) = \frac{supp(X \Rightarrow Y)}{supp(X)} = \frac{supp(X \cup Y)}{supp(X)} = \frac{n_{XY}}{n_X} = \frac{P(X \cap Y)}{P(X)} = P(Y|X)$$

Confidence is defined as the proportion of transactions that contain Y in the set of transactions that contain X . This proportion is an estimate for the probability of seeing the rule's consequent under the condition that the transactions also contain the antecedent.

Confidence is directed and gives different values for the rules $X \Rightarrow Y$ and $Y \Rightarrow X$. Association rules have to satisfy a minimum confidence constraint, $conf(X \Rightarrow Y) \geq \gamma$.

Confidence is not downward closed and was developed together with support by Agrawal et al. (the so-called support-confidence framework). Support is first used to find frequent (significant) itemsets exploiting its downward closure property to prune the search space. Then confidence is used in a second step to produce rules from the frequent itemsets that exceed a min. confidence threshold.

A problem with confidence is that it is sensitive to the frequency of the consequent Y in the database. Caused by the way confidence is calculated, consequents with higher support will automatically produce higher confidence values even if there exists no association between the items.

Range: $[0, 1]$

Added Value

Alias: AV, Pavillon Index, Centered Confidence

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Quantifies how much the probability of Y increases when conditioning on the transactions that contain X
Defined as

$$AV(X \Rightarrow Y) = conf(X \Rightarrow Y) - supp(Y) = P(Y|X) - P(Y)$$

Range: $[-.5, 1]$

Casual Confidence

Reference: Yves Kodratoff. Comparing machine learning and knowledge discovery in databases: An application to knowledge discovery in texts. In Georgios Paliouras, Vangelis Karkaletsis, and Constantine D. Spyropoulos, editors, Machine Learning and Its Applications, Advanced Lectures, volume 2049 of Lecture Notes in Computer Science, pages 1-21. Springer, 2001.

Confidence reinforced by negatives given by

$$\text{casual-conf} = \frac{1}{2}[\text{conf}(X \Rightarrow Y) + \text{conf}(\overline{X} \Rightarrow \overline{Y})] = \frac{1}{2}[P(Y|X) + P(\overline{Y}|\overline{X})]$$

Range: $[0, 1]$

Casual Support

Reference: Yves Kodratoff. Comparing machine learning and knowledge discovery in databases: An application to knowledge discovery in texts. In Georgios Paliouras, Vangelis Karkaletsis, and Constantine D. Spyropoulos, editors, Machine Learning and Its Applications, Advanced Lectures, volume 2049 of Lecture Notes in Computer Science, pages 1-21. Springer, 2001.

Support improved by negatives given by

$$\text{casual-supp} = \text{supp}(X \cup Y) + \text{supp}(\overline{X} \cup \overline{Y}) = P(X \cap Y) + P(\overline{X} \cap \overline{Y})$$

Range: $[0, 2]$

Centered Confidence

Alias: relative accuracy, gain

Reference: N Lavrac, PA Flach, and B Zupan, Rule Evaluation Measures: A Unifying View, Springer-Verlag, in Proceedings of the 9th International Workshop on Inductive Logic Programming, ILP '99 (London, UK, 1999), pp. 174-185.

$$CC(X \Rightarrow Y) = \text{conf}(X \Rightarrow Y) - \text{supp}(Y)$$

Range: $[-1, 1 - 1/n]$

Certainty Factor

Alias: CF, Loevinger

Reference: F. Galiano, I. J. Blanco, D. Sánchez, and M. Vila. Measuring the accuracy and interest of association rules: A new framework. Intelligent Data Analysis, 6:221-235, 2002.

The certainty factor is a measure of the variation of the probability that Y is in a transaction when only considering transactions with X . An increasing CF means a decrease in the probability that Y is not in a transaction that X is in. Negative CFs have a similar interpretation.

$$CF(X \Rightarrow Y) = \frac{\text{conf}(X \Rightarrow Y) - \text{supp}(Y)}{\text{supp}(\overline{Y})} = \frac{P(Y|X) - P(Y)}{1 - P(Y)}$$

Range: $[-1, 1]$ (0 indicates independence)

Chi-Squared

Reference: Sergey Brin, Rajeev Motwani, and Craig Silverstein. Beyond market baskets: Generalizing association rules to correlations. In SIGMOD 1997, Proceedings ACM SIGMOD International Conference on Management of Data, pages 265-276, Tucson, Arizona, USA, May 1997.

For the analysis of 2×2 contingency tables, the chi-squared test statistic is a measure of the relationship between two binary variables (X and Y). The chi-squared test statistic is used as a test for independence between X and Y . The chi-squared test statistic is:

$$\text{chi-squared}(X \Rightarrow Y) = \sum_i \frac{(O_i - E_i)^2}{E_i} = \frac{(n_{XY} - \frac{n_X n_Y}{n})^2}{\frac{n_X n_Y}{n}} + \frac{(n_{\bar{X}Y} - \frac{n_{\bar{X}} n_Y}{n})^2}{\frac{n_{\bar{X}} n_Y}{n}} + \frac{(n_{X\bar{Y}} - \frac{n_X n_{\bar{Y}}}{n})^2}{\frac{n_X n_{\bar{Y}}}{n}} + \frac{(n_{\bar{X}\bar{Y}} - \frac{n_{\bar{X}} n_{\bar{Y}}}{n})^2}{\frac{n_{\bar{X}} n_{\bar{Y}}}{n}} = n \frac{P(X \neq Y)^2}{P(X)P(Y)}$$

O_i is the observed count of contingency table cell i and E_i is the expected count given the marginals. The statistic has approximately a χ^2 distribution with 1 degree of freedom (for a 2×2 contingency table). The critical value for $\alpha = 0.05$ is 3.84; higher chi-squared values indicate that the null-hypothesis of independence between LHS and the RHS should be rejected (i.e., the rule is not spurious). Larger chi-squared values indicate stronger evidence that the rule represents a strong relationship. The statistic can be converted into a p-value using the χ^2 distribution.

Notes: The contingency tables for some rules may contain cells with low expected values (less than 5) and thus Fisher's exact test might be more appropriate. Each rule represents a statistical test, and correction for multiple comparisons may be necessary.

Range: $[0, \infty]$

Collective Strength

Reference: C. C. Aggarwal and P. S. Yu. A new framework for itemset generation. In PODS 98, Symposium on Principles of Database Systems, pages 18-24, Seattle, WA, USA, 1998.

$$S(X) = \frac{1 - v(X)}{1 - E[v(X)]} \frac{E[v(X)]}{v(X)} = \frac{P(X \cap Y) + P(\bar{Y}|\bar{X})}{P(X)P(Y) + P(\bar{X})P(\bar{Y})}$$

where $v(X)$ is the violation rate and $E[v(X)]$ is the expected violation rate for independent items. The violation rate is defined as the fraction of transactions that contain some of the items in an itemset but not all. Collective strength gives 0 for perfectly negative correlated items, infinity for perfectly positive correlated items, and 1 if the items co-occur as expected under independence.

Problematic is that for items with medium to low probabilities, the observations of the expected values of the violation rate is dominated by the proportion of transactions that do not contain any of the items in X . For such itemsets, collective strength produces values close to one, even if the itemset appears several times more often than expected together.

Range: $[0, \infty]$

Confidence Boost

Reference: José L. Balcázar. Formal and computational properties of the confidence boost of association rules. ACM Trans. Knowl. Discov. Data, 7(4), December 2013.

Confidence boost is the ratio of the confidence of a rule to the confidence of any more general rule (i.e., a rule with the same consequent but one or more items removed in the LHS).

$$\text{confidence-boost}(X \Rightarrow Y) = \frac{\text{conf}(X \Rightarrow Y)}{\max_{X' \subset X} (\text{conf}(X' \Rightarrow Y))} = \frac{\text{conf}(X \Rightarrow Y)}{\text{conf}(X \Rightarrow Y) - \text{improvement}(X \Rightarrow Y)}$$

Values larger than 1 mean the new rule boosts the confidence compared to the best, more general rule. The measure is related to the improvement measure.

Range: $[0, \infty]$ (> 1 indicates a rule with confidence boost)

Conviction

Reference: Sergey Brin, Rajeev Motwani, Jeffrey D. Ullman, and Shalom Turk. Dynamic itemset counting and implication rules for market basket data. In SIGMOD 1997, Proceedings ACM SIGMOD International Conference on Management of Data, pages 255-264, Tucson, Arizona, USA, May 1997.

$$\text{conviction}(X \Rightarrow Y) = \frac{1 - \text{supp}(Y)}{1 - \text{conf}(X \Rightarrow Y)} = \frac{P(X)P(\bar{Y})}{P(X \cap \bar{Y})}$$

where $\bar{Y} = E_{\neg Y}$ is the event that Y does not appear in a transaction. Conviction was developed as an alternative to confidence which was found to not capture the direction of associations adequately. Conviction compares the probability that X appears without Y if they were dependent on the actual frequency of the appearance of X without Y . In that respect, it is similar to lift (see the section about lift on this page). However, in contrast to lift, it is a directed measure since it also uses the information of the absence of the consequent. An interesting fact is that conviction is monotone in confidence and lift.

Range: $[0, \infty]$ (1 indicates independence; rules that always hold have ∞)

Cosine

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Cosine is a null-invariant measure of correlation between the items in X and Y defined as

$$\text{cosine}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\sqrt{(\text{supp}(X)\text{supp}(Y))}} = \frac{P(X \cap Y)}{\sqrt{P(X)P(Y)}} = \sqrt{P(X|Y)P(Y|X)}$$

Range: $[0, 1]$ (0.5 means no correlation)

Coverage

Alias: LHS Support

It measures the probability that a rule $X \Rightarrow Y$ applies to a randomly selected transaction. It is estimated by the proportion of transactions that contain the antecedent of the rule $X \Rightarrow Y$. Therefore, coverage is sometimes called antecedent support or LHS support.

$$\text{cover}(X \Rightarrow Y) = \text{supp}(X) = P(X)$$

Range: $[0, 1]$

Descriptive Confirmed Confidence

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Confidence confirmed by the confidence of the negative rule.

$$\text{confirmed-conf} = \text{conf}(X \Rightarrow Y) - \text{conf}(X \Rightarrow \bar{Y}) = P(Y|X) - P(\bar{Y}|X)$$

Range: $[-1, 1]$

Difference of Confidence

Alias: DOC, Difference of Proportions

Reference: Heike Hofmann and Adalbert F. X. Wilhelm. Visual comparison of association rules. *Comput. Stat.*, 16(3):399-415, 2001.

The difference of confidence is the difference of the proportion of transactions containing Y in the two groups of transactions that do and do not contain X . For the analysis of 2×2 contingency tables, this measure of the relationship between two binary variables is typically called the difference of proportion. It is defined as

$$\text{doc}(X \Rightarrow Y) = \text{conf}(X \Rightarrow Y) - \text{conf}(\bar{X} \Rightarrow Y) = P(Y|X) - P(Y|\bar{X}) = n_{XY}/n_X - n_{\bar{X}Y}/n_{\bar{X}}$$

Range: $[-1, 1]$ (0 means statistical independence)

Example and Counter-Example Rate

Example rate reduced by the counter-example rate.

Defined as

$$\text{ecr}(X \Rightarrow Y) = \frac{n_{XY} - n_{X\bar{Y}}}{n_{XY}} = \frac{P(X \cap Y) - P(X \cap \bar{Y})}{P(X \cap Y)} = 1 - \frac{1}{\text{sebag}(X \Rightarrow Y)}$$

The measure is related to the Sebag-Schoenauer Measure.

Range: $[0, 1]$

Fisher's Exact Test

Reference: Michael Hahsler and Kurt Hornik. New probabilistic interest measures for association rules. *Intelligent Data Analysis*, 11(5):437-455, 2007.

If X and Y are independent, then the n_{XY} is a realization of the random variable C_{XY} which has a hypergeometric distribution with n_Y draws from a population with n_X successes and $n_{\bar{X}}$ failures. The p-value for Fisher's one-sided exact test giving the probability of observing a contingency table with a count of at least n_{XY} given the observed marginal counts is

$$\text{p-value} = P(C_{XY} \geq n_{XY})$$

The p-value is related to hyper-confidence. Compared to the Chi-squared test, Fisher's exact test also applies when cells have low expected counts. Note that each rule represents a statistical test, and correction for multiple comparisons may be necessary.

Range: $[0, 1]$ (p-value scale)

Gini Index

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. *Information Systems*, 29(4):293-313, 2004.

Measures quadratic entropy as

$$\text{gini}(X \Rightarrow Y) = P(X)[P(Y|X)^2 + P(\bar{Y}|X)^2] + P(\bar{X})[P(Y|\bar{X})^2 + P(\bar{Y}|\bar{X})^2] - P(Y)^2 - P(\bar{Y})^2$$

Range: $[0, 1]$ (0 means that the rule does not provide any information for the dataset)

Hyper-Confidence

Reference: Michael Hahsler and Kurt Hornik. New probabilistic interest measures for association rules. Intelligent Data Analysis, 11(5):437-455, 2007.

The confidence level for observation of too high/low counts for rules $X \Rightarrow Y$ using the hypergeometric model. Since the counts are drawn from a hypergeometric distribution (represented by the random variable C_{XY} with known parameters given by the counts n_X and n_Y , we can calculate a confidence interval for the observed counts n_{XY} stemming from the distribution. Hyper-confidence reports the confidence level as

$$\text{hyper-conf}(X \Rightarrow Y) = 1 - P[C_{XY} \geq n_{XY} | n_X, n_Y]$$

A confidence level of, e.g., > 0.95 indicates that there is only a 5% chance that the high count for the rule has occurred randomly. Hyper-confidence is equivalent to the statistic used to calculate the p-value in Fisher's exact test. Note that each rule represents a statistical test and correction for multiple comparisons may be necessary.

Hyper-Confidence can also be used to evaluate that X and Y are complementary (i.e., the count is too low to have occurred randomly).

$$\text{hyper-conf}_{\text{complement}}(X \Rightarrow Y) = 1 - P[C_{XY} < n_{XY} | n_X, n_Y]$$

Range: $[0, 1]$

Hyper-Lift

Reference: Michael Hahsler and Kurt Hornik. New probabilistic interest measures for association rules. Intelligent Data Analysis, 11(5):437-455, 2007.

Adaptation of the lift measure where instead of dividing by the expected count under independence ($E[C_{XY}] = n_X/n \times n_Y/n$) a higher quantile of the hypergeometric count distribution is used. This is more robust for low counts and results in fewer false positives when hyper-lift is used for rule filtering. Hyper-lift is defined as:

$$\text{hyper-lift}_{\delta}(X \Rightarrow Y) = \frac{n_{XY}}{Q_{\delta}[C_{XY}]}$$

where n_{XY} is the number of transactions containing X and Y and $Q_{\delta}[C_{XY}]$ is the δ -quantile of the hypergeometric distribution with parameters n_X and n_Y .

δ is typically chosen to use the 99 or 95% quantile.

Range: $[0, \infty]$ (1 indicates independence)

Imbalance Ratio

Alias: IR

Reference: Tianyi Wu, Yuguo Chen, and Jiawei Han. Re-examination of interestingness measures in pattern mining: a unified framework. Data Mining and Knowledge Discovery, January 2010.

Measures the degree of imbalance between two events that the LHS and the RHS are contained in a transaction. The ratio is close to 0 if the conditional probabilities are similar (i.e., very balanced) and close to 1 if they are very different. It is defined as

$$IB(X \Rightarrow Y) = \frac{|P(X|Y) - P(Y|X)|}{P(X|Y) + P(Y|X) - P(X|Y)P(Y|X)} = \frac{|supp(X) - supp(Y)|}{supp(X) + supp(Y) - supp(X \cup Y)}$$

Range: $[0, 1]$ (0 indicates a balanced, typically uninteresting rule)

Implication Index

Reference: Gras R., L'implication statistique. Nouvelle methode exploratoire de donnees. La Pensee Sauvage, Grenoble. 1996.

A variation of the Lerman similarity defined as

$$gras(X \Rightarrow Y) = \sqrt{N} \frac{supp(X \cup \bar{Y}) - supp(X)supp(\bar{Y})}{\sqrt{supp(X)supp(\bar{Y})}}$$

Range: $[0, 1]$

Importance

Reference: MS Analysis Services: Microsoft Association Algorithm Technical Reference.

In the Microsoft Association Algorithm Technical Reference, confidence is called “probability,” and a measure called importance is defined as the log-likelihood of the right-hand side of the rule, given the left-hand side of the rule:

$$importance(X \Rightarrow Y) = \log_{10}(L(X \Rightarrow Y)/L(X \Rightarrow \bar{Y}))$$

where L is the Laplace corrected confidence.

Range: $[-\infty, \infty]$

Improvement

Reference: R. Bayardo, R. Agrawal, and D. Gunopulos. Constraint-based rule mining in large, dense databases. Data Mining and Knowledge Discovery, 4(2/3):217-240, 2000.

The improvement of a rule is the minimum difference between its confidence and the confidence of any proper sub-rule with the same consequent. The idea is that we only want to extend the LHS of the rule if this improves the rule sufficiently.

$$improvement(X \Rightarrow Y) = \min_{X' \subset X} (conf(X \Rightarrow Y) - conf(X' \Rightarrow Y))$$

Range: $[0, 1]$

Jaccard Coefficient

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

A null-invariant measure for dependence using the Jaccard similarity between the two sets of transactions that contain the items in X and Y , respectively. Defined as

$$jaccard(X \Rightarrow Y) = \frac{supp(X \cup Y)}{supp(X) + supp(Y) - supp(X \cup Y)} = \frac{P(X \cap Y)}{P(X) + P(Y) - P(X \cap Y)}$$

Range: $[0, 1]$

J-Measure

Padhraic Smyth and R. Goodman. Rule induction using information theory. In Knowledge Discovery in Databases, 1991.

The J-measure is a scaled version of cross entropy to measure the information content of a rule.

$$J(X \Rightarrow Y) = P(X \cap Y) \log \left(\frac{P(Y|X)}{P(Y)} \right) + P(X \cap \bar{Y}) \log \left(\frac{P(\bar{Y}|X)}{P(\bar{Y})} \right)$$

Range: $[0, 1]$ (0 means that X does not provide information for Y)

Kappa

Alias: Cohen's κ

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Cohen's kappa of the rule (seen as a classifier) given as the rules observed rule accuracy (i.e., confidence) corrected by the expected accuracy (of a random classifier). Kappa is defined as

$$\kappa(X \Rightarrow Y) = \frac{P(X \cap Y) + P(\bar{X} \cap \bar{Y}) - P(X)P(Y) - P(\bar{X})P(\bar{Y})}{1 - P(X)P(Y) - P(\bar{X})P(\bar{Y})}$$

Range: $[-1, 1]$ (0 means the rule is not better than a random classifier)

Klosgen

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Defined as a scaled version of the added value measure.

$$klosgen(X \Rightarrow Y) = \sqrt{supp(X \cup Y)} (conf(X \Rightarrow Y) - supp(Y)) = \sqrt{P(X \cap Y)} (P(Y|X) - P(Y)) = \sqrt{P(X \cap Y)} AV(X \Rightarrow Y)$$

Range: $[-1, 1]$ (0 for independence)

Kulczynski

Reference: Tianyi Wu, Yuguo Chen, and Jiawei Han. Re-examination of interestingness measures in pattern mining: a unified framework. Data Mining and Knowledge Discovery, January 2010. Based on Kulczynski, S. 1927. Die Pflanzenassoziationen der Pieninen. Bulletin International de l'Academie Polonaise des Sciences et des Lettres, Classe des Sciences Mathematiques et Naturelles B, 57-203.

Calculate the null-invariant Kulczynski measure with a preference for skewed patterns.

$$kulc(X \Rightarrow Y) = \frac{1}{2} (conf(X \Rightarrow Y) + conf(Y \Rightarrow X)) = \frac{1}{2} \left(\frac{supp(X \cup Y)}{supp(X)} + \frac{supp(X \cup Y)}{supp(Y)} \right) = \frac{1}{2} (P(X|Y) + P(Y|X))$$

Range: $[0, 1]$ (0.5 means neutral and typically uninteresting)

Lambda

Alias: Goodman-Kruskal's λ , Predictive Association

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Goodman and Kruskal's lambda assesses the association between the LHS and RHS of the rule.

$$\lambda(X \Rightarrow Y) = \frac{\sum_{x \in X} \max_{y \in Y} P(x \cap y) - \max_{y \in Y} P(y)}{n - \max_{y \in Y} P(y)}$$

Range: $[0, 1]$

Laplace Corrected Confidence

Alias: Laplace Accuracy, L

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

$$L(X \Rightarrow Y) = \frac{n_{XY} + 1}{n_X + k},$$

where k is the number of classes in the domain. For association rule k is often set to 2. It is an approximate measure of the expected rule accuracy representing 1 - the Laplace expected error estimate of the rule. The Laplace corrected accuracy estimate decreases with lower support to account for estimation uncertainty with low counts.

Range: $[0, 1]$

Least Contradiction

Reference: Aze, J. and Y. Kodratoff (2004). Extraction de pepites de connaissances dans les donnees: Une nouvelle approche et une etude de sensibilite au bruit. In Mesures de Qualite pour la fouille de donnees. Revue des Nouvelles Technologies de l'Information, RNTI.

$$\text{least-contradiction}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y) - \text{supp}(X \cup \bar{Y})}{\text{supp}(Y)} = \frac{P(X \cap Y) - P(X \cap \bar{Y})}{P(Y)}$$

Range: $[-1, 1]$

Lerman Similarity

Reference: Lerman, I.C. (1981). Classification et analyse ordinaire des donnees. Paris.

Defined as

$$\text{lerman}(X \Rightarrow Y) = \frac{n_{XY} - \frac{n_X n_Y}{n}}{\sqrt{\frac{n_X n_Y}{n}}} = \sqrt{n} \frac{\text{supp}(X \cup Y) - \text{supp}(X) \text{supp}(Y)}{\sqrt{\text{supp}(X) \text{supp}(Y)}}$$

Range: $[0, 1]$

Leverage

Alias: Piatetsky-Shapiro, PS

Reference: Piatetsky-Shapiro, G., Discovery, analysis, and presentation of strong rules. Knowledge Discovery in Databases, 1991: p. 229-248.

$$PS(X \Rightarrow Y) = leverage(X \Rightarrow Y) = supp(X \Rightarrow Y) - supp(X)supp(Y) = P(X \cap Y) - P(X)P(Y)$$

Leverage measures the difference of X and Y appearing together in the data set and what would be expected if X and Y were statistically dependent. The rationale in a sales setting is to find out how many more units (items X and Y together) are sold than expected from the independent sells.

Using minimum leverage thresholds incorporates at the same time an implicit frequency constraint. E.g., for setting a min. leverage thresholds to 0.01% (corresponds to 10 occurrences in a data set with 100,000 transactions) one first can use an algorithm to find all itemsets with min. support of 0.01% and then filter the found item sets using the leverage constraint. Because of this property, leverage also can suffer from the rare item problem.

Leverage is a unnormalized version of the phi correlation coefficient.

Range: $[-1, 1]$ (0 indicates independence)

Lift

Alias: Interest, interest factor

Reference: S. Brin, R. Motwani, J. D. Ullman, and S. Tsur. Dynamic itemset counting and implication rules for market basket data. In Proc. of the ACM SIGMOD Int'l Conf. on Management of Data (ACM SIGMOD '97), pages 265-276, 1997.

Lift was originally called interest by Brin et al. Later, lift, the name of an equivalent measure popular in advertising and predictive modeling became more common. Lift is defined as

$$\text{lift}(X \Rightarrow Y) = \text{lift}(Y \Rightarrow X) = \frac{\text{conf}(X \Rightarrow Y)}{\text{supp}(Y)} = \frac{P(Y|X)}{P(Y)} = \frac{P(X \cap Y)}{P(X)P(Y)} = n \frac{n_{XY}}{n_X n_Y}$$

Lift measures how many times more often X and Y occur together than expected if they were statistically independent. A lift value of 1 indicates independence between X and Y . For statistical tests, see the Chi-squared test statistic, Fisher's exact test, and hyper-confidence.

Lift is not downward closed and does not suffer from the rare item problem. However, lift is susceptible to noise in small databases. Rare itemsets with low counts (low probability), which by chance occur a few times (or only once) together, can produce enormous lift values.

Range: $[0, \infty]$ (1 means independence)

MaxConfidence

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Symmetric, null-invariant version of confidence defined as

$$\text{maxConf}(X \Rightarrow Y) = \max\{\text{conf}(X \Rightarrow Y), \text{conf}(Y \Rightarrow X)\} = \max\{P(Y|X), P(X|Y)\}$$

Range: $[0, 1]$

Mutual Information

Alias: Uncertainty

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Measures the information gain for Y provided by X.

$$M(X \Rightarrow Y) = \frac{\sum_{i \in \{X, \bar{X}\}} \sum_{j \in \{Y, \bar{Y}\}} \frac{n_{ij}}{n} \log \frac{n_{ij}}{n_i n_j}}{\min(-\sum_{i \in \{X, \bar{X}\}} \frac{n_i}{n} \log \frac{n_i}{n}, -\sum_{j \in \{Y, \bar{Y}\}} \frac{n_j}{n} \log \frac{n_j}{n})} = \frac{\sum_{i \in \{X, \bar{X}\}} \sum_{j \in \{Y, \bar{Y}\}} P(i \cap j) \log \frac{P(i \cap j)}{P(i)P(j)}}{\min(-\sum_{i \in \{X, \bar{X}\}} P(i) \log P(i), -\sum_{j \in \{Y, \bar{Y}\}} P(j) \log P(j))}$$

Range: $[0, 1]$ (0 means that X does not provide information for Y)

Odds Ratio

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

For the analysis of 2×2 contingency tables, the odds ratio is a measure of the relationship between two binary variables. It is defined as the ratio of the odds of a transaction containing Y in the groups of transactions that do and do not contain X.

$$OR(X \Rightarrow Y) = \frac{\frac{P(Y|X)}{1-P(Y|X)}}{\frac{P(Y|\bar{X})}{1-P(Y|\bar{X})}} = \frac{\frac{conf(X \Rightarrow Y)}{1-conf(X \Rightarrow Y)}}{\frac{conf(\bar{X} \Rightarrow Y)}{1-conf(\bar{X} \Rightarrow Y)}} = \frac{n_{XY}n_{\bar{X}\bar{Y}}}{n_{X\bar{Y}}n_{\bar{X}Y}}$$

A confidence interval around the odds ratio can be calculated (see Li et al., 2014) using a normal approximation.

$$\omega = z_{\alpha/2} \sqrt{\frac{1}{n_{XY}} + \frac{1}{n_{X\bar{Y}}} + \frac{1}{n_{\bar{X}Y}} + \frac{1}{n_{\bar{X}\bar{Y}}}}$$

$$CI(X \Rightarrow Y) = [OR(X \Rightarrow Y) \exp(-\omega), OR(X \Rightarrow Y) \exp(\omega)]$$

where $\alpha/2$ is the critical value for a confidence level of $1 - \alpha$.

Range: $[0, \infty]$ (1 indicates that Y is not associated with X)

Phi Correlation Coefficient

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

The correlation coefficient between the transactions containing X and Y represented as two binary vectors. Phi correlation is equivalent to Pearson's Product Moment Correlation Coefficient ρ with 0-1 values and related to the chi-squared test statistics for 2×2 contingency tables.

$$\phi(X \Rightarrow Y) = \frac{nn_{XY} - n_X n_Y}{\sqrt{n_X n_Y n_{\bar{X}} n_{\bar{Y}}}} = \frac{P(X \cap Y) - P(X)P(Y)}{\sqrt{P(X)(1-P(X))P(Y)(1-P(Y))}} = \sqrt{\frac{\chi^2}{n}}$$

Range: $[-1, 1]$ (0 when X and Y are independent)}

Ralambondrainy

Reference: Jean Diatta, Henri Ralambondrainy, and André Totohasina. Towards a Unifying Probabilistic Implicative Normalized Quality Measure for Association Rules, pages 237-250. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007.

Defined as the support of the counter examples.

$$\text{ralambondrainy}(X \Rightarrow Y) = \frac{n_{X\bar{Y}}}{n} = \text{supp}(X \Rightarrow Y) = P(X \cap \bar{Y})$$

Range: $[0, 1]$ (smaller is better)

Relative Linkage Disequilibrium

Reference: Ron Kenett and Silvia Salini. Relative linkage disequilibrium: A new measure for association rules. In Petra Perner, editor, Advances in Data Mining. Medical Applications, E-Commerce, Marketing, and Theoretical Aspects, pages 189-199, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.

RLD is an association measure motivated by indices used in population genetics. It evaluates the deviation of the support of the whole rule from the support expected under independence given the supports of X and Y .

$$D = \frac{n_{XY}n_{\bar{X}\bar{Y}} - n_{X\bar{Y}}n_{\bar{X}Y}}{n}$$
$$\text{RLD} = \begin{cases} D/(D + \min(n_{X\bar{Y}}, n_{\bar{X}Y})) & \text{if } D > 0 \\ D/(D - \min(n_{XY}, n_{\bar{X}\bar{Y}})) & \text{otherwise.} \end{cases}$$

Range: $[0, 1]$

Relative Risk

Reference: Siström C.L. and Garvan, C.W. (2004), Proportions, Odds, and Risk, Radiology, 230:1, 12-19.

For the analysis of 2×2 contingency tables, relative risk is a measure of the relationship between two binary variables. It is defined as the ratio of the proportion of transactions containing Y in the two groups of transactions the do and do not contain X . In epidemiology, this corresponds to the ratio of the risk of having disease Y in the exposed (X) and unexposed (\bar{X}) groups.

$$\text{RR}(X \Rightarrow Y) = \frac{n_{XY}/n_X}{n_{\bar{X}Y}/n_{\bar{X}}} = \frac{P(Y|X)}{P(Y|\bar{X})} = \frac{\text{conf}(X \Rightarrow Y)}{\text{conf}(\bar{X} \Rightarrow Y)}$$

Range: $[0, \infty]$ ($RR = 1$ means X and Y are unrelated)

Rule Power Factor

Reference: Suresh Ochin and Nisheeth Joshi Kumar. Rule power factor: A new interest measure in associative classification. In 6th International Conference On Advances In Computing and Communications, ICACC 2016, Cochin, India, 2016.

Weights the confidence of a rule by its support. This measure favors rules with high confidence and high support at the same time.

Defined as

$$\text{rpf}(X \Rightarrow Y) = \text{supp}(X \Rightarrow Y) \text{conf}(X \Rightarrow Y) = \frac{P(X \cap Y)^2}{P(X)}$$

Range: $[0, 1]$

Right-Hand-Side Support

Alias: RHS support, consequent support

Support of the right-hand-side of the rule.

$$\text{RHSSupp}(X \Rightarrow Y) = \text{supp}(Y) = P(Y)$$

Range: $[0, 1]$

Sebag-Schoenauer

Reference: M. Sebag and M. Schoenauer. Generation of rules with certainty and confidence factors from incomplete and incoherent learning bases. In Proceedings of the European Knowledge Acquisition Workshop (EKAW'88), 1988.

Defined as

$$\text{sebag}(X \Rightarrow Y) = \frac{\text{conf}(X \Rightarrow Y)}{\text{conf}(X \Rightarrow \bar{Y})} = \frac{P(Y|X)}{P(\bar{Y}|X)} = \frac{\text{supp}(X \cup Y)}{\text{supp}(X \cup \bar{Y})} = \frac{P(X \cap Y)}{P(X \cap \bar{Y})}$$

Range: $[0, 1]$

Standardized Lift

Reference: P.D. McNicholas, T.B. Murphy, and M. O'Regan. Standardising the lift of an association rule. Computational Statistics & Data Analysis, 52(10):4712–4721, 2008.

Standardized lift uses the minimum and maximum lift that can reach for each rule to standardize lift between 0 and 1. The possible range of lift is given by the minimum

$$\lambda = \frac{\max\{P(X) + P(Y) - 1, 1/n\}}{P(X)P(Y)}.$$

and the maximum

$$v = \frac{1}{\max\{P(X), P(Y)\}}$$

The standardized lift is defined as

$$\text{stdLift}(X \Rightarrow Y) = \frac{\text{lift}(X \Rightarrow Y) - \lambda}{v - \lambda}.$$

The standardized lift measure can be corrected for minimum support and minimum confidence used in rule mining by replacing the minimum bound λ with

$$\lambda^* = \max \left\{ \lambda, \frac{4s}{(1+s)^2}, \frac{s}{P(X)P(Y)}, \frac{c}{P(Y)} \right\}.$$

Range: $[0, 1]$

Varying Rates Liaison

Reference: Bernard, Jean-Marc and Charron, Camilo (1996). L'analyse implicative bayesienne, une methode pour l'etude des dependances orientees. II : modele logique sur un tableau de contingence Mathematiques et Sciences Humaines, Volume 135 (1996), p. 5-18.

Defined as the lift of a rule minus 1 (0 represents independence).

$$VRL(X \Rightarrow Y) = lift(X \Rightarrow Y) - 1$$

Range: $[-1, \infty]$ (0 for independence)

Yule's Q

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Defined as

$$Q(X \Rightarrow Y) = \frac{\alpha - 1}{\alpha + 1}$$

where $\alpha = OR(X \Rightarrow Y)$ is the odds ratio of the rule.

Range: $[-1, 1]$

Yule's Y

Reference: Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right objective measure for association analysis. Information Systems, 29(4):293–313, 2004.

Defined as

$$Y(X \Rightarrow Y) = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$$

where $\alpha = OR(X \Rightarrow Y)$ is the odds ratio of the rule.

Range: $[-1, 1]$