

Rule mining

Oleg Shpynov

JetBrains Biolabs

November 19, 2014

Agenda

- Association rule
- Support-Confidence framework
- Lift, Conviction
- Comparison¹
- Application to Epigenetics

¹Comparing Rule Measures for Predictive Association Rules <http://www.di.uminho.pt/~pja/ps/conviction.pdf>

Notions

- T - set of items.
- Transaction d - subset of T
- DataBase D - list of transaction $d_i \subseteq T$
- Association Rule $X \rightarrow Y$.
 $X \subseteq T, Y \subseteq T, X \cap Y = \emptyset$
- $supp(X) = \frac{\#\{d_i | X \subseteq d_i\}}{\#\{D\}} = P(X)$
- $conf(X \rightarrow Y) = \frac{supp(X \rightarrow Y)}{supp(X)} = \frac{P(X \wedge Y)}{P(X)} = P(Y|X)$
- $lift(X \rightarrow Y) = lift(Y \rightarrow X) = \frac{conf(X \rightarrow Y)}{supp(Y)} = \frac{P(X \wedge Y)}{P(X)P(Y)}$
- $conviction(X \rightarrow Y) = \frac{1 - supp(Y)}{1 - conf(X \rightarrow Y)} = \frac{P(X)P(\neg Y)}{P(X \wedge \neg Y)}$

Measurements

| Measure | Definition | Range |
|---------------|--|---------------|
| confidence | $conf(A \rightarrow C) = \frac{sup(AUC)}{sup(A)}$ | [0, 1] |
| Laplace | $lapl(A \rightarrow C) = \frac{sup(AUC)+1}{sup(A)+2}$ | [0, 1[|
| lift | $lift(A \rightarrow C) = \frac{conf(A \rightarrow C)}{sup(C)}$ | [0, +∞[|
| conviction | $conv(A \rightarrow C) = \frac{1-sup(C)}{1-conf(A \rightarrow C)}$ | [0.5, +∞[|
| leverage | $leve(A \rightarrow C) = sup(A \cup C) - sup(A) \times sup(C)$ | [-0.25, 0.25] |
| χ^2 | $\chi^2(A \rightarrow C) = N \times \sum_{X \in \{A, \neg A\}, Y \in \{C, \neg C\}} \frac{(sup(XUY) - sup(X) \cdot sup(Y))^2}{sup(X) \times sup(Y)}$ | [0, +∞[|
| Jaccard | $jacc(A \rightarrow C) = \frac{sup(AUC)}{sup(A)+sup(C)-sup(AUC)}$ | [0, 1] |
| cosine | $cos(A \rightarrow C) = \frac{sup(AUC)}{\sqrt{sup(A) \times sup(C)}}$ | [0, 1] |
| ϕ -coeff | $\phi(A \rightarrow C) = \frac{leve(A \rightarrow C)}{\sqrt{(sup(A) \times sup(C)) \times (1-sup(A)) \times (1-sup(C))}}$ | [-1, 1] |
| mutual inf. | $MI(A \rightarrow C) = \frac{\sum_i \sum_j sup(A_i \cup C_j) \times \log(\frac{sup(A_i \cup C_j)}{sup(A_i) \times sup(C_j)})}{\min(\sum_i -sup(A_i) \times \log(sup(A_i)), \sum_j -sup(C_j) \times \log(sup(C_j)))}$ | [0, 1] |

Epigenetics

DataBases

- Binned
- Gene Locus (TSS in BivalentHCPExperiment)
- any locations (in this terms any kind of LocationsRule can be reproduced)

Predicates

- CpG predicate
- Enrichment predicate max vote within (DZip)? PoissonHMM model
- ChromHMM predicate location is marked if $> 50\%$ intersection with state

Hypotheses

- Cpg content \Rightarrow Histone enrichment
- Histone enrichment \Rightarrow Histone enrichment
- ChromHMM State \Rightarrow Histone enrichment
- Histone enrichment \Rightarrow ChromHMM State

Hypothesis *Condition* \Rightarrow *Target* on DataBase *D*.

DATABASE RULE

- What if *Target* is property of *D* itself?
- Compute fold change $\text{confidence}(TRUE \Rightarrow \text{Target})$ on *D* vs *D.sampleGenomic()*.
- If $\text{foldChange} > \text{threshold}$ then *Target* is a *D* property!

Hypothesis $Condition \Rightarrow Target$ on DataBase D .

RULE

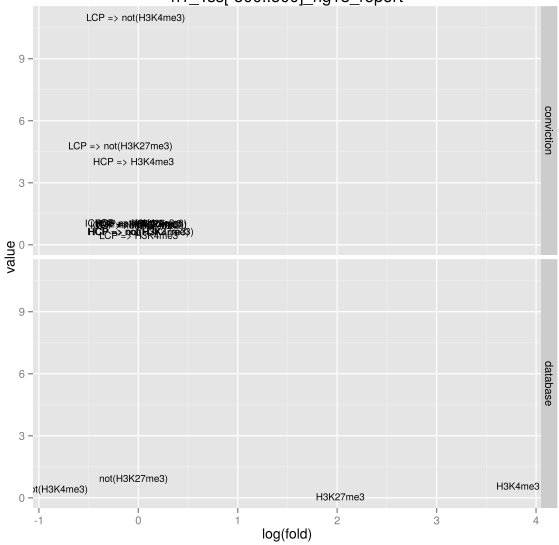
- Measure $conviction(Condition \Rightarrow Target)$ on D .
- What if we get good value by chance?
- We can sample and get confidence intervals for $conviction$ - check how stable it is.
- Check $foldChange$ on D vs $D.sample$
- If $foldChange < threshold$ then we found a RULE!

BivalentHCPExperiment

- FOUND DATASET RULE: $Tss[-500..500]_{protein}$: H3K4me3 is 43.49834364520371 more often vs genomic [0.5695912512755348]
- FOUND RULE: $HCP \Rightarrow H3K4me3$
DATA: $Tss[-500..500]_{protein}$ [51939]
CONDITION: HCP 0.45285816053447314 [23521]
TARGET: $H3K4me3$ 0.5695912512755348 [29584]
CONFIDENCE: $H3K4me3|HCP$: 0.8931167892521577 [21007]
LIFT: 1.5679959747487768
CONVICTION: 4.026906992342143
- FOUND RULE: $LCP \Rightarrow not(H3K4me3)$
DATA: $Tss[-5000..-2500]_{protein}$ [51939]
CONDITION: LCP 0.6139317276035349 [31887]
TARGET: $not(H3K4me3)$ 0.9702920733937889 [50396]
CONFIDENCE: $not(H3K4me3)|LCP$: 0.996173989400069 [31765]
LIFT: 1.0266743558109808
CONVICTION: 7.7647266860020165

● ...

h1_Tss[-500..500]_hg18_report



ChromHMMChipSeqExperiment

- *ChipSeq* \Rightarrow *ChipSeq*
H3K4me3 \Rightarrow *H3K4me2*
- *ChromHMM* \Rightarrow *ChipSeq*
H3K4me2 as result for many rules with high conviction
- *ChipSeq* \Rightarrow *ChromHMM*
No rules found.

References

- Comparing Rule Measures for Predictive Association Rules
<http://www.di.uminho.pt/~pja/ps/conviction.pdf>
- Measures overview http://michael.hahsler.net/research/association_rules/measures.html
- Web Data mining - Bing Liu
<http://link.springer.com/book/10.1007/978-3-540-37882-2>