



Biomarker Discovery: Computational Approaches

Seungchan Kim

Center for Computational Systems Biology Prairie View A&M University CCSB@PVAMU [https://ccsb.pvamu.edu]







PRAIRIE VIEW



Welcome to CCSB@PVAMU

https://ccsb.pvamu.edu/

The Center for Computational Systems Biology at the Prairie View A&M University (CCSB@PVAMU) was established in 2017, by a funding from the Texas A&M University Systems Chancellor's Research Initiative (CRI) to become a nationally recognized computational systems biology research center.

As a newly established center, we are rapidly developing research programs to study 1) biomedical questions of high translational significance and impact and 2) challenging questions in plant and agricultural science, utilizing various genomics and computational tools. The center aims to achieve:

- · Recruitment of high-quality faculty and research scientists
- Promoting interdisciplinary research across Roy G. Perry College of Engineering, College of Agriculture and Human Sciences, College of Arts and Sciences, and other research centers/labs within PVAMU campus
- Development of active collaboration with leading biomedical research institutes as well as bioinformatics research centers
- Development of training programs for bioinformatics and computational systems biology at both the undergraduate and the graduate levels.



.

Recent posts





CCSB Collaboration: SU2C







What is a Biomarker?





Biomarkers can be

- Pulse
- Blood pressure
- Basic chemistries or more complex laboratory tests of blood and other tissues





Molecular Biomarkers

- Gene mutations
- Single Nucleotide Polymorphisms
- Gene expressions
- Small molecules such as miRNAs
- Emerging from analysis of high throughput molecular measurements such as DNA and RNA sequencing data





What is a Biomarker?

- NIH [1998] a characteristic that is objectively measured and evaluated as an indicator of normal biological processes, pathogenic processes, or pharmacologic responses to a therapeutic intervention
- WHO [2001] any substance, structure, or process that can be measured in the body or its products and influence or predict the incidence of outcome or disease, or
- WHO [2001] almost any measurement reflecting an interaction between a biological system and a potential hazard, which may be chemical, physical, or biological.

Curr Opin HIV AIDS. 2010 November ; 5(6): 463–466





What is a Biomarker?

- NIH [1998] a characteristic that is objectively <u>measured</u> and evaluated as an <u>indicator</u> of normal biological processes, pathogenic processes, or pharmacologic responses to a therapeutic intervention
- WHO [2001] any substance, structure, or process that can be <u>measured</u> in the body or its products and influence or <u>predict</u> the incidence of outcome or disease, or
- WHO [2001] almost any <u>measurement reflecting</u> an interaction between a biological system and a potential hazard, which may be chemical, physical, or biological. <u>The measured response may be</u> <u>functional and physiological, biochemical at the cellular level, or a</u> <u>molecular interaction.</u>

Curr Opin HIV AIDS. 2010 November ; 5(6): 463–466



Biomarkers

PRAIRIE VIEW A&M UNIVERSITY

- classification and prediction,
- as surrogate outcomes in clinical trials,
- as measures of toxic or preventive exposures, or
- as a guide to individual treatment choice





Biomarker Discovery: Workflow











Bottle necks/Challenges

- Sample size
 - How to estimate statistical power
- Preprocessing
 - Normalization
 - Batch correction
- Internal validation
 - Statistical validation
 - Error estimation



Sample Size

PRAIRIE VIEW

- Is the study sufficiently statistically powered to test hypothesis?
- Traditional method to estimate sample size is often based on a simple, irrelevant hypothesis test
 - Effect size often arbitrarily chosen
- Simulation-based approach is more desirable, where
 - Synthetic data are generated from a plausible data model,
 - Classifier design and evaluation should performed with varying size of cases and controls to determine sample size





Sample Size



 The study is adequately powered (>0.9) to detect a 2 fold difference in as few as 25 miRNAs out of 7,000 while adjusting for multiple tests using an FDR correction of 0.05





Sample Size: effect size vs. FDR







Sample Size: simulation-based





Normalization

PRAIRIE VIEW

- Compensate for technical and/or biological covariates such as (in RNAseq):
 - Sequencing depth
 - Transcript length
- Also, done to transform the measurements so that they fit into a mathematical model – often called standardization
- Aggregate vs single-sample normalization







sample



Batch correction

PRAIRIE VIEW

- Minimize systematic, undesired difference in measurements between batches of samples
 - Difference between different sites where the samples are collected
 - Difference between different times when the samples are collected
 - Standardized protocols to collect samples should mitigate this issue, but ...
- Some of the conditions are not controllable.



Site 1

Site 2

Batch: Multiple sites

Before Batch Correction













Another Example



After







Another Example



cancer

•

•











Error Estimation

- TPR or sensitivity
- FPR or 1 specificity
- ROC

PRAIRIE VIEW A&M UNIVERSITY

- AUC
- Error or accuracy



• How well (accurately) can we estimate these?





Error Estimation

$$\hat{\varepsilon} = E\left[\left(\hat{\varphi}_{N,k}(\boldsymbol{X}_{i}) - Y_{i}\right)^{2}\right] \qquad \qquad \text{Estimated Error}$$

$$= E\left[\left(\varphi(\boldsymbol{X}_{i}) + \left(\hat{\varphi}_{k}(\boldsymbol{X}_{i}) - \varphi(\boldsymbol{X}_{i})\right) + \left(\hat{\varphi}_{N,k}(\boldsymbol{X}_{i}) - \hat{\varphi}_{k}(\boldsymbol{X}_{i})\right) - Y_{i}\right)^{2}\right]$$

$$= E\left[\left(\left(\varphi(\boldsymbol{X}_{i}) - Y_{i}\right) + \left(\hat{\varphi}_{k}(\boldsymbol{X}_{i}) - \varphi(\boldsymbol{X}_{i})\right) + \left(\hat{\varphi}_{N,k}(\boldsymbol{X}_{i}) - \hat{\varphi}_{k}(\boldsymbol{X}_{i})\right)\right)^{2}\right]$$





Interval Validation



CV: cross-validation – find λ (lambda) for Lasso





Network-based Approach







EDDY: Evaluation of Differential DependencY







EDDY + CTRP-CCLE



- Speyer et al., Pac Symp Biocomput. 2017; 22: 497–508.
- Identifies pathways enriched with differential dependency between sensitive and non-sensitive cancer cell lines, as in DDNs
- Discover mediators of drug sensitivity, i.e. potential targets?





Fighting Cancer, Cell by Cell







Fighting Cancer, Cell by Cell







What about Bladder Cancer!



Translational Team Science Award (DoD-USAMRMC-CDMRP-TTSA), "Development of Classifiers for Novel Bladder Cancer Subtypes"

Woonyoung Choi (JHMI) Seungchan Kim (PVAMU)



















Pathways enriched with differential dependency between BIE and BIS

Pathway	# genes	p-val	Rewiring	Mediators
SYNTHESIS OF PIPS AT THE LATE ENDOSOME MEMBRANE	10	0.0025	0.68	PIKFYVE
MTORC1 MEDIATED SIGNALLING	11	0.0183	0.88	EIF4EBP1, MLST8
PHOSPHORYLATION OF CD3 AND TCR ZETA CHAINS	16	0.0225	0.27	CD4, CD3E, CD3D, CD3G, PAG1, CSK
SEMA3A PAK DEPENDENT AXON REPULSION	15	0.0236	0.65	LIMK1 PLXNA1
ELEVATION OF CYTOSOLIC CA2 LEVELS	10	0.0246	0.60	TRPC3
SYNTHESIS OF PIPS AT THE EARLY ENDOSOME MEMBRANE	12	0.0261	0.73	PI4K2A MTMR2
CELL EXTRACELLULAR MATRIX INTERACTIONS	14	0.0327	0.55	PARVA FERMT2
SYNTHESIS SECRETION AND INACTIVATION OF GLP1	19	0.0390	0.87	CDX2 PAX6
CD28 DEPENDENT VAV1	11	0.0394	0.50	PAK2 FYN
GRB2 SOS PROVIDES LINKAGE TO MAPK SIGNALING FOR INTERGRINS	15	0.0454	0.61	ITGB3 SOS1 TLN1
THE ROLE OF NEF IN HIV1 REPLICATION AND DISEASE PATHOGENESIS	28	0.0479	0.69	AP2S1, CD8B, AP1S2, ELMO1, AP1S1, B2M









Center for Computational Systems Biology



Prairie View A&M University CCSB@PVAMU https://ccsb.pvamu.edu





Q&A