

Detecting and modeling time shifts in microarray time series data applying Gaussian processes

Max Zwießele

Eberhard Karls Universität Tübingen

September 13, 2010



MAX-PLANCK-GESELLSCHAFT



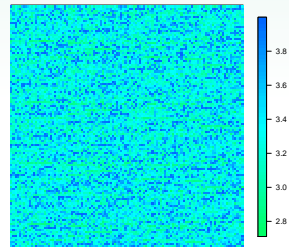
BIOLOGISCHE KYBERNETIK

Why Time Shifts in Microarray Time Series?

Why Microarrays?

Microarray Advantages:

- Huge capacity (thousands of gene expression levels at a time).
- Multiple experiments over time (and replicates):
 - Time Series of experiment.
 - Time Series with several replicates.
- Different conditions:
 - Detect differentially expressed genes.

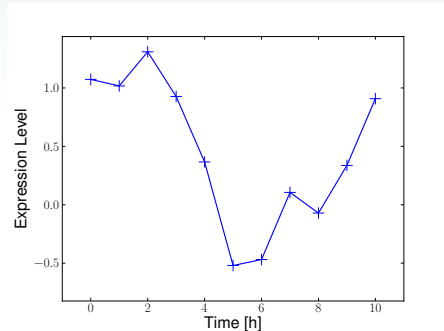


Why Time Shifts in Microarray Time Series?

Why Microarrays?

Microarray Advantages:

- Huge capacity (thousands of gene expression levels at a time).
- Multiple experiments over time (and replicates):
 - Time Series of experiment.
 - Time Series with several replicates.
- Different conditions:
 - Detect differentially expressed genes.

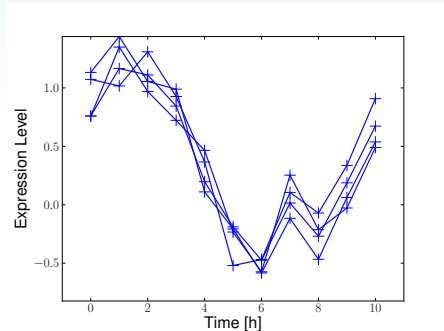


Why Time Shifts in Microarray Time Series?

Why Microarrays?

Microarray Advantages:

- Huge capacity (thousands of gene expression levels at a time).
- Multiple experiments over time (and replicates):
 - Time Series of experiment.
 - Time Series with several replicates.
- Different conditions:
 - Detect differentially expressed genes.

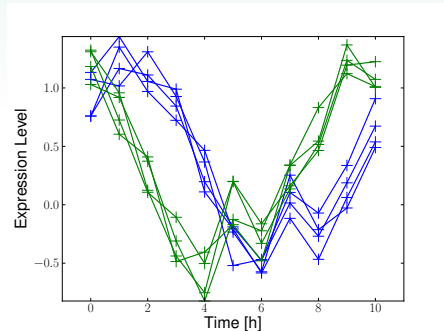


Why Time Shifts in Microarray Time Series?

Why Microarrays?

Microarray Advantages:

- Huge capacity (thousands of gene expression levels at a time).
- Multiple experiments over time (and replicates):
 - Time Series of experiment.
 - Time Series with several replicates.
- Different conditions:
 - Detect differentially expressed genes.

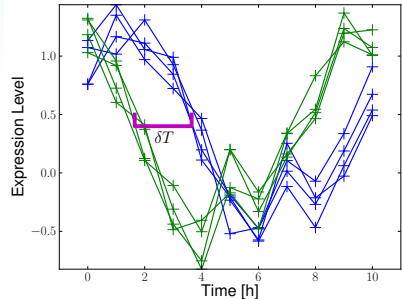


Why Time Shifts in Microarray Time Series?

Why Time Shifts?

- ① *Technical Problems:* [3]
 - Varying efficiencies of intermediary steps:
 - Imaging scanner (Laser).
 - Amount of oligos per cell.
- ② *Biological Problems:*
 - Dependency to environmental conditions.
 - Synchronization of replicates is hard. (Cell-Cycle)

Assumption:
Time shifts in data



Outline

- 1 Why Time Shifts in Microarray Time Series?
- 2 Gaussian Process Regression w.r.t. Time Shifts
 - Dealing with Time Shifts in Data
 - Model Time Shifts with Gaussian Processes
- 3 Applications
 - Detecting Differentially Expressed Genes
 - Correlation of Transcription Factors and Time Shifts
- 4 Summary and Discussion



Outline

- 1 Why Time Shifts in Microarray Time Series?
- 2 Gaussian Process Regression w.r.t. Time Shifts
 - Dealing with Time Shifts in Data
 - Model Time Shifts with Gaussian Processes
- 3 Applications
 - Detecting Differentially Expressed Genes
 - Correlation of Transcription Factors and Time Shifts
- 4 Summary and Discussion



Dealing with Time Shifts in Data

Stegle *et al*'s Approach

- Used Gaussian Processes (GPs) to search functions

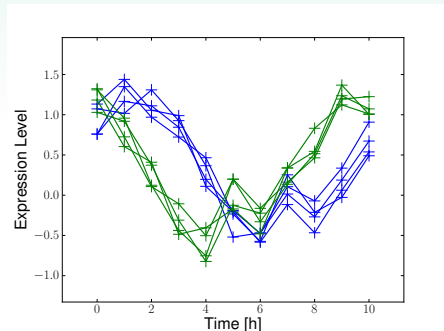
$$f : \mathbf{x} \mapsto f(\mathbf{x}) = \mathbf{y} . \quad (1)$$

- For differential expression detection, they compared:

■, ■ individual model against
 ■ shared model.

- Providing the Bayes factor:

$$BF = \log \frac{p(A|\mathcal{H}_{GP}, \hat{\theta}_I) p(B|\mathcal{H}_{GP}, \hat{\theta}_I) p(\hat{\theta}_I)}{p(A, B|\mathcal{H}, \hat{\theta}_S) p(\hat{\theta}_S)}$$



Dealing with Time Shifts in Data

Stegle *et al*'s Approach

- Used Gaussian Processes (GPs) to search functions

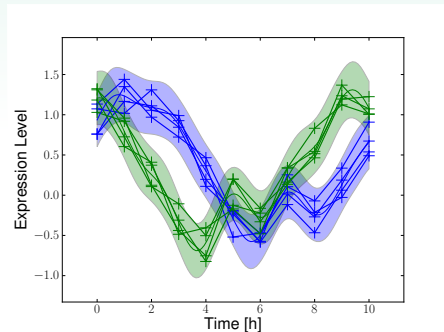
$$f : \mathbf{x} \mapsto f(\mathbf{x}) = \mathbf{y} . \quad (1)$$

- For differential expression detection, they compared:

■, ■ individual model against
■ shared model.

- Providing the Bayes factor:

$$BF = \log \frac{p(A|\mathcal{H}_{GP}, \hat{\theta}_I) p(B|\mathcal{H}_{GP}, \hat{\theta}_I) p(\hat{\theta}_I)}{p(A, B|\mathcal{H}, \hat{\theta}_S) p(\hat{\theta}_S)}$$






Dealing with Time Shifts in Data

Stegle *et al*'s Approach

- Used Gaussian Processes (GPs) to search functions

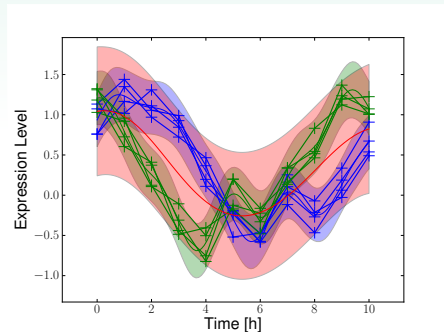
$$f : \mathbf{x} \mapsto f(\mathbf{x}) = \mathbf{y} \quad (1)$$

- For differential expression detection, they compared:

  individual model against
 shared model.

- Providing the Bayes factor:

$$BF = \log \frac{p(A|\mathcal{H}_{GP}, \hat{\theta}_I) p(B|\mathcal{H}_{GP}, \hat{\theta}_I) p(\hat{\theta}_I)}{p(A, B|\mathcal{H}, \hat{\theta}_S) p(\hat{\theta}_S)}$$






Dealing with Time Shifts in Data

Stegle *et al*'s Approach

- Used Gaussian Processes (GPs) to search functions

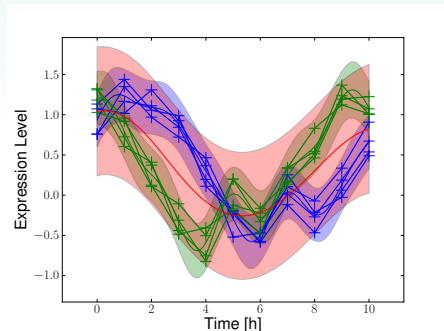
$$f : \mathbf{x} \mapsto f(\mathbf{x}) = \mathbf{y} \quad (1)$$

- For differential expression detection, they compared:

  individual model against
 shared model.

- Providing the Bayes factor:

$$BF = \log \frac{p(A|\mathcal{H}_{GP}, \hat{\theta}_I) p(B|\mathcal{H}_{GP}, \hat{\theta}_I) p(\hat{\theta}_I)}{p(A, B|\mathcal{H}, \hat{\theta}_S) p(\hat{\theta}_S)}$$



Dealing with Time Shifts in Data

Regression w.r.t. Time Shifts

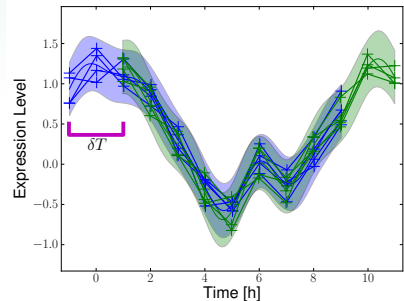
- We allow functions f to shift time line

$$f:\mathbf{x} \mapsto f(\mathbf{x} - \delta T) = \mathbf{y} \quad , \quad (3)$$

→ Matching time series .

- Now compare

  individual model against
 shared model.



Dealing with Time Shifts in Data




Regression w.r.t. Time Shifts

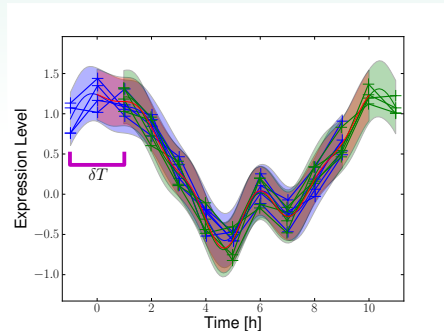
- We allow functions f to shift time line

$$f:\mathbf{x} \mapsto f(\mathbf{x} - \delta T) = \mathbf{y} \quad , \quad (3)$$

→ Matching time series .

- Now compare

  individual model against
 shared model.



Gaussian Process Regression w.r.t. Time Shifts

Probability for unseen data

Two main results of GP regression [1] of interest:

- 1 Probability for previous unseen outputs \mathbf{y}^* :

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{D}, \theta) = \mathcal{N}(\mathbf{y}^*|\underline{\mathbf{y}}^*, \text{cov}(\mathbf{y}^*)), \text{ where} \quad (4)$$

$$\underline{\mathbf{y}}^* = \mathbf{K}_{\mathbf{x}^*, \mathbf{x}} [\mathbf{K}_{\mathbf{x}} + \sigma^2 \mathbf{E}]^{-1} \mathbf{y}, \quad (5)$$

$$\text{cov}(\mathbf{y}^*) = \mathbf{K}_{\mathbf{x}^*} - \mathbf{K}_{\mathbf{x}^*, \mathbf{x}} [\mathbf{K}_{\mathbf{x}} + \sigma^2 \mathbf{E}]^{-1} \mathbf{K}_{\mathbf{x}, \mathbf{x}^*}, \quad (6)$$

and \mathbf{E} denotes the identity matrix.

- 2 Probability for hyperparameters to describe the training data:

$$\log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{1}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi, \quad (7)$$

$$\mathbf{K} = \mathbf{K}_{\mathbf{x}} + \sigma^2 \mathbf{E}. \quad (8)$$



Gaussian Process Regression w.r.t. Time Shifts

Novel Covariance Function with Time Shift Parameter

Requirements:

- Covariance matrix $K_{\mathbf{x}, \mathbf{x}^*} = [k(x_r, x'_{r'}) : \forall x_r, x'_{r'} \in \mathbf{x}, \mathbf{x}^*]$.
- Time shift $T_r \in \mathbf{T}$ per replicate $(\mathbf{x}, \mathbf{y})_r$.

New covariance function with time shift per replicate:

$$k(x_r, x'_{r'}) := A^2 \cdot \exp\left(-\frac{d^2}{2L^2}\right) + \sigma, \quad (9)$$

$$d = \|(x_r - T_r) - (x'_{r'} - T_{r'})\| . \quad (10)$$

Important: The most probable set of hyperparameters $\hat{\theta} = (A, L, \mathbf{T}, \sigma)$ are learned from data! (Appendix 1)

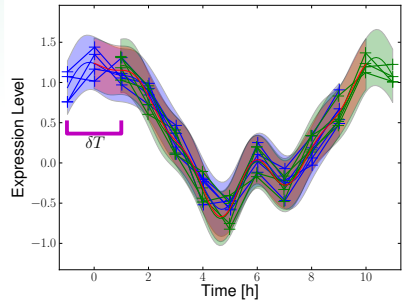


Gaussian Processes Regression w.r.t. Time Shifts

Bayes Factors [2]

Three Bayes factors for hypothesis testing:

- ① Individual vs. shared regression, included \mathbf{T} :
- ② Individual vs. shared regression, where for shared $\mathbf{T} = (0, \dots, 0)$:
- ③ Shared with, vs. shared without time shift:



$$\mathcal{S}_{\delta T}^{\mathcal{I}} := \log \frac{p(\mathbf{C}_k | \mathcal{H}_{GP}, \hat{\theta}_{\mathcal{I}}^{\mathbf{T}}) \cdot p(\mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_{\mathcal{I}}^{\mathbf{T}}) \cdot p(\hat{\theta}_{\mathcal{I}}^{\mathbf{T}})}{p(\mathbf{C}_k \cup \mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_{\mathcal{S}}^{\mathbf{T}}) \cdot p(\hat{\theta}_{\mathcal{S}}^{\mathbf{T}})} = 2.225$$

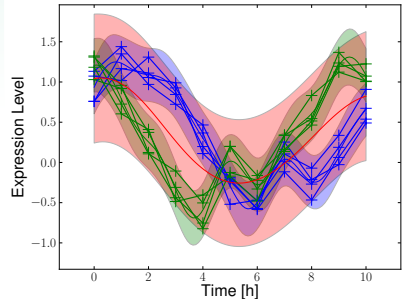


Gaussian Processes Regression w.r.t. Time Shifts

Bayes Factors [2]

Three Bayes factors for hypothesis testing:

- ① Individual vs. shared regression, included \mathbf{T} :
- ② Individual vs. shared regression, where for shared $\mathbf{T} = (0, \dots, 0)$:
- ③ Shared with, vs. shared without time shift:



$$\mathcal{S}_0^{\mathcal{I}} := \log \frac{p(\mathbf{C}_k | \mathcal{H}_{GP}, \hat{\theta}_{\mathcal{I}}^{\mathbf{T}}) \cdot p(\mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_{\mathcal{I}}^{\mathbf{T}}) \cdot p(\hat{\theta}_{\mathcal{I}}^{\mathbf{T}})}{p(\mathbf{C}_k \cup \mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_S^{(0, \dots, 0)}) \cdot p(\hat{\theta}_S^{(0, \dots, 0)})} = 82.780$$

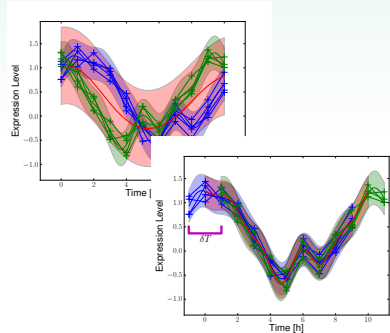


Gaussian Processes Regression w.r.t. Time Shifts

Bayes Factors [2]

Three Bayes factors for hypothesis testing:

- ① Individual vs. shared regression, included \mathbf{T} :
- ② Individual vs. shared regression, where for shared $\mathbf{T} = (0, \dots, 0)$:
- ③ Shared with, vs. shared without time shift:



$$\mathcal{S}^S := \log \frac{p(\mathbf{C}_k \cup \mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_S^{\mathbf{T}}) \cdot p(\hat{\theta}_S^{\mathbf{T}})}{p(\mathbf{C}_k \cup \mathbf{C}_{k'} | \mathcal{H}_{GP}, \hat{\theta}_S^{(0, \dots, 0)}) \cdot p(\hat{\theta}_S^{(0, \dots, 0)})} = 80.556$$



Outline

- 1 Why Time Shifts in Microarray Time Series?
- 2 Gaussian Process Regression w.r.t. Time Shifts
 - Dealing with Time Shifts in Data
 - Model Time Shifts with Gaussian Processes
- 3 Applications
 - Detecting Differentially Expressed Genes
 - Correlation of Transcription Factors and Time Shifts
- 4 Summary and Discussion



Detecting Differentially Expressed Genes

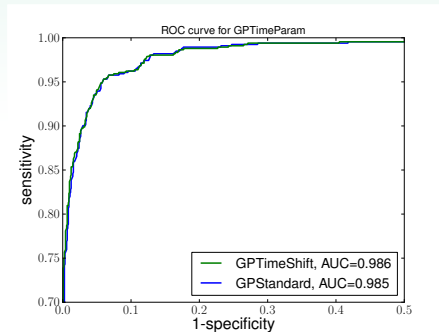
Experiment

- GPTimeShift applied on *Arabidopsis thaliana* data set
 - Plant leafs inoculated by fungal pathogen *Botrytis cinerea*.
 - Harvested every 2h up to 48h post-inoculation.
 - Contains four replicates for each condition of $n = 30,336$ genes.
- Goal: Detect differentially expressed genes.



Detecting Differentially Expressed Genes

Result



Test set of $n = 1,890$ genes, manually classified into either 'differentially expressed' or '**not** differentially expressed'.

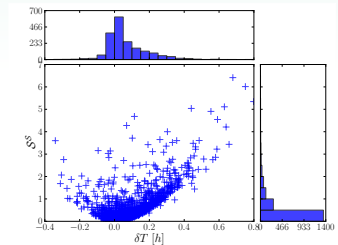


Detecting Differentially Expressed Genes

Conclusion

No significant improvement by including time shifts:

- Microarray problems solved by preprocessing and normalization.
- No significant time shifts in data.



Correlation of Transcription Factors and Time Shifts

Experiment

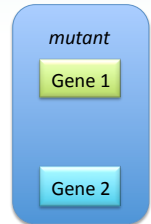
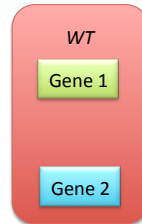
- GPTimeShift applied on *yeast* data set:
 - Knocked out *mei4* (coding for a TF).
 - Harvested every 1h up to 8h post-perturbation.
 - Only one replicate per condition of $n = 4,410$ genes.
- Goal: Detect TFs and corresponding targets by time shift correlations.



Correlation of Transcription Factors and Time Shifts

Results

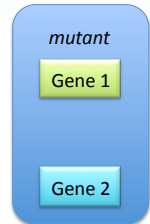
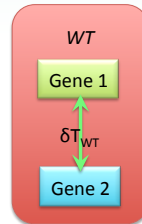
- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^I > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.



Correlation of Transcription Factors and Time Shifts

Results

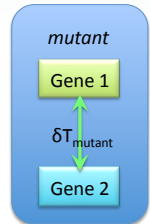
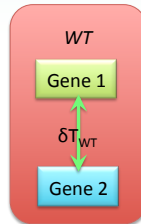
- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^I > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.



Correlation of Transcription Factors and Time Shifts

Results

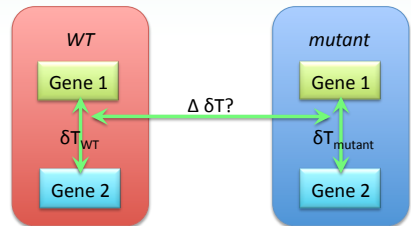
- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^I > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.



Correlation of Transcription Factors and Time Shifts

Results

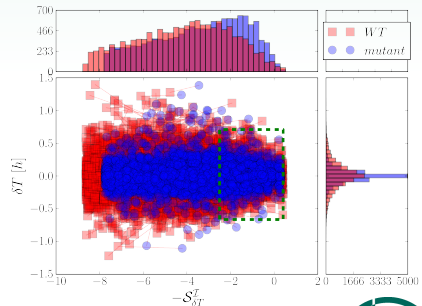
- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^I > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.



Correlation of Transcription Factors and Time Shifts

Results

- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^{\mathcal{I}} > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.

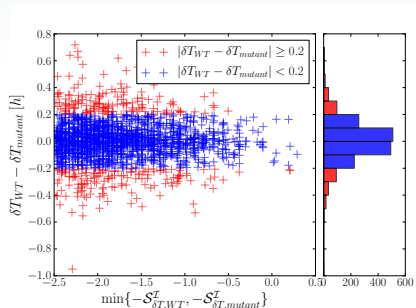




Correlation of Transcription Factors and Time Shifts

Results

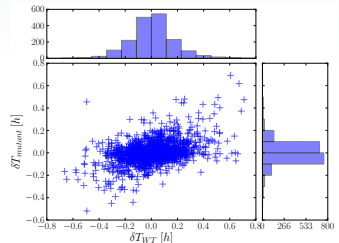
- Compare time shifts within conditions (*WT* and *mutant*) of gene pairs against each other.
- Get correlated gene pairs ($-\mathcal{S}_{\delta T}^I > -2.5$), due to comparing time series of different genes.
- Select only tails of time shift difference distribution.



Correlation of Transcription Factors and Time Shifts

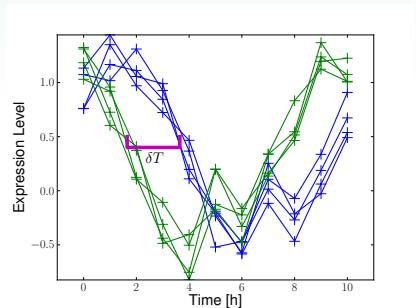
Conclusion

- Unfortunately only a subset of predicted pairs are TFs:
 - Algorithm for searching TFs and targets has to be adapted.
- Proof of principle:
 - TFs lead to time shifts in data.
 - Knocking out a TF leads to less time shifts.



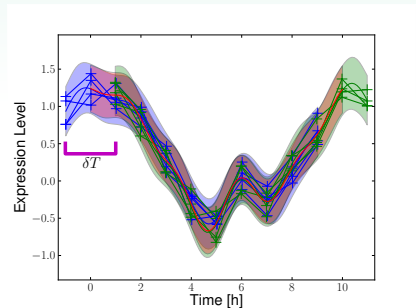
Summary

- Technical and Biological problems in microarray analyzes lead to time shifts.
- We used Gaussian Processes to model these time shifts.
- Modeling time shifts allowed us to
 - improve differential gene expression detection.
 - analyze TF(-target) and time shift correlations.



Summary

- Technical and Biological problems in microarray analyzes lead to time shifts.
- We used Gaussian Processes to model these time shifts.
- Modeling time shifts allowed us to
 - improve differential gene expression detection.
 - analyze TF(-target) and time shift correlations.



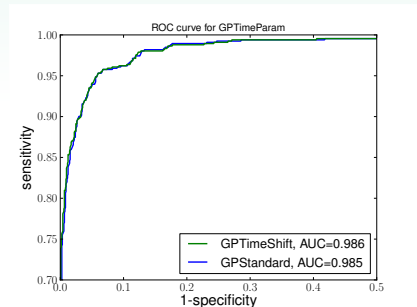
Summary

- Technical and Biological problems in microarray analyzes lead to time shifts.
- We used Gaussian Processes to model these time shifts.
- Modeling time shifts allowed us to
 - improve differential gene expression detection.
 - analyze TF(-target) and time shift correlations.



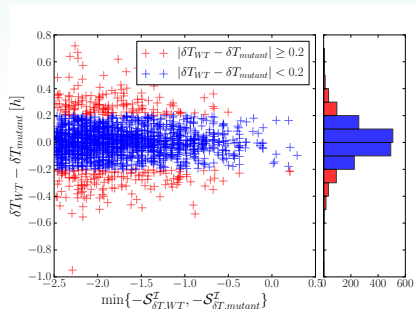
Summary

- Technical and Biological problems in microarray analyzes lead to time shifts.
- We used Gaussian Processes to model these time shifts.
- Modeling time shifts allowed us to
 - improve differential gene expression detection.
 - analyze TF(-target) and time shift correlations.



Summary

- Technical and Biological problems in microarray analyzes lead to time shifts.
- We used Gaussian Processes to model these time shifts.
- Modeling time shifts allowed us to
 - improve differential gene expression detection.
 - analyze TF(-target) and time shift correlations.



Discussion



References



Carl E. Rasmussen and Christopher K. I. Williams.

Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning).

The MIT Press, December 2005.



O. Stegle, K. Denby, D.L. Wild, Z. Ghahramani, and K.M. Borgwardt.

A robust bayesian two-sample test for detecting intervals of differential gene expression in microarray time series.

In Research in Computational Molecular Biology: 13th Annual International Conference, Recomb 2009, Tucson, Arizona, USA, May 18-21, 2009, Proceedings, page 201. Springer, 2009.



A. Zien, B. Schoelkopf, K. Tsuda, and JP Vert.

A primer on molecular biology.

Kernel methods in computational biology, page 3, 2004.



Learning Hyperparameters [1]

With the probability for the outputs being described by the hyperparameters

$$\log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}^{-1}\mathbf{y} - \frac{1}{2}\log |\mathbf{K}| - \frac{n}{2}\log 2\pi \quad , \quad (7) \quad (15)$$

we can denote the most probable set of hyperparameters $\hat{\theta}$ with:

$$\hat{\theta} = \arg \max_{\theta} \{p(\theta|\mathbf{D})\} \quad (16)$$

$$= \arg \max_{\theta} \{p(\mathbf{y}|\mathbf{x}, \theta) \cdot p(\theta)\} \quad (17)$$

$$= \arg \max_{\theta} \{\log p(\mathbf{y}|\mathbf{x}, \theta) + \log p(\theta)\} \quad , \quad (18)$$

where $p(\theta)$ is the hyperprior probability for the hyperparameters θ .

