



# Machine learning-based AI applied to breeding

**Alencar Xavier**

Breeding Analyst at Corteva  
Adjunct professor at Purdue

# Adequate use of



# Outline

## 1. Introduction

- More data
- Branching ML

## 2. Machines

- Filters
- Engines

## 3. Analytics

- Target  $G \times E \times M$
- Validation
- Cases of study

## 4. Conclusion

# 1. Introduction

- More data
- Branching ML

## 2. Machines

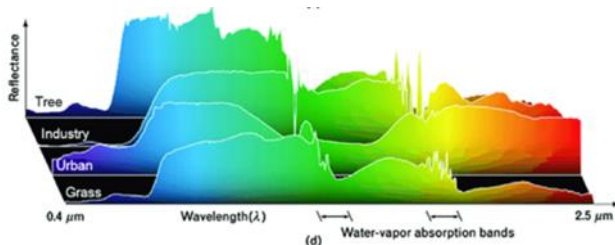
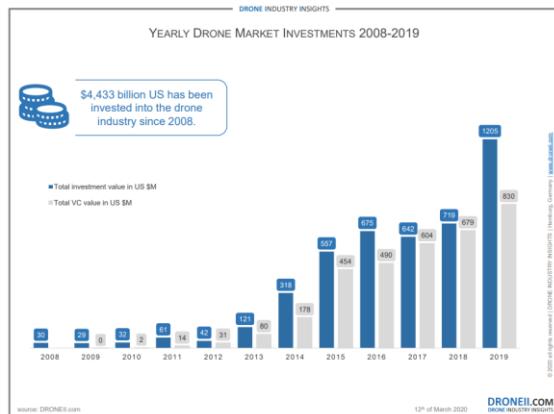
- Filters
- Engines

## 3. Analytics

- Target  $G \times E \times M$
- Validation
- Cases of study

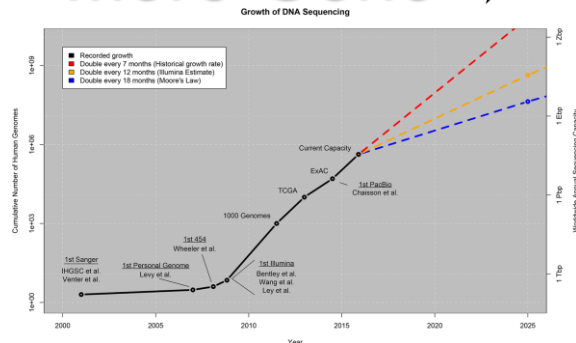
## 4. Conclusion

# More Pheno

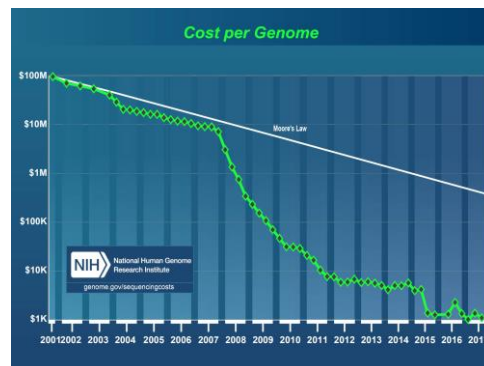


<https://www.mdpi.com/2076-3417/12/5/2570>

# 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 genomic? *PLoS biology*, 13(7), e1002195.

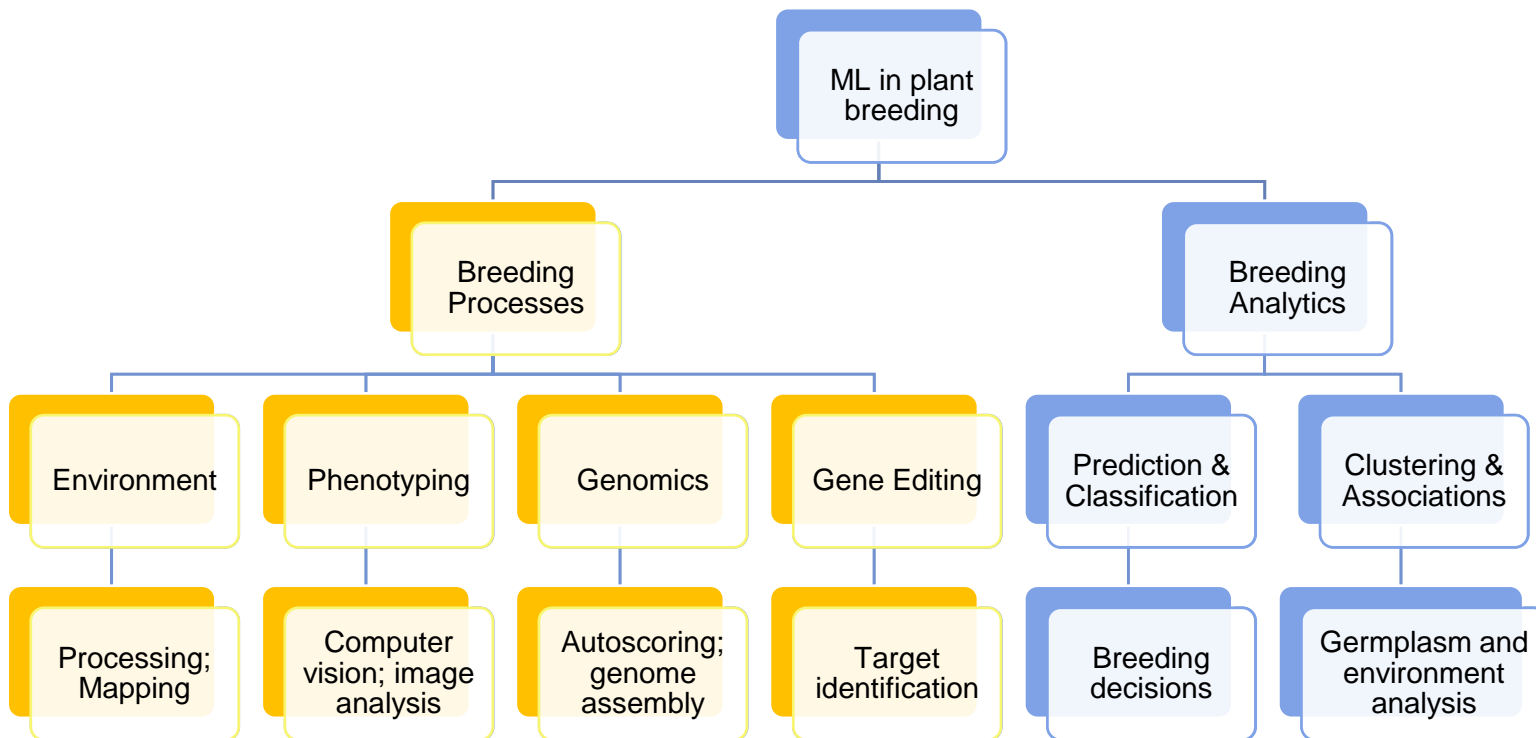
# More Env



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

# More Computing



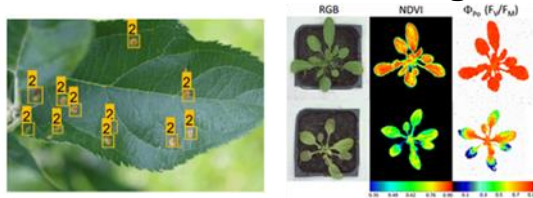


# ML in breeding processes

Enhancing databases, automating lab tasks field work

## phenotyping

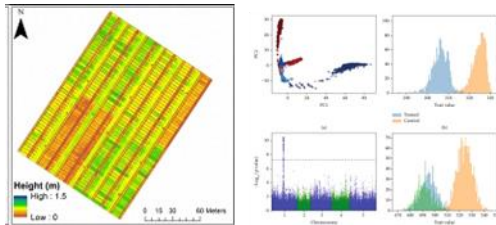
### Disease, stress scoring



<https://www.mdpi.com/2673-2688/2/3/26>  
<https://www.biomedcentral.com/collections/phenomics>

### Phenotype automation

(e.g., plant height, identify new traits)



<https://www.mdpi.com/2072-4292/8/12/1031>  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7706325/>

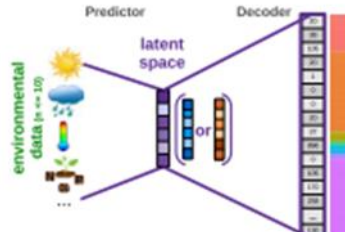
## environment

### Mapping / zoning



<https://www.publish.csiro.au/cp/CP14007>

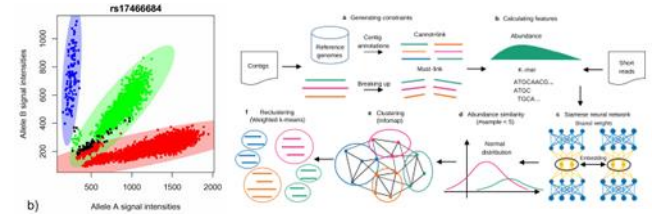
### Latent weather, soil



<https://doi.org/10.1093/bioinformatics/btaa971>

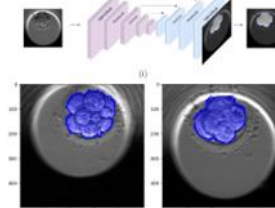
## biotech

### SNP calls, genome assembly



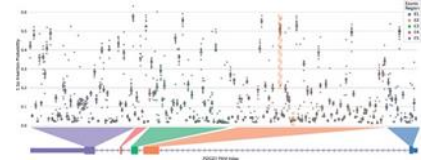
<https://doi.org/10.1186/1753-6561-3-s7-s58>  
<https://www.nature.com/articles/s41467-022-29843-y>

### Embryo rescue DH production

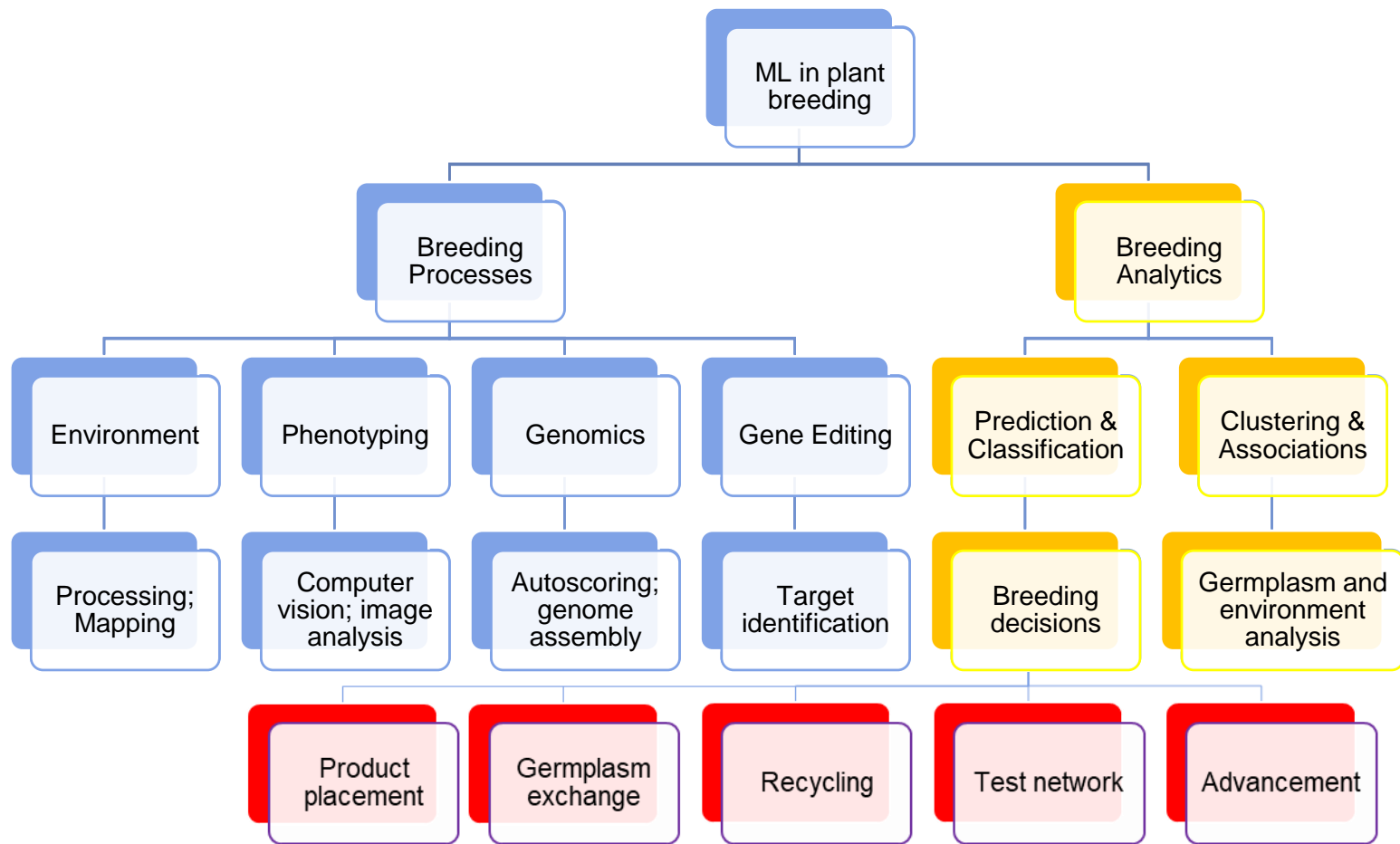


<https://www.nature.com/articles/s41598-022-06336-y>

### Gene editing targets



<https://doi.org/10.1093/bioinformatics/btab268>





## 1. Introduction

- More data
- Branching ML

## 2. Machines

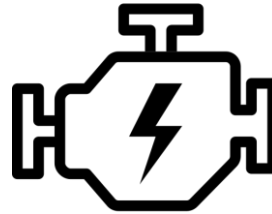
- Filters
- Engines

## 3. Analytics

- Target  $G \times E \times M$
- Validation
- Cases of study

## 4. Conclusion

# Machine Learning Engines



A new approach fits multivariate genomic prediction models efficiently

Alencar Xavier<sup>1,2\*</sup> and David Habier<sup>1†</sup>



Walking through the statistical black boxes of plant breeding

Alencar Xavier<sup>1</sup> · William M. Muir<sup>2</sup> · Bruce Craig<sup>3</sup> · Katy Martin Rainey<sup>1</sup>

Efficient Estimation of Marker Effects in Plant Breeding

Alencar Xavier<sup>1,2\*</sup>  
<sup>1</sup>Corteva Agriscience, 8333 N 132nd Ave, Johnston, IA, and <sup>2</sup>York University, 4700 Keele St, West Lutherville, IN  
DOI: 10.1002/2019.00000.00000.00000

Technical nuances of machine learning: implementation and validation of supervised methods for genomic prediction in plant breeding

Alencar Xavier<sup>1\*</sup>

Impact of Genomic Prediction Model, Selection Intensity, and Breeding Strategy on the Long-Term Genetic Gain and Genetic Erosion in Soybean Breeding

Éder David Borges da Silva<sup>1,2\*</sup>, Alencar Xavier<sup>1,2\*</sup> and Marcos Ventura Faria<sup>1</sup>

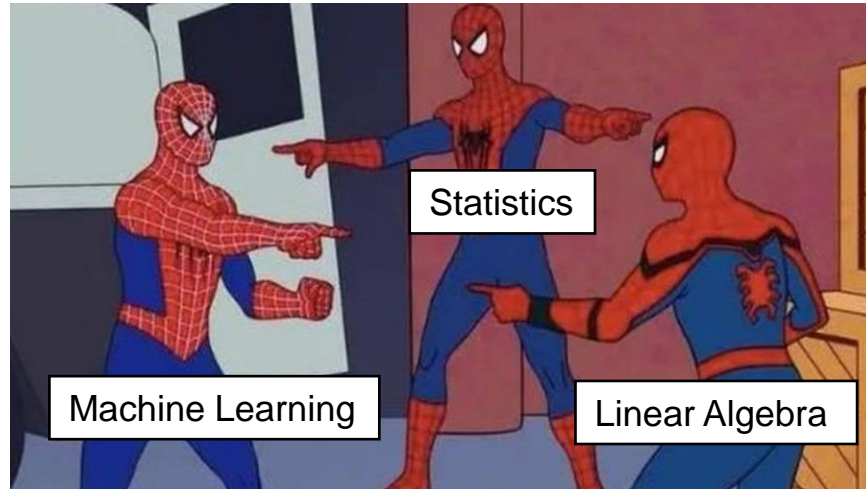
Using unsupervised learning techniques to assess interactions among complex traits in soybeans

Alencar Xavier · Benjamin Hall · Shaun Casteel · William Muir · Katy Martin Rainey

Article

Joint Modeling of Genetics and Field Variation in Plant Breeding Trials Using Relationship and Different Spatial Methods: A Simulation Study of Accuracy and Bias

Éder David Borges da Silva<sup>1,2,\*</sup>, Alencar Xavier<sup>1,2,\*</sup> and Marcos Ventura Faria<sup>1,2</sup>



# Key idea of supervised learning: FILTERING

Simple filter

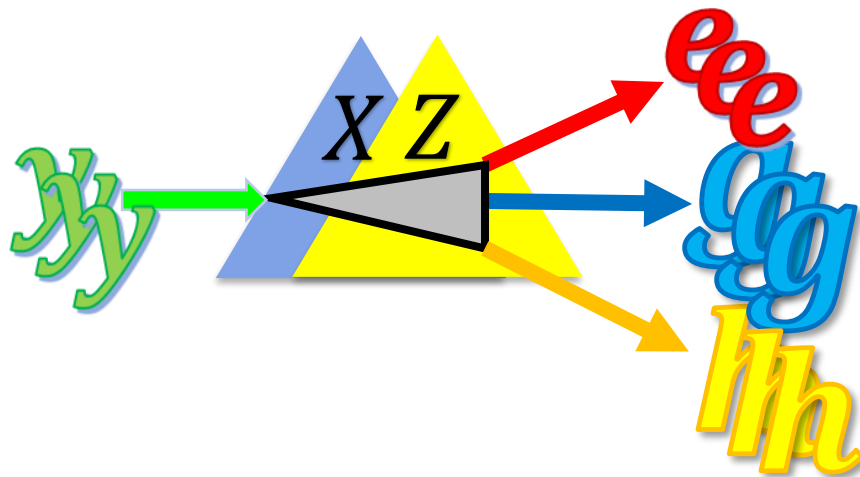
$$y = g + e$$

Multiple filters

$$y = g + h + e$$

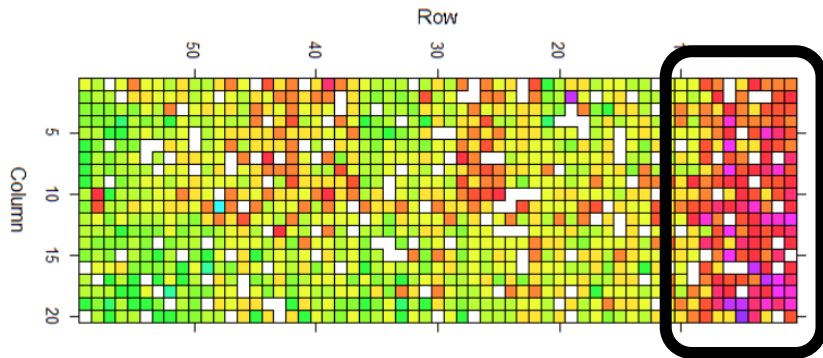
Multi-task filter

$$Y = G + H + E$$

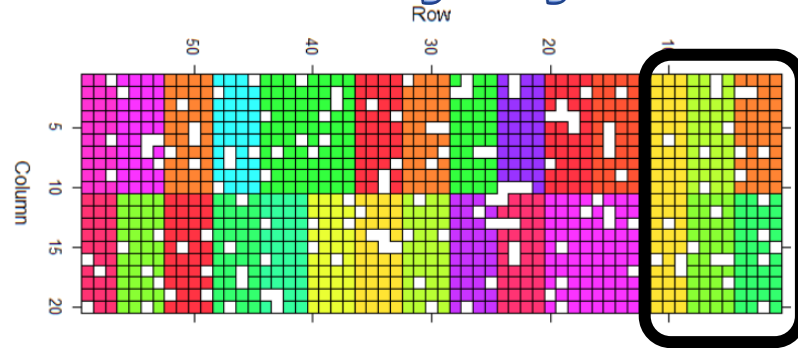


# Why bother with multiple filters?

## Field variation



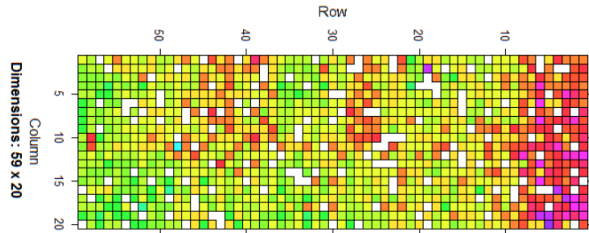
## Family layout



Some families were placed on unfavorable side of the field...

SoyNAM field,  
Indiana 2014

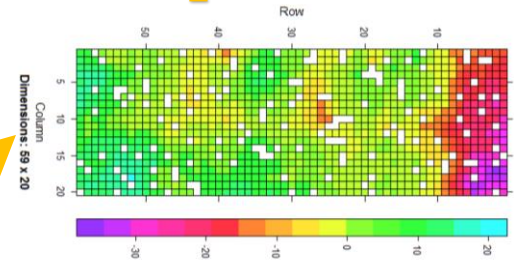
# Pheno



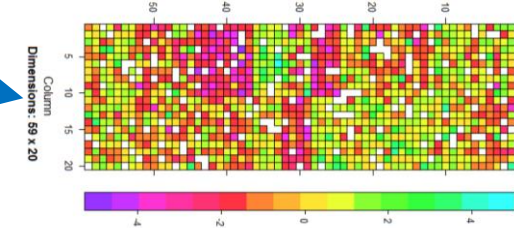
G, R

Separation of  
tangled signals!

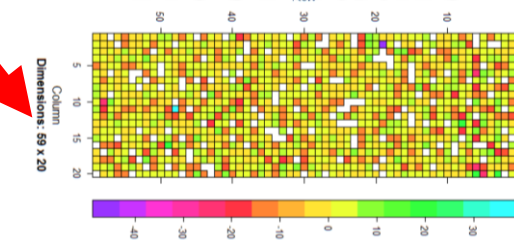
## Spatial



## Genetics



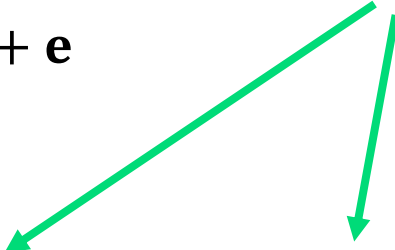
## Residuals



# Why bother with multi-task filters?

Simple (bivariate) model:

**INFORMATION GAIN**

$$y = g + e$$
$$Var \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \sigma_{a_1}^2 & \sigma_{a_{12}} \\ \sigma_{a_{12}} & \sigma_{a_2}^2 \end{bmatrix} + \begin{bmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{bmatrix}$$


# Why bother with multi-task filters?

$$y = Zg + e, \quad y \sim N(0, V)$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} Z_1 & 0 \\ 0 & Z_2 \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$$

- Covariance structure

$$V = G \otimes \Sigma_a + I \otimes \Sigma_e = G \otimes \begin{bmatrix} \sigma_{a_1}^2 & \sigma_{a_{12}} \\ \sigma_{a_{12}} & \sigma_{a_2}^2 \end{bmatrix} + I \otimes \begin{bmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{bmatrix}$$

- Model equation

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

- Univariate vs bivariate

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

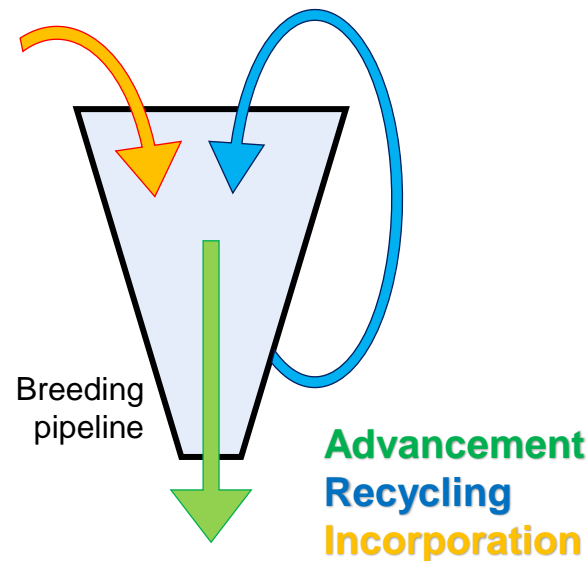
$$g_1 | g_2 = (Z_1' \Sigma_e^{11} Z_1 + G^{-1} \Sigma_a^{11})^{-1} (Z_1' (\Sigma_e^{11} y_1 + \Sigma_e^{12} y_2) - (Z_1' \Sigma_e^{12} Z_2 + G^{-1} \Sigma_a^{12}) g_2)$$

**INFORMATION GAIN**

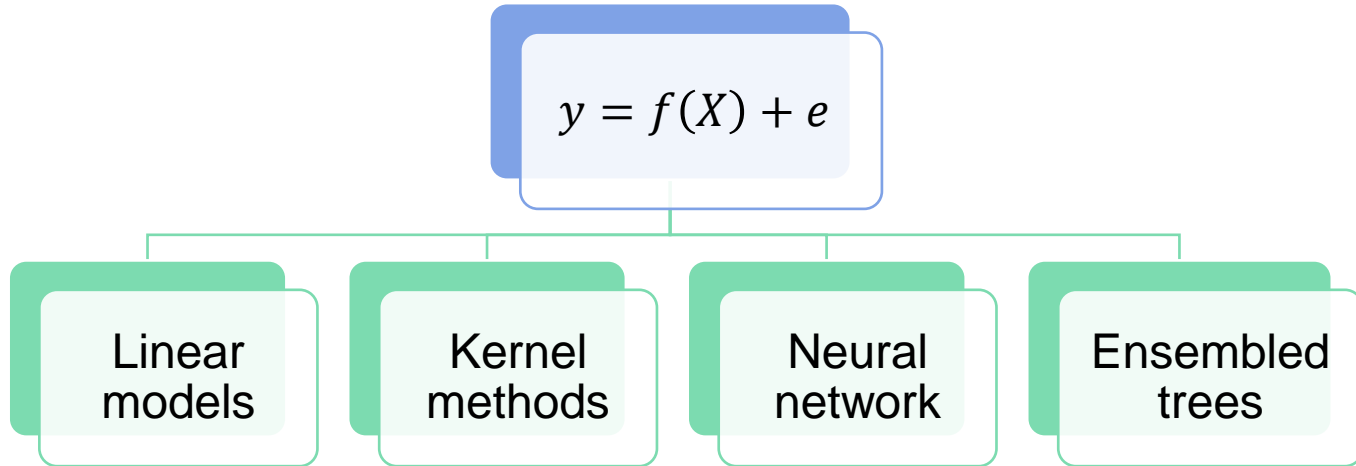


# Does the choice of filter matter?

- **ADDITIVE LINEAR FILTERS** (GEBV)
  - *Pattern:* ADDITIVE GENETICS - **heritable**
  - *Method:* GBLUP, RIDGE, LASSO
  - *Suits:* **RECYCLING**, **ADVANCEMENT**
- **NON-LINEAR FILTERS** (EGV)
  - *Pattern:* **ANY GENETIC SIGNAL**
  - *Method:* RKHS, DNN, Random Forest
  - *Suits:* **ADVANCEMENT**, **PRODUCT PLACEMENT**



# Main classes of learners



# Solving: $y = Xb + e$

*Finding  $\rightarrow \operatorname{argmin}(e'e + \lambda b'b)$*

- Coordinate descent

(Use diagonals of LHS)

$$\hat{b}_j^{t+1} = \frac{x_j'(y - X_{-j}\hat{b}_{-j})}{x_j'x_j + \lambda}$$

- Gradient descent

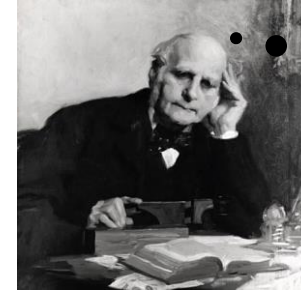
(Does not build LHS)

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

- Second order

(Builds entire LHS)

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



I've created a monster!!

Used for  $p \gg n$  solvers

glmnet, BGLR, bWGR, GS3

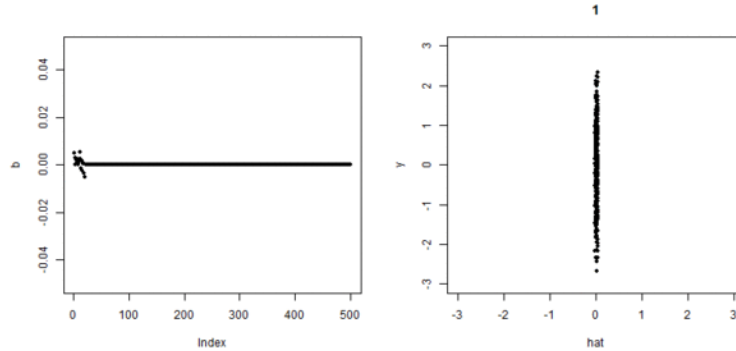
Used for Deep Neural Nets

TensorFlow Keras, PyTorch, MXNet

Used for everything else

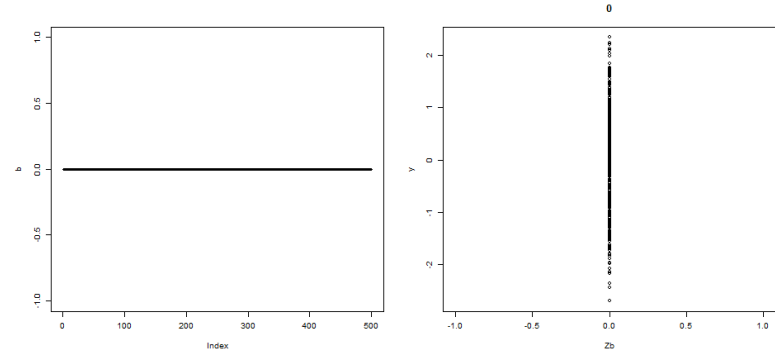
ASREML, lme4, SAS

# Coordinate descent



$$\hat{b}_j^{t+1} = \frac{x_j'(y - X_{-j}\hat{b}_{-j})}{x_j'x_j + \lambda}$$

# Gradient descent



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

What about the deep learning? 🤖

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

i.e., just a “stack of solvers”

## Data > Method

*Unnecessarily complex analysis should not be used as a foil to disguise lower quality datasets*

Kruuk ([2004](#) apud Walsh and Lynch [2018](#))

## 1. Introduction

- More data
- Branching ML

## 2. Machines

- Filters
- Engines

## 3. Analytics

- Target  $G \times E \times M$
- Validation
- Cases of study

## 4. Conclusion

# Analytics



# “Breeding objective”

- Set of traits of interest (**TOI**)

bred into a

- Target population of genotypes (**TPG**)

for a given

- Target population of environments (**TPE**)



# TPE, TPG, TPM

- **Target population of environments (TPE)**

- Influences accuracies via GxE correlation
- Which environments should I be able to predict?

- **Target population of genotypes (TPG)**

- Influences accuracies via genetic relationship
- Which genetics should I be able to predict?

- **Target population of management (TPM)**

- Herein nested in TPE

**From QTLs to Adaptation Landscapes: Using Genotype-To-Phenotype Models to Characterize G×E Over Time**

Daniela Bustos-Korts<sup>1\*</sup>, Marcos Malosetti<sup>1</sup>, Karine Chenu<sup>2</sup>, Scott Chapman<sup>3,4</sup>, Martin P. Boer<sup>1</sup>, Bangyou Zheng<sup>2</sup> and Fred A. van Eeuwijk<sup>1\*</sup>

**What Should Students in Plant Breeding Know About the Statistical Aspects of Genotype × Environment Interactions?**

Fred A. van Eeuwijk,\* Daniela V. Bustos-Korts, and Marcos Malosetti

**An Equation to Predict the Accuracy of Genomic Values by Combining Data from Multiple Traits, Populations, or Environments** 

Yvonne C J Wientjes , Piter Bijma, Roel F Veerkamp, Mario P L Calus

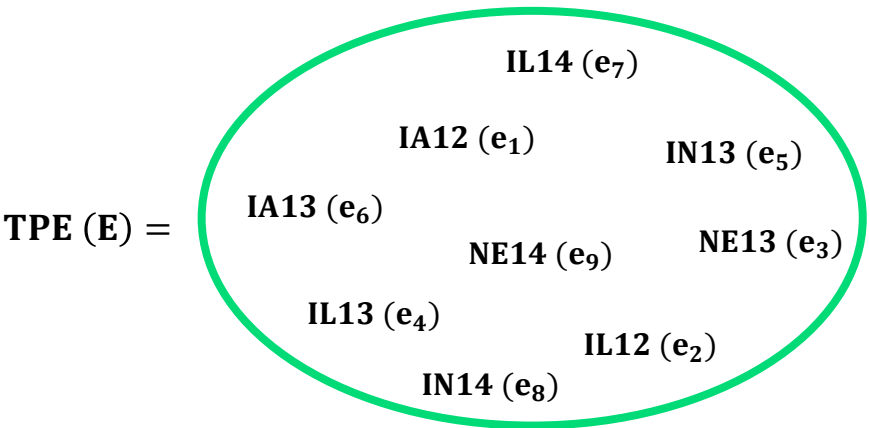
Genetics, Volume 202, Issue 2, 1 February 2016, Pages 799–823,  
<https://doi.org/10.1534/genetics.115.183269>

# TPE

- Any given trial happens in each environment-management combination, that is sample of much larger population:

$$e_i \in E$$

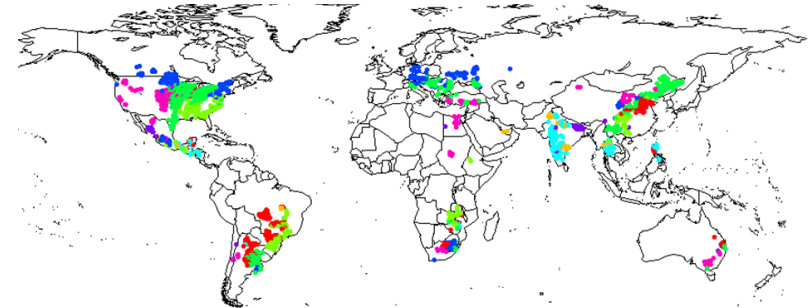
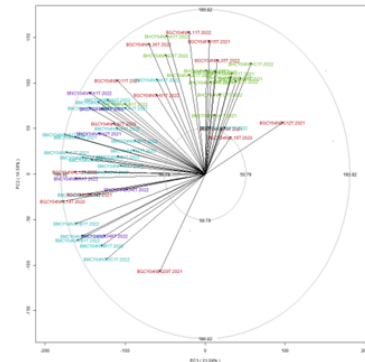
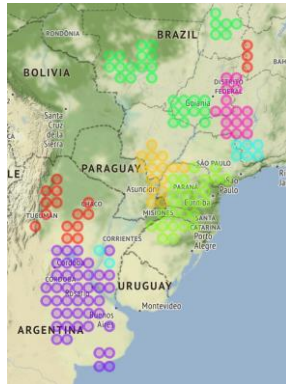
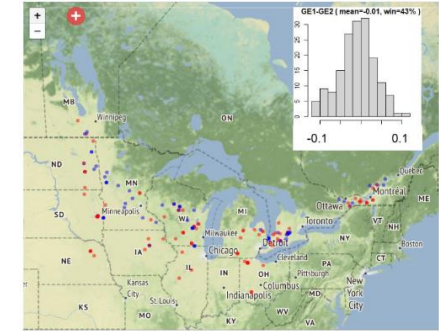
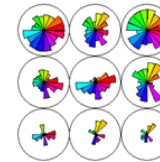
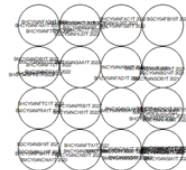
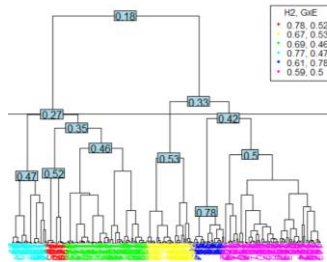
That is:



$$E[MET] = E[TPE]?$$

$$\begin{bmatrix} y_{e_i} \\ y_{e_j} \\ \underline{g_E} \end{bmatrix} = \begin{bmatrix} \sigma_{g(e_i)}^2 + \sigma_{\epsilon(e_i)}^2 & \sigma_{g(e_i, e_j)} & \sigma_{g(e_i, E)} \\ \sigma_{g(e_j, e_i)} & \sigma_{g(e_j)}^2 + \sigma_{\epsilon(e_j)}^2 & \sigma_{g(e_j, E)} \\ \sigma_{g(E, e_i)} & \sigma_{g(E, e_j)} & \sigma_{g(E)}^2 \end{bmatrix}$$

# NOTE: GxExM patterns within TPE are largely assessed using different methods of ML



# TPG + TPE

- Accuracy ([Wientjes et al 2016](#)) = correlation( true signal, estimated signal ),
- It is a function of heritability, GxE, representativeness of the calibration set
- For:

$$y = g + e,$$
$$\text{var}(y) = V, \quad \text{var}(g) = G$$

Then accuracy is

$$a_i = \text{cor}(g_i, \hat{g}_i) = \frac{\text{cov}(g_i, \hat{g}_i)}{\text{var}(g_i)\text{var}(\hat{g}_i)} = \frac{\text{var}(\hat{g}_i) r_{\text{GxE}}^2}{\text{var}(g_i)\text{var}(\hat{g}_i)} = r_{\text{GxE}}^2 \sqrt{\frac{G_{i,y} V^{-1} G_{y,i}}{G_{i,i}}}$$

Thus, we **know** how much signal to expect in any given prediction

# Validation schemes

## 1) CV type – Test intent

- **Random CV** = Upper-bound predictive potential
- **Leave-one-out** = Assess structured scenarios (e.g., geography-out, year-out)
- **Holdout** = Reproduce true applications (e.g., predict individuals from upcoming)

## 2) TPE/TPG relation

	Genotype	Environment	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

## 3) Signal availability

Genetic information available in different cross-validation setups

- ***Intra-family***: Linkage\*
- ***Within-family***: Linkage and LD
- ***Across-family***: Relationships\*\*, Linkage and LD
- ***Leave-family-out***: Relationships and LD
- ***Untested environments***: Same as above x ( GxE )

# Validation metrics

- **Correlations**

- Most common metrics in breeding (e.g., predictability)
- Pertinent to **ranking** and selection of complex traits

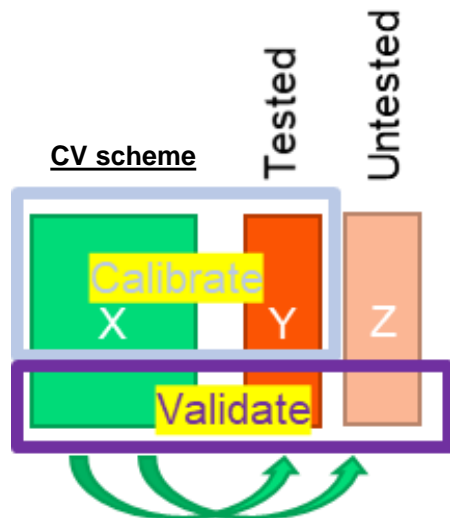
- **Prediction error**

- Utilized when the predicted values must be as close as possible to original scale
- Pertinent to risk prediction (e.g., disease risk)

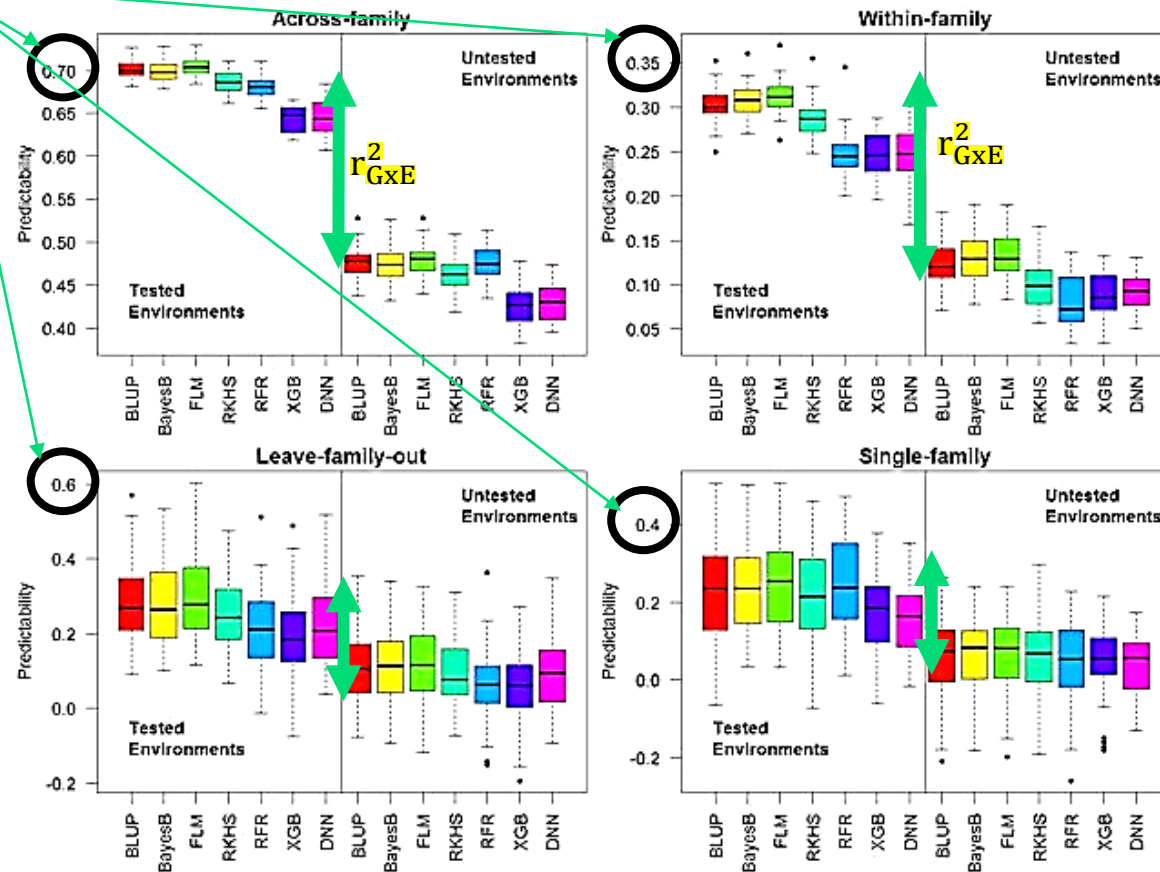
- **Success**

- Accommodate complex or subjective criteria, independent or otherwise
- Pertinent to decision involving data from multiple sources (e.g., advancement)

Amount of signal that can be captured in different structures



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

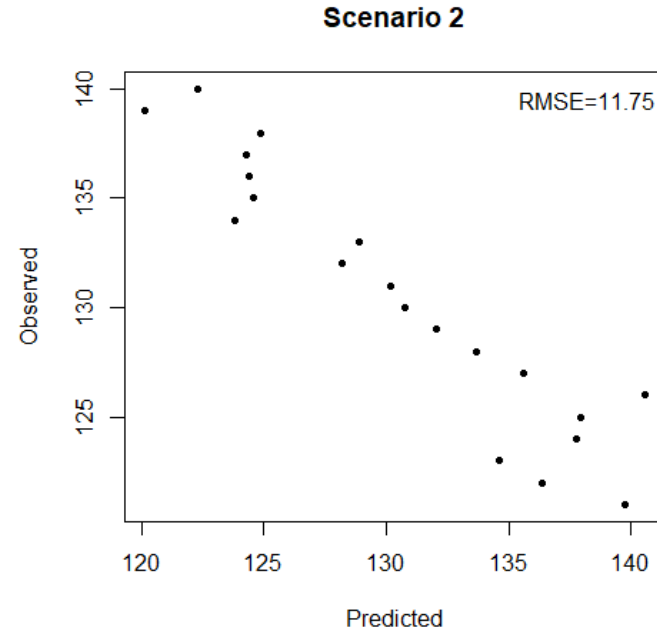
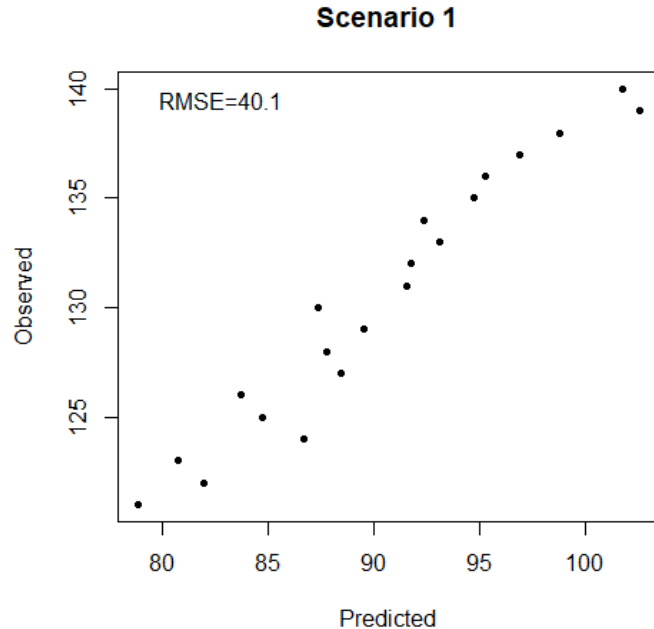


# Case of study





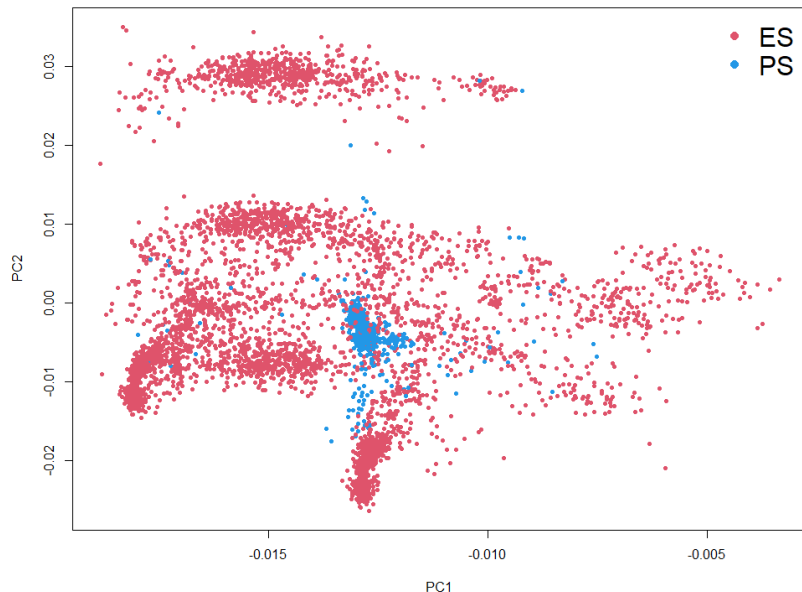
# Evaluation criterion



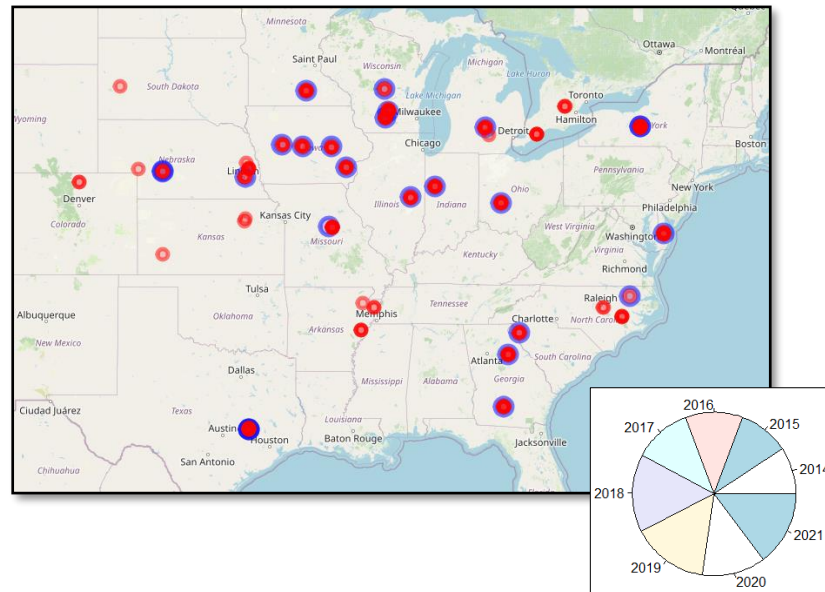
# 2022 G2F GxE prediction competition

TPG

Population structure



TPE



# What was modeled?

$$y|E_i = \mu_i + g|E_i \quad (\text{Two FILTERS})$$

Phenotype @  $i^{\text{th}}$  Loc =  $i^{\text{th}}$  Loc Mean + Genetic effect @  $i^{\text{th}}$  Loc

- The winning approach:
  - Predict location means using mixed model and random forest
  - Predict genetic performance with index from multi-response based on TPE/TPG

# 2022 G2F GxE prediction competition

## Realized results

Team Name	Within RMSE
CLAC	2.329
igorkf	2.345
phenomaize	2.374
UCD_MegaLMM	2.387
CGM	2.391
breedingteam	2.398
Purdue	2.402
SmAL	2.425
ML_APT	2.472
MPB_Group	2.544

## Ranking with alternative metrics

Team Name	Cor Within Loc	Team Name	Cor Across Loc
CLAC	0.357	breedingteam	0.650
CGM	0.353	DataJanitors	0.644
MPB_Group	0.342	CLAC	0.631
UCD_MegaLMM	0.338	Purdue	0.631
SmAL	0.285	UCD_MegaLMM	0.628
DeepCropVision	0.281	phenomaize	0.617
CropEnthusiast	0.279	igorkf	0.600
AllModelsAreWrong	0.272	CGM	0.587
DataJanitors	0.256	SmAL	0.586
supermanwasd	0.243	AllModelsAreWrong	0.575

Source: Jacob Washburn, Jose Ignacio Varela, Alencar Xavier

## 1. Introduction

- More data
- Branching ML

## 2. Machines

- Filters
- Engines

## 3. Analytics

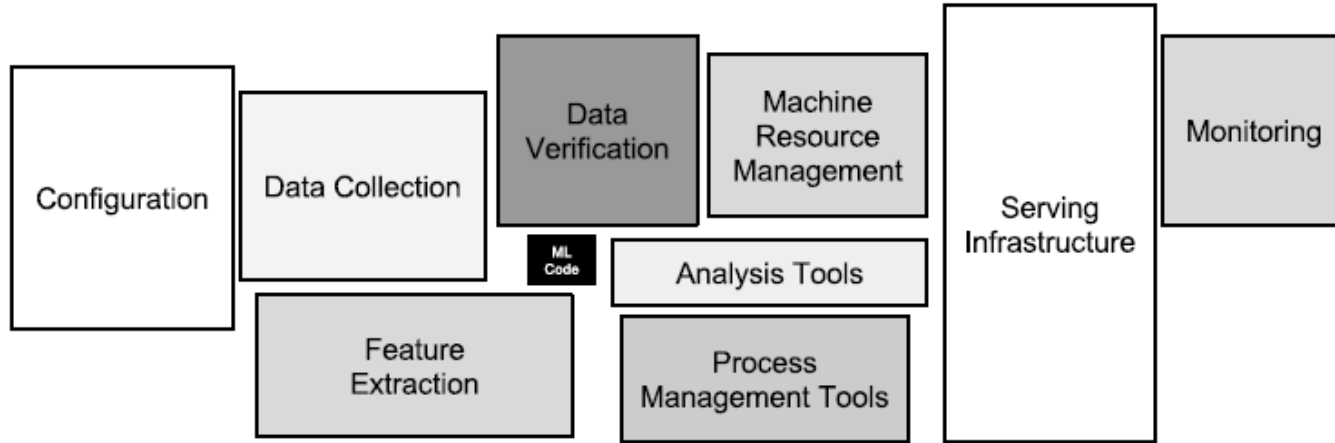
- Target  $G \times E \times M$
- Validation
- Cases of study

## 4. Conclusion

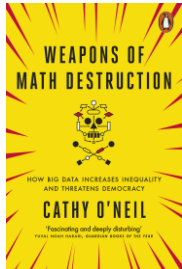
# There is more to ML than proof of concepts using cross-validations

## Hidden Technical Debt in Machine Learning Systems

doi/10.5555/2969442.2969519



- 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?



# Thank you for your attention!

## Final remarks:

- 1) Plant breeding uses machine learning for multiple purposes in processes and analytics
- 2) Filter settings are important to maximize signal, but it is less important than data
- 3) Validation metrics and validation schemes matter to design meaningful models

## Questions??

***Alencar Xavier***

[Alencar.Xavier@Corteva.com](mailto:Alencar.Xavier@Corteva.com)