

Distributed recipe mining to combine historic and social media knowledge

Motivation



Social Media

develops during the time

redundant

adaptions

Manhattan

6 cl Rye

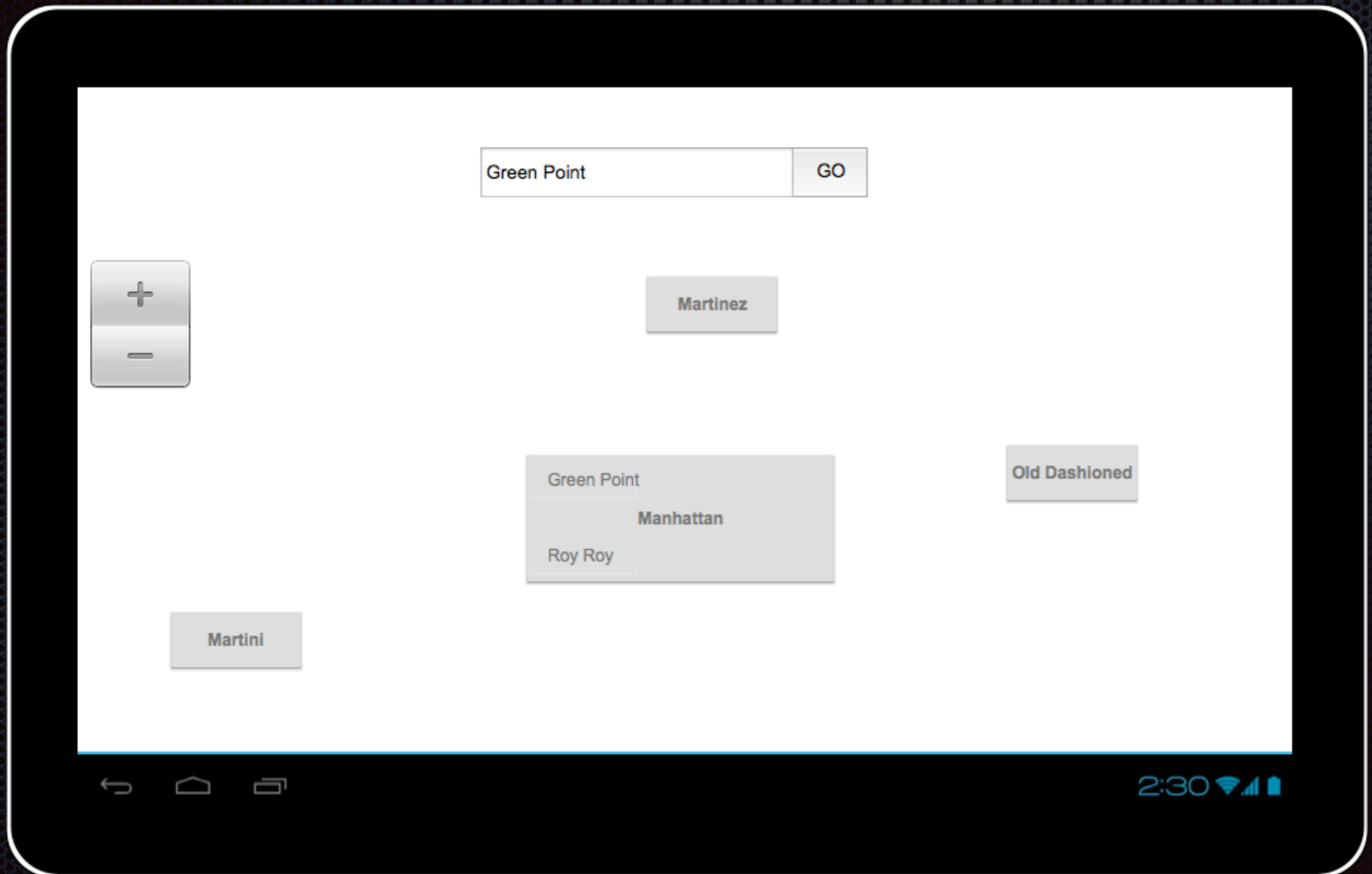
3 cl roter Wermut

stir, strain in cocktail glass

Historic
Books

What is the main
recipe?

Search on cluster map



Inner Cluster


Back to Overview

Recipes, Ingerediants or Tools

Main Recipe

Manhattan

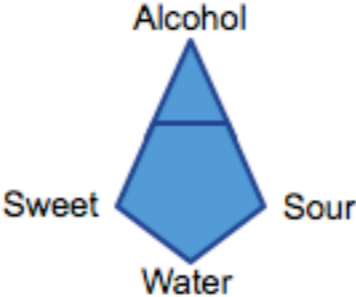
6 cl Rye
3 cl roter Wermut
rühren, im Cocktailspitz



Statistics

How many findings?	122
How old is it?	140 year

Alcohol



Sweet Sour

Water

Adaptions

Main Recipe
Green Point
Rob Roy

Timeline

1878 Jerry Thomas S. 122	1922 Savoy S. 43	...
--------------------------	------------------	-----

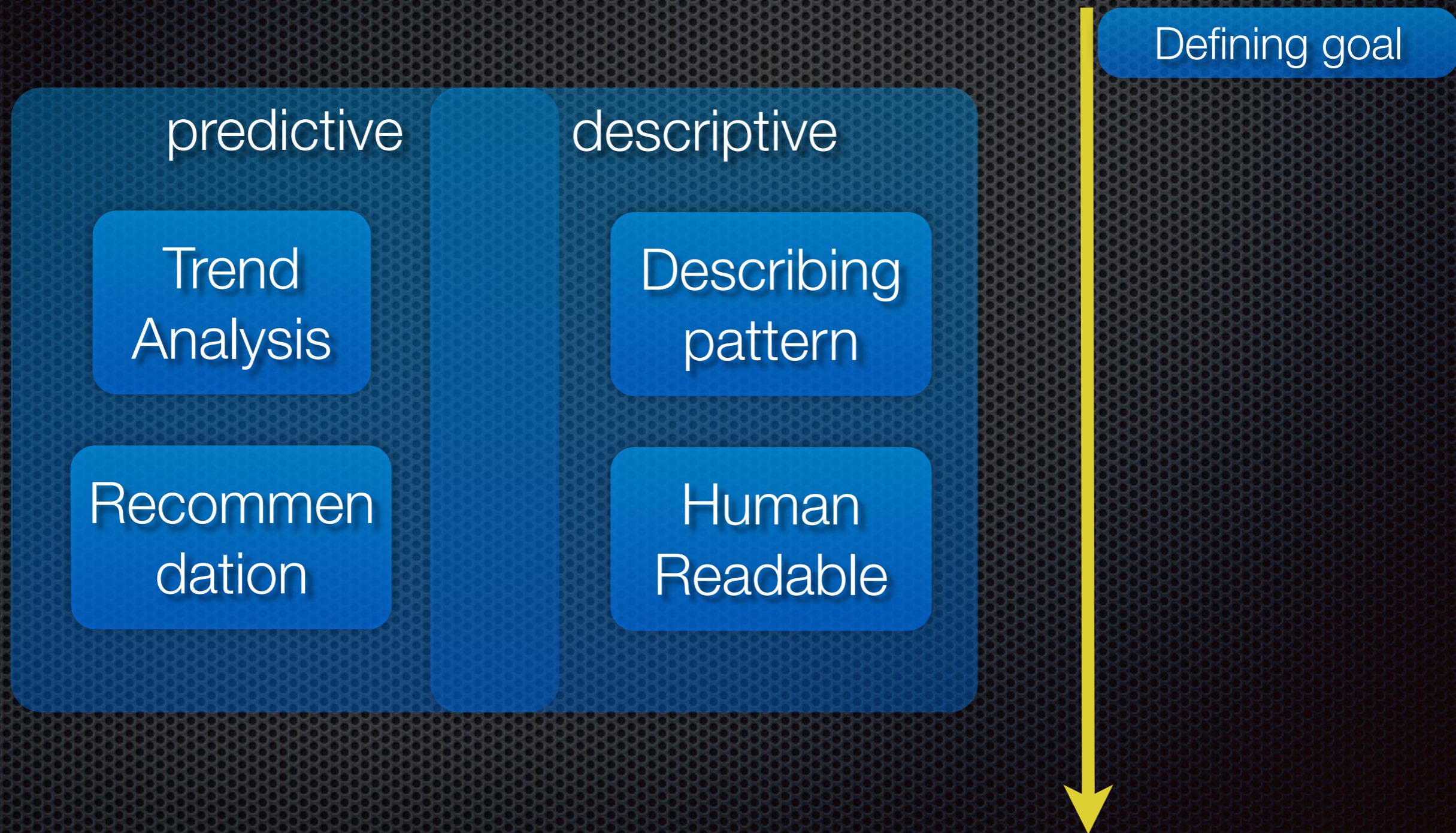
2:30

KDD process [FPSS96]

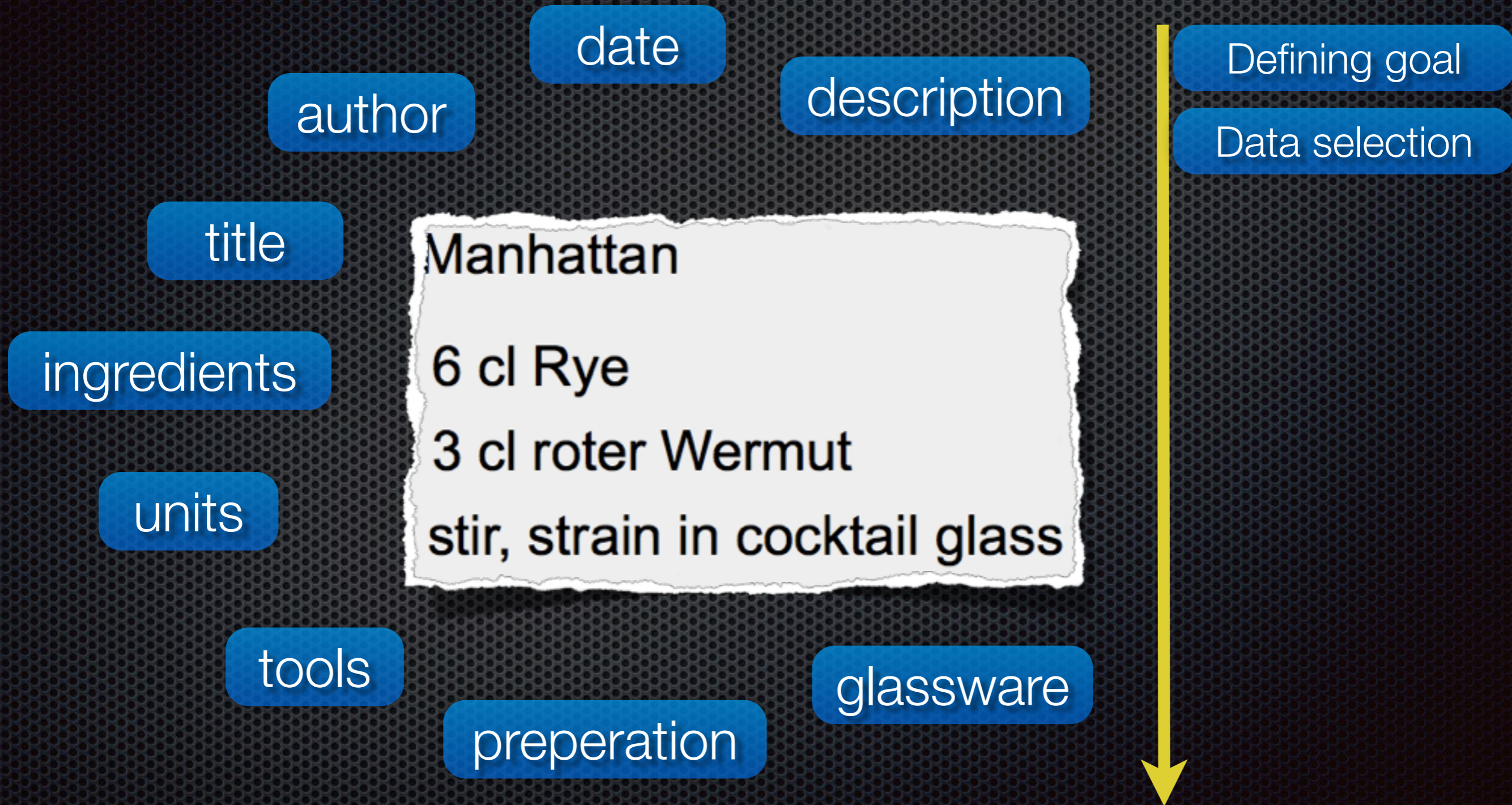
Knowledge discovery in database and data mining



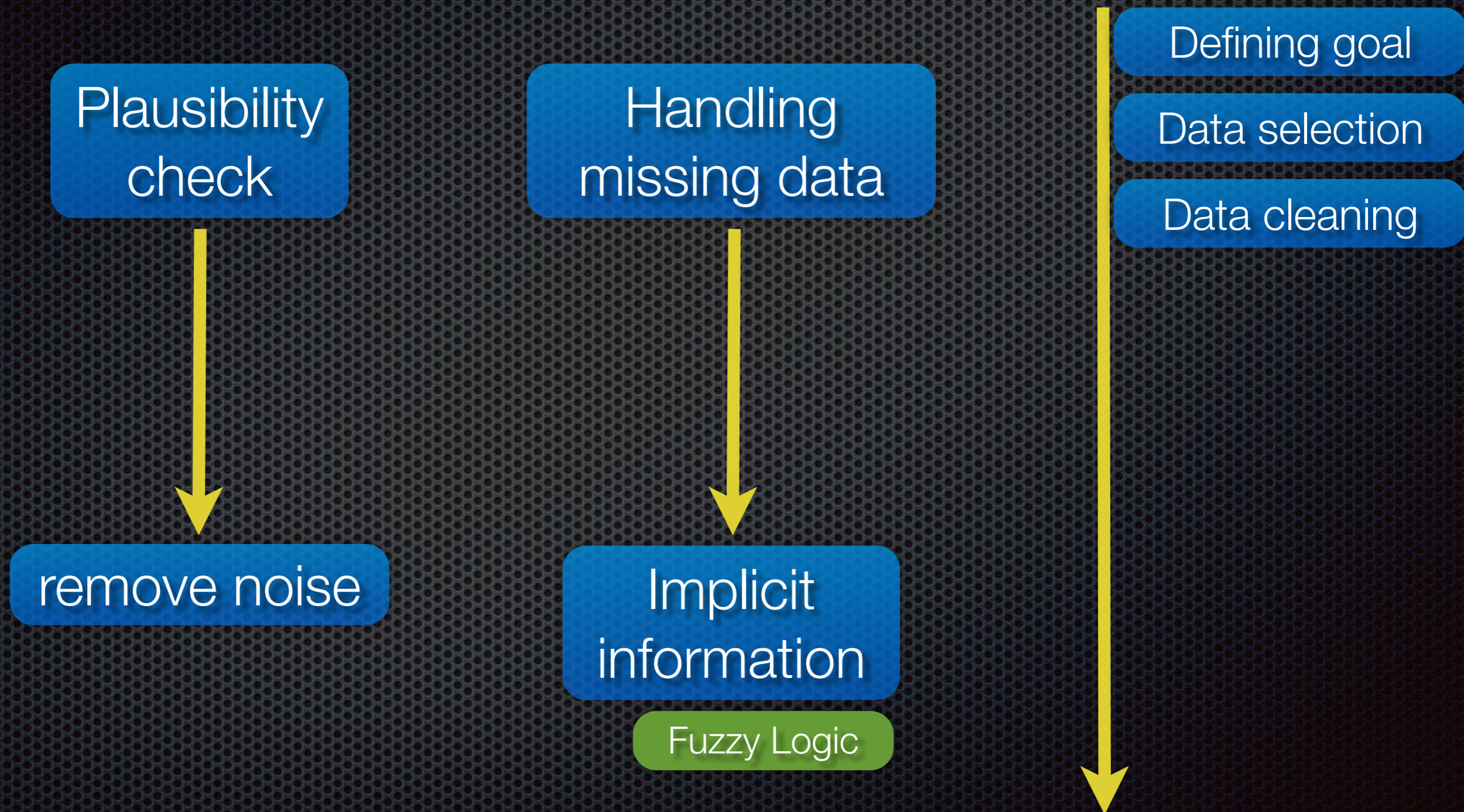
Defining goal domain and feature selection



Feature selection

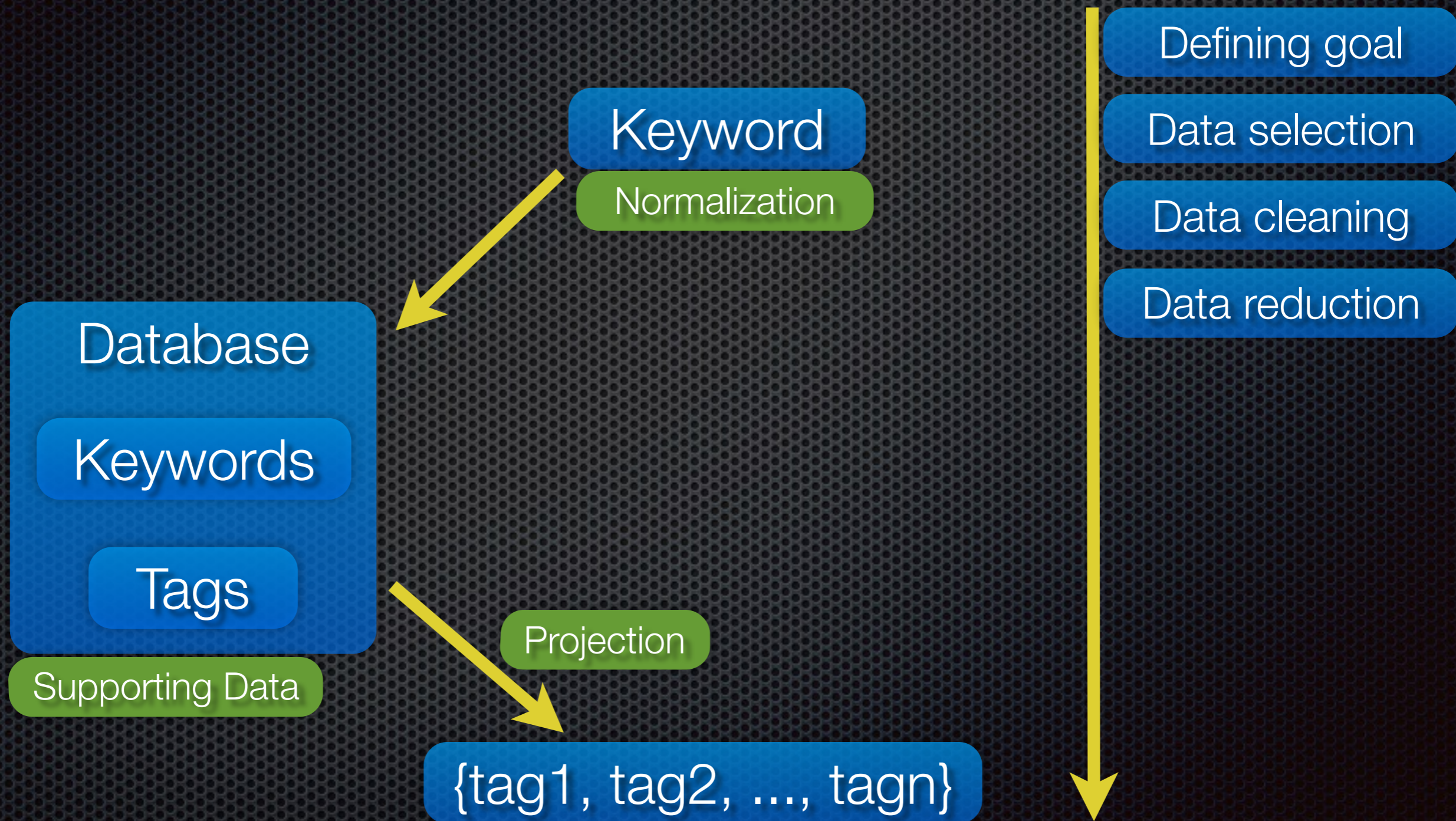


Data cleaning



Data reduction & projection

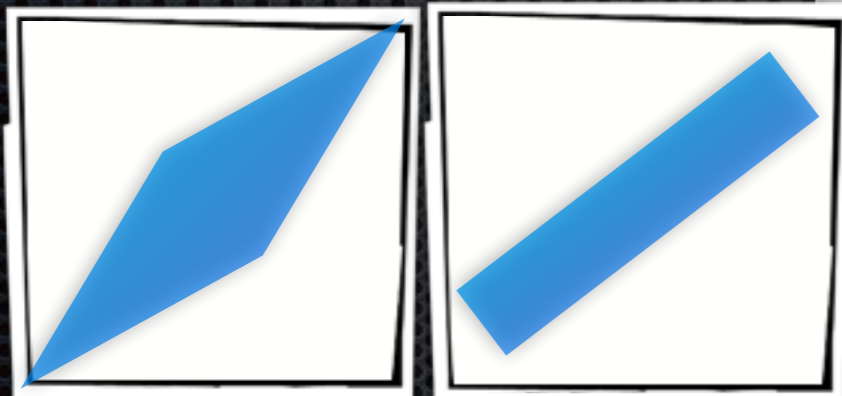
Feature extraction



Distance function

Euklid distance [JMF99]

$$d_2(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{k=1}^d (x_{i,k} - x_{j,k})^2 \right)^{1/2}$$



$$\text{lev}_{a,b}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

Dynamic-Time-Warping

[HLWG08]

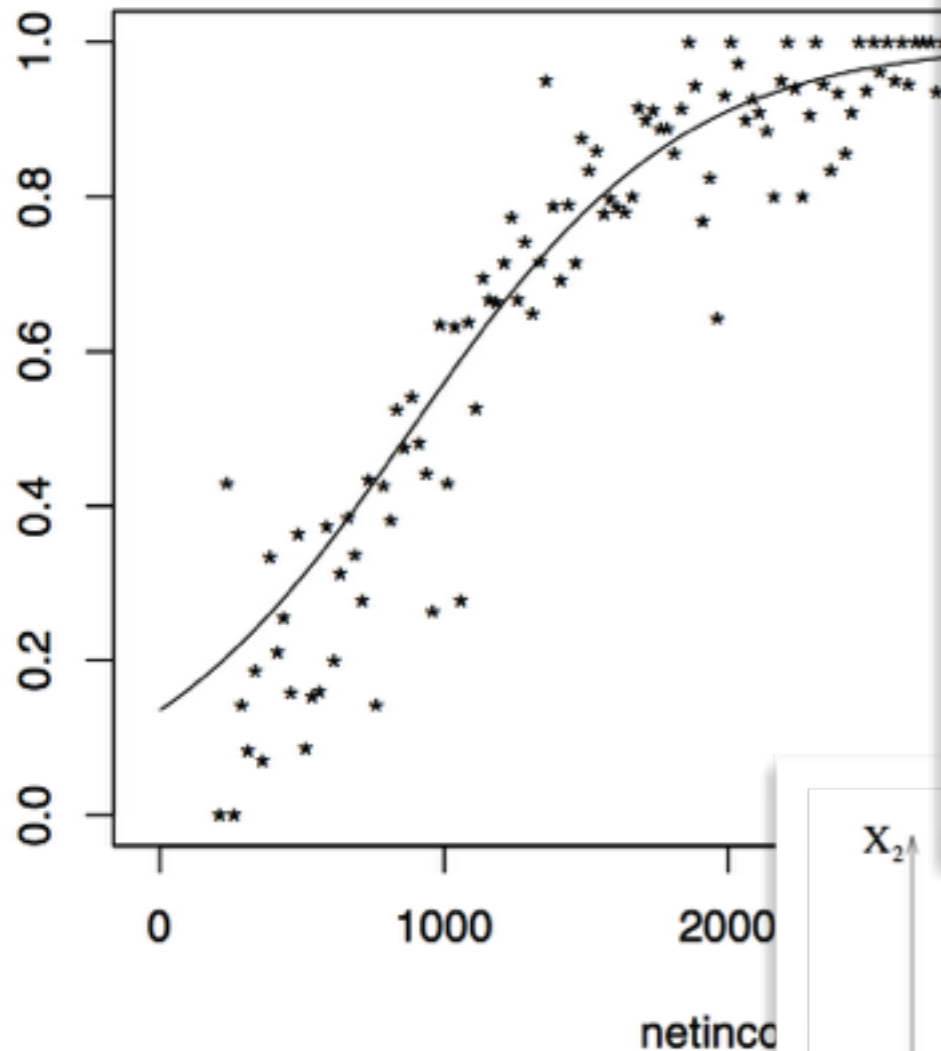
Speech Recognition

Levenshtein-Distance

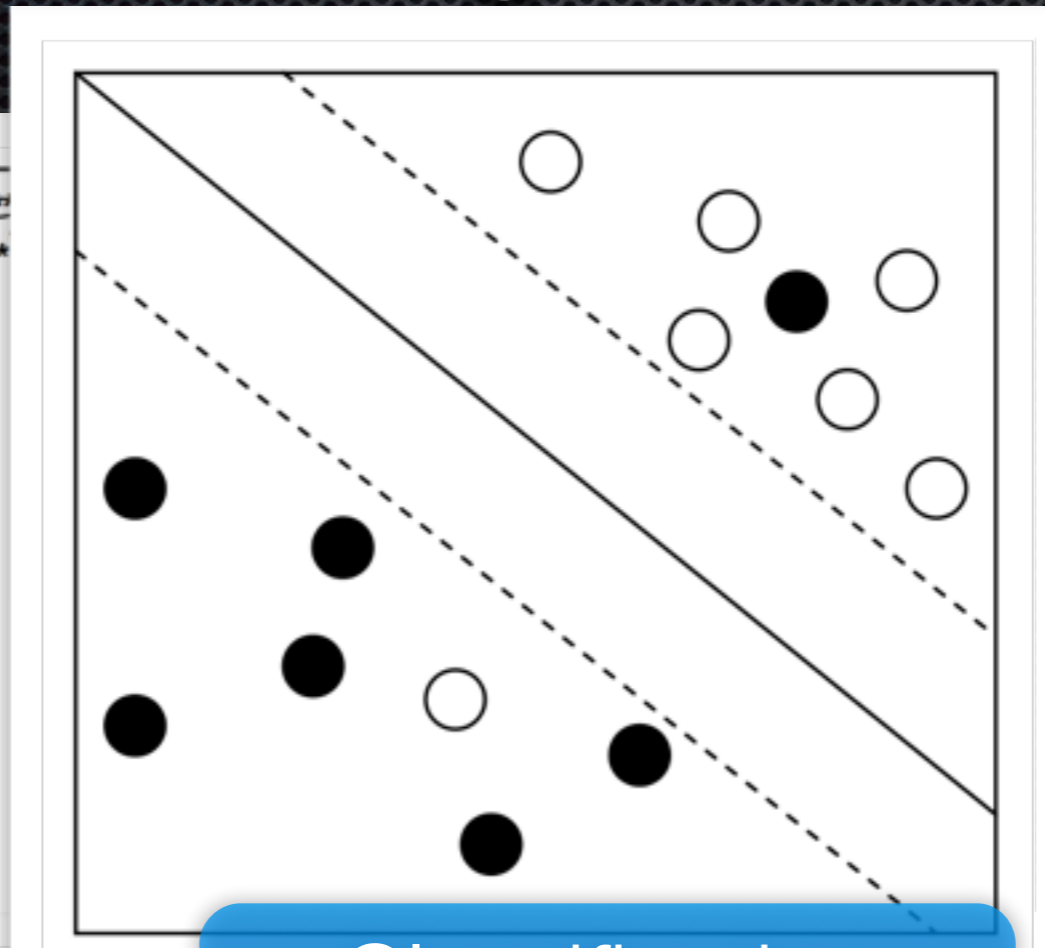
String recognition

Choosing mining function

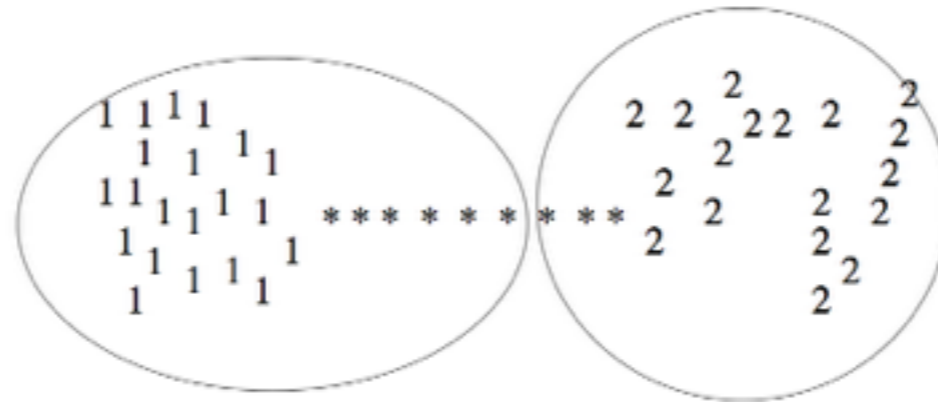
[FKPT07, S. 512]



Regression



Classification



Clustering

[LW06, S. 6]

Defining goal

Data selection

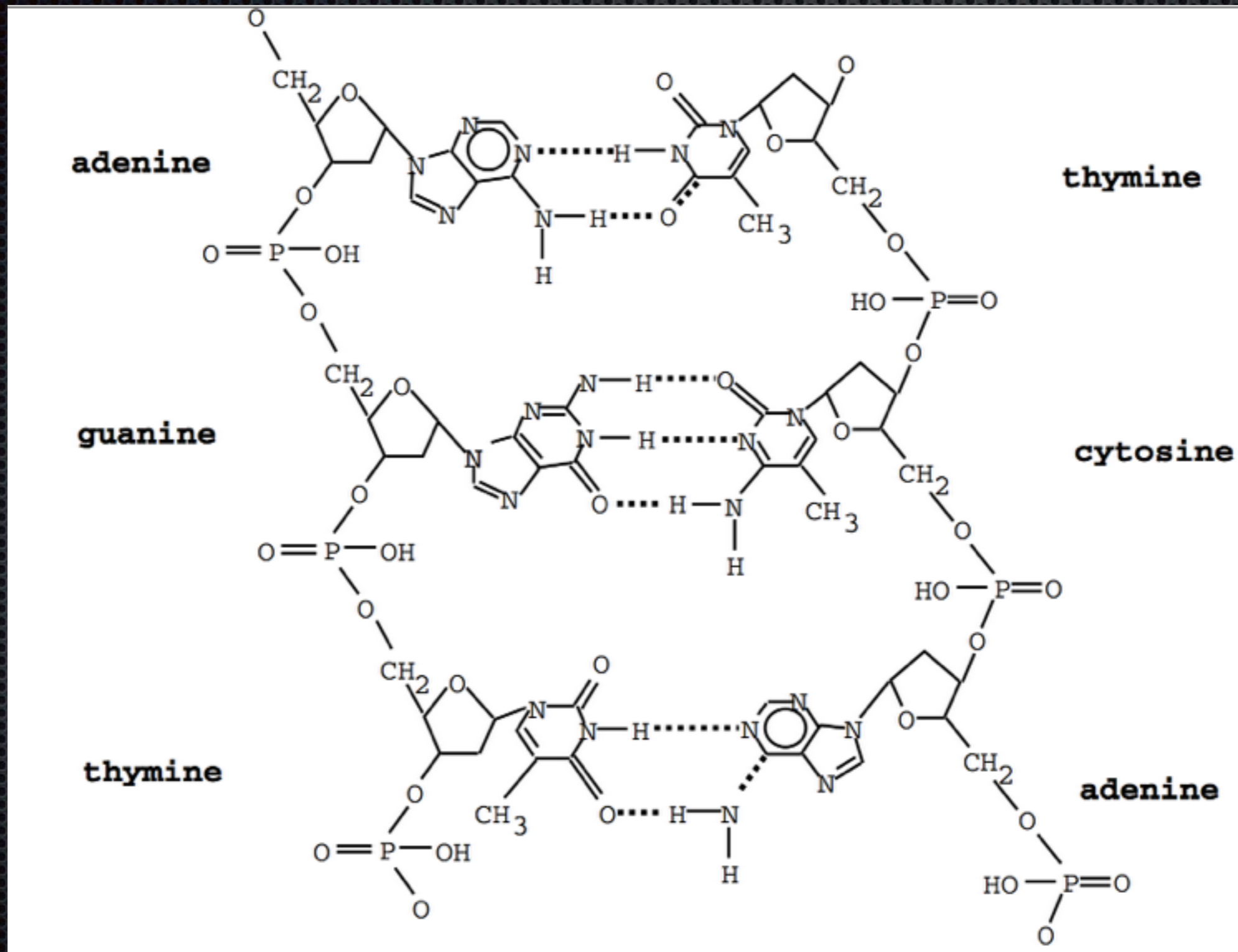
Data cleaning

Data reduction

Mining function

[JMF99, S. 277]

Chemical Model [MHC06, S. 80]



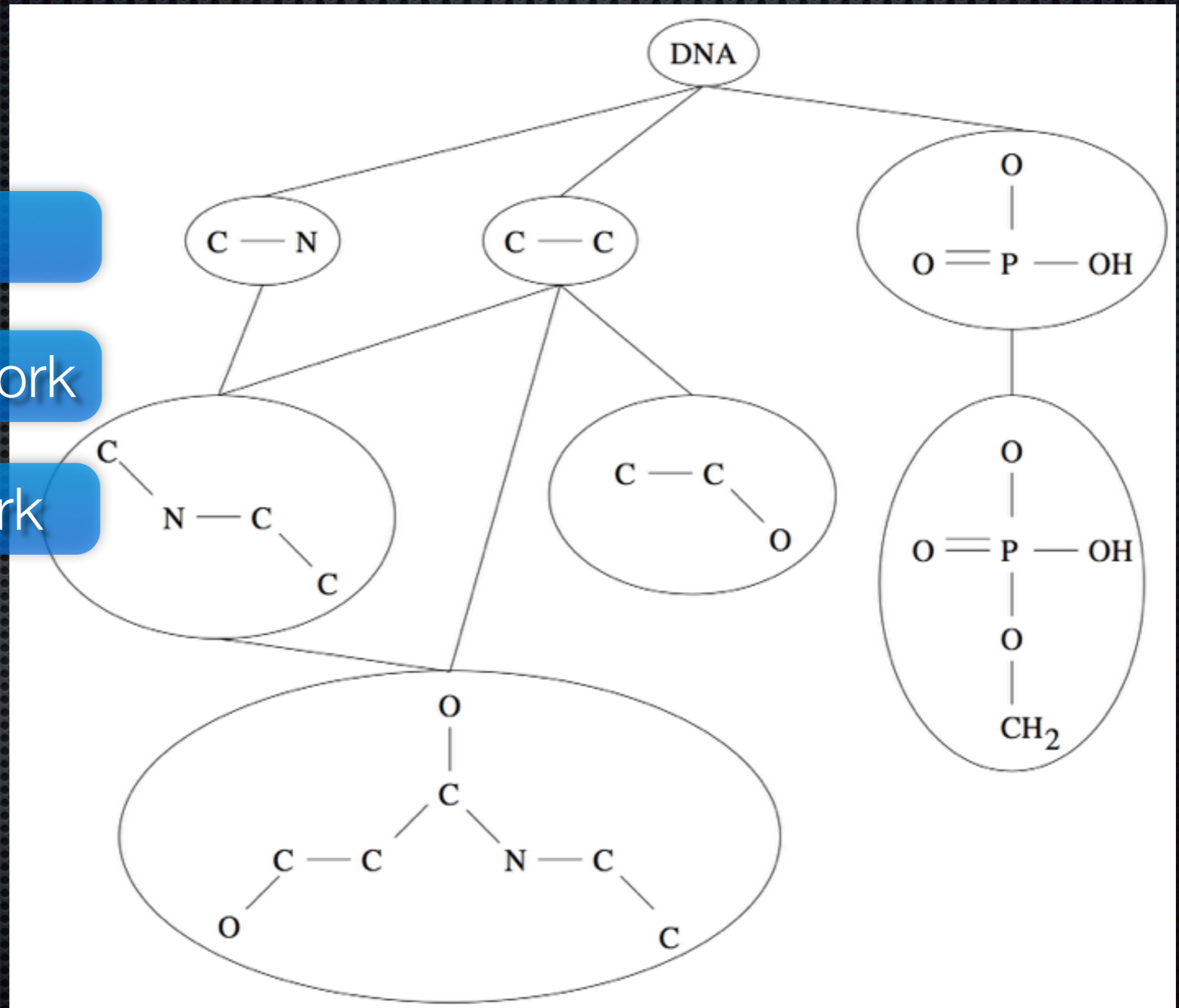
(Graph) Representation

[MHC06, S. 80]

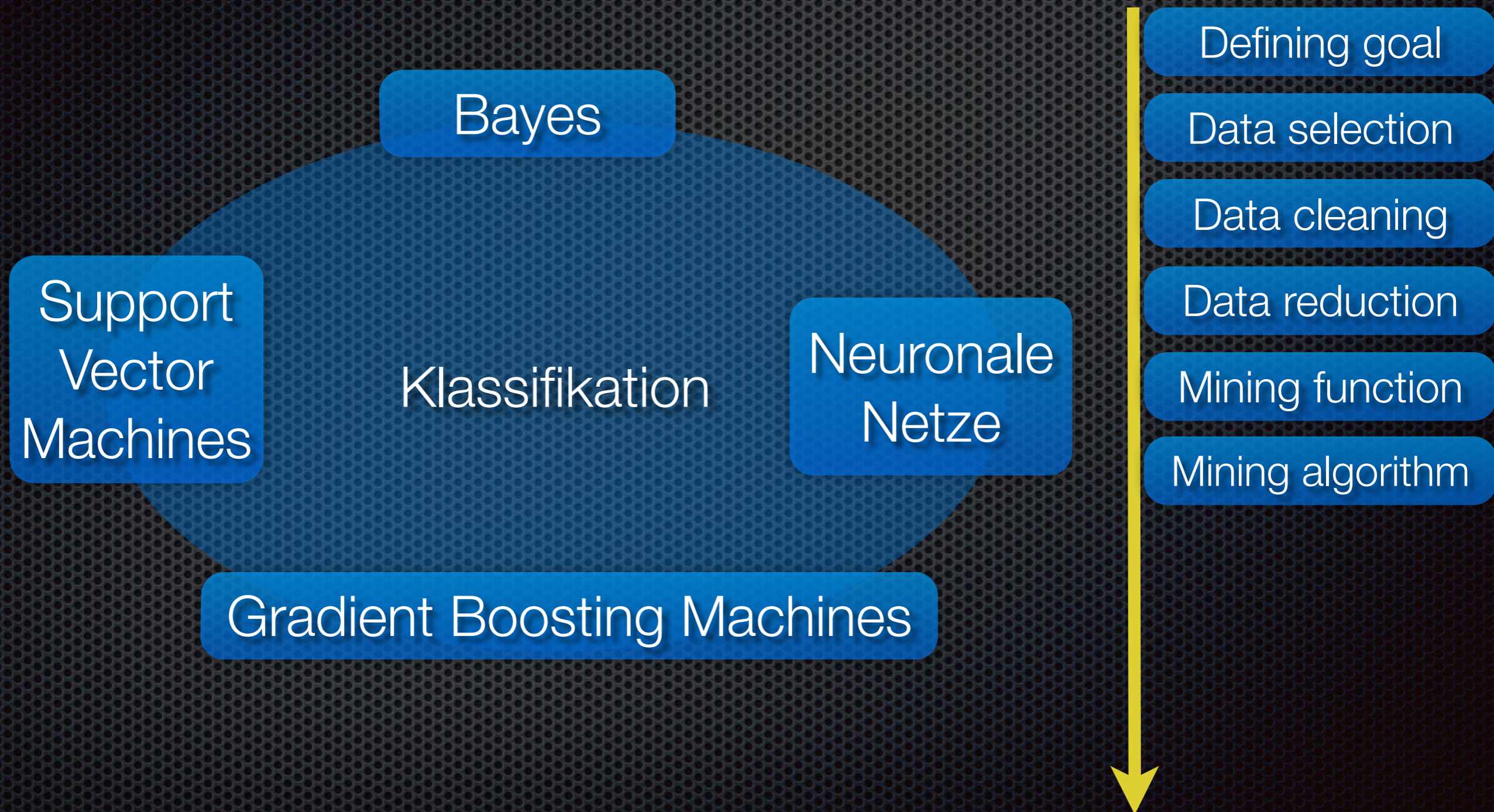
Trees

Bayesian network

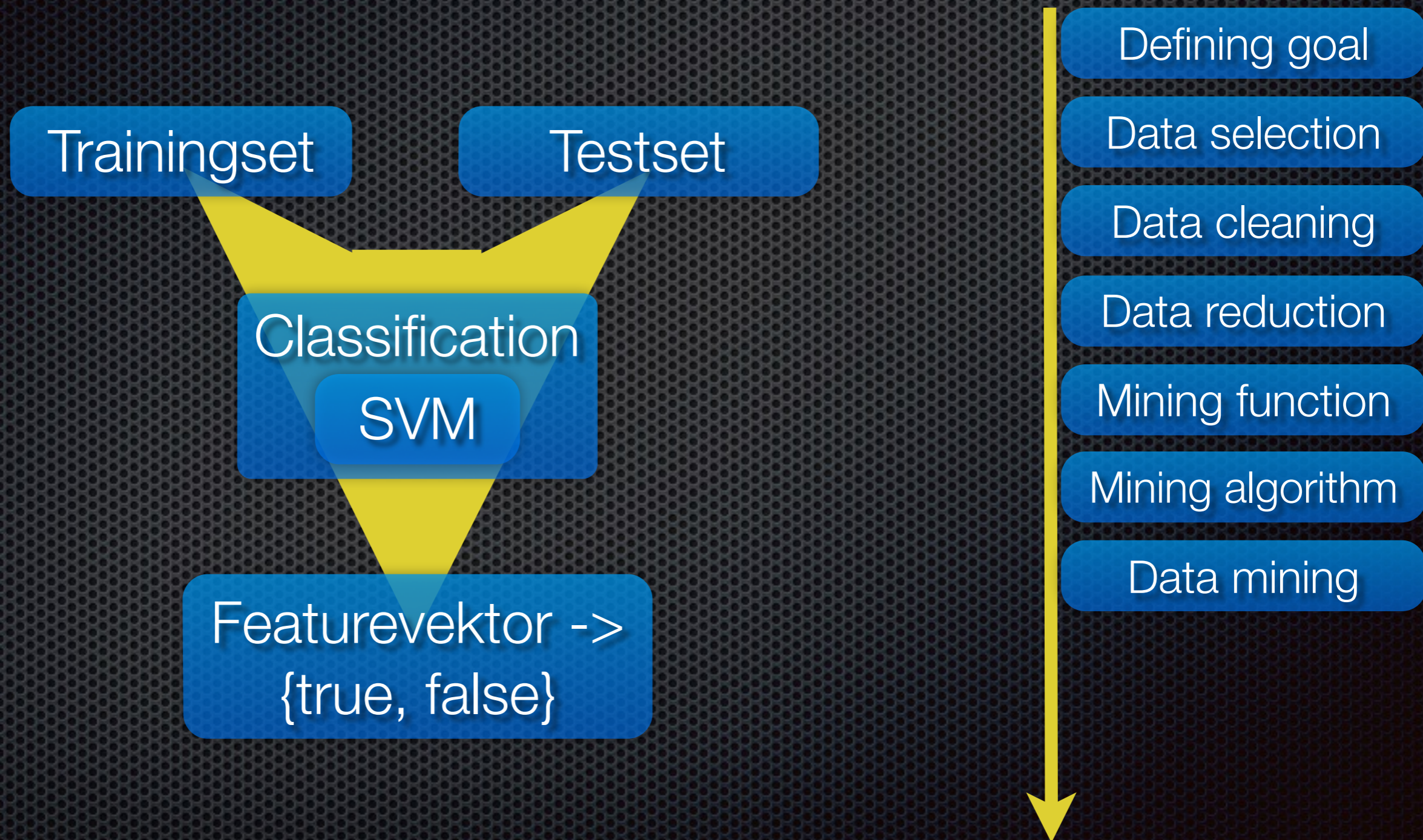
Neural network



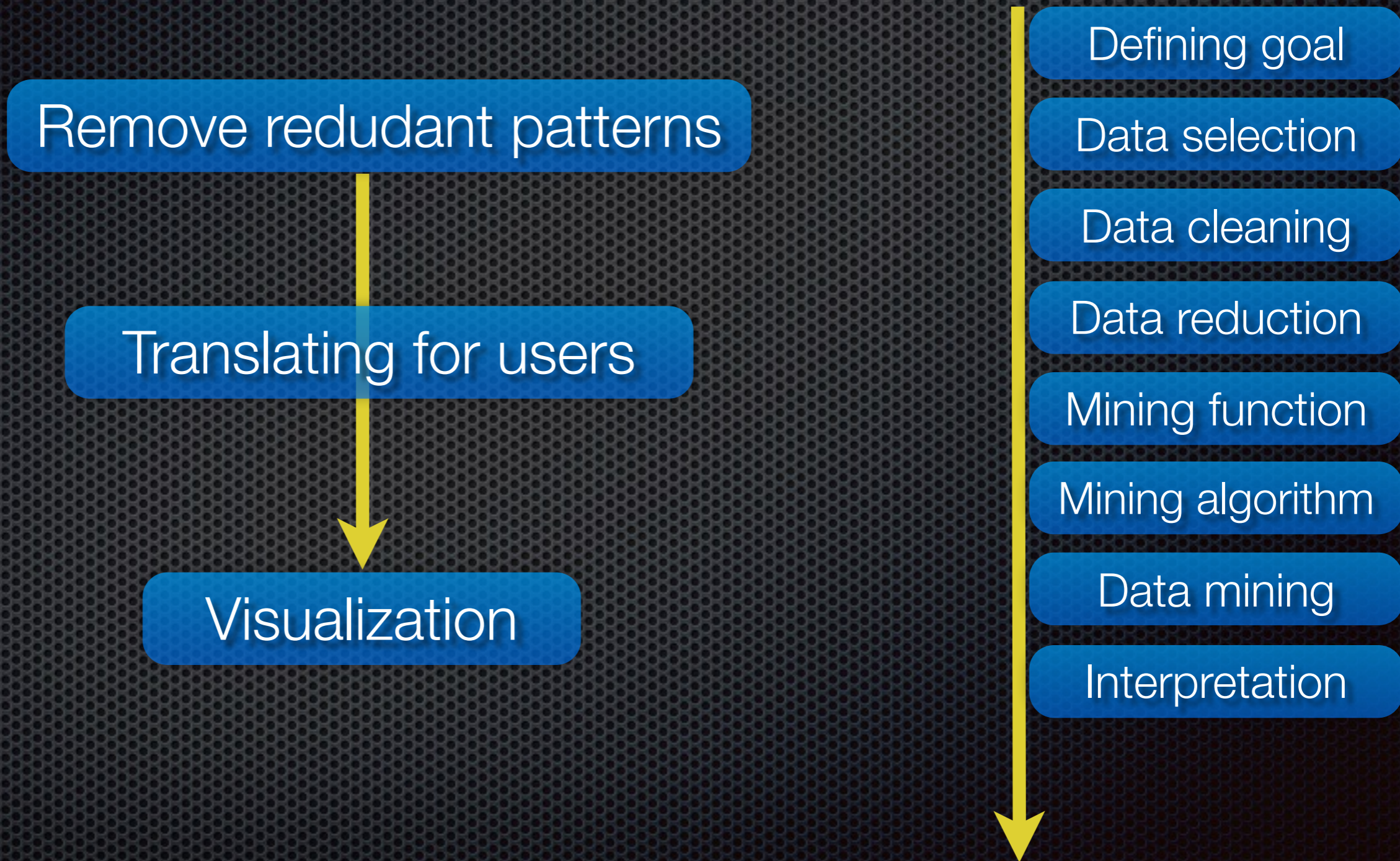
Choosing mining algorithm



Data mining: Classification



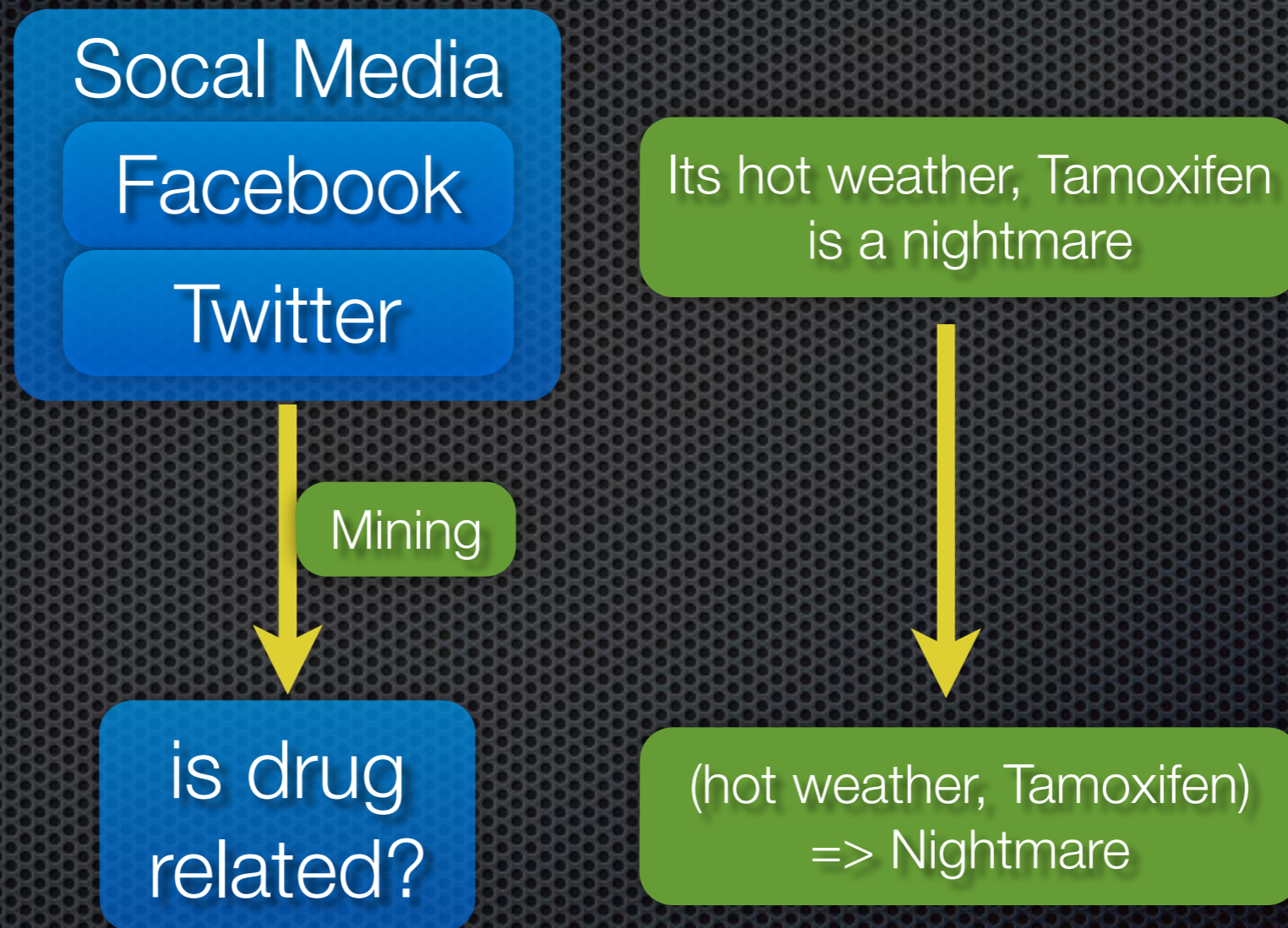
Interpretation



Optimizing



Large scale Twitter mining of drug related adverse events [BTY12]



Classification Training

User Training set

Potential drug users

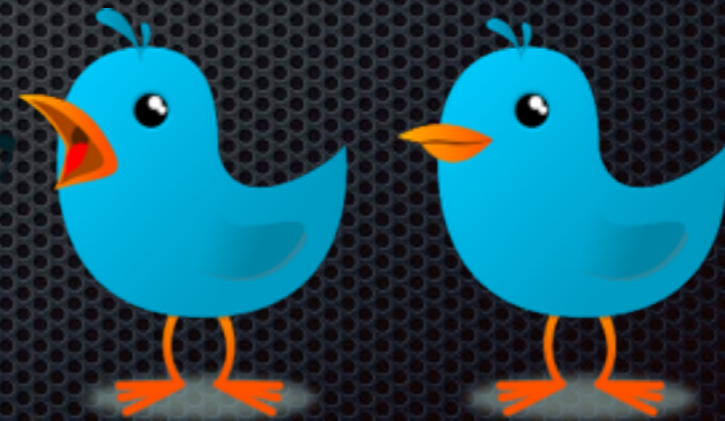
Corresponding Keywords in Timeline



Tweet Training set

Medicine Ontology

Unified Medicine Language system

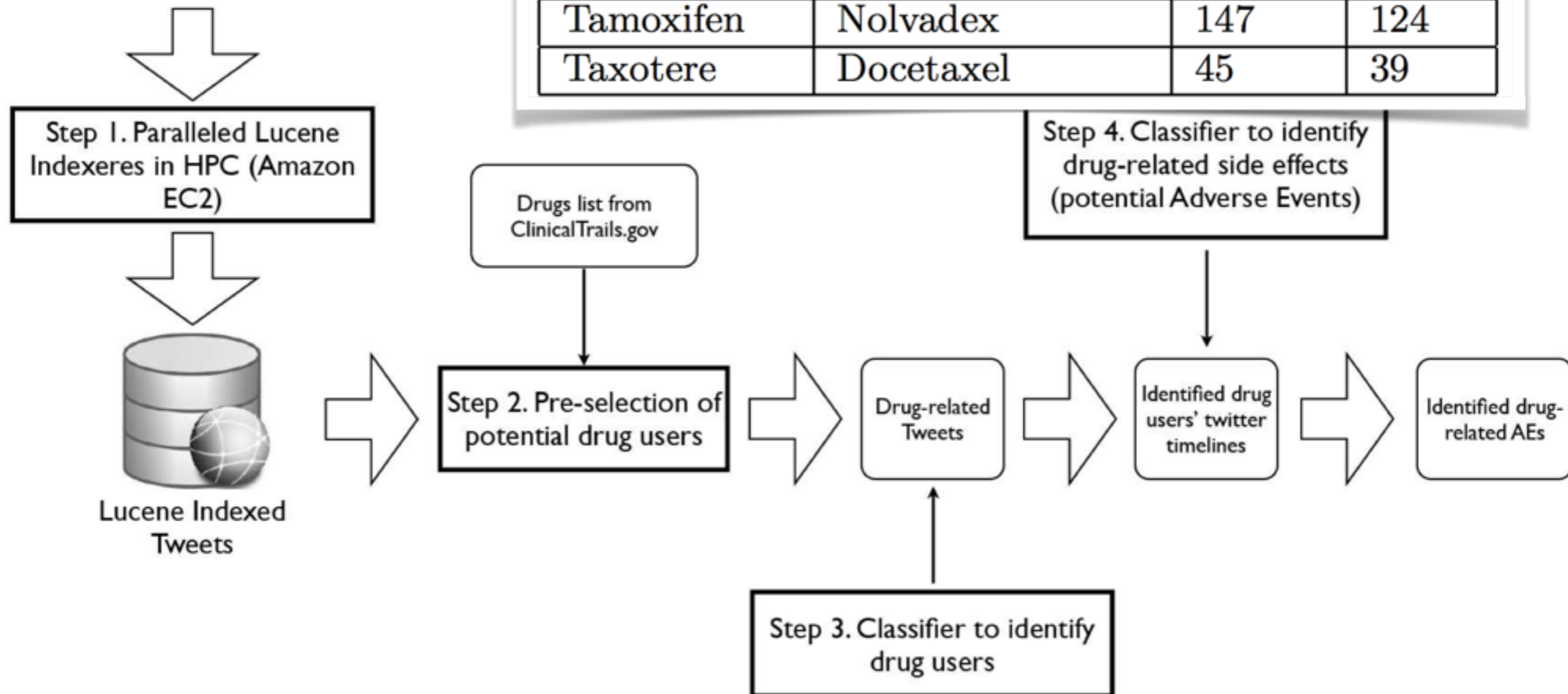


Mining Process to find user

[BTY12, S. 27]

Drug name	Synonym(s)	# of tweets	# of users
Avastin	Bevacizumab	264	236
Melphalan	ALKERAN	23	15
Rupatadin	Rupafin, Urtimed	10	10
Tamoxifen	Nolvadex	147	124
Taxotere	Docetaxel	45	39

2,102,176,189 Tweets



Feature Extraction of Tweets

Bag of words

Action

State of
Drug Using

Count

Hash Tags

Reply Tags

URLs

Pronouns

Drug names

Relevance

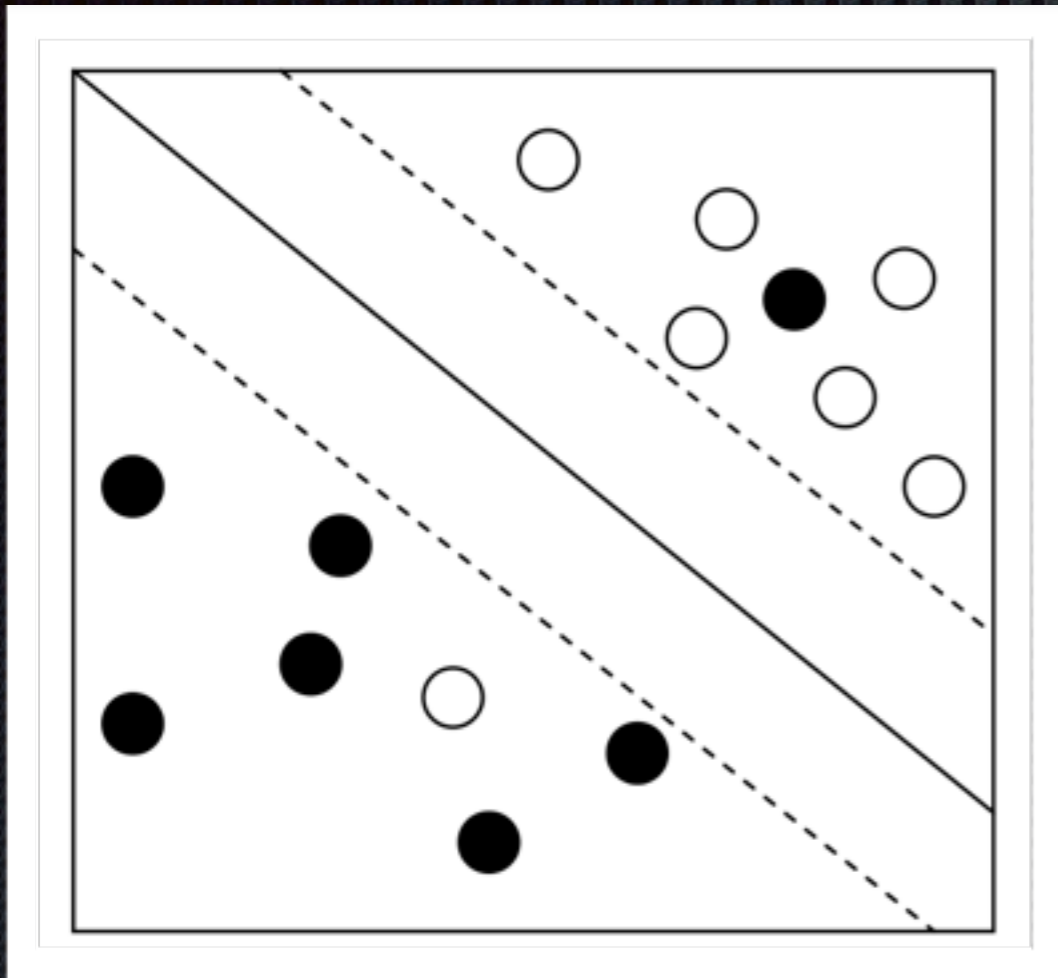
Count

Semantic Types

Semantic Groups

Relevance in Unified
Medicine Language system

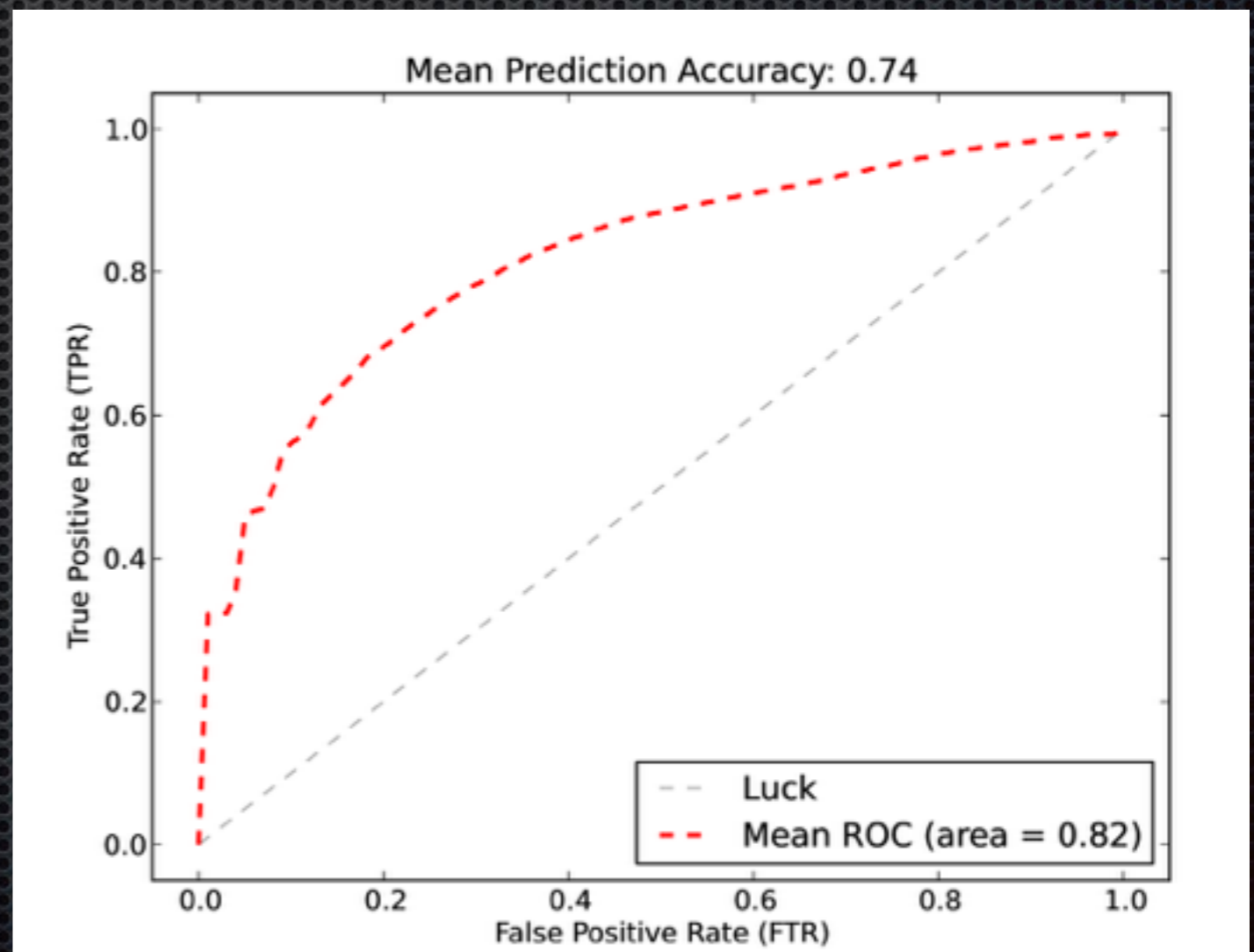
Support Vector Machine



[LW06, S. 6]

decision function $g(x) = \text{sign}(f(x))$

hyper plane $\langle w, x \rangle + b = 0$



[BTY12, S. 30]

Recipe recommendation [TLA12]



Preparation

Dishes



Ingredient combination

removed quantities

removed temperature

key ingredients

modification options

nutritions

Gradient Boosting Machine

Classification

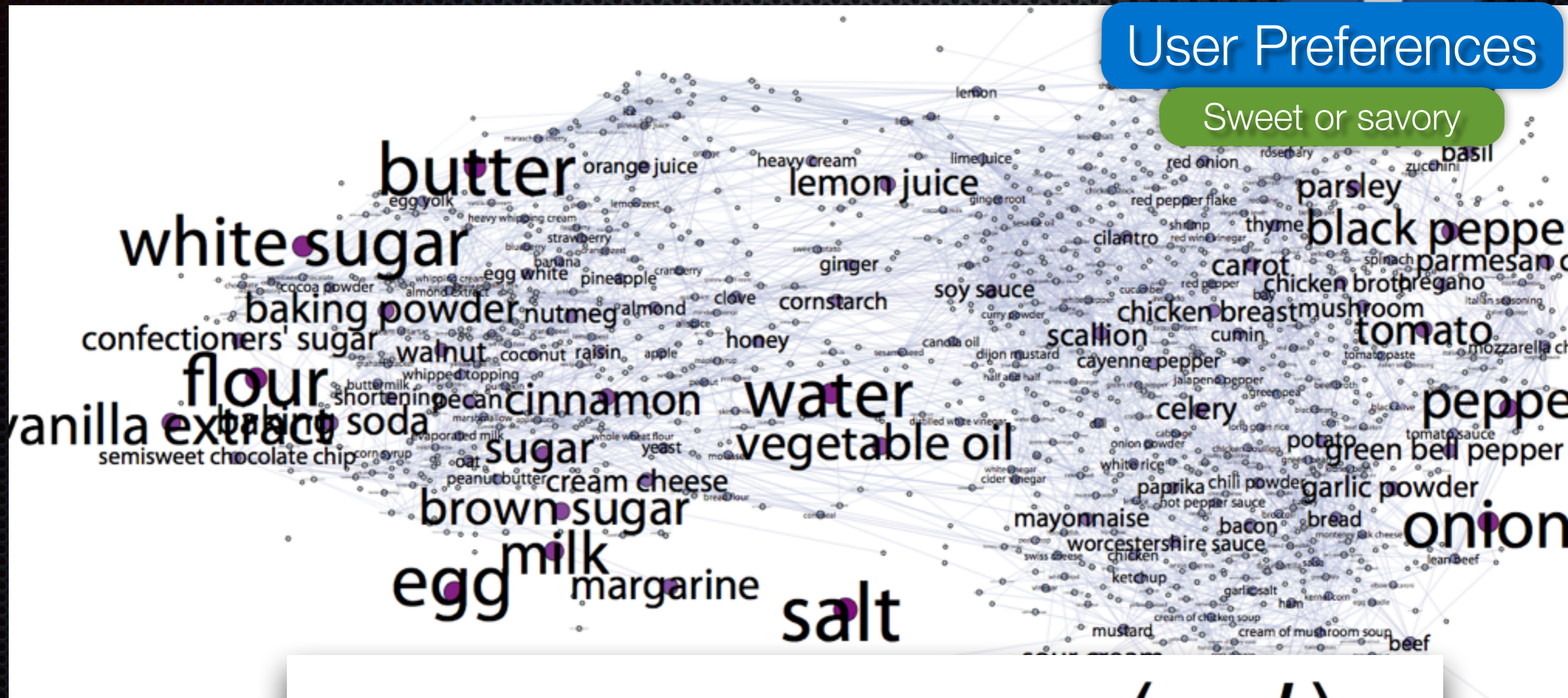
[Fri01]

Co-occurrence network



User Preferences

Sweet or savory



[TLA12, S. 301]

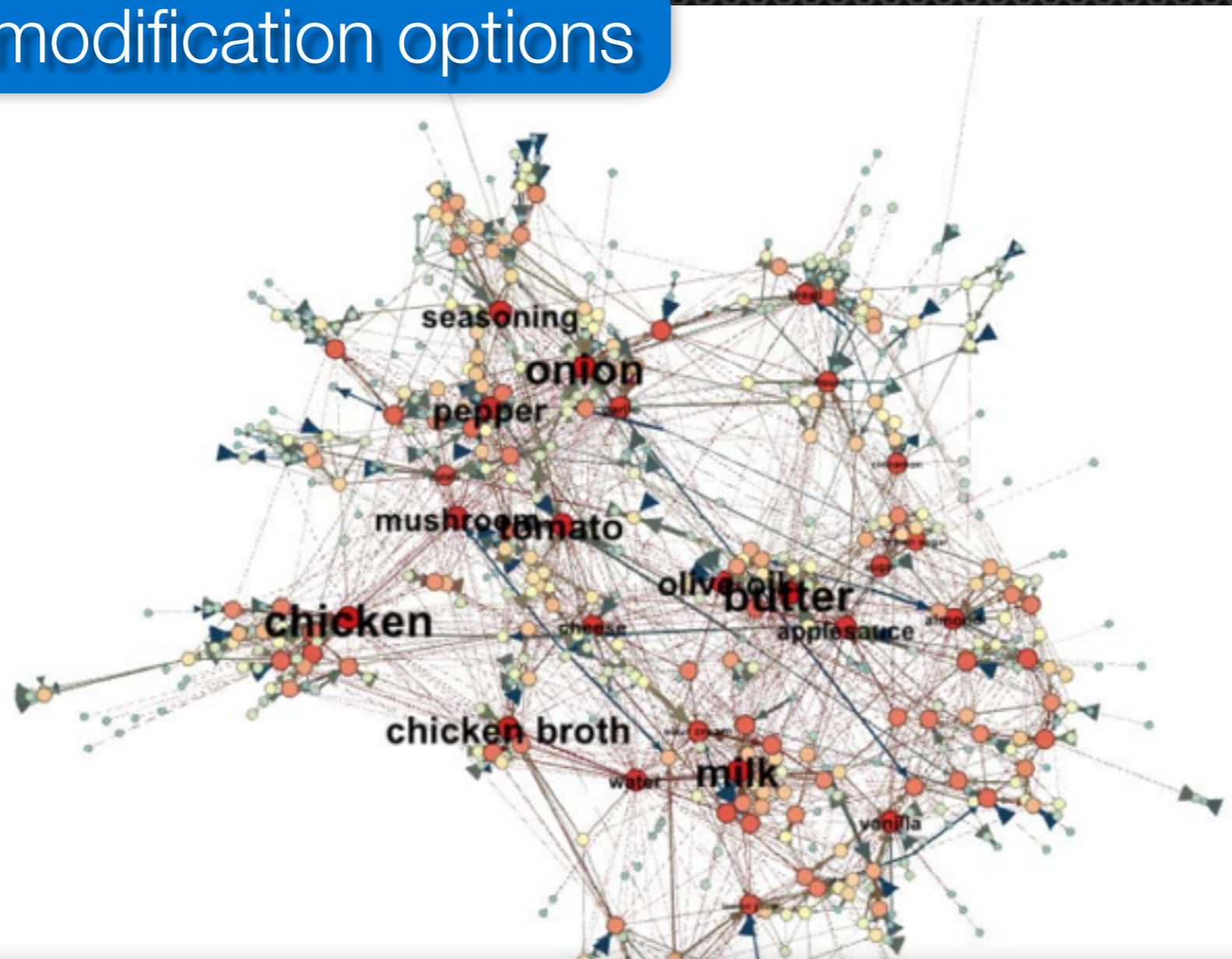
$$PMI(a, b) = \log \frac{p(a, b)}{p(a)p(b)},$$

pointwise mutual information

Substitution network

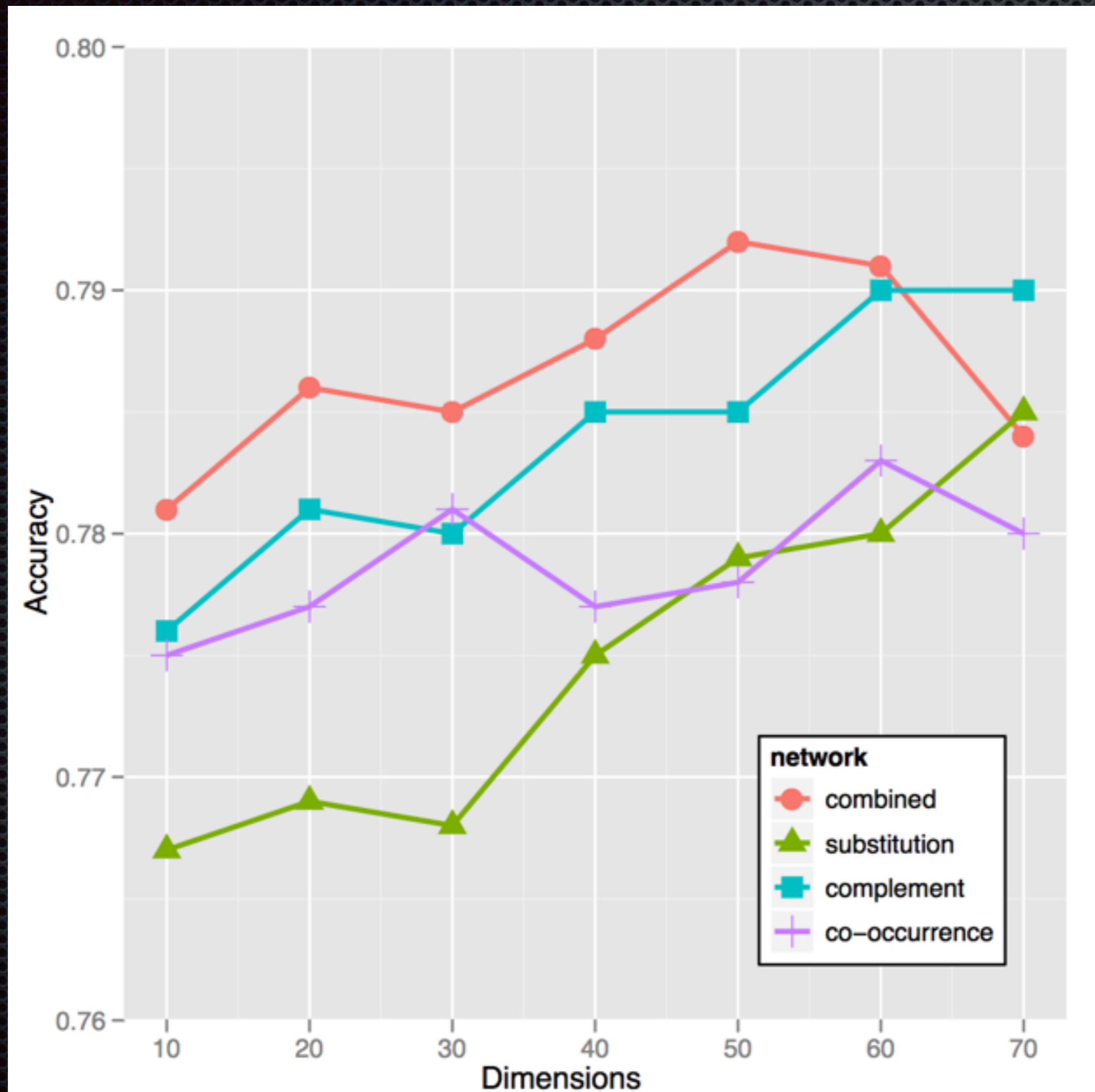
[TLA12, S. 302]

based on modification options



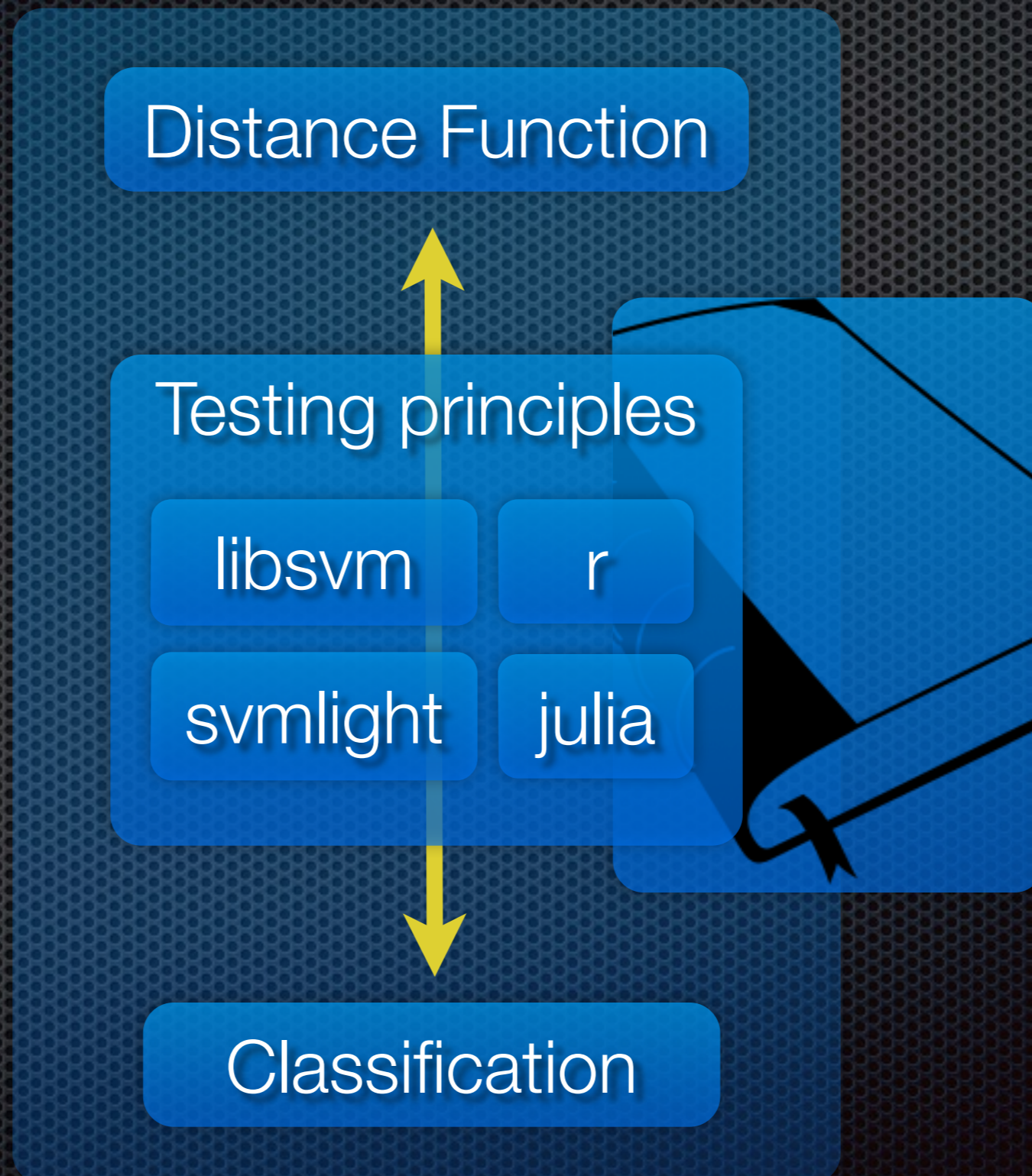
$$\text{PMI}(a \rightarrow b) = \log \frac{p(a \rightarrow b)}{p(a)p(b)},$$

Prediction performance



[TLA12, S. 307]

Roadmap Project 1



Sources

- [BTY12] Bian, Jiang ; Topaloglu, Umit ; Yu, Fan: Towards Large-scale Twitter Mining for Drug- related Adverse Events. In: Proceedings of the 2012 International Workshop on Smart Health and Wellbeing. New York, NY, USA : ACM, 2012 (SHB '12). – ISBN 978–1–4503– 1712–2, 25–32
- [FKPT07] Fahrmeir, Ludwig ; Künstler, Rita ; Pigeot, Iris ; Tutz, Gerhard: Statistik. Springer- Verlag Berlin Heidelberg, 2007 (Springer-Lehrbuch). <http://books.google.de/books?id=ZinjP103iRcC>. – ISBN 9783540697398
- [FPSS96] Fayyad, Usama ; Piatetsky-Shapiro, Gregory ; Smyth, Padhraic: The KDD Pro- cess for Extracting Useful Knowledge from Volumes of Data. In: Commun. ACM 39 (1996), November, Nr. 11, 27–34. <http://dx.doi.org/10.1145/240455.240464>. – DOI 10.1145/240455.240464. – ISSN 0001–0782
- [Fri01] Friedman, Jerome H.: Greedy function approximation: a gradient boosting machine. In: Annals of Statistics (2001), S. 1189–1232
- [JMF99] Jain, A. K. ; Murty, M. N. ; Flynn, P. J.: Data Clustering: A Review. In: ACM Com- put. Surv. 31 (1999), September, Nr. 3, 264–323. <http://dx.doi.org/10.1145/331499.331504>. – DOI 10.1145/331499.331504. – ISSN 0360–0300
- [LW06] Lovell, Brian C. ; Walder, Christian J.: Support vector machines for business applica- tions. (2006)
- [MHC06] Maulik, U. ; Holder, L.B. ; Cook, D.J.: Advanced Methods for Knowledge Discovery from Complex Data. Springer, 2006 (Advanced Information and Knowledge Processing). http://books.google.de/books?id=OOOSx1X2-_sC. – ISBN 9781846282843
- [TLA12] Teng, Chun-Yuen ; Lin, Yu-Ru ; Adamic, Lada A.: Recipe Recommendation Using Ingredient Networks. In: Proceedings of the 3rd Annual ACM Web Science Conference. New York, NY, USA : ACM, 2012 (WebSci '12). – ISBN 978–1–4503–1228–8, 298–307