# Adequate use of



INFORMATION IS BOTH POWER— AND PROFIT.

# Outline

Three faces on machine learning

1. Good

2. Bad

3. Ugly

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

1. **Good**

   • Intro and motivation

   • Filters in PB

2. **Bad**

   • Metrics of success

   • Adequate validation

3. **Ugly**

   • Optimization
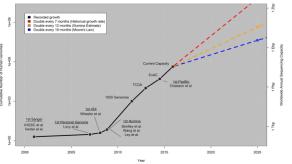
   • From RR to NN

# Part 1 – Good learners

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group
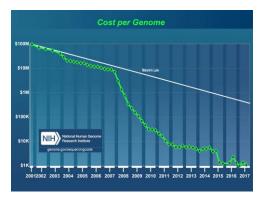
# More Pheno



# More Geno



The Cost of Sequencing a Human Genome. NIH.
https://www.genome.gov/27565109/the-cost-of-sequencing-a-human-genome/



Stephens, Z. D.et al. (2015). Big data: astronomical or genomical? *PLoS biology*, *13*(7), e1002195.

# More Env

- **UC Merced - GridMET**
- **NWS - NOAA**
- **NASA - GISS**
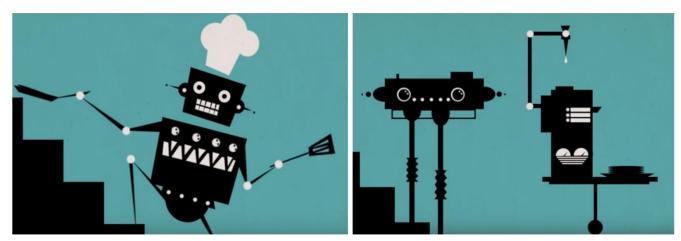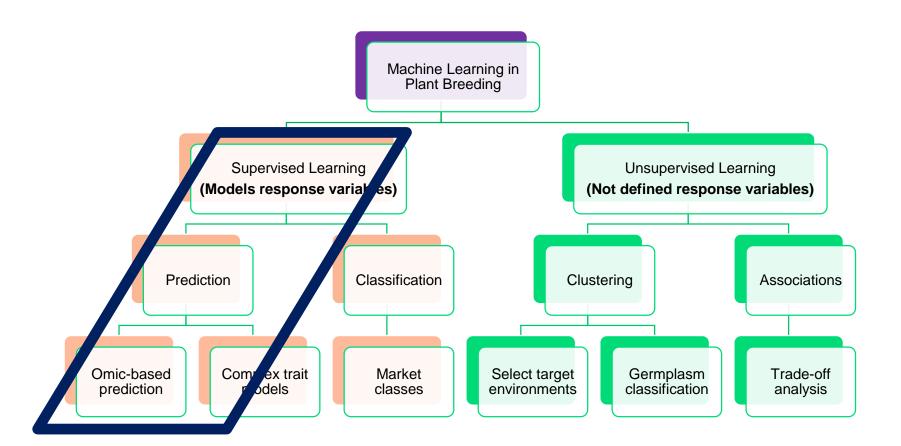- **Harmonized SoilDB**
- **USDA - SSURGO**

# More Computing

# Why is machine learning good for?

Good for solving **single** **well-defined** problem



Source: https://www.youtube.com/watch?v=MPR3o6Hnf2g

Machine Learning in Plant Breeding

- Supervised Learning **(Models response variables)**
  - Prediction
    - Omic-based prediction
    - Complex trait models
  - Classification
    - Market classes
- Unsupervised Learning **(Not defined response variables)**
  - Clustering
    - Select target environments
    - Germplasm classification
  - Associations
    - Trade-off analysis

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

CORTEVA agriscience

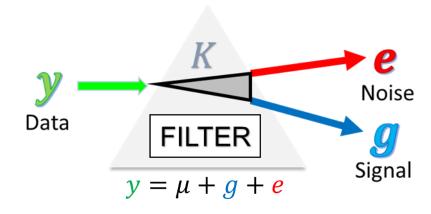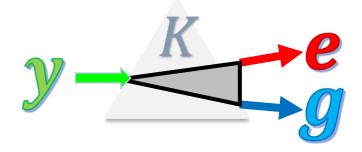# Key of supervised learning: FILTERING
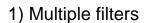
# Generalizations of simple filters



1) Multiple filters

2) Multi-response filters

# 1) When bother with multiple filters?



Some families were placed on unfavorable side of the field…

SoyNAM field,
Indiana 2014

# Pheno

# Spatial

# Genetics

# Residuals

**G, R**

# Multiple filters benefit separability of signals



Improvement in genetic accuracy

+0.12

Improvement in spatial accuracy

+0.02

Results derived from simulation of field with1500 plots, 1275 non-replicated entries + checks, from single trials

# 2) When bother with multi-response filters?

Filter one input at a time
(parallelizable)

Filter passing multiple input at once



$$\text{Cor}(g_A, g_B, ) \neq 0$$

**New information**

Genetic correlation table

|       | $g_1$      | $g_2$      | $g_3$      |
|-------|------------|------------|------------|
| $g_1$ | 1          | $\rho_{12}$ | $\rho_{13}$ |
| $g_2$ | $\rho_{21}$ | 1          | $\rho_{23}$ |
| $g_3$ | $\rho_{31}$ | $\rho_{32}$ | 1          |

EXTRA INFO

- **Bivariate model**

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}, \qquad y \sim N(0, G \otimes \Sigma_a + I \otimes \Sigma_e)$$

- Model equation

$$\begin{bmatrix} Z_1' \Sigma_e^{11} Z_1 + G^{-1} \Sigma_a^{11} & Z_1' \Sigma_e^{12} Z_2 + G^{-1} \Sigma_a^{12} \\ Z_2' \Sigma_e^{12} Z_1 + G^{-1} \Sigma_a^{12} & Z_2' \Sigma_e^{22} Z_2 + G^{-1} \Sigma_a^{22} \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} Z_1' (\Sigma_e^{11} y_1 + \Sigma_e^{12} y_2) \\ Z_2' (\Sigma_e^{22} y_2 + \Sigma_e^{12} y_1) \end{bmatrix}$$

- **Univariate vs bivariate solution for a given predictor**

EXTRA INFORMATION

$$g_1 = (Z_1' \Sigma_e^{11} Z_1 + G^{-1} \Sigma_a^{11})^{-1} (Z_1' \Sigma_e^{11} y_1)$$

$$g_1 | g_2 = (Z_1' \Sigma_e^{11} Z_1 + G^{-1} \Sigma_a^{11})^{-1} (Z_1' (\Sigma_e^{11} y_1 + \Sigma_e^{12} y_2) - (Z_1' \Sigma_e^{12} Z_2 + G^{-1} \Sigma_a^{12}) g_2)$$

# Sparse & Directed Filtering

Genetic correlation table

|       | $g_1$     | $g_2$     | $g_3$     |
|-------|-----------|-----------|-----------|
| $g_1$ | 1         | $\rho_{12}$ | $\rho_{13}$ |
| $g_2$ | $\rho_{21}$ | 1         | $\rho_{23}$ |
| $g_3$ | $\rho_{31}$ | $\rho_{32}$ | 1         |

## Sparse

|       | $g_1$     | $g_2$     | $g_3$     |
|-------|-----------|-----------|-----------|
| $g_1$ | 1         | $\rho_{12}$ | $-$       |
| $g_2$ | $\rho_{21}$ | 1         | $\rho_{23}$ |
| $g_3$ | $-$       | $\rho_{32}$ | 1         |

## Directed

|       | $g_1$     | $g_2$     | $g_3$     |
|-------|-----------|-----------|-----------|
| $g_1$ | 1         | $-$       | $-$       |
| $g_2$ | $\rho_{21}$ | 1         | $-$       |
| $g_3$ | $-$       | $\rho_{32}$ | 1         |

# Modern multi-response machinery

MVRF feature splits occur
for all responses together

Sikdar et al. 2019. Price Dynamics on Amazon
Marketplace: A Multivariate Random Forest
Variable Selection approach
http://dx.doi.org/10.2139/ssrn.3518690

Multi-response comes natural for DNNs

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

# Modern multi-response machinery

**MegaLMM** (SEM-like approach)

$$Y = F\Lambda + XB + U + E$$
$$F = X_F B_F + U_F + E_F$$
$$Y \sim N(0, \Psi_G + \Psi_R + \Lambda'\Lambda)$$
$$U \sim N(0, K, G)$$
$$G = \Psi_G + \Lambda'\Psi_{FG}\Lambda$$
$$E \sim N(0, I, R)$$
$$R = \Psi_R + \Lambda'\Psi_{FE}\Lambda$$

*Solved using Gibbs sampling

**A** Estimated prediction accuracy vs # traits

**B** Hours vs # traits (Rate 3, 2, 1, 0.5)

Method: MegaLMM, MCMCglmm, MTG2, phenix

"Classical" mixed model framework has been evolving in this area

**Tilde-Hat and Gauss-Seidel (THAGS)**

Statistical model
- $y_k = \mu_k + X_k \beta_k + e_k$, for trait k
- $\beta \sim N(0, I \otimes \Sigma_\beta)$
- $e \sim N(0, I \otimes \Sigma_e)$, $cov(e_i, e_j) = 0$
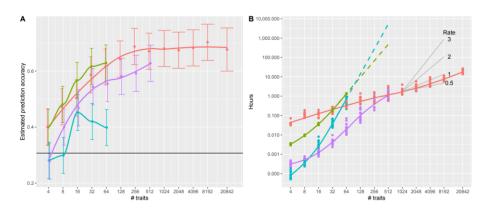
Coefficient updates: Multivariate Gauss-Seidel
- $\beta_j^{t+1} = (X_j'\Sigma_e^{-1}X_j + \Sigma_\beta^{-1})^{-1}\Sigma_e^{-1}X_j'(X_j\beta_j^t + e^t)$
- $e^{t+1} = e^t - X_j'(\beta_j^{t+1} - \beta_j^t)$

Covariance updates: Tilde-Hat method
- $\Sigma_{\beta(A,B)} = \frac{(y_A - \hat{\mu}_A)'(X_A\hat{\beta}_B) + (y_B - \hat{\mu}_B)'(X_B\hat{\beta}_A)}{n_A\sum_{j=1}^{J}\sigma_{X(A)j}^2 + n_B\sum_{j=1}^{J}\sigma_{X(B)j}^2}$
- $\Sigma_{e(i,i)} = \frac{(y_i - \hat{\mu}_i)'\hat{e}_i}{n_i - 1}$

**a) Computation time** — Bivariate, Trivariate; Time (log minutes) vs Number of observations (1000, 2000, 3000); THAGS, AIREML

**b) Estimator bias and accuracy** — H2 bias, GenCor bias, BV accuracy; THAGS, AIREML
$-0.02 (0.15)$  $0.01 (0.15)$  $0 (0.32)$  $0 (0.18)$  $0.91 (0.09)$  $0.91 (0.08)$
$Acc = Cor(\hat{g}, g)$
$Bias = \hat{\theta} - \theta$

Number of observations: 1000, 2000, 3000, 4000, 5000; Time (minutes) vs Number of traits

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

1. **Good**
   - Intro and motivation
   - Filters in PB

2. **Bad**
   - Metrics of success
   - Adequate validation

3. **Ugly**
   - Optimization
   - From RR to NN

(avoiding)

# Part 2 – Bad learning

**Hidden Technical Debt in Machine Learning Systems**

- How easily can an entirely new algorithmic approach be tested at full scale?
- What is the transitive closure of all data dependencies?
- How precisely can the impact of a new change to the system be measured?
- Does improving one model or signal degrade others?
- How quickly can new members of the team be brought up to speed?

# Chasing the right signal

- Breeding value (**GEBV**)
  - *Pattern*: <u>ADDITIVE</u> GENETICS
  - *Method*: GBLUP, BayesABC, LASSO
  - *Suits*: **RECYCLING**, **ADVANCEMENT**

- Genomic value (**EGV**)
  - *Pattern*: ANY GENETICS
  - *Method*: RKHS, DNN, Random Forest
  - *Suits*: **ADVANCEMENT**

Breeding pipeline

**Advancement**
**Recycling**
**Incorporation**

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

# Defining the problem: <u>Metrics for success</u>

1. Scientist (<u>**why**</u>): to define the problem mathematically (<span style="color:red">easy to get it wrong</span>)

2. Metric (<u>**what**</u>): Correlations, MSPE, Jaccard index, Accuracy, Success (<span style="color:red">1|0</span>)

3. Testing (<u>**how**</u>): Simulation or cross-validation (CV) on real data? WF vs AF?

   <span style="color:red">How to design an adequate cross-validation???</span>

**Alencar.Xavier@Corteva.com**
**Corteva Biostatistics, Methods group**

**CORTEVA**
agriscience

# Two metrics, two different answers



Grain yield on the SoyNAM population: Populations 1:8 predicting populations 9:15

Alencar.Xavier@Corteva.com
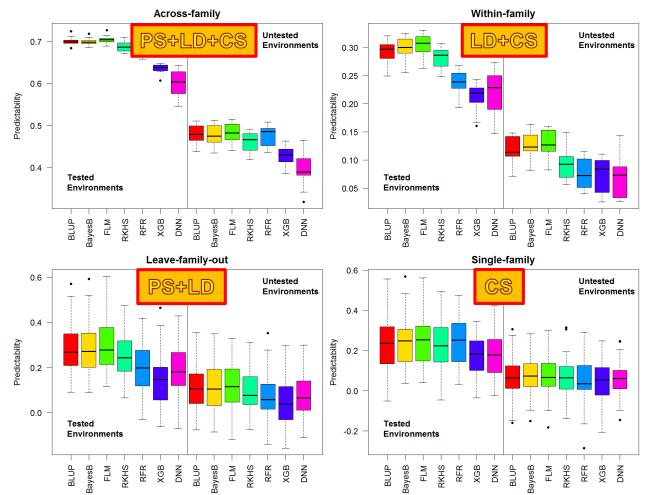Corteva Biostatistics, Methods group

# Testing machines for different scenarios of genomic prediction

| | Genotype | Environment | Prediction Difficulty |
|---|---|---|---|
| **CV00** | New | New | ***** |
| **CV0** | Observed | New | *** |
| **CV1** | New | Observed | *** |
| **CV2** | Observed | Observed | * |

Adapted from Crossa et al. (2017) doi.org/10.1016/j.tplants.2017.08.011

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

CORTEVA agriscience

**CV scheme**

Type of information captured by SNP
- Population structure (PS)
- Linkage disequilibrium (LD)
- Cosegregation / Haplotype (CS)

**Across-family**

PS+LD+CS

Untested Environments

Tested Environments

BLUP BayesB FLM RKHS RFR XGB DNN BLUP BayesB FLM RKHS RFR XGB DNN

**Within-family**

LD+CS

Untested Environments

Tested Environments

BLUP BayesB FLM RKHS RFR XGB DNN BLUP BayesB FLM RKHS RFR XGB DNN

**Leave-family-out**

PS+LD

Untested Environments

Tested Environments

BLUP BayesB FLM RKHS RFR XGB DNN BLUP BayesB FLM RKHS RFR XGB DNN

**Single-family**

CS

Untested Environments

Tested Environments

BLUP BayesB FLM RKHS RFR XGB DNN BLUP BayesB FLM RKHS RFR XGB DNN

**SoyNAM data**
ES: 2012 (7 loc)
PS: 2013 (4 loc)
#Fam = 40
Genos = 5600
SNPs = 4300
Obs: 3k-5k obs/loc

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

1. **Good**

   - Intro and motivation

   - Filters in PB

2. **Bad**

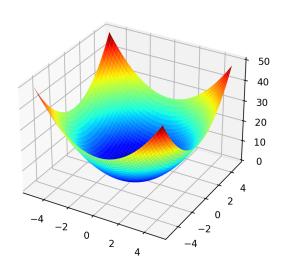   - Metrics of success

   - Adequate validation

3. **Ugly**

   - Optimization

   - From RR to NN

# Part 3 – Ugly learning

# Divergency in philosophy

- In quantitative genetics:
  - Parameters: Variances + Regression coefficients
  - Function: Likelihood (complex and convex)
  - Tuning: Generally not needed
  - **Method**: First order (EM, MCMC), second order (AI, MIVQUE)

- In machine learning:
  - Parameters: Regression coefficients
  - Function: MSE, L2 (simple)
  - Tuning: Cross validations, need secondary objective function
  - **Method:** First order: coordinate & gradient descent

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

# Solving: $y = Xb + e$

$Finding \rightarrow argmin(e'e + \lambda b'b)$

I've created a monster!!

- Coordinate descent

  (Use diagonals of LHS)

$$b_j^{t+1} = \frac{x_j'(y - X_{-j}b_{-j})}{x_j'x_j + \lambda}$$

**Used for WGR (RR, BayesA)**

glmnet, BGLR, bWGR, GS3

- Gradient descent

  (Does not build LHS)

$$b^{t+1} = b^t - \frac{2r}{n}[X'(y - Xb^t) + \lambda b^t]$$

**Used for Deep Neural Nets**

TF/Keras, PyTorch, MXNet, h2o

- Second order

  (Builds entire LHS)

$$b = (X'X + \lambda)^{-1}(X'y)$$

**Used for underline{everything else}**

ASREML, lme4, SAS, BLUPF90

**Alencar.Xavier@Corteva.com**
**Corteva Biostatistics, Methods group**

**CORTEVA** agriscience™

# Coordinate descent of a RRBLUP

$b_0 \rightarrow b_1$

$y = Xb + e$

$b_1 \rightarrow b_2$

$\dots \rightarrow b_3$

$$f(b_{RR}) \rightarrow \text{argmin}\left(\sum_{i=1}^{n} e^2 + \lambda \sum_{j=1}^{p} b^2\right)$$

$$SS = e'e + \lambda(b'b)$$
$$SS = (y - xb)'(y - xb) + \lambda(b'b)$$
$$SS = y'y + (xb)'(xb) - 2(y'xb) + \lambda(b'b)$$
$$SS = y'y + b'(x'x)b - 2(y'xb) + \lambda(b'b)$$

$$\frac{\partial SS}{\partial b} = y'y + 2(x'x)b^2 - 2(y'x)b + 2\lambda b^2$$

$$0 = 2b(x'x) - 2(y'x) + 2b\lambda$$

$$-2b(x'x + \lambda) = -2(y'x)$$

$$b = \frac{-2(y'x)}{-2(x'x + \lambda)} = \frac{y'x}{x'x + \lambda}$$

**Marker effects**

**Fitness**

Alencar.Xavier@Corteva.com
Corteva Biostatistics, Methods group

CORTEVA
agriscience

# Gradient descent of RRBLUP

Scatter plot between true solution and solution at iteration x

**Small** learning rate

**Large** learning rate

- Model (Ridge)

$$y = X\beta + e$$

- Gradient descent

$$\beta^{t+1} = \beta^t - \alpha\nabla$$

- Where
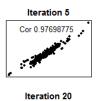
$$\alpha = \text{Learning rate}$$

$$f = (y - Xb)'(y - Xb) + \lambda b'b$$

$$\nabla = f' = \frac{\partial f}{\partial b} = -2X'(y - Xb) + 2\lambda b$$

- Thus

$$\nabla = n^{-1}(-2X'e + 2\lambda\beta^t) = -2n^{-1}(X'e - \lambda\beta^t)$$

$$\beta^{t+1} = \beta^t + \frac{2\alpha}{n}(X'e - \lambda\beta^t)$$

| | | | | | |
|---|---|---|---|---|---|
| **Iteration 5** Cor 0.97698775 | **Iteration 10** Cor 0.99172954 | **Iteration 15** Cor 0.99781079 | **Iteration 5** Cor 0.98381857 | **Iteration 10** Cor 0.99978843 | **Iteration 15** Cor 0.99999994 |
| **Iteration 20** Cor 0.99959892 | **Iteration 25** Cor 0.99995334 | **Iteration 30** Cor 0.99999695 | **Iteration 20** Cor -0.14664758 | **Iteration 25** Cor 0.32280437 | **Iteration 30** Cor -0.35354069 |
| **Iteration 35** Cor 0.99999991 | **Iteration 40** Cor 1 | **Iteration 45** Cor 1 | **Iteration 35** Cor 0.38514267 | **Iteration 40** Cor -0.41341605 | **Iteration 45** Cor 0.4252109 |

Deep Neural Network

# From Ridge Regression to Deep Neural Net

Models illustrated without intercept

Linear model

$$y = Xb + e$$

PLS/PCR model

$$y = (XB_1)b_2 + e$$

NN model

$$y = \alpha(XB_1)b_2 + e$$

Deep NN model

$$y = \alpha(\alpha(XB_1)B_2)b_3 + e$$

$\alpha$ = activation function

**Alencar.Xavier@Corteva.com**
**Corteva Biostatistics, Methods group**

# The ~~Scary~~ Deep Neural Network

$$Y = \alpha(\alpha(XB_1)B_2)B_3 + E$$

$H_1$   $H_2$   $H_3$

- **Fit layers**
  - $H_1 = \alpha(XB_1)$
  - $H_2 = \alpha(H_1B_2)$
  - $H_3 = H_2B_3$

- **Compute gradients** (aka. get residuals)
  - $E_3 = Y - H3$
  - $E_2 = \alpha(E_3B_3')$
  - $E_1 = \alpha(E_2B_2')$

- **Update coefficient**
  - $B_1 = B_1 - \gamma\left(\frac{2r_1}{n}\left[X'E_1 - \lambda B_1\right]\right)$
  - $B_2 = B_2 - \gamma\left(\frac{2r_2}{n}\left[H_1'E_2 - \lambda B_2\right]\right)$
  - $B_3 = B_3 - \gamma\left(\frac{2r_3}{n}\left[H_2'E_3 - \lambda B_3\right]\right)$

**Top-to-bottom code ~ 30 lines**

**Alencar.Xavier@Corteva.com**
**Corteva Biostatistics, Methods group**

CORTEVA agriscience

# Thank you for your attention!

**Remarks**:

1) ML can be a powerful tool for plant breeding

2) Cross-validation must be carefully design to understand the machinery

3) The scary black boxes are usually simple methods

# Questions??

*Alencar Xavier*

Alencar.Xavier@Corteva.com

# Some free lit to look up!