

Hierarchical Bayesian Modelling Identifies Shared Gene Function

Peter Sykacek^{1,2} & Gos Micklem^{1,3}

Departments of Genetics¹, Pathology² & CCBI³

University of Cambridge

ps408@cam.ac.uk

Data & Biology:

Richard Clarkson, Cris Print

Methodological Discussions:

David J. C. MacKay & Inference Group

Problem Statement

- Assumption: Several microarray experiments are obtained such that slides can be mapped to a biological state of interest.
- Shared genetic function: Interesting genes are **across experiments** informative about these biological states.
- Task: find those genes! Actually two problems:
 - Cross annotation of genes (potentially different species)
 - Calculate a measure across experiments

This talk shows how we may obtain such a measure using a probabilistic approach.

Biological States of Experiments

Mammary Gland tc. (lact. day & hours of involution)

biol. state	L ₀	L ₅	L ₁₀	I ₁₂	I ₂₄	I ₄₈	I ₇₂	I ₉₆
Type 1 Apoptosis	-	-	-	+	+	?	-	-
Type 2 Apoptosis	-	-	-	-	-	?	+	+
Apoptosis	-	-	-	+	+	+	+	+
Differentiation	+	+	+	?	-	-	-	-
Inflammation	?	-	-	+	+	?	-	-
Remodelling	-(?)	-	-	-	-	?	+	+
Acute Phase	+	-	-	-	+	+	+	+

Serum Deprived Apoptosis (duration in hours)

biol. state	t ₀	t ₂₈	t ₄₈
Type 2 Apoptosis	-	+	+
Apoptosis	-	+	+
Differentiation	+	-	-

Probabilistic Approach



Thomas Bayes (1701 - 1763)
Learning from data based on a
decision theoretic framework

Probabilistic Approach



Thomas Bayes (1701 - 1763)
Learning from data based on a
decision theoretic framework

$$p(I|\mathcal{D}) = \frac{p(\mathcal{D}|I)p(I)}{p(\mathcal{D})}$$

First consequence: we
must revise beliefs ac-
cording to Bayes theorem

Probabilistic Approach



Thomas Bayes (1701 - 1763)
Learning from data based on a
decision theoretic framework

$$p(I|\mathcal{D}) = \frac{p(\mathcal{D}|I)p(I)}{p(\mathcal{D})}$$

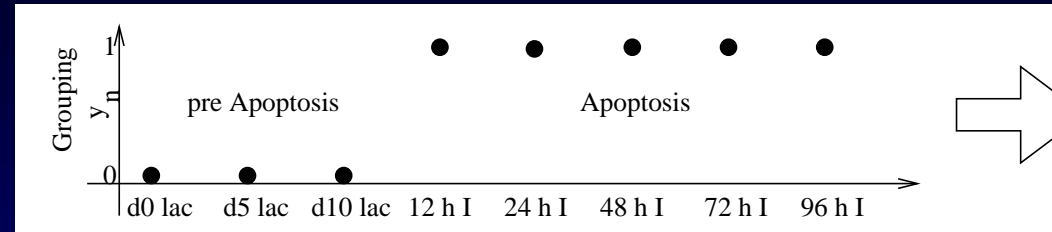
First consequence: we must revise beliefs according to Bayes theorem

$$\alpha_{opt} = \operatorname{argmax}_{\alpha} \langle u(\alpha) \rangle, \text{ where } \langle u(\alpha) \rangle = \int_G u(\alpha, I)p(I|\mathcal{D})dI.$$

Second consequence: Decisions by maximising expected utilities

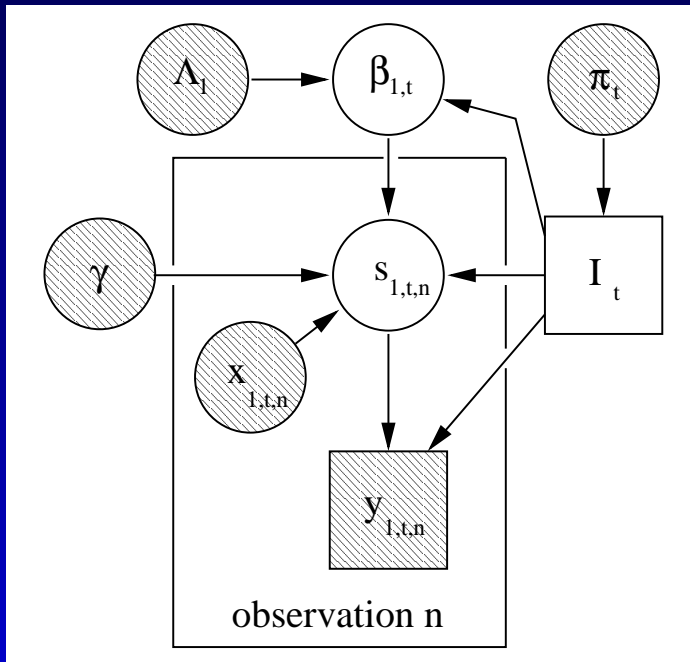
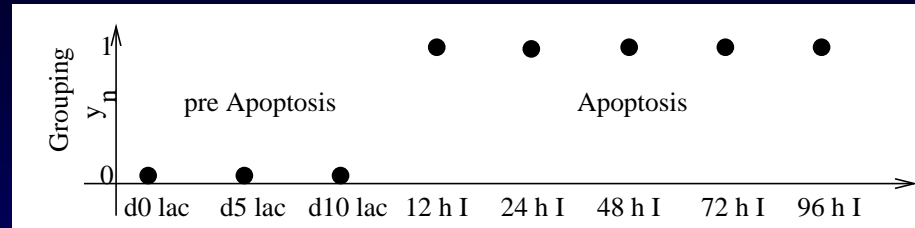
Probabilistic Gene Ranking

Apoptosis (lac. vs. inv.!) in the Mouse Mammary Gland



Probabilistic Gene Ranking

Apoptosis (lac. vs. inv.!) in the Mouse Mammary Gland



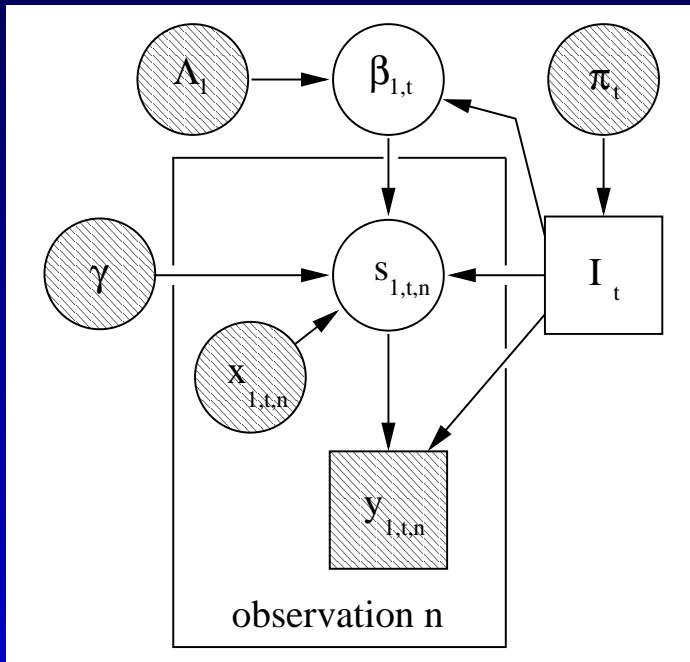
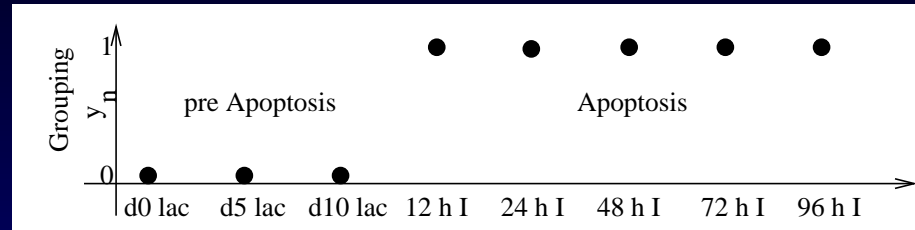
Latent variable probit GLM.

$$\text{if } I_t = \begin{cases} 1 : s_{1,t,n} \sim 1 + x_{t,n} \\ 0 : s_{1,t,n} \sim 1 \end{cases}$$

$s_{1,t,n}$ is a one dimensional Gaussian random variable with mean $\beta_{t,1}^T x_{t,n}$ and precision γ .

Probabilistic Gene Ranking

Apoptosis (lac. vs. inv.!) in the Mouse Mammary Gland



Latent variable probit GLM.

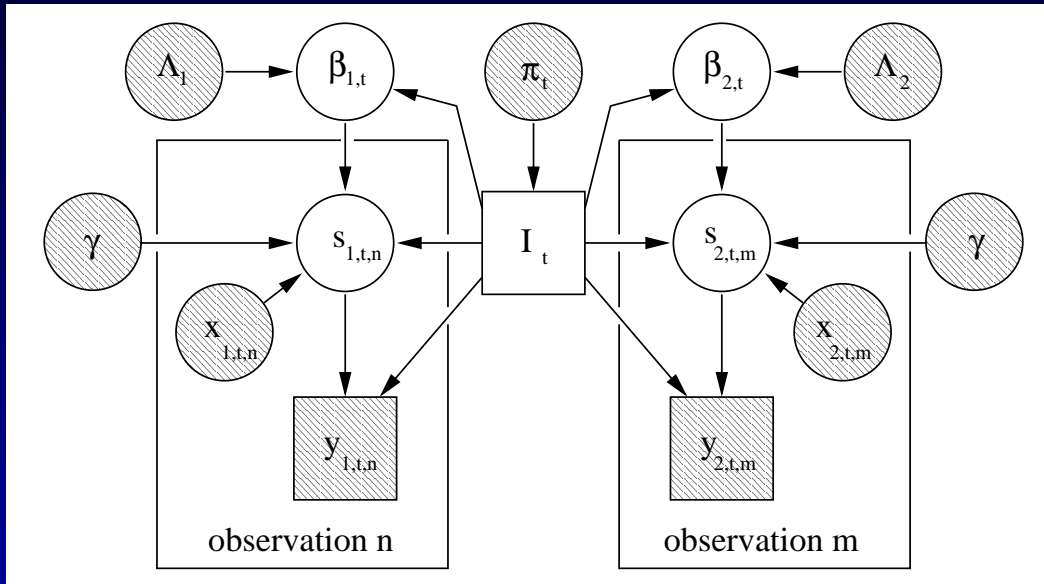
$$\text{if } I_t = \begin{cases} 1 : s_{1,t,n} \sim 1 + x_{t,n} \\ 0 : s_{1,t,n} \sim 1 \end{cases}$$

$s_{1,t,n}$ is a one dimensional Gaussian random variable with mean $\beta_{t,1}^T x_{t,n}$ and precision γ .

As an alternative to p-values, the posterior $P(I_t | \mathcal{D}_1)$, serves as a probabilistic rank measure. (VB-eqns.)

Shared Gene Function

Include Information about Endothelial Cell Death

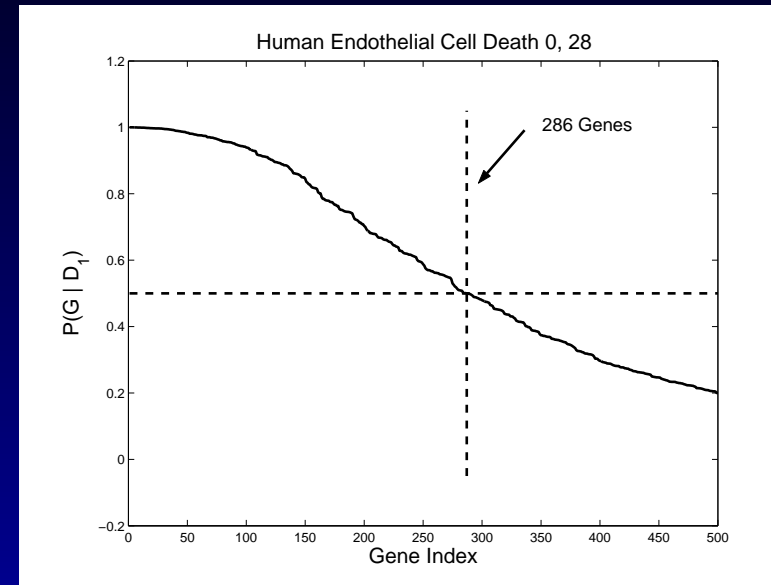
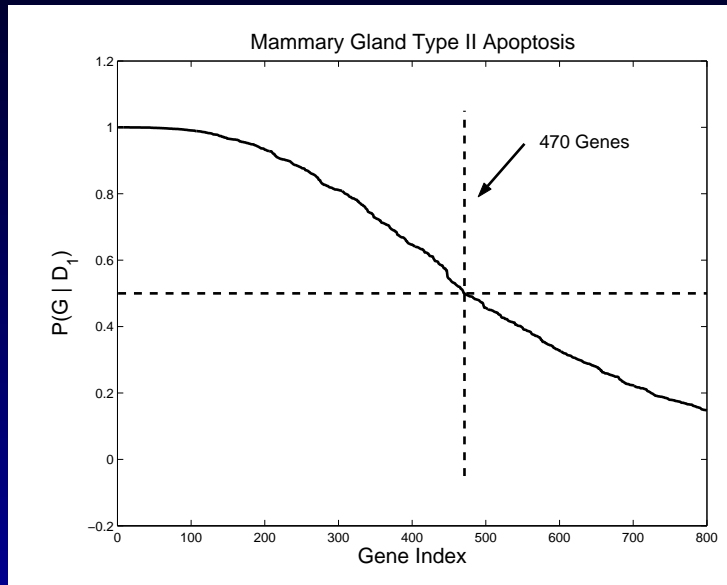


Model 0 hrs. vs. 28 hrs. as latent variable probit GLM. Calculate $P(\mathcal{D}_2|I_t)$, the marginal likelihood.

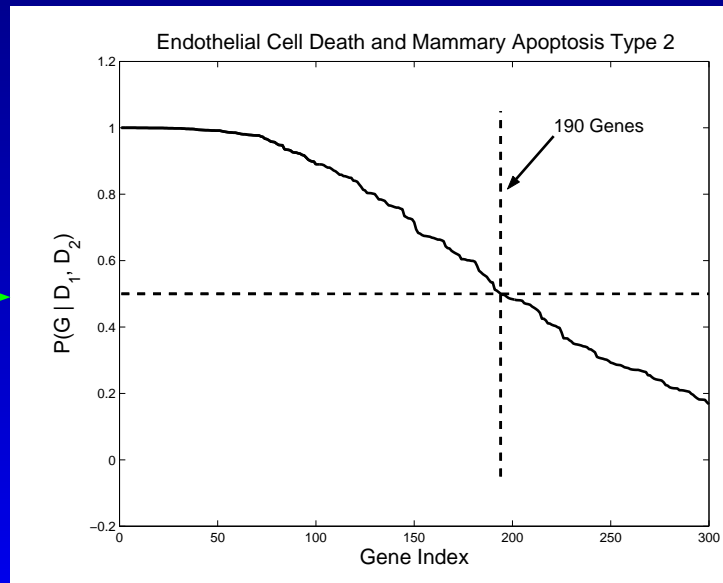
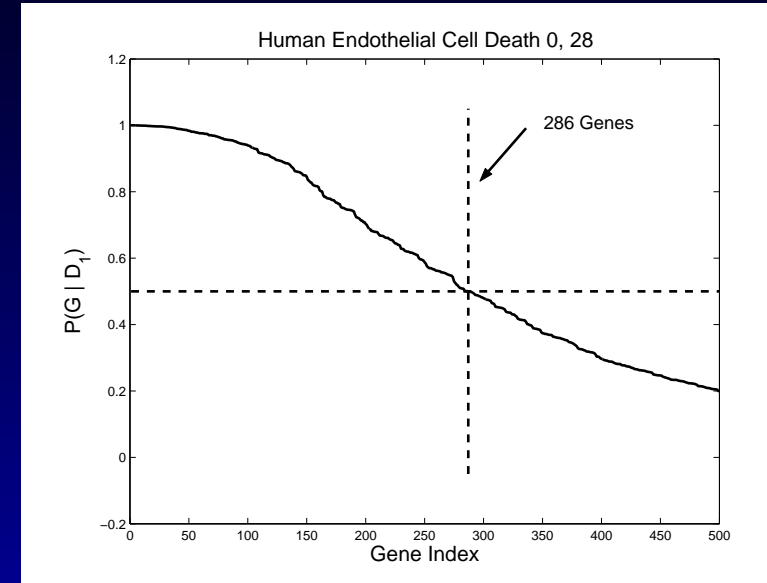
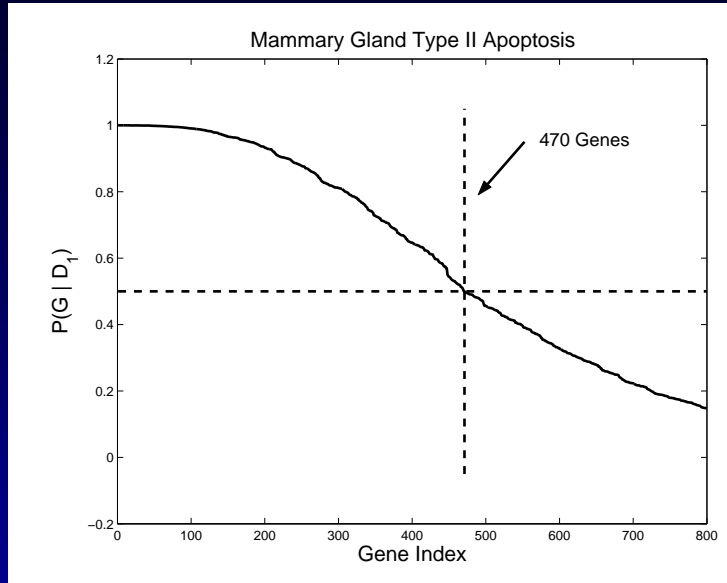
Bayes theorem gives a *principled* measure for ranking

$$P(I_t|\mathcal{D}_1, \mathcal{D}_2) = \frac{P(I_t|\mathcal{D}_1)p(\mathcal{D}_2|I_t)}{p(\mathcal{D}_2|\mathcal{D}_1)}$$

Don't Do That at Home!



Don't Do That at Home!



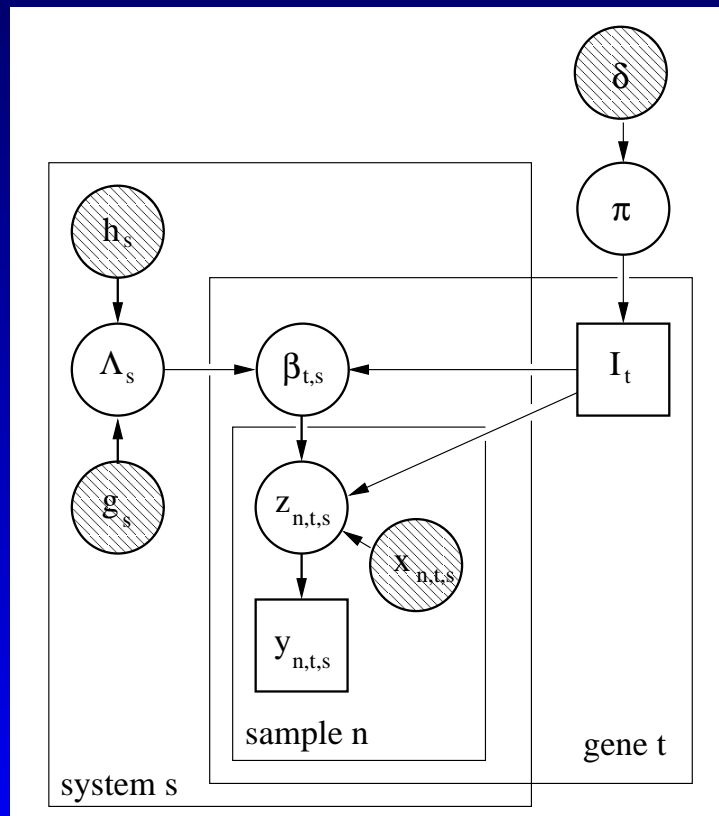
Improving on Previous Model

- Hyper parameters (π_t , Λ_1 and Λ_2) influence probability measure $P(I_t|\mathcal{D}_1, \mathcal{D}_2)$.
- Less critical for $P(I_t = 1|\pi_t)$ (e.g. 0.5 for ignorance). However even a pragmatic approach for adjusting Λ like $\min_t p(\hat{\beta}_t|\Lambda) = 0.95 p(\mathbf{0}|\Lambda)$ is not convincing. (Why 0.95 ?)

Improving on Previous Model

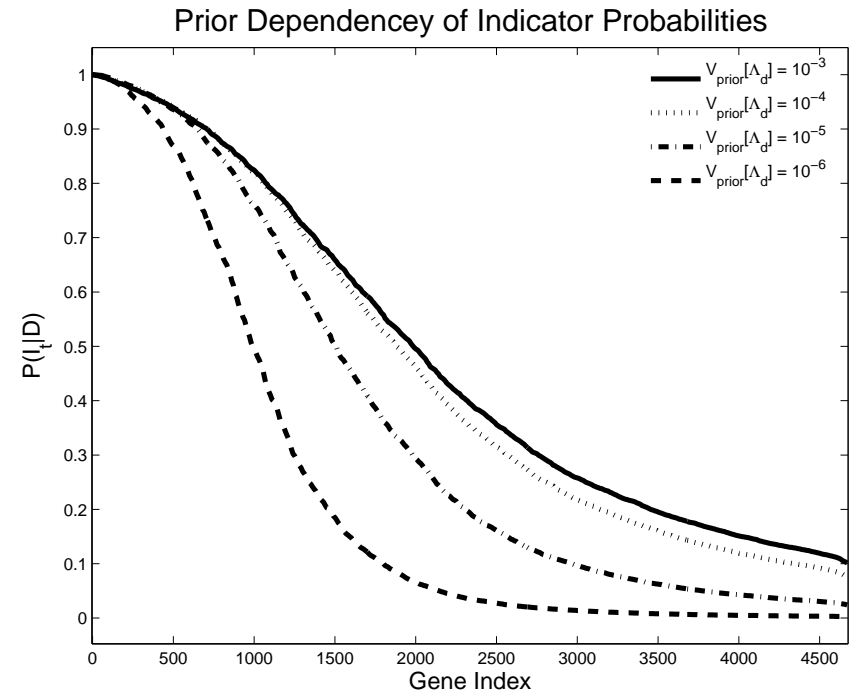
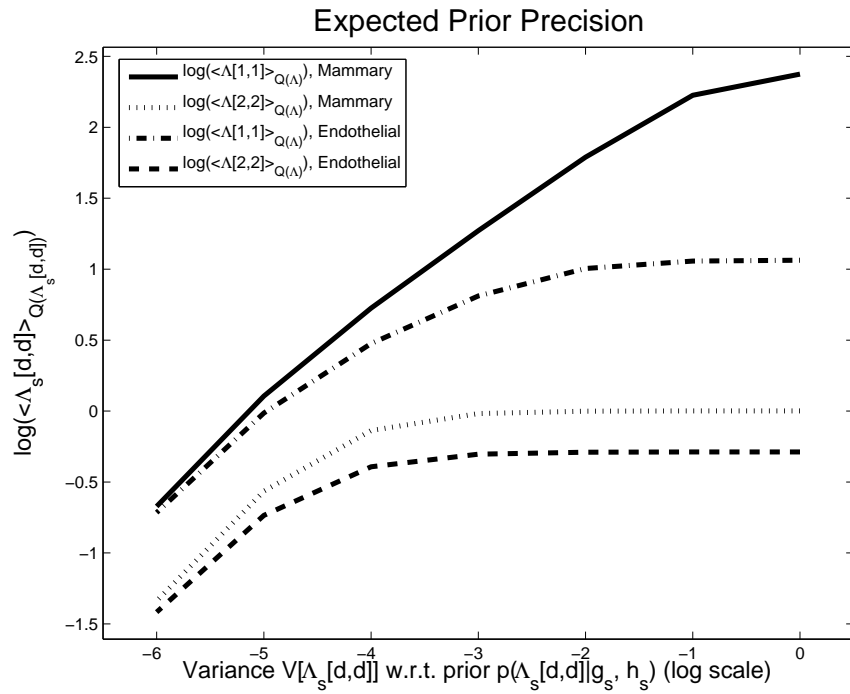
- Hyper parameters (π_t , Λ_1 and Λ_2) influence probability measure $P(I_t|\mathcal{D}_1, \mathcal{D}_2)$.
- Less critical for $P(I_t = 1|\pi_t)$ (e.g. 0.5 for ignorance). However even a pragmatic approach for adjusting Λ like $\min_t p(\hat{\beta}_t|\Lambda) = 0.95 p(\mathbf{0}|\Lambda)$ is not convincing. (Why 0.95 ?)

Better solution uses hierarchical priors



- all genes contribute to inference of Λ_s
- hierarchical priors for sensitivity analysis
- $Q(I_t)$ approximates gene measure
- using *one* model gets all marginals right

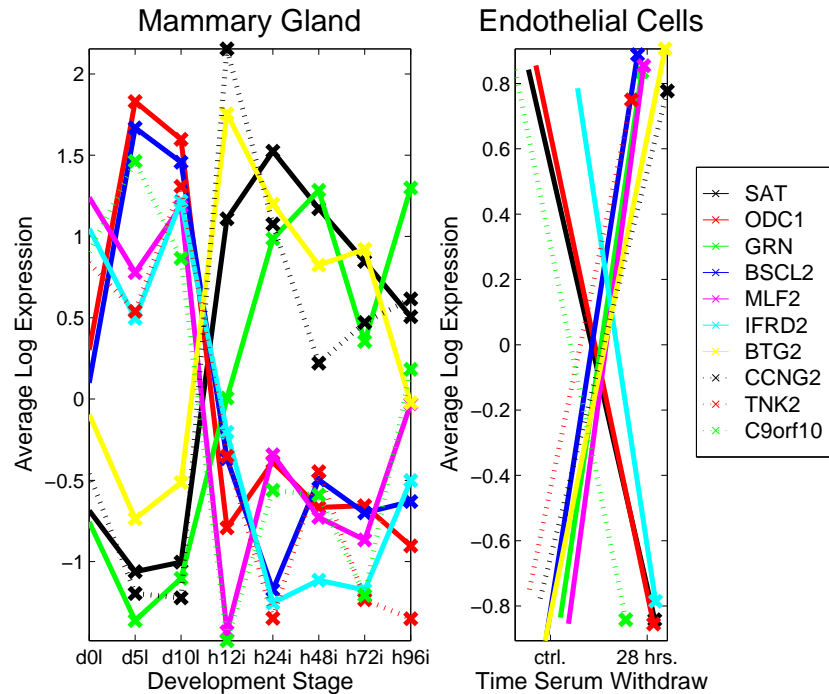
Sensitivity Check



For the hyper parameters this suggests $g \leq 0.01$ and $h \leq 1$.

We also conclude that equal cost results in many potential candidate genes.

Top Ten



Top 10 $P(I_t = 1 | \mathcal{D}_1, \mathcal{D}_2)$ for Mammary lactation vs. involution *and* Endothelial cell death (result updated 01 2007).

Gene Symbol	$P(I_t \mathcal{D})$
SAT	0.99951
ODC1	0.99921
GRN	0.99921
BSCL2	0.99919
MLF2	0.99884
IFRD2	0.99867
BTG2	0.99843
CCNG2	0.99826
TNK2	0.99789
C9orf10	0.99783

Summary

- A relatively straight forward approach allows inference of shared gene function.
- Beware of non hierarchical models - arbitrary gene measures can be adjusted for using the “right” prior.
- Variational methods provide a rather efficient tool to explore models before deciding for a final possibly MCMC based implementation.

Table of Contents

- Problem Statement
- Biological States of Experiments
- Probabilistic Concepts
- Probabilistic Gene Ranking
- Shared Genetic Function
- Discussion of Priors
- Combined Analysis
- Combined Results
- Summary

Variational Bayes

Mean field ansatz plus Jensens inequality. For all pdfs $Q(\theta)$:

$$\begin{aligned} \log \left(\int_{\theta} p(D|\theta)p(\theta)d\theta \right) &\geq \\ &\int_{\theta} (\log(p(D|\theta)) + \log(p(\theta)) - \log(Q(\theta)))Q(\theta)d\theta \\ &= \log(p(D)) + \int_{\theta} (\log(p(\theta|D)) - \log(Q(\theta)))Q(\theta)d\theta \end{aligned}$$

the last integral is a negative Kullback Leibler divergence and thus smaller or equal zero.

+ easy to compute; - systematic error as only an approximation.

[back](#)

Variational Bayes II

Joint Distribution implied by the previous DAG

$$p(I_t, \beta_{1,t}, S_{1,t}, D_{1,t} | \Lambda_1, \pi_t, \gamma, X_{1,t}) = P(I_t | \pi_t) p(\beta_{1,t} | \Lambda_1, I_t) \\ \times \prod_n \left(p(s_{1,t,n} | \beta_{1,t}, \mathbf{x}_{1,t,n}, I_t, \gamma) P(y_{1,t,n} | s_{1,t,n}, I_t) \right)$$

where $S_{1,t} = \{s_{1,t,1}, \dots, s_{1,t,N}\}$ and $D_{1,t} = \{y_{1,t,1}, \dots, y_{1,t,N}\}$.

- Approximate posterior by a mean field expansion $Q(\beta_{1,t} | I_t) \prod_n Q(s_{1,t,n} | I_t)$.
- Write down negative free energy and maximize the functional iteratively w.r.t. all Q-distributions.
- The negative free energy $F_{\max}(Q)$ approximates the log marginal likelihood and thus $P(I_t | D_{1,t}, \Lambda_1, \pi_t, \gamma, X_{1,t})$.

[back](#)