

Biologically motivated constraints for machine learning models in biomedicine

This manuscript ([permalink](#)) was automatically generated from [greenelab/biopriors-review@982a997](#) on August 27, 2019.

Authors

- **Jake Crawford**

 [0000-0001-6207-0782](#) ·  [jjc2718](#) ·  [jjc2718](#)

Graduate Group in Genomics and Computational Biology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA; Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA · Funded by Grant XXXXXXXX

- **Jane Roe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [janeroe](#)

Department of Something, University of Whatever; Department of Whatever, University of Something

Abstract

Introduction

1. Why is this necessary? Where do traditional models fail?
 - Biological data is often “wide” ($n \ll p$), so it can be difficult for models to learn feature covariance structure directly from the data.
 - For “tall” data ($n \gg p$) as is often used in non-biological ML, sufficiently flexible models (NNs, etc) can effectively learn covariance structure directly from the data, and incorporating structured priors is often unnecessary.
 - Frequently, there is no natural way to encode biological knowledge as real-valued features: knowledge can take many forms (pathways, networks, hierarchies, ontologies/knowledge graphs). Most ML models operate on real-valued features, so modelers must be creative when deciding how to constrain the model based on the structure of the relevant biological knowledge.
2. What is a “biologically motivated constraint”? How is this different from regularization?
 - a. Regularization refers specifically to adding a smoothness or penalty term to the model’s loss function to reflect a biological desideratum.
 - For example, a model using the expression of genes as predictors could have a penalty requiring genes in the same biological pathway to have a similar role in the predictive model.
 - Or, a model using eQTL expression levels to predict a phenotype could have a penalty requiring genes that are nearby on the chromosome to have a similar role in the predictive model. (could give explicit examples of loss functions if it makes things clearer)
 - b. We use “constraint” as a more general term that can take many forms, depending on the choice of machine learning model.
 - Regularization is a type of biological constraint, but there exist constraints that do not fit the definition of regularization given above.
 - Example: DCell (Ma et al 2018). Not really an explicit “regularizer” mathematically speaking, but the structure of the NN is biologically inspired and constrained using biological knowledge.
 - c. Note that the term “constraint” has a specific meaning in mathematical optimization (one might minimize some function $f(x)$ subject to a constraint $g(x) > c$). Herein, we use the term more generally, and none of the constraints we discuss take this form. (should we come up with a different term to use for this?)

Sequence-Constrained Models

Applications in regulatory biology

1. Early examples
 - 10.1126/science.1254806 (can cite as an early example of using hand-engineered features as input to a 2-layer MLP, preceded most CNN examples on sequence data -> splicing exon inclusion prediction)
 - 10.1038/nbt.3300 (CNN for prediction of DNA/RNA binding protein sequence specificity)
 - 10.1038/nmeth.3547 (CNN for multi-task variant effect prediction from noncoding genome, one of the earlier examples of CNN on genomic sequence data)
 - 10.1101/gr.200535.11 (CNN for prediction of chromatin accessibility from DNA sequence, evaluated to predict effects of SNPs)
2. More recent examples
 - 10.1101/737981 (CNN for prediction of TF binding probability and binding strength from sequence, used to identify and interpret cis-regulatory code i.e. order and spacing of TF binding motifs, no need for strong modeling assumptions + allow lots of context ~1000 base windows)

- 10.1016/j.cell.2019.04.046 (prediction of alternative polyadenylation from sequence, used to predict probability of poly-A cleavage and engineer new PASs)
- DNA methylation models typically operate on microarray beta values (proportion of DNA methylated at each CpG locus), so these wouldn't be considered sequence-constrained. An exception: DeepSignal (10.1093/bioinformatics/btz276) - CNN that works on electrical signals around CpG sites from Nanopore sequencing reads, to predict 5mC or 6mA methylation status
- 10.1101/329334 and 10.1093/bioinformatics/bty612: prediction of gene expression (or differential gene expression) from histone modification signals. Signals are discretized using bins (no raw sequence input) but bidirectional LSTM considers context in both directions and allows for learning complex histone mark dependencies.

Applications in variant calling and mutation detection

- 10.1038/nbt.4235 (actually works on images of read pileups, but uses sequence context)
- 10.1038/s41467-019-09025-z (single molecule sequencing: PacBio/ONT, works on featurization of sequence defined in paper, outperforms \wedge on these)
- 10.1101/097469 (similar to DeepVariant, but works on raw sequences rather than images, + less preprocessing/claims to learn more info directly from sequences)
- 10.1101/601450 (calling of short indels using sequence information)
- 10.1038/s41467-019-09027-x (mutation detection in tumors)

Applications in CRISPR guide selection

- 10.1371/journal.pcbi.1005807 (shows that a wide variety of sequence attributes are useful for predicting cleavage efficiency, didn't use contextual info?)
- 10.1038/nbt.3437, 10.1038/s41551-017-0178-6 (models for predicting on-target efficacy and off-target effects, uses contextual info about sequence i.e. order 2+ features)
- 10.1038/nbt.4061 (CNN-based model for Cpf1 activity prediction, convolutions of length 5 on one-hot encoded sequence)
- 10.1186/s13059-018-1459-4 (similar to \wedge but for Cas9, also uses ChIP-seq data)
- 10.1101/636472 (similar to \wedge)
- 10.1101/505602 (uses STRING -> gene network-based features, in addition to sequence features, for on-target efficacy prediction)

Network- and Pathway-Constrained Models

Applications in bulk transcriptomics

- 10.1016/j.cels.2019.04.003 (pathways -> transfer learning for rare diseases)
- 10.1038/s41540-019-0086-3 (PPI + pathways -> phenotype prediction/cancer subtyping)
- <https://arxiv.org/pdf/1906.10670> (HumanBase -> AML drug response prediction)
- 10.1093/bioinformatics/bty429 (HINT PPI)
- 10.1038/s41598-017-03141-w (PPI -> mutated cancer gene detection)
- 10.1371/journal.pcbi.1006657 (PPI + network optimization -> cancer outcome prediction)
- 10.1186/s12859-018-2500-z (pathways -> NN structure, for predicting patient outcomes in GBM)
- 10.1186/s12859-017-1984-2 (gene regulatory network -> NN structure, for patient outcomes in kidney transplantation and ulcerative colitis)
- 10.1038/s41598-018-19635-0 (gene regulatory network -> group lasso, same datasets as \wedge)
- 10.1093/bioinformatics/bty945 (gene regulatory network -> NN structure, for predicting gene expression given gene knockouts)

Applications in single cell transcriptomics

- 10.1038/s41592-019-0456-1 (pathways -> cell type deconvolution/pathway-level eQTLs)

- 10.1093/nar/gkx681 (PPI; protein-DNA interactions)
- 10.1101/544346 (PPI)

Applications in genetics

1. Genotype-phenotype associations (GWAS/eQTL)
 - 10.1093/bioinformatics/btu293
 - 10.1093/bioinformatics/btx677 (multivariate version of ^)
 - 10.1111/biom.13072 (phenotype prediction using variant and GE data)
2. Using pathways as a prior to study SNP-SNP or gene-gene interactions (not sure if this is truly a constraint on the model, more just a strategy to reduce the number of hypothesis tests)
 - 10.1101/182741, 10.1371/journal.pgen.1006973 (pathways + GWAS data -> GIs)
 - 10.1371/journal.pgen.1006516 (pathways + GWAS data -> GIs)
3. Using SNP-SNP or gene-gene networks to reduce hypothesis testing burden for detecting GIs
 - can find some examples if we decide this is a direction worth going

Other constrained models

Ontology-constrained models

1. Phenotype prediction
 - 10.1038/nmeth.4627 (gene ontology -> NN structure, to predict effects of mutations on growth in yeast)
2. Function prediction
 - 10.1093/bioinformatics/btx624 (GO hierarchy relationships -> NN structure, for function prediction)
 - 10.1093/bioinformatics/btx252 (PPI network + tissue ontology -> function prediction)

Otherwise constrained models

1. Cell cycle information
 - 10.1038/nbt.3102
 - 10.1101/526848 (in principle the denoising method could be generalized to other gene sets, but here they used cell cycle-relevant gene sets and emphasized the utility of this)
2. Circadian rhythms: 10.1073/pnas.1619320114 (circular node autoencoder for modeling periodic gene expression)
3. TAD/3D chromatin structure information? Can probably find some examples of this

Conclusions

1. What is outside of the scope of this review? (this can also go in introduction?)
 - a. Biological "constraints" vs. feature selection or feature extraction from heterogeneous biological data (e.g. network embedding approaches)
 - Example: one could use a network-based feature extraction method (e.g. Node2Vec) to convert each gene in a PPI network into a set of real-valued features, then use those + gene expression as input to a model
 - For purposes of keeping this review short enough, I'm trying to stay away from papers like ^, but still unclear to me where exactly the line should be drawn. Almost any ML model that operates on sequence data can be viewed as having a feature extraction component, for example.
 - b. Could kind-of consider many single-cell dimension reduction methods biologically constrained (e.g. dropout/zero inflation modeling approaches, etc), but this is way too broad for this review - maybe refer the reader to other recent reviews of these methods.

- c. Could also consider omics integration methods (combining, for example, gene expression and epigenetic data) to be biologically constrained, but we refer the reader to [10.1016/j.inffus.2018.09.012](https://doi.org/10.1016/j.inffus.2018.09.012) for further detail on these methods.

References
