

Automated machine learning for large-scale multi-modal biomarker studies.

Mike A. Nalls, PhD
NIA / NIH c/o
Data Tecnica Int'l, LLC
mike@datacnica.com
 @mike_nalls

Quick intro to key concepts in machine learning (ML).

Images from this great review I had no part of.

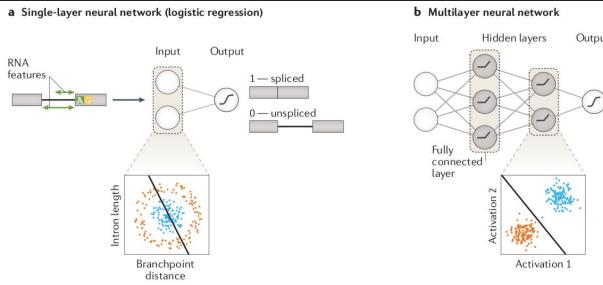
REVIEWS

Deep learning: new computational modelling techniques for genomics

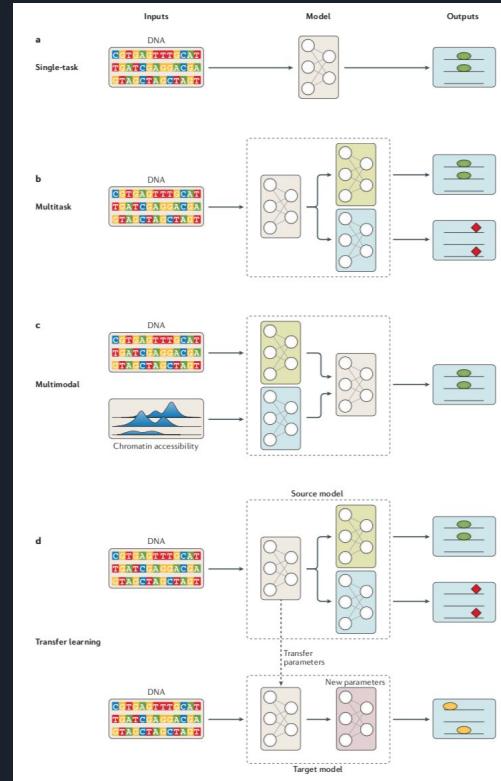
Gökçen Eraslan^{1,2,5}, Žiga Avsec^{3,5}, Julien Gagneur^{3*} and Fabian J. Theis^{1,2,4*}

Abstract | As a data-driven science, genomics largely utilizes machine learning to capture dependencies in data and derive novel biological hypotheses. However, the ability to extract new insights from the exponentially increasing volume of genomics data requires more expressive machine learning models. By effectively leveraging large data sets, deep learning has transformed fields such as computer vision and natural language processing. Now, it is becoming the method of choice for many genomics modelling tasks, including predicting the impact of genetic variation on gene regulatory mechanisms such as DNA accessibility and splicing.

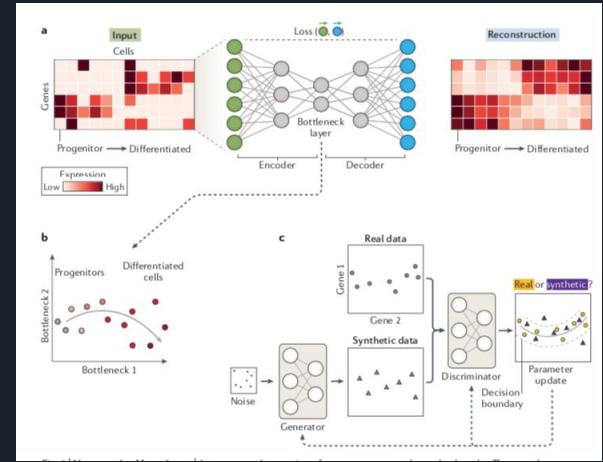
Linear models (regression) are ML.



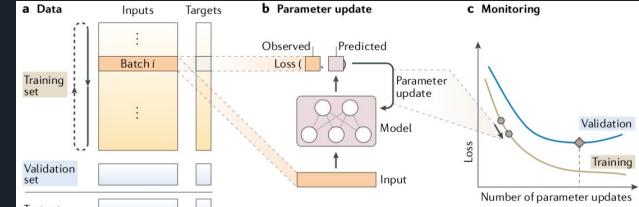
Supervised learning



Unsupervised learning



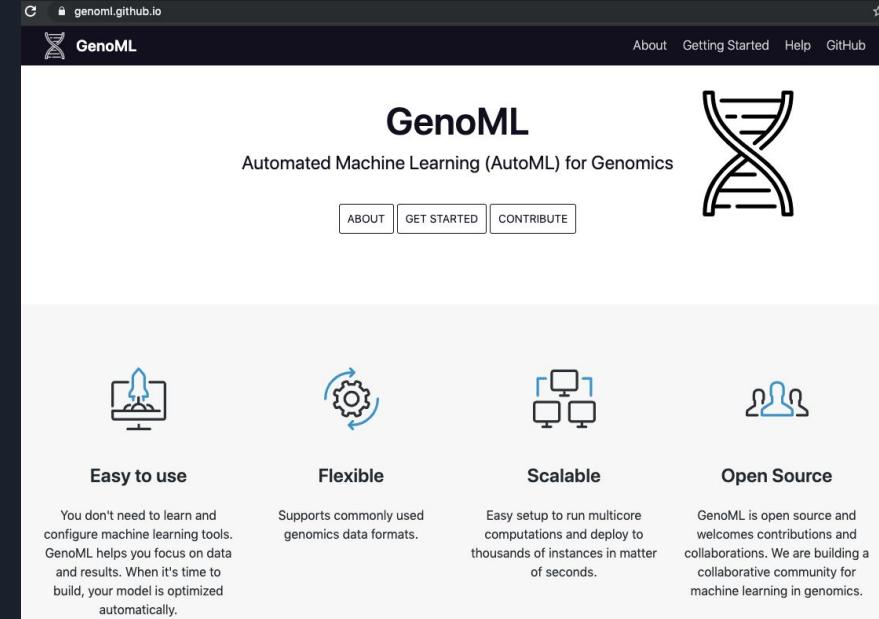
Cross-validation plus tuning



Democratizing and automating ML

GenoML - genoml.github.io.

- Older beta version available (R)
- Current version (python) in October '19
- Automated ML tuned for clinical and genomics applications
- Feature selection
- Compete dozens of most popular algorithms for classification, regression and clustering
- Longitudinal and meta/federated learning methods under development
- Model zoo
- Terra ready!



Current python version is quite fast and memory efficient from preprocessing to training and tuning... competing 12 models on 1,000 samples @ 100,000 features takes ~1 hour using 6 cores and 64GB RAM.



Why use ML as an analytic framework for biomarker studies?

- PRS are great but you need a huge GWAS, got one applicable to your trait or even your population of interest?
- Maybe problems you are dealing with can't be solved by a line?
- Hey, are you trying to include things that aren't SNPs?
- Need to run production scale analytics?
- Did you know business is really far ahead of science in terms of analytics?
- Maybe just looking at one predictor at a time isn't a great idea?

Multimodal ML in AMP-PD: pilot work.

- 872 samples from PPMI passing QC
 - “Idiopathic PD” with no other NDD
 - Controls
- Genome sequencing data
 - Reduced to ~50K variants
 - GWAS hits (90 SNPs) - Nalls et al 2019
 - LD pruned tag SNPs at MAF > 5%
- Clinical candidate predictors
 - UPSIT
 - Family history
 - Age
 - Genetic sex
- Transcriptomic data
 - ~50k normalized transcripts
 - Pre-release
 - Jensen and Craig labs
 - Thanks David and Kendall!

nih.gov/research-training/accelerating-medicines-partnership-amp

ACCELERATING MEDICINES PARTNERSHIP (AMP)

Accelerating Medicines Partnership (AMP)

Alzheimer's Disease
Type 2 Diabetes
Rheumatoid Arthritis and Lupus
Parkinson's Disease

On this page

AMP Partners | Budget | Opportunity | Challenge | Impact | Governance

Overview

The Accelerating Medicines Partnership (AMP) is a public-private partnership between the National Institutes of Health (NIH), the U.S. Food and Drug Administration (FDA), multiple biopharmaceutical and life science companies and non-profit organizations to transform the current model for developing new diagnostics and treatments by jointly identifying and validating promising biological targets for therapeutics. The ultimate goal is to increase the number of new diagnostics and therapies for patients and reduce the time and cost of developing them.

AMP was launched in February 2014, with projects in three disease areas:

- Alzheimer's disease
- type 2 diabetes
- autoimmune disorders of rheumatoid arthritis and systemic lupus erythematosus (lupus)

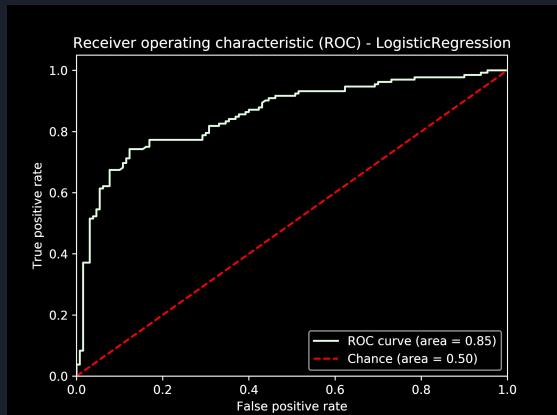
In January 2018, an AMP project on Parkinson's disease was launched with nine partners.



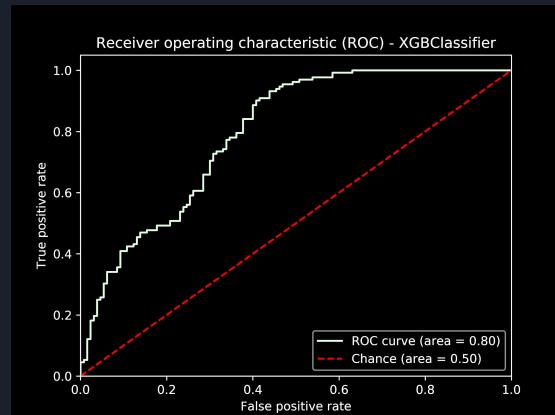
clinical VS genotype VS omic

Performance in 30% test samples

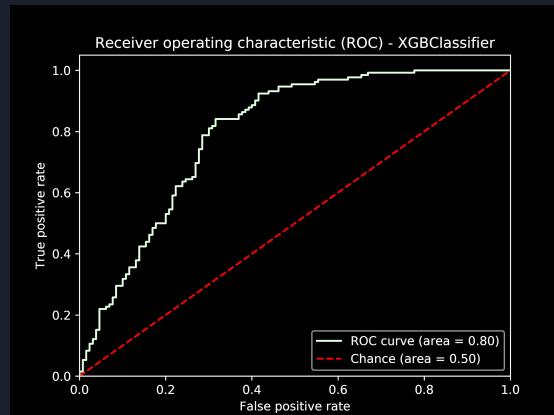
| Dataset | Algorithm | AUC_Percent | Accuracy_Percent | Balanced_Accuracy_Percent | Log_Loss | Sensitivity | Specificity | PPV | NPV | Runtime_Seconds |
|----------|--------------------|-------------|------------------|---------------------------|----------|-------------|-------------|-------|-------|-----------------|
| Omic | XGBClassifier | 79.639 | 75.191 | 75.152 | 0.540 | 0.803 | 0.700 | 0.731 | 0.778 | 248.344 |
| Genotype | XGBClassifier | 79.557 | 70.611 | 70.594 | 0.522 | 0.727 | 0.685 | 0.701 | 0.712 | 83.112 |
| Clinical | LogisticRegression | 85.487 | 79.389 | 79.452 | 0.477 | 0.712 | 0.877 | 0.855 | 0.750 | 0.012 |



clinical



genotype

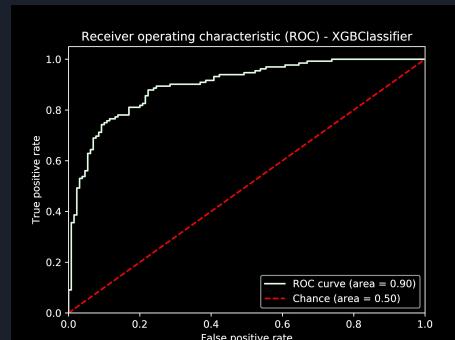
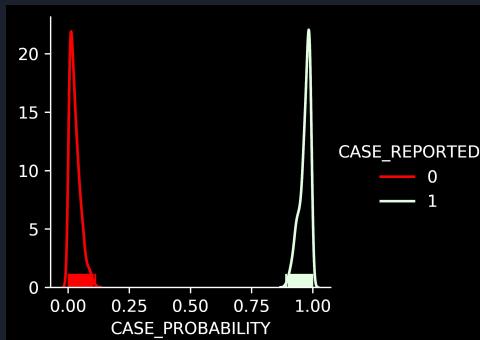


omic

Multiple modalities in the same model perform best.

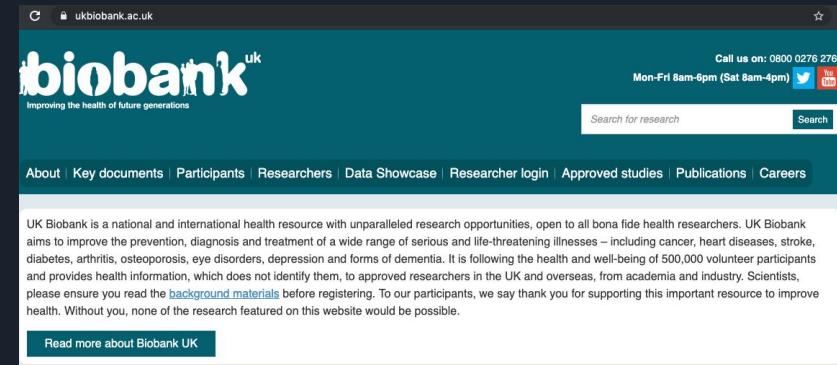
Performance in 30% test samples

| Algorithm | AUC_Percent | Accuracy_Percent | Balanced_Accuracy_Percent | Log_Loss | Sensitivity | Specificity | PPV | NPV | Runtime_Seconds |
|-------------------------------|-------------|------------------|---------------------------|----------|-------------|-------------|--------|--------|-----------------|
| XGBClassifier | 89.8776 | 80.5344 | 80.5536 | 0.4078 | 0.7803 | 0.8308 | 0.8240 | 0.7883 | 308.8823 |
| GradientBoostingClassifier | 89.7786 | 82.8244 | 82.8613 | 0.4103 | 0.7803 | 0.8769 | 0.8655 | 0.7972 | 482.1020 |
| BaggingClassifier | 84.3590 | 75.1908 | 75.2331 | 0.6044 | 0.6970 | 0.8077 | 0.7863 | 0.7241 | 158.2600 |
| AdaBoostClassifier | 83.8054 | 79.7710 | 79.7844 | 0.6553 | 0.7803 | 0.8154 | 0.8110 | 0.7852 | 223.0328 |
| LogisticRegression | 81.2238 | 73.6641 | 73.6364 | 0.7095 | 0.7727 | 0.7000 | 0.7234 | 0.7521 | 19.0327 |
| LinearDiscriminantAnalysis | 78.3625 | 72.1374 | 72.1212 | 0.5645 | 0.7424 | 0.7000 | 0.7153 | 0.7280 | 36.2314 |
| MLPClassifier | 77.9720 | 72.1374 | 72.1620 | 0.8206 | 0.6894 | 0.7538 | 0.7398 | 0.7050 | 229.5861 |
| SVC | 74.4814 | 69.4656 | 69.5280 | 0.5967 | 0.6136 | 0.7769 | 0.7364 | 0.6645 | 360.2292 |
| SGDClassifier | 72.8147 | 72.9008 | 72.8147 | 9.3597 | 0.8409 | 0.6154 | 0.6894 | 0.7921 | 2.9697 |
| RandomForestClassifier | 71.9930 | 60.3053 | 60.4138 | 0.6056 | 0.4621 | 0.7462 | 0.6489 | 0.5774 | 1.5003 |
| KNeighborsClassifier | 63.6946 | 59.9237 | 60.0233 | 1.5562 | 0.4697 | 0.7308 | 0.6392 | 0.5758 | 146.4244 |
| QuadraticDiscriminantAnalysis | 44.6911 | 44.6565 | 44.6911 | 19.1150 | 0.4015 | 0.4923 | 0.4454 | 0.4476 | 26.8305 |

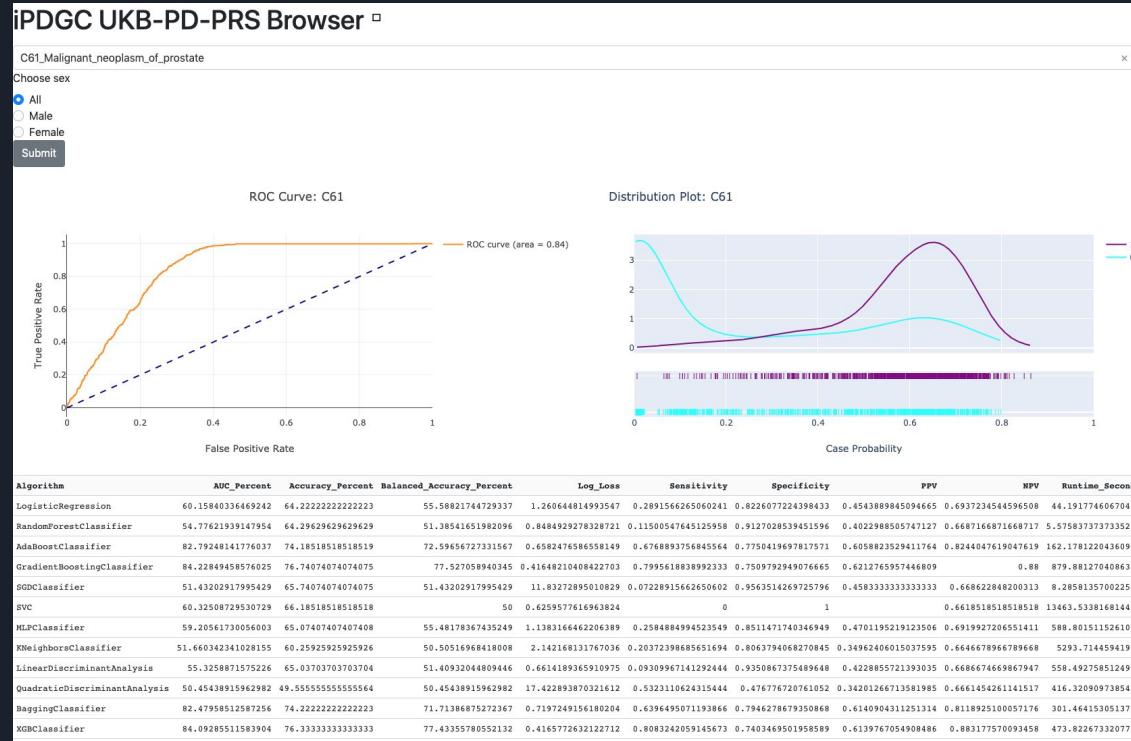


Production scale analyses in the UK Biobank.

- > 400K participants
 - Sampling to reduce resource use
 - 1:1 case to control
 - 5K samples per trait maximum
- > 7M common imputed variants
 - LD pruned to ~31K variants
- Demographics
 - Age
 - Genetic sex
 - Townsend score
 - All analyses stratified by ancestry
 - Sex stratification as well when needed
- ICD10 codes
 - Granular and parent level analyses
 - Only codes with > 1K cases



UKB results browser and model zoo.



Should be up in a couple weeks for multiple ancestry groups and > 2K codes.
Trained and tuned model files will be available for download and application to your data

Application of UKB results to PD studies.

- Progression cohorts are generally small and underpowered
- Predict comorbidities and ancillary outcomes from UKB to PD cohorts
- Predict aspects of progression markers that overlap with ICD10 codes
 - For example ... constipation, depression etc from UPDRS
- Endpoint prediction in trial design

PD Progression Meta-GWAS Browser
Please provide the inputs and hit "Create Table" to begin.

Tables Download all results About

LRRK2

P $<$ 0.05

Phenotype

- Constipation
- Cognitive Impairment
- Daytime Sleepiness
- Depression
- Dyskinesias
- Hypersmia
- Insomnia
- Motor fluctuations
- REM sleep Behavior Disorder
- Restless Legs Syndrome
- Hoehn and Yahr score
- Hoehn and Yahr score of 3 or more

Baseline analysis by logistic regression model

| SNP | RSID | CHR | START | REF | ALT | MAF | BETA | SE | P | N | NSTUDY | Isq | FUNC |
|-----|-------------|-----------|-------|----------|-----|-----|--------|---------|--------|---------|--------|-----|---------------|
| 1 | 12:40562742 | rs7300197 | 12 | 40562742 | G | T | 0.2123 | -0.228 | 0.1111 | 0.04006 | 1350 | 6 | 0 intergenic |
| 2 | 12:40566066 | rs7314494 | 12 | 40566066 | C | A | 0.3097 | -0.216 | 0.1025 | 0.03516 | 1350 | 6 | 0 intergenic |
| 3 | 12:40581231 | | 12 | 40581231 | T | A | 0.339 | -0.2428 | 0.0983 | 0.01353 | 1350 | 6 | 20 intergenic |
| 4 | 12:40592225 | | 12 | 40592225 | G | A | 0.3189 | -0.1935 | 0.0901 | 0.03177 | 1724 | 7 | 0 intergenic |
| 5 | 12:40620815 | | 12 | 40620815 | A | G | 0.3177 | -0.1796 | 0.0902 | 0.04648 | 1724 | 7 | 0 intronic |
| 6 | 12:40621904 | rs4567538 | 12 | 40621904 | C | T | 0.3174 | -0.1841 | 0.0903 | 0.04138 | 1724 | 7 | 0 intronic |
| 7 | 12:40661403 | | 12 | 40661403 | T | A | 0.3152 | -0.1938 | 0.0878 | 0.02732 | 1724 | 7 | 1.2 intronic |
| 8 | 12:40665227 | | 12 | 40665227 | C | G | 0.3309 | -0.1965 | 0.0873 | 0.02446 | 1724 | 7 | 31.6 intronic |

https://pdgenetics.shinyapps.io/pd_progmetagwasbrowser/

- >4k samples
- 25k obs @ ~4 years
- Progression model predictions will be included
 - In the works ;-)
 - Basis for federated learning work

Next steps...

- UKB
 - Models across ancestry groups
 - Fit relevant models to PD data to predict comorbidities
 - Browser and model zoo go live
- PD progression
 - Predict PD specific outcomes
- PD-AMP
 - Integrate more molecular datatypes
 - Include more cohorts
- GenoML
 - Python release
 - Currently optimizing
 - Expand ecosystem
 - LSTM
 - Federated learning
 - Outlier detection
 - Unsupervised learning

Thanks!

If you have any questions, please email mike@data tecnica.com.

Shout outs to:

Andy Singleton - Good idea guy

Hirotaka Iwaka - Lead, clinical data and progression studies

Jeff Kim - Lead, open data and browser development

Mary Makarios - Lead, notebooks for multimodal analysis in AMP-PD

Hampton Leonard - Lead, production scale analyses in UKB (and everywhere else)

Cornelis Blauwendraat - Data management for literally everything

Faraz Faghri - Co-lead, GenoML ecosystem development

Biowulf and Terra contributors - Resources to get all these analyses done!



app.terra.bio/#workspaces

Click to go back, hold to see history

WORKSPACES

SEARCH WORKSPACES

hpc.nih.gov

WORKSPACES

Filter by Tags Access levels Billing project

Create a New Workspace +

AMP PD - Beta 1 - Getting Started - 2019v1beta0220 Last changed: Apr 29, 2019

The purpose of this workspace is to provide getting started information and no

AMP PD - Beta - Getting Started - 2019v1beta0220 Last changed: May 1, 2019

BIOWULF

HIGH PERFORMANCE COMPUTING AT THE NIH