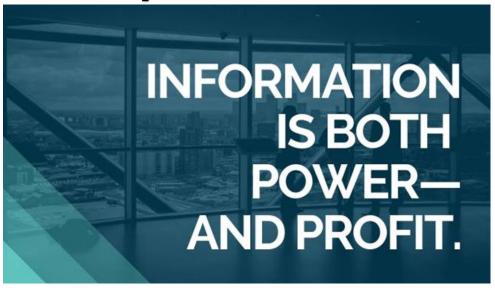


Learning from data:

Technical Nuances of Machine Learning in Plant Breeding

Alencar Xavier, 02/04/2021
Research Scientist at Corteva Biostatistics
Adjunct professor at Purdue University

Adequate use of





Outline

Three faces on machine learning

- 1. Good
- 2. Bad
- 3. Ugly





1. Good

- Intro and motivation
- Filters in PB

2. Bad

- Metrics of success
- Adequate validation

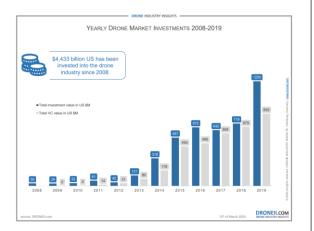
3. Ugly

- Optimization
- From RR to NN

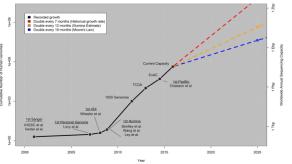
Part 1 – Good learners



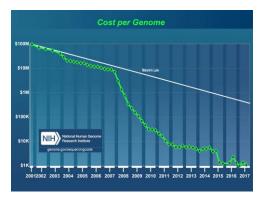
More Pheno



More Geno Growth of DNA Sequencing



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 2

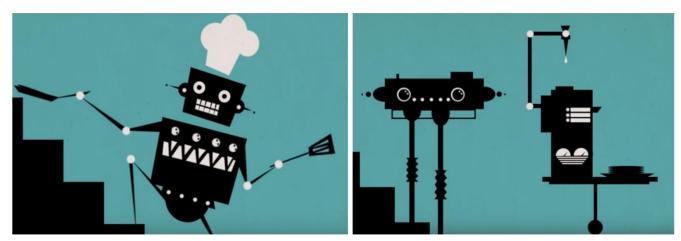






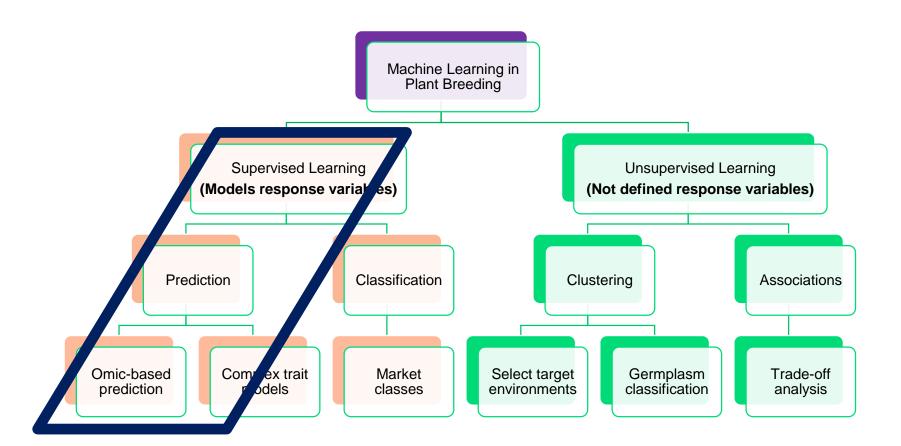
Why is machine learning good for?

Good for solving single well-defined problem



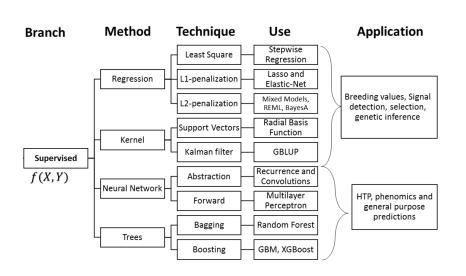
Source: https://www.youtube.com/watch?v=MPR3o6Hnf2g

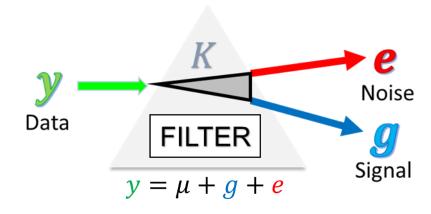






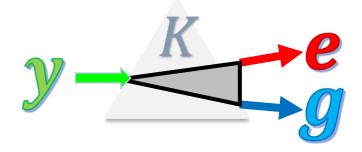
Key of supervised learning: FILTERING

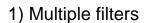


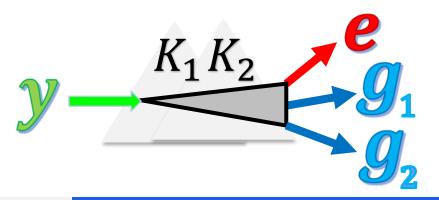




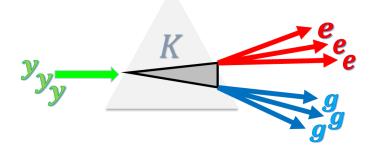
Generalizations of simple filters





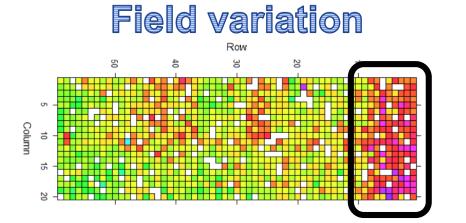


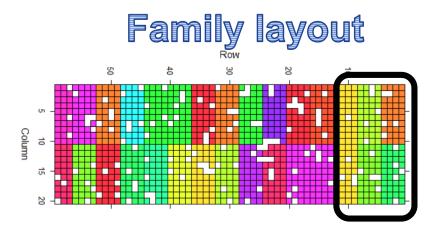






1) When bother with multiple filters?



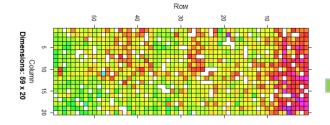


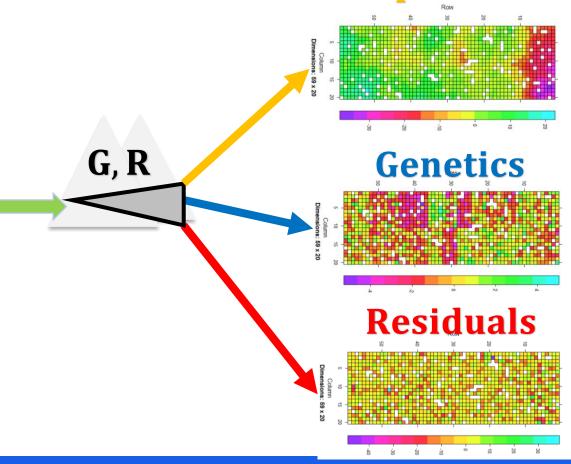
Some families were placed on unfavorable side of the field...

SoyNAM field, Indiana 2014



Pheno



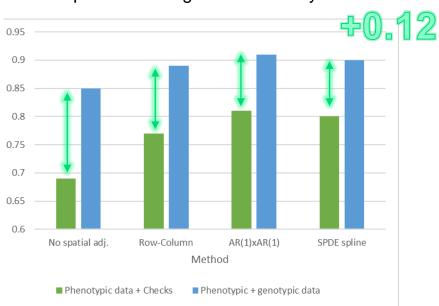


Spatial

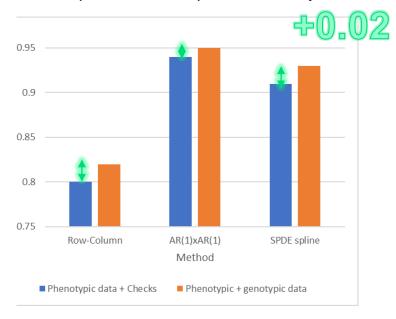


Multiple filters benefit separability of signals

Improvement in genetic accuracy



Improvement in spatial accuracy

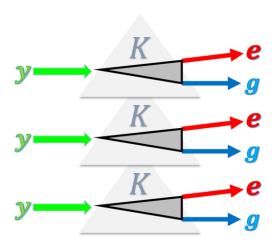


Results derived from simulation of field with 1500 plots, 1275 non-replicated entries + checks, from single trials

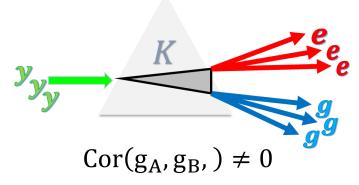


2) When bother with multi-response filters?

Filter one input at a time (parallelizable)



Filter passing multiple input at once



New information

Genetic correlation table

,	stictic correlation tab			
		g_1	g_2	g_3
	g_1	1	$ ho_{12}$	$ ho_{13}$
	g_2	$ ho_{21}$	1	ρ_{23}
	g_3	$ ho_{31}$	ρ_{32}	1





Bivariate model

Objective gain of information!

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

Model equation

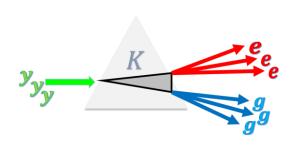
$$\begin{bmatrix} Z_1' \Sigma_e^{11} Z_1 + G^{-1} \Sigma_a^{11} & Z_1' \Sigma_e^{12} Z_2 + G^{-1} \Sigma_a^{12} \\ Z_2' \Sigma_e^{12} Z_1 + G^{-1} \Sigma_a^{12} & Z_2' \Sigma_e^{22} Z_2 + G^{-1} \Sigma_a^{22} \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} Z_1' (\Sigma_e^{11}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

· Univariate vs bivariate solution for a given predictor

$$\begin{split} \mathbf{g}_1 &= (\mathbf{Z}_1' \boldsymbol{\Sigma}_e^{11} \mathbf{Z}_1 + \mathbf{G}^{-1} \boldsymbol{\Sigma}_a^{11})^{-1} (\mathbf{Z}_1' \boldsymbol{\Sigma}_e^{11} \mathbf{y}_1) \\ \\ \mathbf{g}_1 | \mathbf{g}_2 &= (\mathbf{Z}_1' \boldsymbol{\Sigma}_e^{11} \mathbf{Z}_1 + \mathbf{G}^{-1} \boldsymbol{\Sigma}_a^{11})^{-1} (\mathbf{Z}_1' (\boldsymbol{\Sigma}_e^{11} \mathbf{y}_1 + \boldsymbol{\Sigma}_e^{12} \mathbf{y}_2) - (\mathbf{Z}_1' \boldsymbol{\Sigma}_e^{12} \mathbf{Z}_2 + \mathbf{G}^{-1} \boldsymbol{\Sigma}_a^{12}) \mathbf{g}_2) \end{split}$$

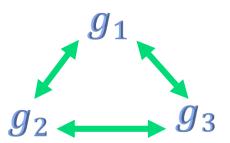


Sparse & Directed Filtering



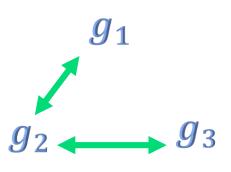
Genetic correlation table

	g_1	g_2	g_3
g_1	1	$ ho_{12}$	$ ho_{13}$
g_2	$ ho_{21}$	1	$ ho_{23}$
g_3	$ ho_{31}$	$ ho_{32}$	1



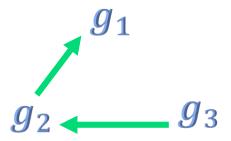
Sparse

	g_1	g_2	g_3
g_1	1	$ ho_{12}$	l
g_2	$ ho_{21}$	1	ρ_{23}
g_3	_	ρ_{32}	1



Directed

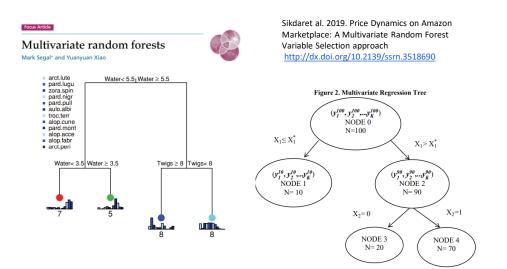
	g_1	g_2	g_3
g_1	1	ı	ı
g_2	$ ho_{21}$	1	_
g_3	_	$ ho_{32}$	1



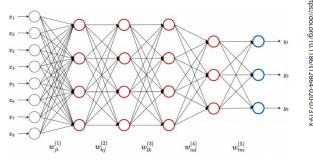


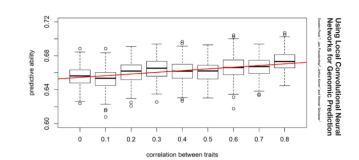
Modern multi-response machinery

MVRF feature splits occur for all responses together



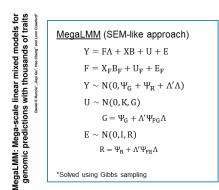
Multi-response comes natural for DNNs

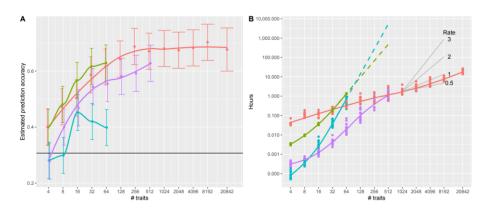






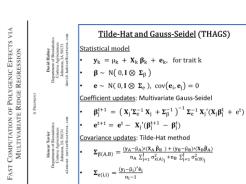
Modern multi-response machinery

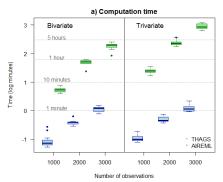


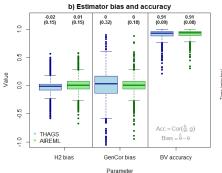


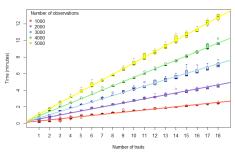
→ MegaLMM → MCMCglmm → MTG2 → phenix

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









1. Good

- Intro and motivation
- Filters in PB

2. Bad

- Metrics of success
- Adequate validation

3. Ugly

- Optimization
- From RR to NN

(avoiding)

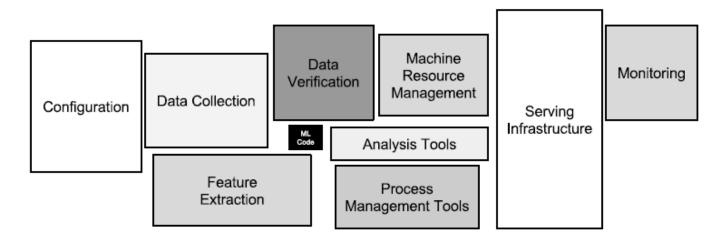
Part 2 – Bad learning



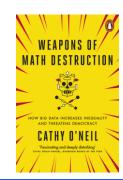
There is usually more to ML than a proof of concept with cross-validations

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?

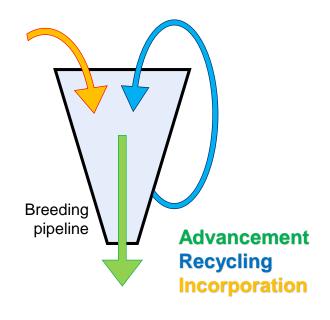




Chasing the right signal

- Breeding value (GEBV)
 - Pattern: ADDITIVE GENETICS
 - Method: GBLUP, BayesABC, LASSO
 - Suits: RECYCLING, ADVANCEMENT

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





Defining the problem: Metrics for success

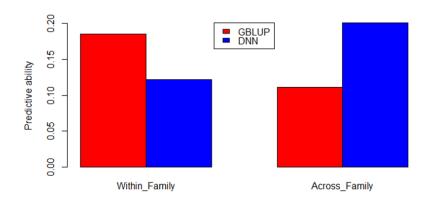
1. Scientist (why): to define the problem mathematically (easy to get it wrong)

2. Metric (what): Correlations, MSPE, Jaccard index, Accuracy, Success (1|0)

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

How to design an adequate cross-validation???

Two metrics, two different answers



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

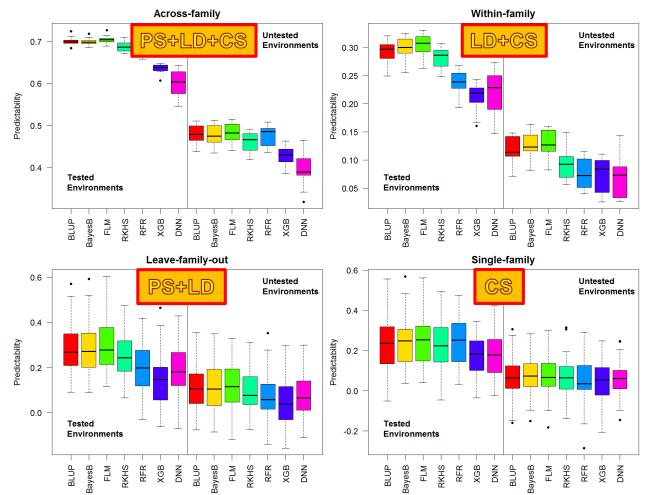


Testing machines for different scenarios of genomic prediction

	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





SoyNAM data

ES: 2012 (7 loc)

PS: 2013 (4 loc)

#Fam = 40

Genos = 5600 SNPs = 4300

Obs: 3k-5k obs/loc



Untested

Tested

- Population structure (PS)
- Linkage disequilibrium (LD)
- Cosegregation / Haplotype (CS)



CV scheme

1. Good

- Intro and motivation
- Filters in PB

2. Bad

- Metrics of success
- Adequate validation

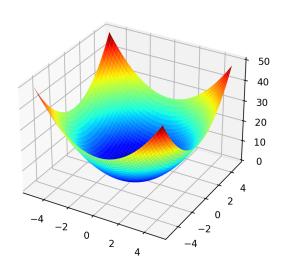
3. Ugly

- 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 <u>secondary objective function</u>
 - Method: First order: coordinate & gradient descent





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)

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

Second order

(Builds entire LHS)

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

Used for Deep Neural Nets

TF/Keras, PyTorch, MXNet, h2o

Used for everything else

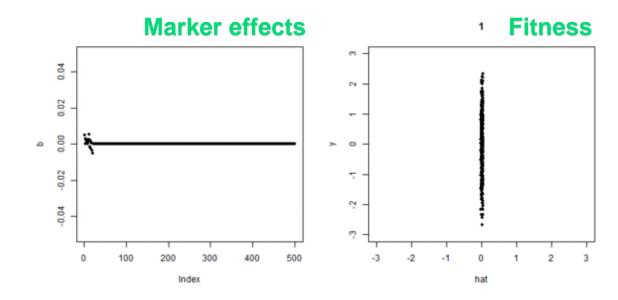
ASREML, Ime4, SAS, BLUPF90



Coordinate descent of a RRBLUP



$$\begin{split} f(b_{RR}) &\to argmin \Big(\sum_{i=1}^{n} e^2 + \lambda \sum_{j=1}^{p} b^2 \Big) \\ SS &= e'e + \lambda(b'b) \\ SS &= (y - xb)'(y - xb) + \lambda(b'b) \\ SS &= y'y + (xb)'(xb) - 2(y'xb) + \lambda(b'b) \\ SS &= y'y + b'(x'x)b - 2(y'xb) + \lambda(b'b) \\ \frac{\partial SS}{\partial b} &= y'y + 2(x'x)b^2 - 2(y'x)b + 2\lambda b^2 \\ 0 &= 2b(x'x) - 2(y'x) + 2b\lambda \\ -2b(x'x + \lambda) &= -2(y'x) \\ b &= \frac{-2(y'x)}{-2(x'x + \lambda)} = \frac{y'x}{x'x + \lambda} \end{split}$$



Gradient descent of RRBLUP

Scatter plot between true solution and solution at iteration x

Small learning rate

Large learning rate



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

• Where

$$\alpha$$
 = Learning rate

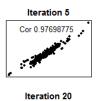
$$f = (y - Xb)'(y - Xb) + \lambda b'b$$

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

Thus

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

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



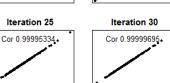
Iteration 35

Cor 0.99999991.



Iteration 40

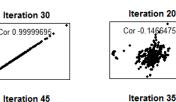
Cor 1

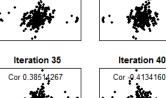


Cor 1

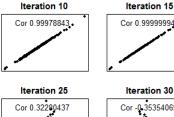
Iteration 15

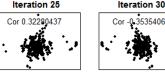
Cor 0.99781079 •

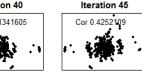




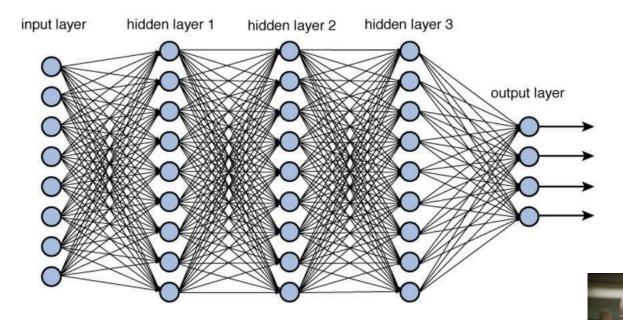
Iteration 5







Deep Neural Network



WHAT IS

THIS NONSENSE?



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

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

 α = activation function



The Scary Deep Neural Network

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

- Fit layers
 - $H_1 = \alpha(XB_1)$
 - $H_2 = \alpha(H_1B_2)$
 - $H_3 = H_2 B_3$
- Compute gradients (aka. get residuals)
 - $E_3 = Y H3$
 - $E_2 = \alpha(E_3 B_3')$
 - $E_1 = \alpha(E_2B_2')$
- Update coefficient
 - $B_1 = B_1 \gamma \left(\frac{2r_1}{n} \left[X' E_1 \lambda B_1 \right] \right)$
 - $B_2 = B_2 \gamma \left(\frac{2r_2}{n} [H_1' E_2 \lambda B_2] \right)$
 - $B_3 = B_3 \gamma \left(\frac{2r_3}{n} [H_2' E_3 \lambda B_3] \right)$

Top-to-bottom code ~ 30 lines

```
dnn = function(y, X, nit=1000, batch=250, RELU=FALSE, Leak=0.1,
                     dropout=0, Lambda=0.1,LrnRate = 1,Nodes1=4, Nodes2=4){
   if(is.null(ncol(v))) v = matrix(v)
  muY = colMeans(y,na.rm=1); sdY = apply(y,2,sd,na.rm=1); y = apply(y,2,scale) ActFun = tanh; if(RELU) ActFun = function(x)\{x[x<0]=x[x<0]\} Leak; return(x)
   DropOut = function(x,prc=dropout) {x[sample(length(x),length(x)*prc)]:return(x)
   n = nrow(X); p = ncol(X); k = ncol(y)
  n1 = Nodes1; n2 = Nodes2; lmb = Lambda; rate = LrnRate/c(p,n1,n2)
 b1 = matrix(rnorm(n1*p,0,1/p),p,n1)
b2 = matrix(rnorm(n1*n2,0,1/n1),n1,n2)
b3 = matrix(rnorm(n2,0,1/n2),n2,k)
 CNV1 = CNV2 = CC
     w = sample(n,batch,replace=T); y0 = y[w,]; X0 = X[w,]
    Hl = ActFun(X0%%bl); H2 = ActFun(H1%%b2); H3 = H2%%b3
e3 = y0-H3; ff(anyMc(e3)) e3[fs.na(e3)] e2
e2 = ActFun(e3 %% t(b3)); e1 = ActFun(e2 %% t(b2))
     \begin{array}{lll} b1 = b1 + DropOut(t(X0)\%\%(e1) - lmb^*b1)^*(2/n)^*rate[\\ b2 = b2 + DropOut(t(H1)\%\%(e2) - lmb^*b2)^*(2/n)^*rate[\\ \end{array}
     b3 = b3 + DropOut(t(H2)%*%(e3))*(2/n)*rate[
     CNV1 = c(CNV1,mean(apply(e3,2,var,na.rm=
   CNV2 = c(CNV2,mean(diag(cor(H3,V0,use='p')))) out = list(af=ActFun,b1=b1,b2=b2,b3=b3,conv_MSE=CNV1,conv_GOF=CNV2,mu=muY,sd=sd)
predict.smalldnn = function(object,newdata)
   x = object; h = x$af(x$af(newdata%*%x$b1)%*%x$b2)%*%x$b3
   for(i in 1:ncol(h)) h[,i] = x[mu[i] + x[sd[i]*h[,i]
```

Thank you for your attention!

Remarks:

- 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!

