ENDOGENOUS PERTURBATION ANALYSIS OF CANCER -

INTEGRATED ANALYSIS OF MRNA EXPRESSION AND

COPY NUMBER VARIATION.

Rebecka Jörnsten

Mathematical Sciences, University of Gothenburg/Chalmers

jornsten@chalmers.se

November 19, 2010

OUTLINH

DATA INTEGRATION NETWORK MODELING

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION *Experimental* Structural Pathway enrichment Prognostic

Conclusion and Future work

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

1 Data Integration - Network modeling

2 Modeling Steady-State MRNA levels

- Regularization
- Alternative representation
- Model Validation

3 Application

- Analysis of glioblastoma
- Experimental and Structural Validation

4 Conclusion and Future work

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and Future work

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

CANCER

ENDOGENOUS PERTURBATION ANALYSIS OF

GOALS

- Construct regulatory network and predictive models for cancer pathways
- Identify disease-specific key regulators and their targets

???Is THIS STATISTICS???

- Challenges in biology: high-throughput technologies (get information on activities of 10000+ genes in different cancer tumors)
- Statistics to help identify which among these 10000+ genes are interacting and how

Outline

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

ENDOGENOUS PERTURBATION ANALYSIS OF CANCER

GOALS

- Construct regulatory network and predictive models for cancer pathways
- Identify disease-specific key regulators and their targets



- Data now available at multiple levels of these biological processes (genome seq, transcription, proteomics)
- By integrating multiple data sources we hope to identify direct, causal interactions.

JÖRNSTEN. EPOC

Outline

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Multiple levels of data



OUTLINE

Data Integration -Network modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and UTURE WORK

▲□▶ ▲□▶ ▲ 臣▶ ▲ 臣▶ 三臣 - のへ⊙

PATHWAY AND NETWORK MODELING

CORRELATION-BASED

- Use mRNA expression data only
- Construct networks based on mRNA-mRNA correlation (relationship)
- Pro: massive amounts of mRNA expression data available
- Con: "passive" observation makes network reconstruction difficult

EXPERIMENTAL APPROACHES

- Knock-out of single genes observe the system
- Pro: targeted perturbation and question
- Con: expensive, limited amounts of data

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

PATHWAY AND NETWORK MODELING

USING NATURALLY OCCURRING GENETIC VARIATION

• Genetic variation: in a population of individuals, different yeast strains, crosses

ENDOGENOUS PERTURBATIONS

- Focus on *acquired* genetic variation, e.g. mutations in tumors.
- View this as a perturbation of the biological system and observe the response

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Outline

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION *Experimental* Structural Pathway enrichment Prognostic

EPoC

ENDOGENOUS PERTURBATION

- We view each tumor's Copy Number Aberration (CNA) profile as a system perturbation that simultaneously affects multiple genes, and
- the mRNA profiles as the steady-state response to that perturbation
- CNAs tend to appear in a patient-specific and multifactorial manner - ideal for network construction
- Ongoing projects like The Cancer Genome Atlas means massive amounts of mRNA/CNA data are available

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL ANE STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

DATA - 186 TUMORS FROM THE CANCER GENOME ATLAS CONSORTIUM (TCGA)

JÖRNSTEN. EPOC

Outline

DATA INTEGRATION -NETWORK MODELING

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION *Experimental* Structural Pathway enrichment Prognostic

Conclusion and Future work



I. Tumor

samples

3. Molecular profiles

- Genetic profile
- \sim 10s of mutations affecting protein seq
- ~100s of copy-number altered genes

Epigenetic profile

~100s of hypermethylated promoters

Transcriptional profile

- ~100s-1000s altered message RNA levels
- ~10s-100s altered micro-RNA levels

DATA - 186 TUMORS FROM THE CANCER GENOME ATLAS CONSORTIUM (TCGA)



JÖRNSTEN. EPOC

Outline

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and Future work

▲ロト ▲御ト ▲ヨト ▲ヨト 三回 - のへの

MODEL WITH MULTIPLICATIVE EFFECTS

$$\frac{dy_i}{dt} = u_i \alpha_i \prod_j y_j^{w_{ij}} - \beta_i \prod_j y_j^{v_{ij}}$$

- where y_i is the mRNA expression of gene i, and u_i the copy number.
- Elements w_{ij} and v_{ij} denote the effect of gene j on i during synthesis/degradation respectively,
- and the α_i and β_i denote "environmental effects (non-CNA perturbations)"



JÖRNSTEN. EPOC

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic CONCLUSION AND

Future work

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

STEADY-STATE SOLUTION

- We take CNA profile u and a baseline reference \tilde{u} , and likewise for mRNA y, \tilde{y} .
- Define $\Delta u_i = \log(u_i) \log(\tilde{u}_i), \ \Delta y_i = \log(y_i) \log(\tilde{y}_i)$
- At steady-state we can write:

$$\Delta u_i + \sum_j (w_{ij} - v_{ij}) \Delta y_i +$$

$$+\overbrace{\left(\log(\alpha_i) - \log(\tilde{\alpha}_i)\right) - \left(\log(\beta_i) - \log(\tilde{\beta}_i)\right)}^{\gamma_i} = 0$$

- We denote the *direct* causal influence of transcript j on i by $a_{ij} = w_{ij} v_{ij}$.
- We want to find the "important" *a*_{ij} (i.e. large/real effects)

For each gene we can now write a linear equation that represents this genes mRNA expression:

JÖRNSTEN. EPOC

Outline

Data Integration -Network modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

STEADY-STATE SOLUTION

$$a_{ii}y_i = \sum_{j \neq i} -a_{ij}y_j - u_i + \gamma_i$$

Note: we dropped the Δy symbol here - for convenience we just use y to denote the baseline adjusted mRNA and similarly for CNA

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and Future work

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

REGRESSION MODELING

•
$$a_{ii}y_i = \overbrace{j \neq i}^{model} -a_{ij}y_j - u_i + \gamma_i$$
 (from now on, set $a_{ii} = 1$)

• We want to use data from the 186 tumors to *estimate* the *a*_{ii}

Data

- Observe mRNA and CNA for tumors $t = 1, \dots, T = 186$ and genes $i = 1, \dots, n = 10000+$.
- We assume the same *a_{ij}* describe the mRNA expression for all tumors and
- the term γ_i captures the tumor specific (unmeasured) variation = referred to as "noise" or "error".
- Data: y_{it}, u_{it}.
- Find the a_{ij}s that best match up the ys and the us.

Outline

Data Integration Network Modeling

Modeling Steady-State MRNA levels

Regression modeling

Regularization Alternative Representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

REGRESSION MODELING

• Least Squares: Find aiis to minimize

$$\sum_{t=1}^{T} (y_{it} - (\sum_{j \neq i} -a_{ij}y_{jt} - u_{it}))^2 = \sum_{t=1}^{T} \gamma_{it}^2$$

• That is: we minimize the average "noise" or errors over all tumors, i.e. how the observed mRNA y_i deviates from the model a_{ij}



OUTLINE

DATA INTEGRATION · NETWORK MODELING

Modeling Steady-State MRNA levels

Regression modeling

REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and Future work

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

REGRESSION MODELING

- When we minimize the squared errors we get an *estimate* of the *true* a_{ij}
- Which interactions are real (a_{ij} > 0 or < 0), and which are not (a_{ij} = 0)?
- There are statistical tests for deciding if an estimated a_{ij} is large enough to be declared "real" (depends on the magnitude of the estimated a_{ij} , the number of samples(tumors), the amount of noise, ...).





OUTLINE

Data Integration Network modeling

Modeling Steady-State MRNA levels

Regression modeling

Regularization Alternative Representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic CONCLUSION AND

Conclusion and Future work

Multivariate response regression modeling

• For the whole system we can write

$$Ay + u = \Gamma$$

- That is, row i in A contains the interactions that gene i has with the other genes a_{ij}, etc.
- We are solving 10000+ regression problems!!!

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling

REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic CONCLUSION AND

FUTURE WORK

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

NETWORK SOLUTION

 $A\mathbf{y} + \mathbf{u} = \Gamma$

- A is called the *network matrix*, our parameters of interest.
- Γ includes all the tumor-specific differences, unmeasured environmental effects, SNPs etc.



JÖRNSTEN. EPOC

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling

REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic CONCLUSION AND

TUTURE WORK

三 のへで

OH-OH! ANOTHER COMPLICATION

Let's go back to the original regression model for gene i

$$a_{ii}y_{it} = \underbrace{\sum_{j \neq i}^{model} -a_{ij}y_{jt} - u_{it}}_{j \neq i} + \gamma_{it}$$

- We have data from 186 tumors
- We have p = 10000+ unknowns in the model (all the a_{ij} 's)
- This is called the p > T problem we are asking more questions about the data than we have observations.
- Consequence: there is an infinite number of solutions a_{ij} 's that fit the data equally well!!!

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression Modeling

REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION *Experimental* Structural Pathway enrichment Prognostic

p > T problem



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで



(日本)(同本)(日本)(日本)(日本)

• Each gene takes its turn to be *y*, and the other genes form columns in *X*.

JÖRNSTEN, EPOC

LEAST SQUARES

We want to minimize

10

$$Q(a) = \sum_{t=1}^{T} \gamma_t^2 = (\gamma_1 \gamma_2 \cdot \gamma_T) \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \cdot \\ \gamma_T \end{pmatrix} = \gamma^T \gamma = (y - Xa)^T (y - Xa)^{\text{Model:}}_{\text{Request}}$$

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ → ヨ → の Q @

$$\frac{dQ}{da} = -2X^{T}(y - Xa) = 0 \rightarrow (X^{T}X)a = X^{T}y$$

When p > T (here 10000 vs 186) the inverse of $X^T X$ does not exist! No unique solution a

JÖRNSTEN, EPOC

RIZATION

SOLUTION 1 - SUBSET OF GENES Pick max T genes to be in the model. Then

 $\hat{a} = (X_S^T X_S)^{-1} X_S^T y$

where \hat{a} means we have estimated the interactions a, and X_S means we are only using a subset S of the p = 10000 columns in X.

- PRO: Simple and interpretable model â "genes in set S affect gene y"
- CON: Which set S to use? How large max T but perhaps better to use fewer genes? There are 2¹⁰⁰⁰⁰ possible combinations of genes we can put in set S!!!

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression Modeling

REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Solution 2 - Regularization

We "fix" the inverse of $X^T X$:

$$\hat{\boldsymbol{a}} = (\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X} + \lambda \boldsymbol{I})^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{y}$$

That is, we add something to the diagonal of $X^T X$. Now we can take the inverse!

Note, $\hat{a} \simeq a_{original}/(1 + \lambda)$ - That is, we are *shrinking* the estimates of *a* toward 0.

- PRO: Simple fix to the inverse problem
- CON: We made all the *a*s smaller real effects and the *a*'s that are truly zero What should λ be?

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Outline

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression Modeling

REGULARIZATION

Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Solution 2 - Regularization

This is also called "regularized regression" because we can arrive at this solution by solving problem:

$$\min_{a} Q(a) = (y - Xa)^{T} (y - Xa) + \lambda \sum_{j=1}^{10000} a_{j}^{2} = (y - Xa)^{T} (y - Xa) + \lambda a_{\text{STEADY-STATE}}^{\text{MODELING}}$$

That is, we *penalize* large values of a (err on the side of caution). Solution:

$$\frac{dQ}{da} = -2X^{T}(y - Xa) + 2\lambda a = 0 \rightarrow (X^{T}X + \lambda I)a = X^{T}y$$

We "fix" the inverse of $X^T X$:

$$\hat{\boldsymbol{a}} = (\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{y}$$

REGULARIZATION

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Solution 3 - Lasso regression

This idea on penalizing the large *a* values can be generalized to other types of penalties.

Last 10 years, lots of statistics research into penalizing the *absolute values* of *a*: "Lasso regression"

$$\min_{a} Q(a) = (y - Xa)^{T} (y - Xa) + \lambda \sum_{j=1}^{10000} |a_{j}|$$

- - - - -

Turns out the solution looks something like this:

•
$$\hat{a}_j = a_{original} - 1/\lambda$$
 if $\hat{a} > 1/\lambda$

•
$$\hat{a}_j = a_{\textit{original}} + 1/\lambda$$
 if $\hat{a} < -1/\lambda$

• $\hat{a}_i = 0$ otherwise

That is, small a's are set to 0, large a's are shrunk by $1/\lambda$

- PRO: Combines the best of both worlds: selection of non-zero a's, and simple procedure involving just λ rather than the 2¹⁰⁰⁰⁰ combinatorics
- CON: What should λ be?

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression Modeling

REGULARIZATION ALTERNATIVE REPRESENTATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

FINALLY! BACK TO NETWORK SOLUTION

$A\mathbf{y} + \mathbf{u} = \Gamma$

- We solve for A one row (gene) at a time
- For each gene we solve a Lasso regression problem
- Still need to choose λ we'll come back to that in a minute.



JÖRNSTEN. EPOC

OUTLINE

Data Integration Network modeling

Modeling Steady-State MRNA levels

Regression modeling

REGULARIZATION

Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

ALTERNATIVE REPRESENTATION

$$\mathbf{y} = G\mathbf{u} + \Gamma'$$

- $G = -A^{-1}$ is called the *system matrix*.
- *G* represents the *system gain*: where the genetic variation (system input) shows up as amplified (positive or negative) signal in the mRNA expression (system output).

OUTLINE

Data Integration Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative

REPRESENTATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic CONCLUSION AND

FUTURE WORK

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

ALTERNATIVE REPRESENTATION



$$y_1 = a_{12}y_2 + u_1 + \gamma_1$$

$$y_2 = u_2 + \gamma_2$$

$$y_3 = a_{31}y_1 + u_3 + \gamma_3$$

$$y_4 = a_{42}y_2 + u_4 + \gamma_4$$



$$y_{1} = u_{1} + \overbrace{a_{12}}^{g_{12}} u_{2} + \gamma'_{1}$$

$$y_{2} = u_{2} + \gamma'_{2}$$

$$y_{3} = \overbrace{a_{31}}^{g_{31}} u_{1} + \overbrace{a_{31}a_{12}}^{g_{32}} u_{2} + u_{3} + \gamma'_{3}$$

$$y_{4} = \overbrace{a_{42}}^{g_{42}} u_{2} + u_{4} + \gamma'_{4}$$

JÖRNSTEN. EPOC

Outline

Data Integration Network modeling

Modeling Steady-State MRNA levels

Regression Modeling Regularization

ALTERNATIVE REPRESENTATION

Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

ALTERNATIVE REPRESENTATION - WHY?

$$A\mathbf{y} + U = \Gamma$$
, $\mathbf{y} = G\mathbf{u} + \Gamma$

- The genetic variation is thought to be (in part) the *disease cause*, whereas the mRNA expression is the symptom.
- Looking at G is therefore more informative for understanding which part of the biological network is disease specific.
- Practical reason: A is difficult to estimate due to strong correlations (lots of models fit the data equally well).



JÖRNSTEN. EPOC

Outline

Data Integration · Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Choosing λ - the network size

- Remember: λ controls which a's (g's) are zero and which are non-zero.
- Large λ : we risk missing some important interactions in our model
- Small λ: we risk including "false" interactions in our model.
- How do we validate the results?

VALIDATION STATISTICS

We consider two different evaluation measures;

- Network structure consistency Kendall's W
- mRNA prediction minimum average prediction error

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

NETWORK CONSISTENCY

- We randomly split the data into two halves, and estimate the network on each of the two sets of data.
- We compare the network agreement between the two data sets using Kendall's W: measures how well the non-zero and zero a's (or g's) agree in terms of location and magnitude.
- Most consistent networks contain \sim 400 edges.



structure prediction

JÖRNSTEN. EPOC

Outline

Data Integration -Network modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

MRNA PREDICTION

- We split the data into random halves, and estimate a network on each of the two sets of data.
- We use the network model from one half to predict the mRNA levels in the other half, and vice versa.
- Minimum mRNA prediction errors for networks with \sim 10000 edges.



JÖRNSTEN. EPOC

OUTLINE

Data Integration -Network modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Conclusion and Future work

GLIOMA NETWORK ANALYSIS



OUTLINE

Data Integration Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

Analysis of glioblastoma

Experimental and Structural Validation

Experimental Structural

Pathway enrichment Prognostic

GLIOMA NETWORK ANALYSIS





cell differentiation (GO:0030154, 2.29e-09)
 nervous system development (GO:0007399, 1.35e-24)
 cell-cell signaling (GO: GO:0007267, 8.39e-25







GLIOBLASTOMA EXPERIMENTAL STRUCTURAL VALIDATION

Experimenta Structural Pathway enrichment Prognostic

ANALYSIS OF

Conclusion and Future work

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

GLIOMA NETWORK ANALYSIS



(日)、(四)、(E)、(E)、(E)

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

Analysis of glioblastoma

EXPERIMENTAL AND STRUCTURAL VALIDATION

Structural Pathway

Prognostic

EXPERIMENTAL VALIDATION



0

DATA INTEGRATION NETWORK MODELING

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

Analysis of glioblastoma Experimental and

Structural Validation

Experimental Structural Pathway enrichment

Prognostic

Conclusion and Future work

▲□▶ ▲□▶ ▲ 臣▶ ▲ 臣▶ 二臣 - のへで

For 146 tumors, we have mRNA and CNA data from

- Two different labs
- Two different platforms (agilent, affy)

We compare network consistency across network sizes and between competing methods.



Outline

Data Integration -Network modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION

Structural

Pathway enrichment Prognostic

OVERLAP WITH PATHWAY DATABASES





- We map interactions found by different methods to molecular links in pathway repositories HPRD, Reactome, Intact, and NCI-nature.
- Compare pathway links to the shortest paths in networks.
- EPoC-G is clearly enriched for short or direct paths.
 EPoC-A, GLasso, ARACNE and GeneNet are associated with longer paths.

OUTLINE

Data Integration -Network modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment

Prognostic

PREDICTING PATIENT SURVIVAL

DECOMPOSITION OF G

- We argue that G represents the disease mechanism.
- The SVD decomposition of $G = C\Lambda D^T$ has the following meaning:
 - leading columns of *D* are directions of CNA perturbations that are amplified by the system.
 - leading columns of *C* are directions of mRNA transcripts most affected by the directions in *D*.

• Write
$$Y = GU = C\Lambda D^T U$$

- $\rightarrow C^T Y = \Lambda D^T U$
- Projecting mRNA onto columns of C = output, Projecting CNA onto D = input, $\Lambda =$ amplification.

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment **Prognostic**

PREDICTING PATIENT SURVIVAL

DECOMPOSITION OF G



- $C^T Y = \Lambda D^T U$
- Consider the leading projections (first columns of *C* and *D*).
- mRNA profiles of individual patients are projected onto
 C: Z_y = C^TY and CNA profiles are projected by
 Z_u = D^TU
- Compare the survival of the patients using these projected scores.

JÖRNSTEN. EPOC

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment **Prognostic**

Decomposition of G vs A and data

Survival curve based on singular vector score



JÖRNSTEN. EPOC

Outline

Data Integration -Network modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment **Prognostic** CONCLUSION AND

JUTURE WORK

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

- EPoC scales to 10000 genes and produces stable network estimates
- Attained network models exhibit good agreement with pathway databases
- Experimental validation of novel hubs identify interesting therapeutic targets
- The EPoC network provides clinical stratification into long- and short-term survival, whereas competing methods do not.

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Outline

Data Integration · Network Modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

- Extend EPoC to tumor subtype identification (promising results already)
- Common and subtype specific network modules (ongoing work with PhD students in my group)
- Include multiple data sources (e.g. miRNA, methylation)
- Methodological work supervised prognostic network estimation

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION *Experimental* Structural Pathway enrichment Prognostic

Conclusion and Future work

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

- Sven Nelander
- Tobias Abenius
- Teresia Kling, Linnea Schmidt, Bodil Nordlander, Erik Johansson, Torbjörn Nordling, Chris Sander, Björn Nilsson, Peter Gennemark, Keiko Funa, Linda Lindahl
- Cancerfonden, Barncancerfonden, Vetenskapsradet, BioCare, Sahlgrenska-CMR, NB-CNS

OUTLINE

Data Integration -Network Modeling

Modeling Steady-State MRNA levels

Regression modeling Regularization Alternative representation Model Validation

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL ANE STRUCTURAL VALIDATION Experimental Structural Pathway enrichment Prognostic

Outline

Data Integration -Network modeling

Modeling Steady-State MRNA levels

REGRESSION MODELING REGULARIZATION ALTERNATIVE REPRESENTATION MODEL VALIDATION

Application

ANALYSIS OF GLIOBLASTOMA EXPERIMENTAL AND STRUCTURAL VALIDATION Experimental Structural Pathway enrichment

Prognostic

Conclusion and Future work

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへぐ