

Bayes statisztika és Bayes hálók az asztma/allergia kutatásban

Antal Péter

BME

Méréstechnika és Információs Rendszerek Tanszék
Mesterséges Intelligencia Csoport

Tartalom

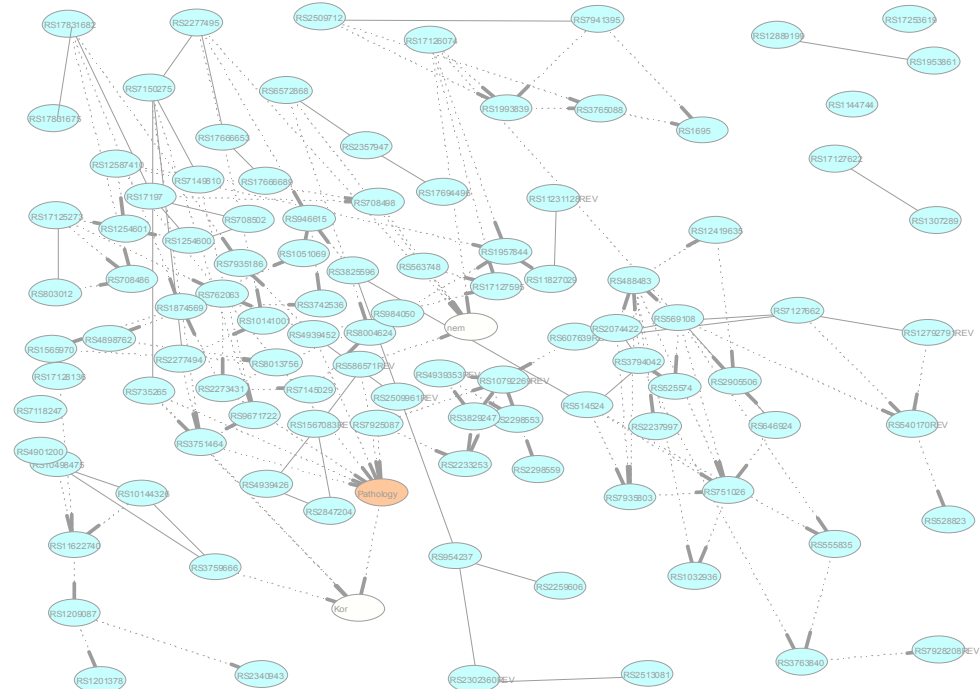
- Bayes hálók (grafikus modellek)
 - Kronológia
 - Az általános Bayes hálós modell
 - Értelmezési szintek
 - Bayes hálók tanulása
- Bayes statisztika
- Bayes hálók és Bayes statisztika orvosbiológiai kutatásban
 - Genetikai asszociációs kísérletek
 - Fúziós kihívások
- Eddigi alkalmazások asthma/allergia kutatásban

Bayes-*

- **Thomas Bayes (c. 1702 – April 17, 1761)**,
brit matematikus és presbiteriánus lelkész

- Bayes tétel
- Bayes statisztika
- Bayesi döntésmélet
- Bayes hálók
- Bayes faktor
- Bayes hiba
- Bayes rizikó
- Bayes döntés
- Bayes becslő
- ...

$$p(\text{Model} | \text{Data}) \propto p(\text{Data} | \text{Model}) p(\text{Model})$$
$$\hat{f} \approx \bar{f} = E_{p(\theta, S | \text{Data})} [f(\theta, S)]$$



Történeti áttekintés

- [1921]: S. Wright's diagrams, „Correlation and causation”
- [1970] The first (medical) applications (Dombal).
- [1979] The axiomatization of independencies (Dawid)
- [1982] The decomposition of a distribution (Lauritzen).
- [1988] Representation of independencies (Pearl).
- [1988] Inference methods in polytrees (Pearl).
- [1989] Exact general inference methods (Lauritzen).
- [1990] Closed form for the structure posterior (Cooper, Herskovits)
- [1995] Bayesian inference over BNs (Madigan)
- [1998] Bayesian inference about causation (Heckerman)
- [2000] Bayesian inference about relevance (Friedman)
- [2000] J.Pearl: Causality, Cambridge University Press
-

Bayes tétel

Egy algebrai trivialis

$$p(X | Y) = \frac{p(Y | X) p(X)}{p(Y)} = \frac{p(Y | X) p(X)}{\sum_x p(Y | X) p(X)}$$

Egy tudománytörténeti paradigmaváltás

$$p(\textit{Model} | \textit{Data}) \propto p(\textit{Data} | \textit{Model}) p(\textit{Model})$$

És pszichológiai konzekvenciák, megerősítés



$P(\text{betegség} | \text{tünet}) = ???$



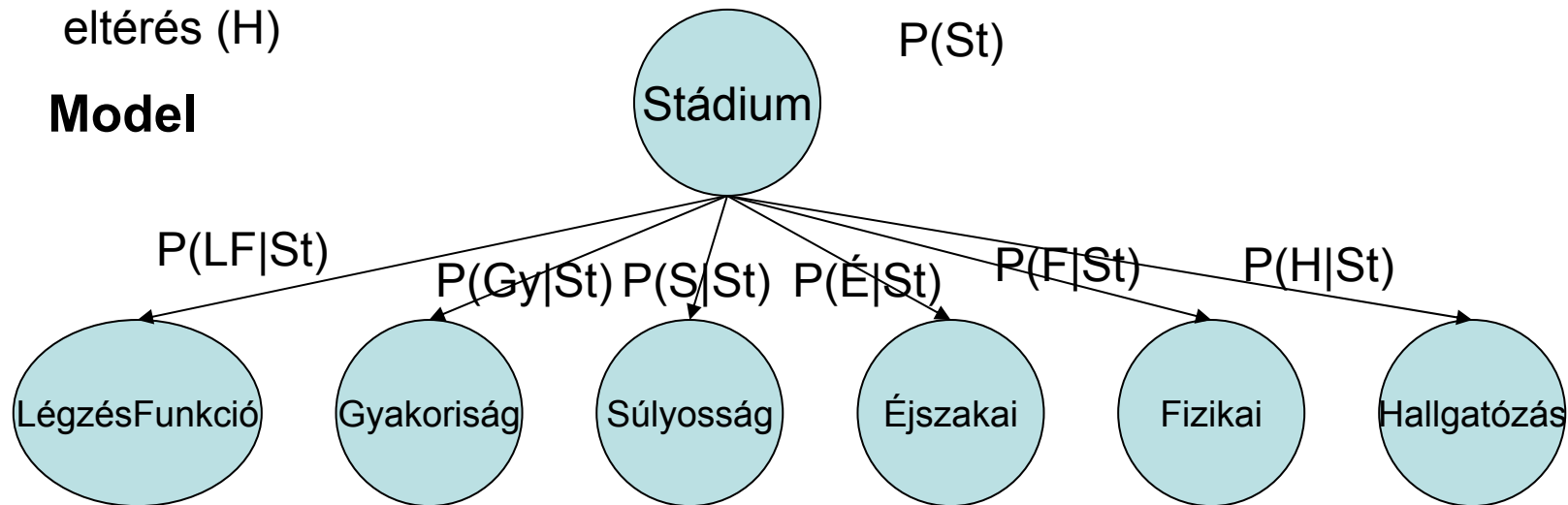
$p(\text{tünet} | \text{betegség}) = 0.79$
 $p(\text{betegség}) = 0.12$

Bayes tételtől a Bayes hálóig: Naív/idióta Bayes háló

Változók:

Asthma-stádium (St), Légzésfunkció(LF), TünetGyakoriság (Gy),
TünetSúlyosság (S), ÉjszakaiTünet (É), FizikaiTerhelés (F), Hallgatósági
eltérés (H)

Model

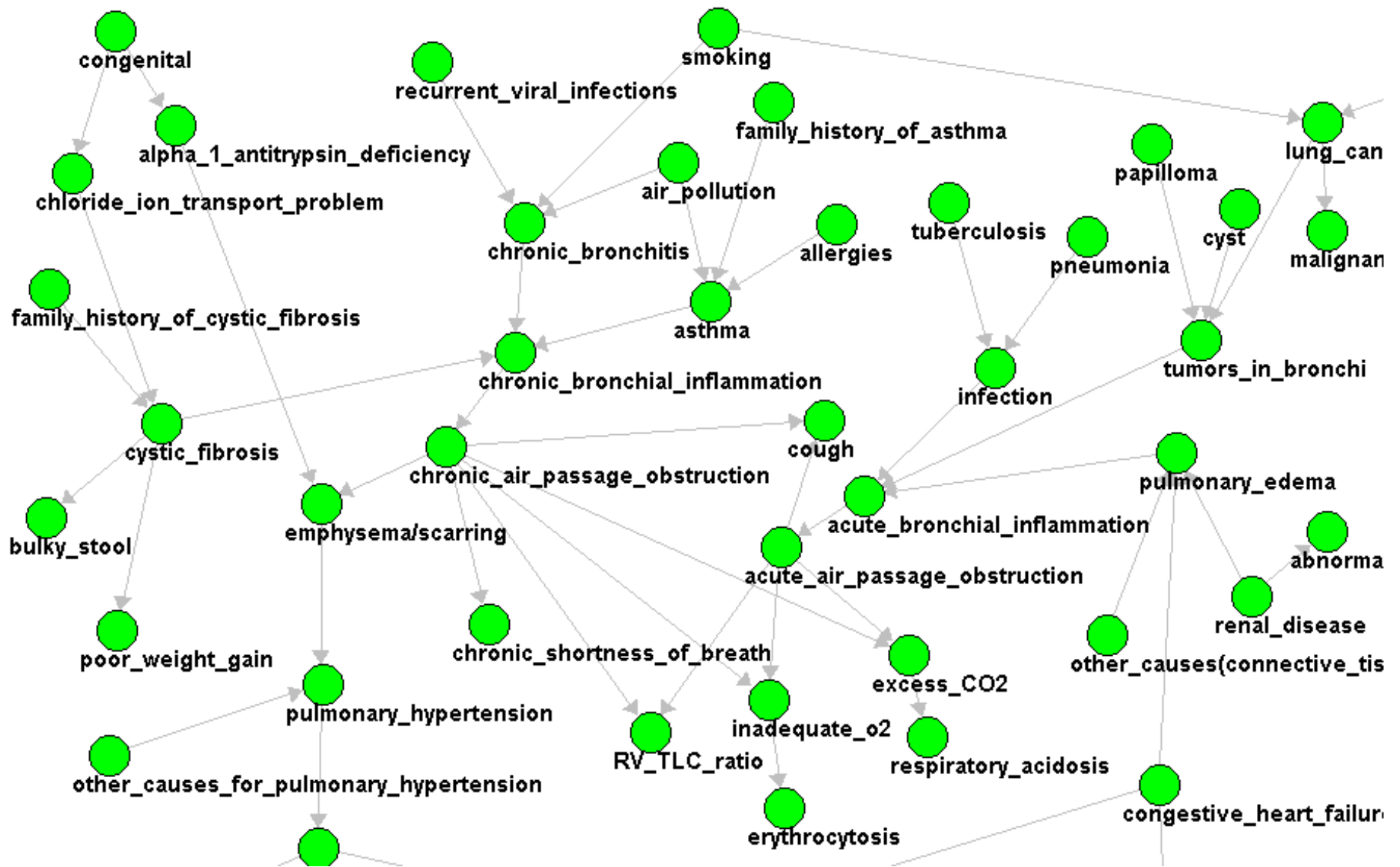


Diagnosztika

$$p(St | LF, Gy, S, É, F, H)$$

$$\propto p(St) p(LF | St) p(Gy | St) p(S | St) p(É | St) p(F | St) p(H | St)$$

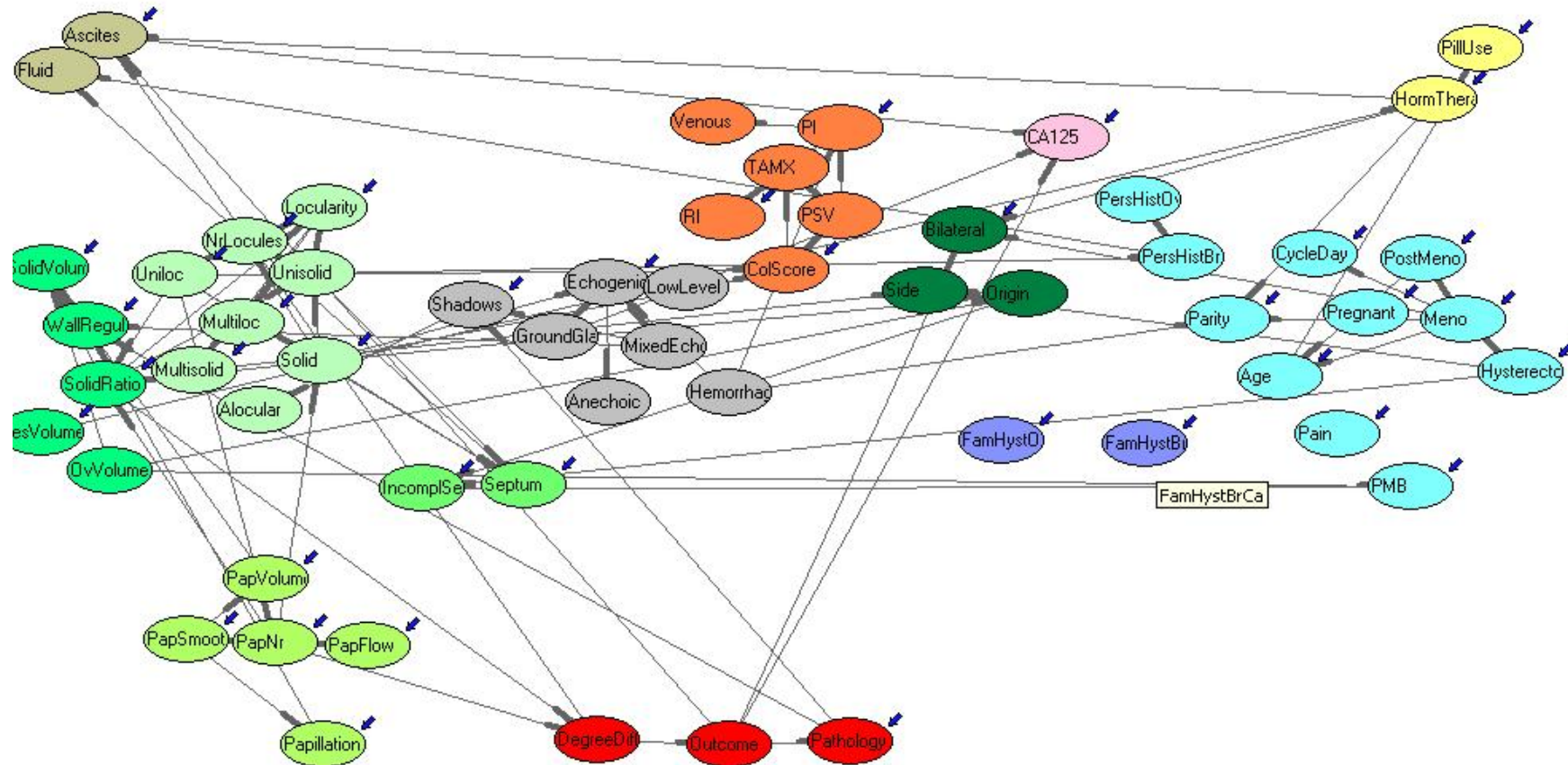
$$= p(St) \prod_i p(Tünet_i | St)$$



Pathfinder

- Pathfinder system. (Heckerman, Probabilistic Similarity Networks, MIT Press, Cambridge MA).
- • Diagnostic system for lymph-node diseases.
- • 60 diseases and 100 symptoms and test-results.
- • 14,000 probabilities.
- • Expert consulted to make net.
- • 8 hours to determine variables.
- • 35 hours for net topology.
- • 40 hours for probability table values.
- • Apparently, the experts found it quite easy to invent the causal links and probabilities.
- Pathfinder is now outperforming the world experts in diagnosis. Being extended to several dozen other medical domains.

És még egy model



Bayes háló értelmezési szintek

- Együttes eloszlás hatékony reprezentálása
 - Hatékony következtetés és tanulás
- Függetlenségi állítások reprezentálása
 - Direkt függések
- Oksági modell reprezentálása
 - „A csomópontok változók,
 - Az élek direkt, okozati relációk”

Eloszlások faktorizációja

- Általános (→ Markov hálók)
 - Potenciálok

$$p(X_1, \dots, X_n) \propto \prod_{j=1, \dots, L} \varphi_j(X_{j_1}, \dots, X_{j_{L_j}})$$

- Speciális (→ Bayes hálók)
 - Feltételes eloszlások

$$\begin{aligned} p(X_1, \dots, X_n) &= \prod_{i=1, \dots, n} p(X_i | X_1, \dots, X_{i-1}) \\ &= \prod_{i=1, \dots, n} p(X_i | Pa(X_i)) \end{aligned}$$

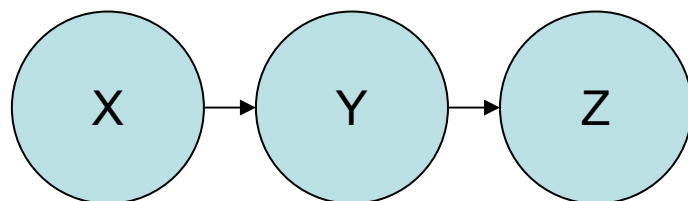
Eloszlások függetlenségeinek gráfos reprezentálása

„A valószínűségszámítás nem a számokról szól”

- Jelölje $I_p(X/Z/Y)$, hogy X, Y függetlenek Z feltétellel p eloszlásban
- Az M_p függetlenségi modell tartalmazza az összes p -beli függetlenségi állítást
- Kérdések
 - Lehet-e minden M_p -t egy gráffal reprezentálni? (nem)
 - Pontosán melyiket lehet? (nem ismert)
 - Mi a legjobb „közelítés”? (problémafüggő)

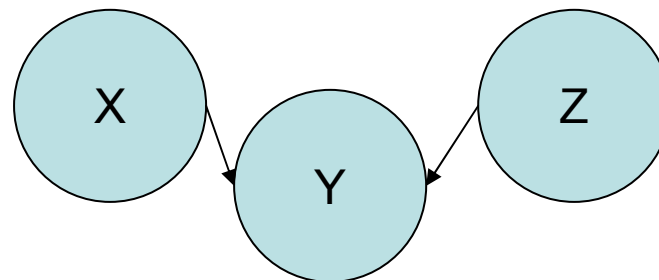
Gráffal nem reprezentálható eloszlások/függetlenségi modellek

- Intranszitivív függés (pl. X és Z lehet független)



- Alsóbbrendű függetlenségből nem következik magasabbrendű

(X,Z)-től együtt függhet Y, mégha páronként független is



A függetlenségi modell vs gráfok II.

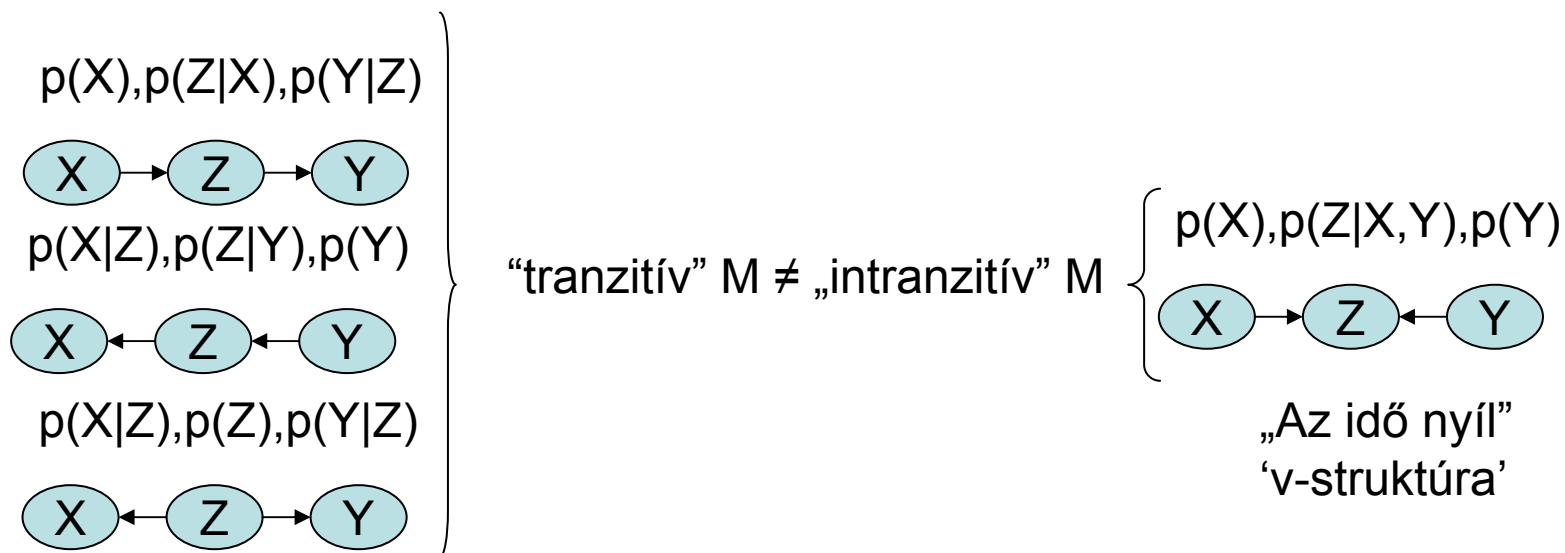
- Jó hírek:
 - Véges karakterizáció egzakt Markov hálós reprezentációra
 - D-elválasztottság helyes és teljes kalkulus
- Sejtések
 - Nincs véges karakterizáció egzakt Bayes hálós reprezentációra
- Rossz hírek
 - Léteznek egzakt módon nem reprezentálható eloszlások (sem Markov, sem Bayes hálókkal)
 - A reprezentáció redundáns → oksági értelmezés???

Principles of causality

- strong **association**,
- X precedes **temporally** Y,
- plausible **explanation** without alternative explanations,
- **necessity** (generally: if cause is removed, effect is decreased or actually: y would not have been occurred with that much probability if x had not been present),
- **sufficiency** (generally: if exposure to cause is increased, effect is increased or actually: y would have been occurred with larger probability if x had been present).

Függetlenségi vs. oksági modell

BAJ: egy függetlenségi modellhez több Bayes háló is tartozhat

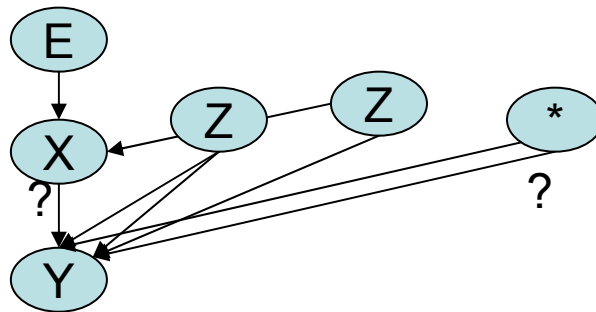


Oksági Bayes hálók

„Correlation \neq causation”, „no cause in, no cause out”



„Csak oksági kapcsolatok” feltevés: több változó esetén bizonyos élek egyértelműek
Több változó esetén bizonyos esetekben még rejtett változók is kizárhatók



KONKLÚZIÓ1: tabula rasa oksági következtetés lehetséges

- K2: oksági következtetések spektruma a tárgyterület független feltevések függvényében
 - Csak oksági relációk
 - Rejtett változók
 - Modell minimalitás
 - Változószám

„Az okozatiság matematizálása”
„a statisztikai idő iránya”

Feladatok Bayes hálóknál

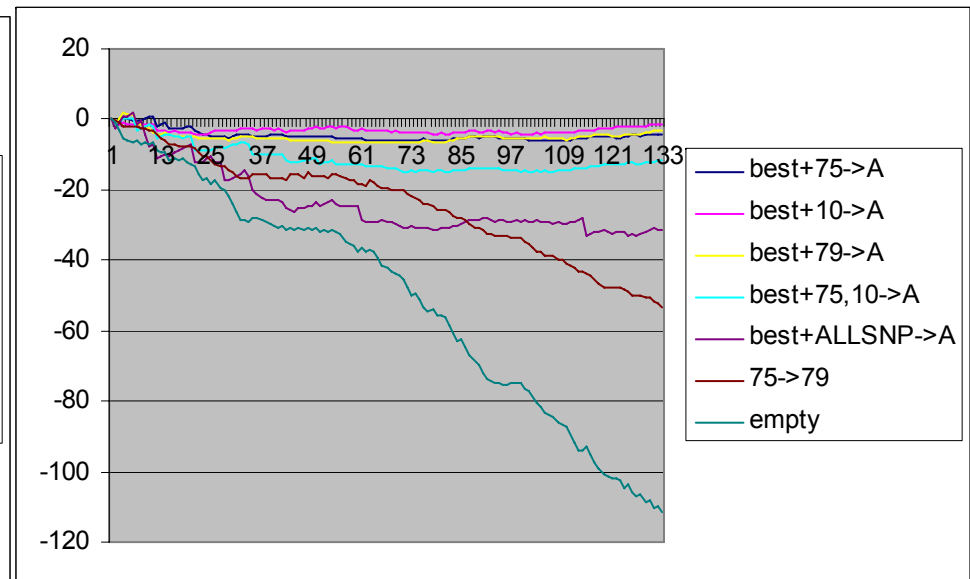
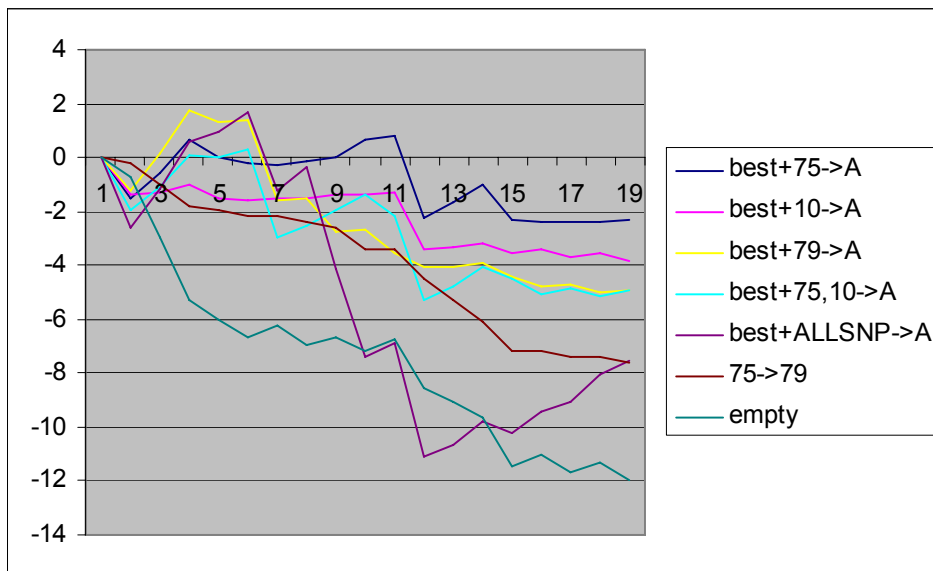
- (Passzív megfigyeléses) következtetés
 - $P(\text{Cél-változók}|\text{Megfigyelt-változók})$
 - Legvalószínűbb érték konfiguráció
- Beavatkozási következtetés
 - $P(\text{Cél-változók}|\text{Megfigyelt-változók}, \text{Beállított-változók})$
- Kontrafaktuális következtetés
 - $P(\text{Cél-változók}|\text{Megfigyelt-változók}, \text{Átállított-változók})$
- Döntési háló, hasznosság és beavatkozási változók
 - További információ értéke
 - Maximális hasznú döntés
- Paraméter tanulás
- Struktúra tanulás
 - → BNET 1995-

Bayes hálók tanulása

- Hipotézisteszteléssel
 - Függetlenségi tesztek diktálják a hálót
- Bayes statisztika esetén

$$p(\text{Model} \mid \text{Data}) \propto p(\text{Data} \mid \text{Model}) p(\text{Model})$$

Változók: SNP75, SNP10, SNP79, Asthma („best model”: 79<-75->10)



Bayes háló tulajdonságok

- Ötlet: hipotézisek, mint Bayes háló jegyek
- Ötlet²: (modellalapú) posterior a hipotézisekre

$$p(\text{ModelTulajdonság} \mid \text{Data})$$

$$= \sum_m p(m \mid \text{Data}) \mathbb{1}(\text{Tulajdonság} - \text{igaz} - m - \text{ben})$$

$$= E_{p(m|\text{Data})}[\mathbb{1}(\text{Tulajdonság} - \text{igaz} - m - \text{ben})]$$

- Génregularizációs kísérletben
 - Páronkénti oksági kapcsolatok: független, közös ok, oksági kapcsolat (direkt/indirekt)
- Asszociációs kísérletben
 - Relevancia???

Feature learning

Feature posterior

$$P(F = f) = \sum_{G:F^G = f} P(G)$$

- Goal: approximate the full-scale summation (integral)
- Practical methods: MCMC (Markov Chain Monte Carlo) sampling

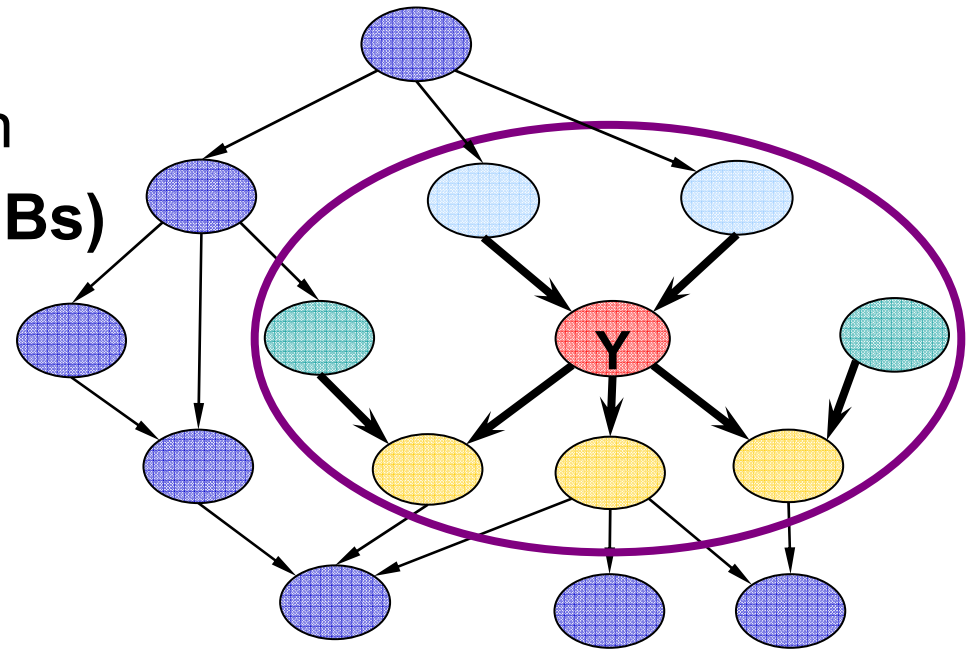
A kapcsolt BN jegyek relevanciára

– Markov Blanket (sub)Graphs (MBGs)

- (1) parents of the node
- (2) its children
- (3) parents of the children

Markov Blanket Sets (MBs)

- the set of nodes which probabilistically isolate the target from the rest of the model



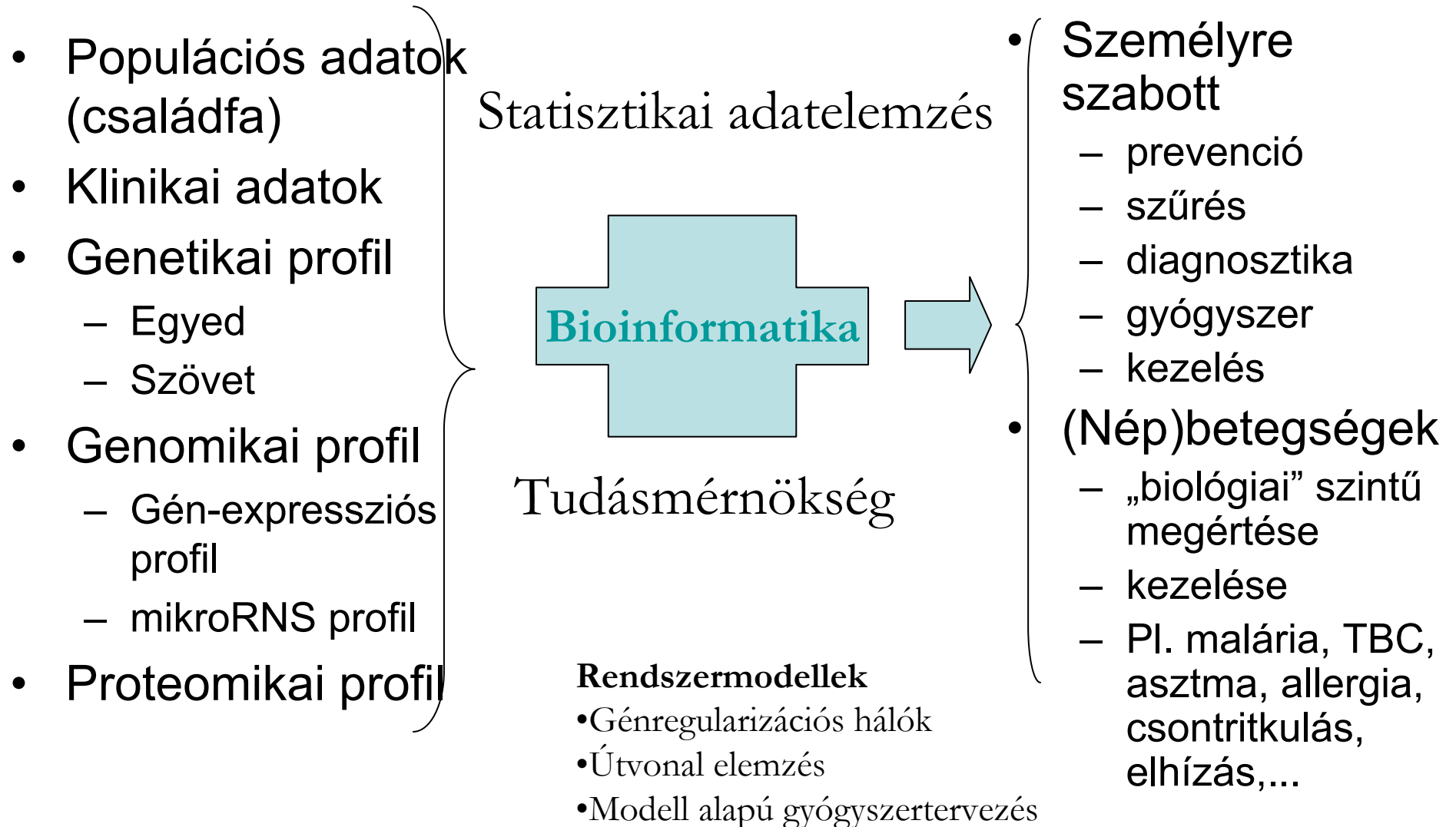
– Markov Blanket Membership (MBM)

- pairwise relationship

Kapcsolt, eltérő absztrakciós szintű hipotézisosztályok

- **Az asszociációs tanulmányok kérdései Bayes hálók strukturális jegyeivel formalizálhatók**
 - Páronkénti asszociáció: Markov Blanket Memberships (MBM)
 - Többváltozós/multifaktoriális asszociáció: Markov Blanket sets (MB)
 - Multifaktoriális asszociáció interakciókkal: Markov Blanket Subgraphs (MBG, C-RPDAG)
 - Teljes interakciós modell: Részlegesen irányított Bayes háló (PDAG)
 - Teljes oksági modell: Bayes háló (BN)
- **Kapcsolt, absztrakciós szinteket alkotnak ezek az asszociációs (feltételes) jegyek**
 - $DAG \Rightarrow PDAG \Rightarrow MBG \Rightarrow C-RPDAG \Rightarrow MB \Rightarrow MBM$

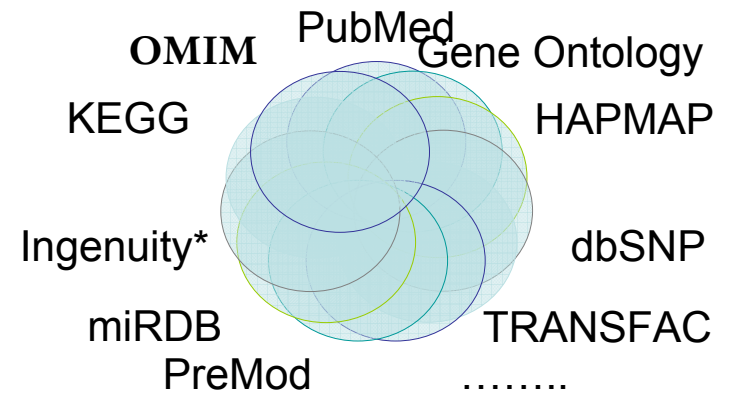
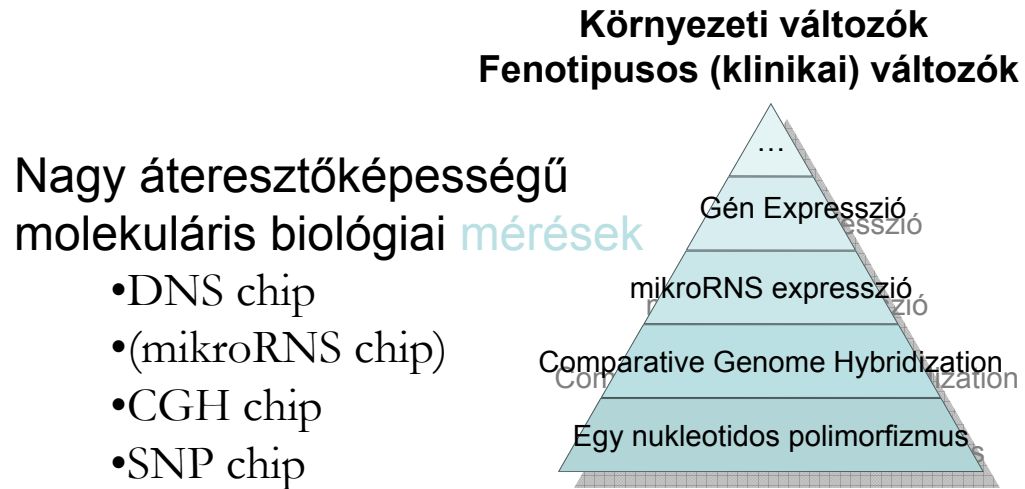
Adatoktól az orvosbiológiai ígéretekig



Bioinformatika: tudás és adat fúzió

Nagy dimenziójú **adatok**

Nagy mennyiségű, elektronikus **tudás**

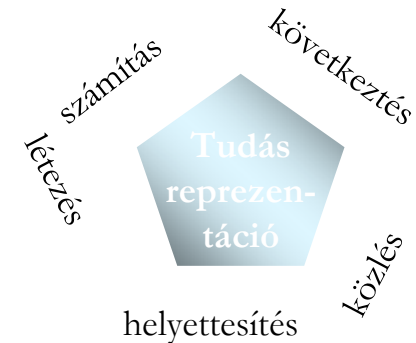


Komplex, hierarchikus, tudásgazdag hipotézisek

- Asszociációk
- Multifaktoriális interakciók
- Okozati relációk
- Alrendszerek, rendszerek

Tudásmérnöki trendek (fúzióhoz)

- Szabadszöveges tudástárak
 - Pl. Pubmed 15 millió absztrakt
 - Cikk állítása formálisan
- Szemi-formális tudásbázisok
- Elektronikus ontológiák, taxonómiák
 - UMLS, MeSH, Gene Ontology, MGED ontológia, SNP ontológia
- Formális tudásbázisok
- **Valószínűségi tudásmodellek/logikák**
- Szemantikus világháló
- Annotált, integrált adatbázisok
- Annotált, integrált web-szolgáltatások



Tudásmérnöki feladatok:

- reprezentáció,
- kinyerés,
- következtetés,
- fenntartás

Enciklopédisták (D.Diderot)

Logikai pozitivizmus/empirizmus

(R. Carnap)

“World Brain” (H.G.Wells,1938)



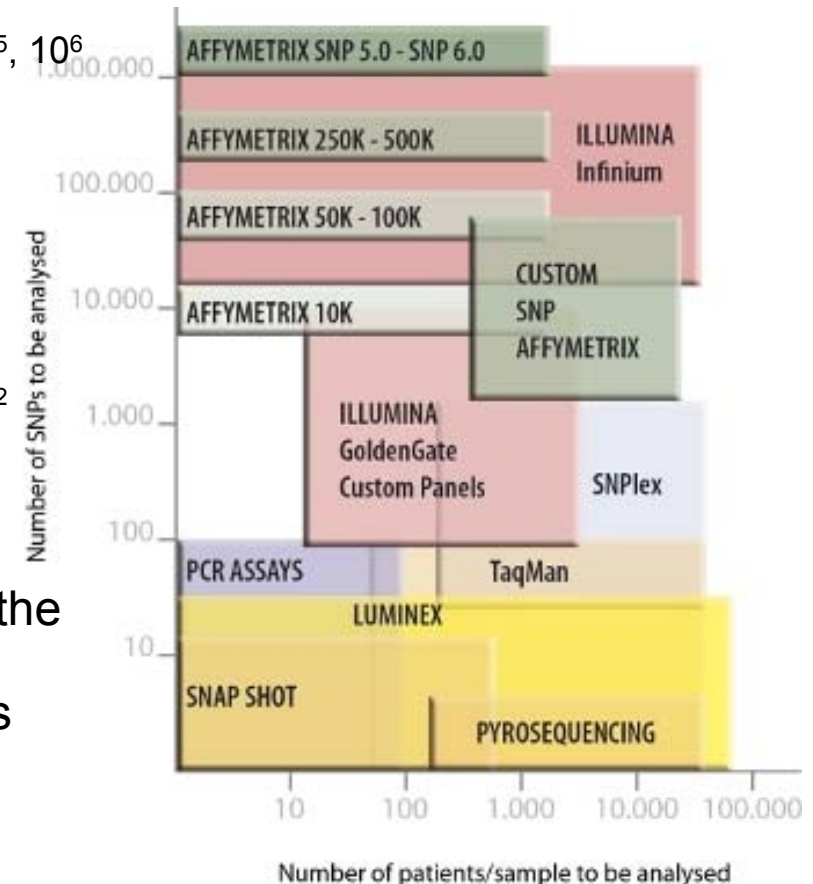
Elektronikus, hatalmas mennyiségű, egységesített, valószínűségi, empirikusan lehorgonyozott, formálisan reprezentált tudás

Indukciós trendek (fúzióhoz)

- Statisztikai, számításelméleti tanulási elméletek
 - Nem-parametrikus konzisztencia eredmények
 - Nem-aszimptotikus konvergencia eredmények
 - ...
 - Bayes statisztika
 - Markov Lánc Monte Carlo módszerek (MCMC)
 - Hibrid/Hierarchikus/Kapcsolt/... MCMC-ek
 - ..
 - Oksági következtetés
 - Oksági kapcsolatok indukciója passzív megfigyelésekből
 - Oksági kapcsolatok formalizálása (kauzális Bayes hálókkal/strukturális egyenletekkel)
 - Megfigyelések és beavatkozások tervezhetősége (kiváltása)
 - (Beavatkozások hatásának valószínűségi modellezése)
 - Kontrafaktuális következtetések)
- W.Ockham
D.Hume
K.R.Popper
- T.Bayes
P.Laplace
- D.Hume
H.Reichenbach

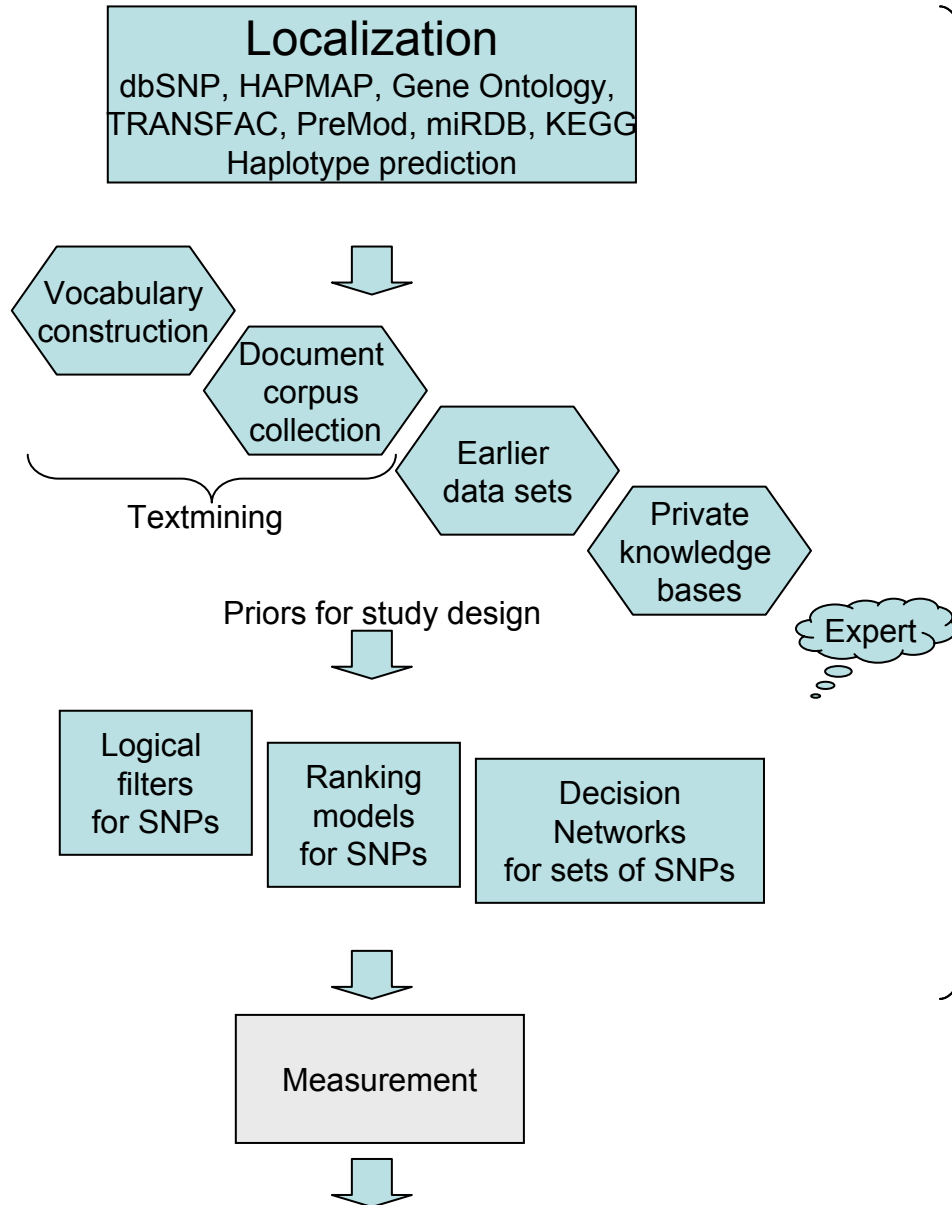
Genetic Association Studies (GAS): a statistical view

- Genome-wide association studies (GWASs)
 - Sample size: ~1000
 - |Response/outcome (target) variables|: ~10
 - Predictor/explanatory:
 - |SNPs| (nominal/numeric, unphased genotype): $\sim 10^5, 10^6$
 - |Environmental variables|: ~10
 - Cost: 1000x1000 Euro
- Partial genome screening studies (PGSSs)
 - Sample size: ~1000
 - |Response/outcome (target) variables|: ~10
 - Predictor/explanatory:
 - |SNPs| (nominal/numeric, unphased genotype): $\sim 10^2$
 - |Environmental variables|: ~10
 - Cost: 10000 Euro
- The goal is the exploration of...
 - the **relevance** of explanatory variables w.r.t. the target variables
 - the **interdependence** of explanatory variables

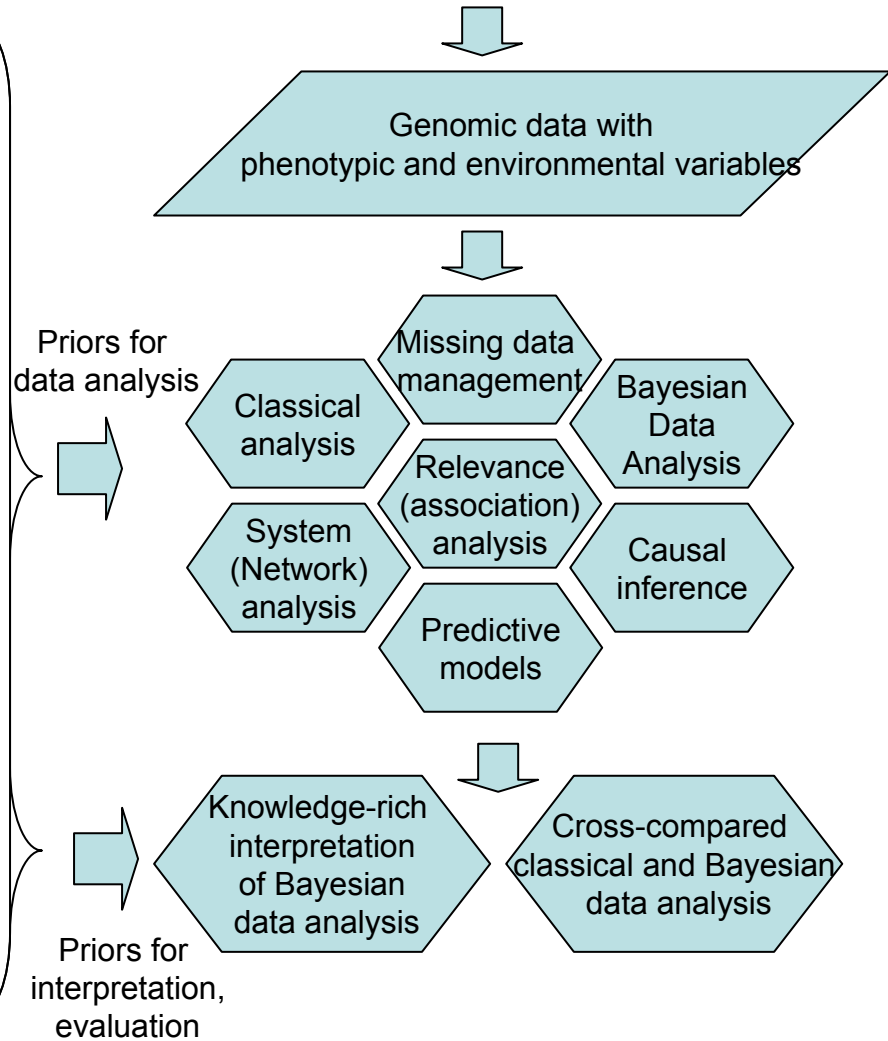


Overview (SNP)

Study design

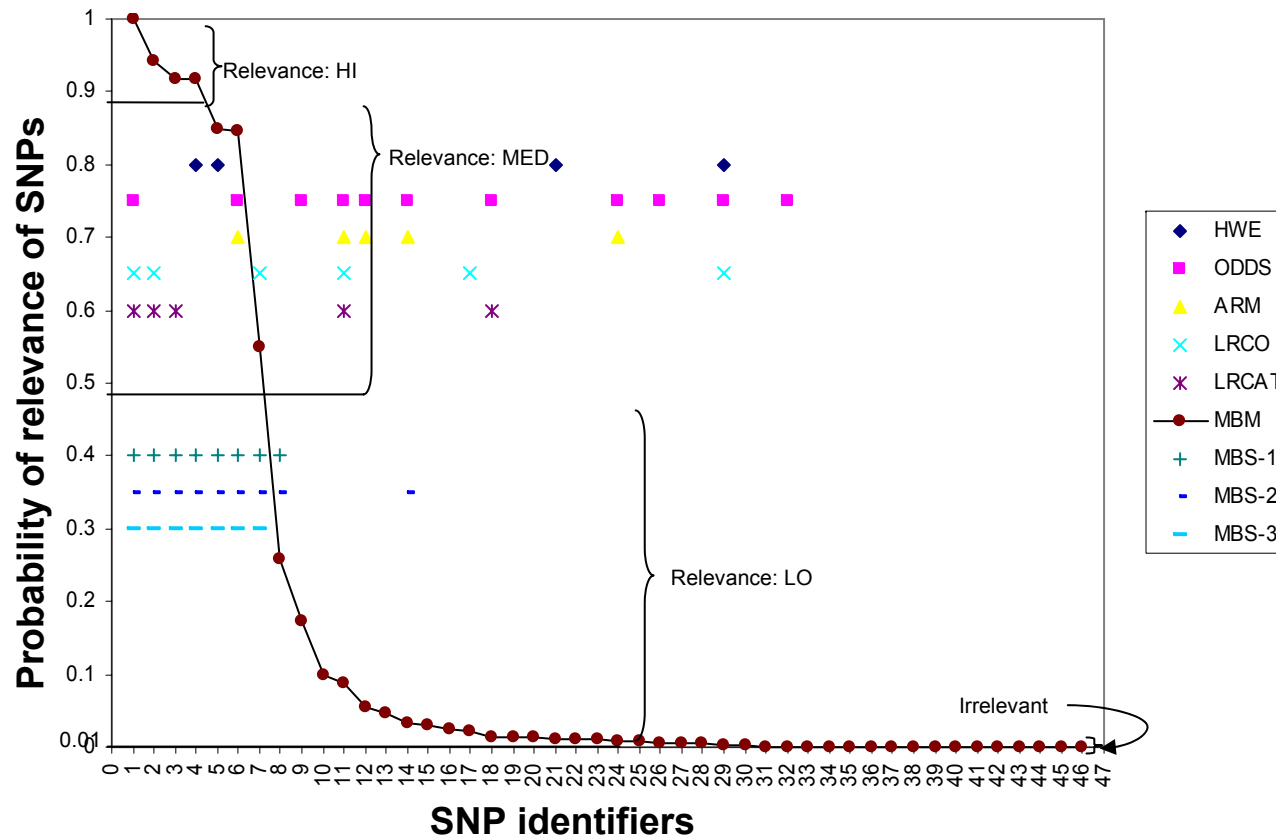


Data analysis



An overview of results

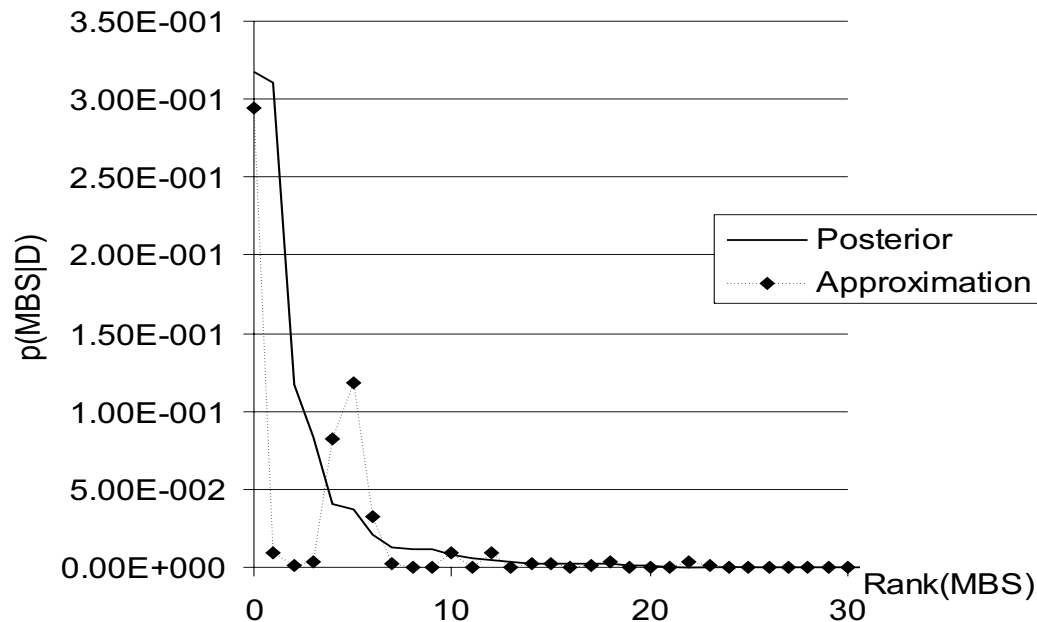
HWE – Hardy-Weinberg equilibrium test, *ODDS* – odds ratio, *ARM* – Cochran-Armitage trend test, *LRCO* – logistic regression (continuous case), *LRCAT* – logistic regression (categorical case), *MBM* – Bayesian pairwise relevance, *MBS-1–9* relevant sets by Bayesian analysis. (only MBM values are numeric, others are arbitrary values for visualization)



Uncertainty/sample complexity in multivariate analysis?

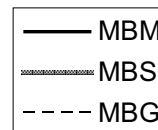
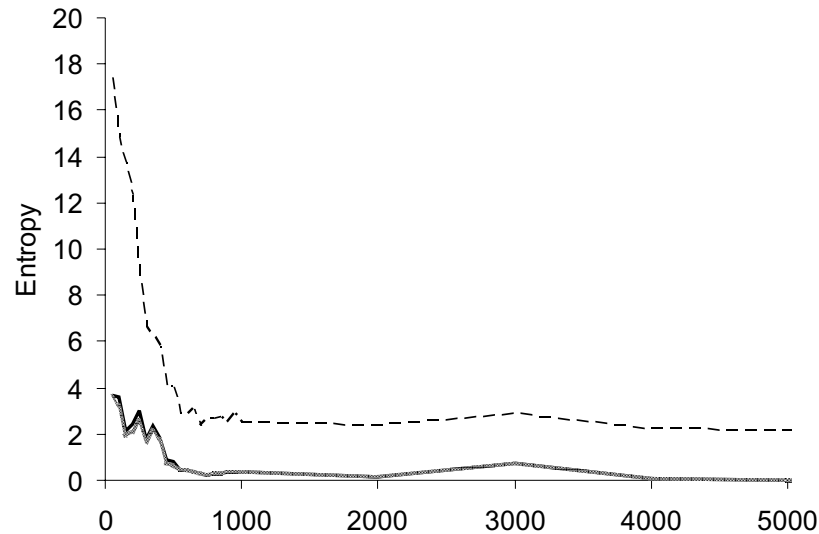
Lessons from the applications of the BMLA

- High-level of uncertainty in multivariate analysis
- There are stable sub-parts (e.g., subset, subgraphs)
- Results for target variables and for certain SNPs could be aggregated

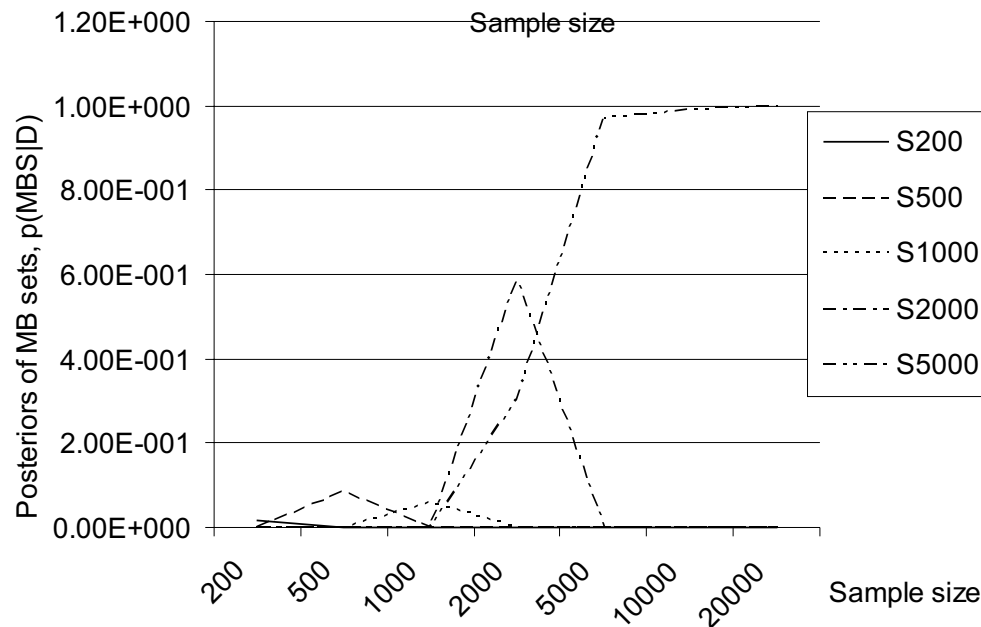


The peakness of the posteriors of the most probable MB sets and their MBM-based approximations. (46 variables, 1000 samples)

“Sample complexity” for FSS



The summed entropies of the $p(\text{MBM}(Y, X_i, G)|D)$ posteriors, and the entropy of the MBS and MBG posteriors for the artificial model (with 7+1 variables).



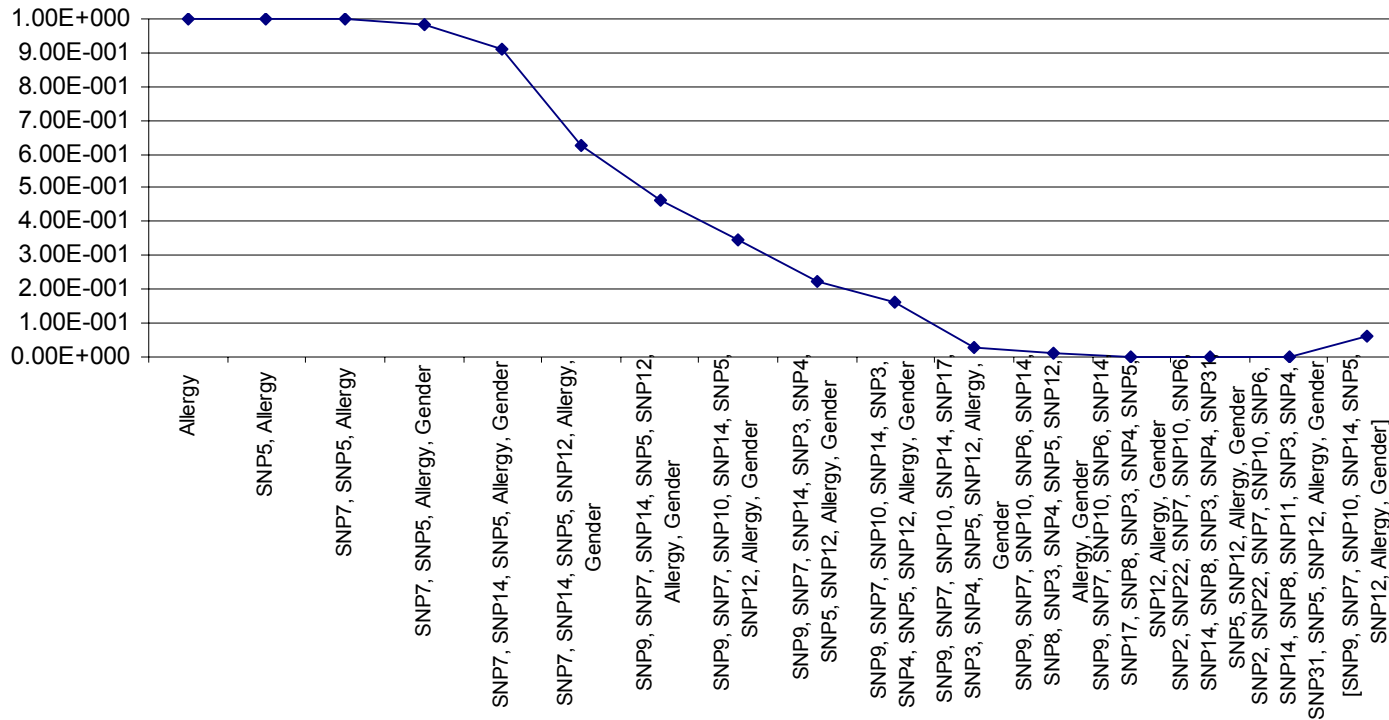
The sequential posteriors of relevant MB sets for the model with 50 variables.

Scalable multivariate analysis

Problem: the “polynomial” gap between simple and complex features
(e.g., MBM (n^2) and MBS (2^n))

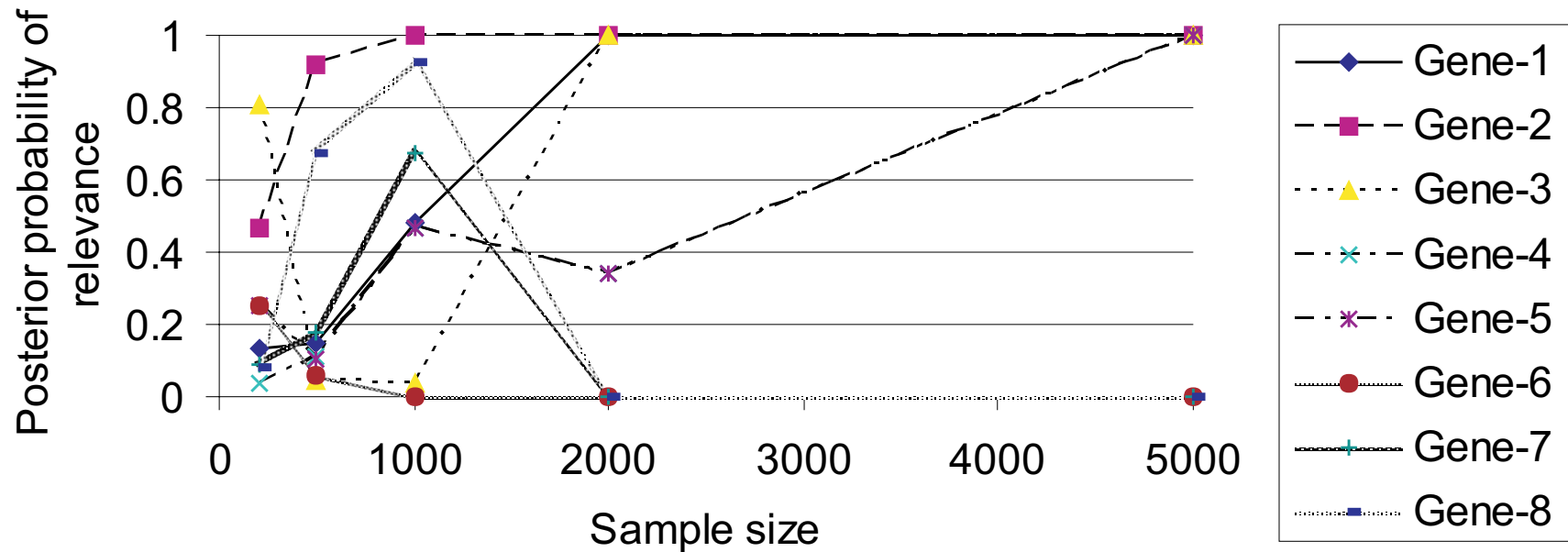
Idea: If all X_i in set S with size k are members of a Markov Boundary set, then S is called a k -ary Markov Boundary subset ($O(n^k)$).

$$p(G : X \subseteq MBS(G) \mid D_N)$$



The maximal posteriors for k -relevance/ k -MBS for Asthma.

SNP-to-gene aggregation



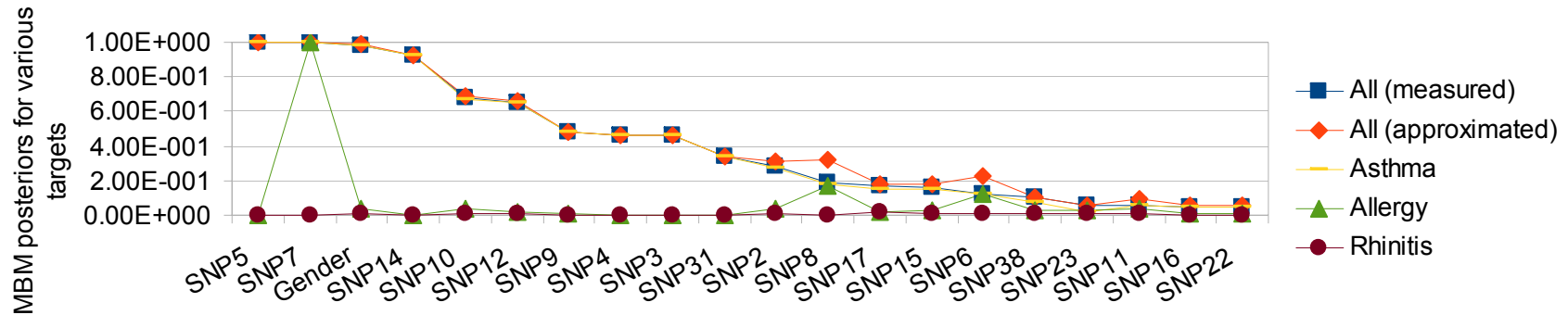
The sequential posteriors that a given gene contains a SNP relevant for asthma

Abstraction levels: SNP, haplo-block, gene,..., pathway

Note that it is different from aggregated multi-variables.

Joint analysis for multiple targets

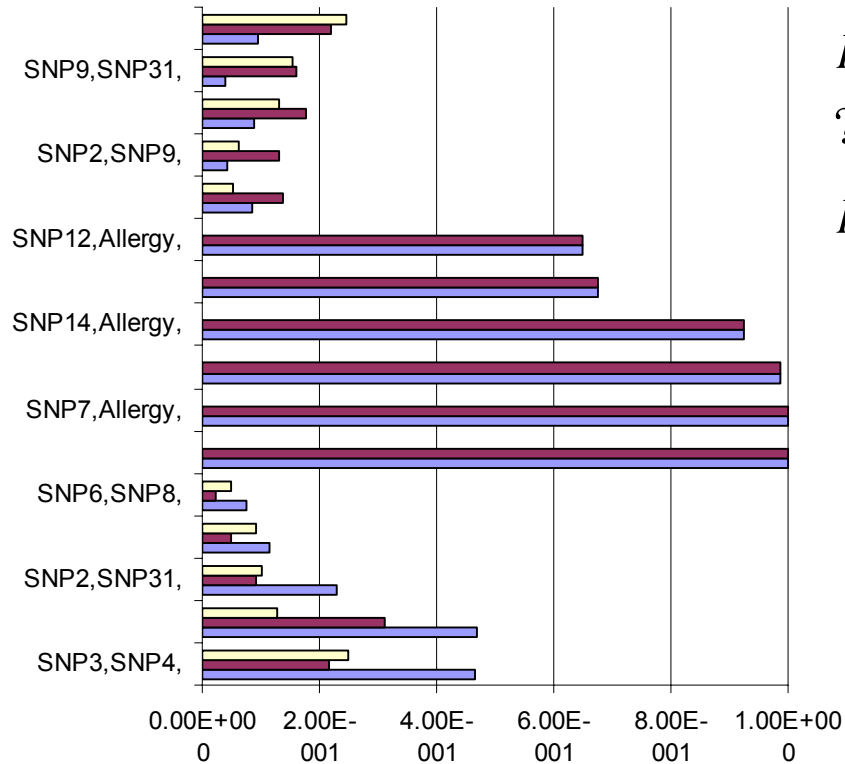
Targets (Outcome): Asthma, Allergy, Rhinitis



The MBM posteriors that a SNP is relevant for asthma, allergy and rhinitis, both separately and jointly.

Note that it cannot be done off-line (contrary to k-MBS/k-MBG and feature aggregation).

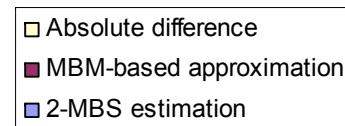
Redundancy and interaction



$$p(G : X_i, X_j \subseteq MBS(G) | D_N)$$

?

$$p(G : X_i \in BS(G) | D_N) p(G : X_j \in MBS(G) | D_N)$$

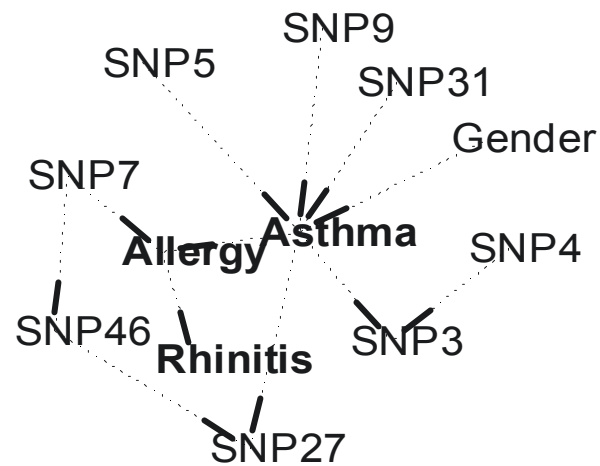


The quantification of interaction and redundancy by the decomposability of the posterior.

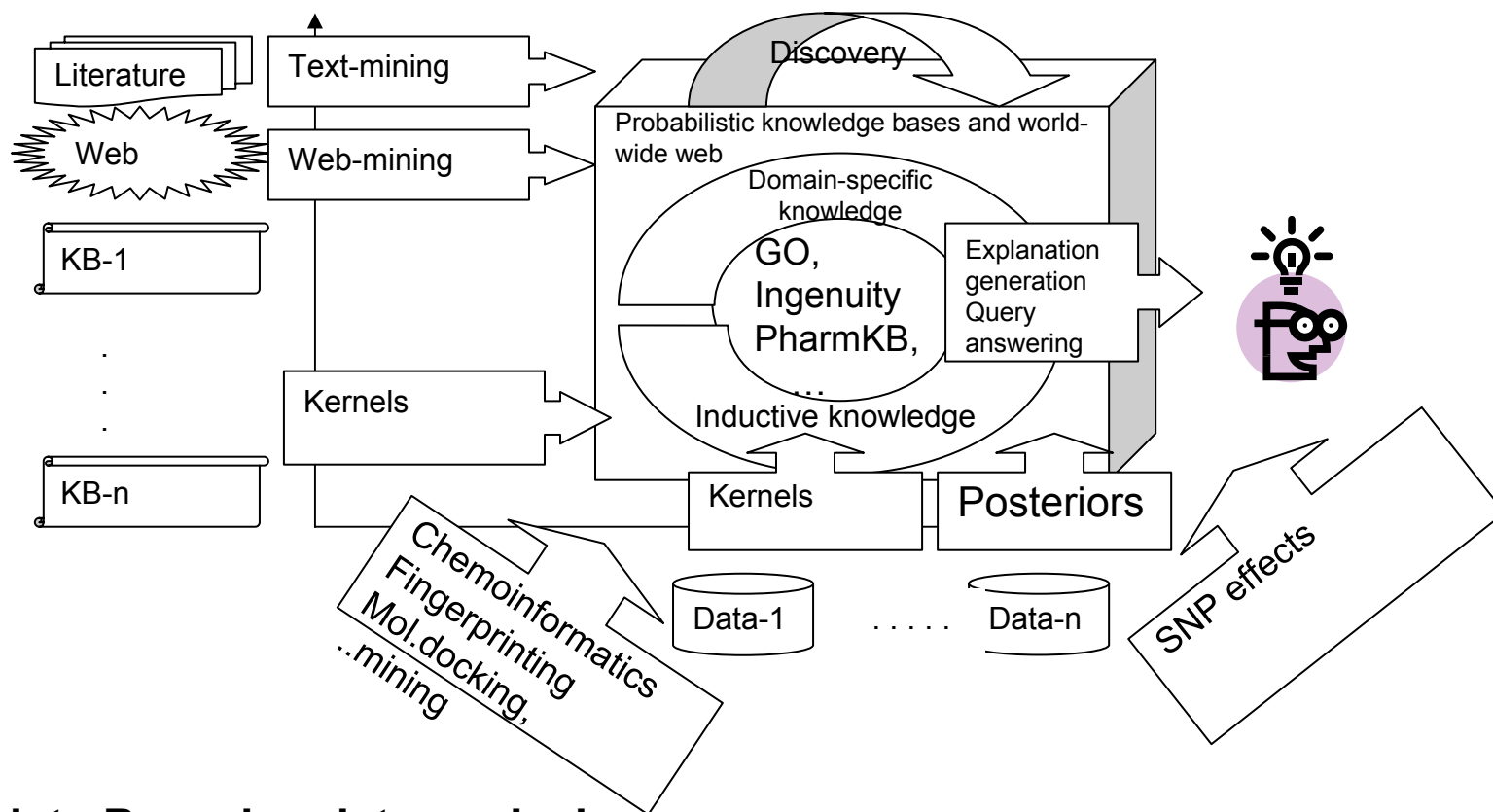
Single measure (instead of many Bayes factors).

Conditional (and contextual) relevance

The maximum a posteriori Markov Blanket subgraph



Fusion using probabilistic logic



- **multivariate Bayesian data analysis,**
- the use of more powerful **logic for representing knowledge,**
- **uncertainty management** in knowledge representation, induction, and inference, syntactic semantics:
$$P(\phi | KB) = \sum_{pKB} P(\phi \text{ is provable} | pKB, KB, kKB) p(pKB)$$

Diákok, kollégák

- Végzett diákok
 - **Ercsényi** Dániel: Bayes hálók strukturális tulajdonságainak vizsgálata MCMC módszerekkel (2005)
 - **Gézi** András: Monte Carlo módszerek Bayes hálók komplex strukturális jegyeinek következtetésére (2006)
 - **Barna** Gergely: Bayes-hálók dekomponált tanulása a priori tudás felhasználásával (2005)
 - **Szabadka** Zoltán (Grolmusz Vince): Fehérje-kismolekula kölcsönhatások vizsgálata molekuláris dokkolással (2006)
 - **Iván** Gábor (Grolmusz Vince): Protein-ligand komplexek kötőhelyeinek vizsgálata adatbányászati algoritmusokkal (2006)
 - **Orbán** György: Kontextuális függések szerepe Bayes hálók tanulásában (2006)
 - **Rohla** Tibor: Monte Carlo módszerek Bayes hálók egyszerű strukturális jegyeinek következtetésére (2007)
 - **Palotai** Robin (Csermely Péter): Modulkeresés, klaszterezés fehérje-fehérje hálózatokban
 - **Temesi** Gergely (Falus András): Információ menedzsment és intelligens adatelemzés az SNP analízis támogatására
 - **Gergely** Balázs: Statisztikus szövegbányászat
- Jelenlegi diplomázó, orvosbio. másoddiplomázó, és doktorandusz diákok
 - **Millinghoffer** András: A first-order probabilistic logic for fusing data and prior knowledge
 - **Hullám** Gábor: Valószínűségi tudásmodellek tanulása
 - **Hajós** Gergely: Döntéstámogató keretrendszer SNP asszociációs kísérletek tervezéséhez
 - **Kecskeméthy** Péter (Oxford): Bayesi relevancia elemzés feltételes modellekkel
 - **Huszár** Ferenc: nem-parametrikus Bayesi modellek a kognitív pszichológiában
 - **Erdélyi Boróka**: Bayesi logisztikus regresszió
 - **Szabados** Péter: orvosbiológiai tudásbázisok
 - **Arany Ádám**: Bayesi térszerkezet predikció

Köszönöm a figyelmet!