Widespread redundancy in -omics profiles of cancer mutation states

ISMB 2022: General Computational Biology track

Jake Crawford

PhD Candidate, Greene Lab University of Pennsylvania

Web: http://jjc2718.github.io Twitter/GitHub/Gmail: @jjc2718

• Background:

- Functional -omics in cancer
- What are we predicting? Why are we predicting it?

• Current work:

- Which cancer -omics is the best -omics? (for our problem)
- Survival prediction
- Multi-omics models

Genetic alterations in cancer



Genetic alterations in cancer



tumor/ malignancy Functional -omics in cancer











Predicting mutation status from -omics data









Can we detect a functional signature...

• for this mutation?



Can we detect a functional signature...

- for this mutation?
- for this data type?

But why?



• Benchmark for cancer-related information content

Question: Does data source or mutated gene contribute most to predictive signal?

But why?



• Benchmark for cancer-related information content

Question: Does data source or mutated gene contribute most to predictive signal?

• Study design

Question: When are data sources redundant? When might they provide unique information?

Performance comparison

Which data type is best for mutation prediction?

Gene expression vs. DNA methylation

- Gene set collected from ¹, ², ³ (~280 cancer-related genes)
- RNA-seq, 27K and 450K methylation array data from TCGA
- Only samples with all three data types measured -> **7,981** total samples



¹ Vogelstein et al. Science 2013, ² Bailey et al. Cell 2018, ³ COSMIC Cancer Gene Census

Performance evaluation



4 folds x 2 random seeds cross-validation, stratified by cancer type



Metric: AUPR (handles continuous output, works well for imbalanced labels)



Baseline: model with permuted mutation labels (controls for cancer type)



Which classifiers beat the baseline?



Each point = results for one gene in one data type





mean difference between true model and baseline with permuted labels



t-test p-value, across 4 cross-validation folds x 2 random seeds



genes

genes

77/272 genes

Expression/methylation summary across genes



Comparing more -omics data types

• Data types from TCGA:



• Use samples where *all data types* are measured -> **5,226** total samples

All data types prediction: comparisons vs. baseline





All data types summary across genes



Survival analysis

Do similar results hold for prognosis prediction?

Survival analysis

- Prediction of TCGA clinical endpoints from -omics data
 - PCA preprocessing for all -omics types
 - Elastic net Cox regression
 - Covariates for age, cancer type, mutation burden
 - Baseline predictor uses only non-omics covariates
- Metric: concordance index (higher = better)

Gene expression vs. methylation for survival



All data types for survival



Multi-omics models

Does combining data types improve performance?

Experimental design

- Can we improve performance by combining more than one data type?
- Gene expression, 27K methylation, 450K methylation
 - Top 5000 PCs for each (results with raw features were similar)
 - Concatenate datasets to form "multi-omics" model
 - All pairs + combination of all 3 data types
- Focused on six "well-predicted" driver genes
 - EGFR, IDH1, KRAS, PIK3CA, SETD2, TP53

Single-omics or multi-omics?



Single-omics or multi-omics?



Takeaways

What have we learned?

Gene, not data type, is the primary source of variability

- Gene expression tends to capture the most signal
- Many genes predictable using one -ome are predictable using multiple -omes
- Mutations with strong functional signatures tend to perturb all the -omes, to some degree



Gene, not data type, is the primary source of variability



Performance by data type, for genes with at least one significant predictor

- Grey dot = significantly better than baseline
- Black dot = same, and statistically equivalent to "best" predictor

Survival prediction is similar between -omes

- Not much difference between -omics types
- Generally fairly difficult to beat the covariate-only baseline



Multi-omics integration doesn't seem to help

- No discernable difference between best single-omics and multi-omics model
- Most "useful" -omics type varies with target gene
- Relatively simple method, maybe room for improvement



Resource for experimental design



Performance by data type, for genes with at least one significant predictor

https://greenelab.github.io/mpmp-manuscript (Figure 6; Supp. Figure 9)

Acknowledgements

Greene Lab:

- Casey Greene
- Natalie Davidson
- Ben Heil
- Ariel Hippen
- Alex Lee
- David Nicholson
- Milton Pividori
- Halie Rando
- Taylor Reiter
- Vince Rubinetti

Collaborators:

- Brock Christensen (Dartmouth)
- Maria Chikina (U of Pittsburgh)

Thesis Committee:

- Marylyn Ritchie
- Mingyao Li
- Kai Tan
- Pablo Cámara
- Donna Slonim (Tufts)

Code and analyses: https://github.com/greenelab/mpmp

Paper in Genome Biology: https://doi.org/10.1186/s13059-022-02705-y

