**Can you open it?**

**Can you open it?**

Yes          No

**Decision trees are made of yes/no questions (nodes)**

How do we know if this is the best branching structure?

# Random Forests

- Made of decision trees

- Uses bootstrapping observations and predictors to make a "forest" of decision trees which are aggregated (Bagging)

- Observations left out (due to bootstrapping – OOB – out-of-bag) are used to estimate error

- Can extract importance of variables "Gini" Score by comparing trees with/without variable included

# Bootstrapping

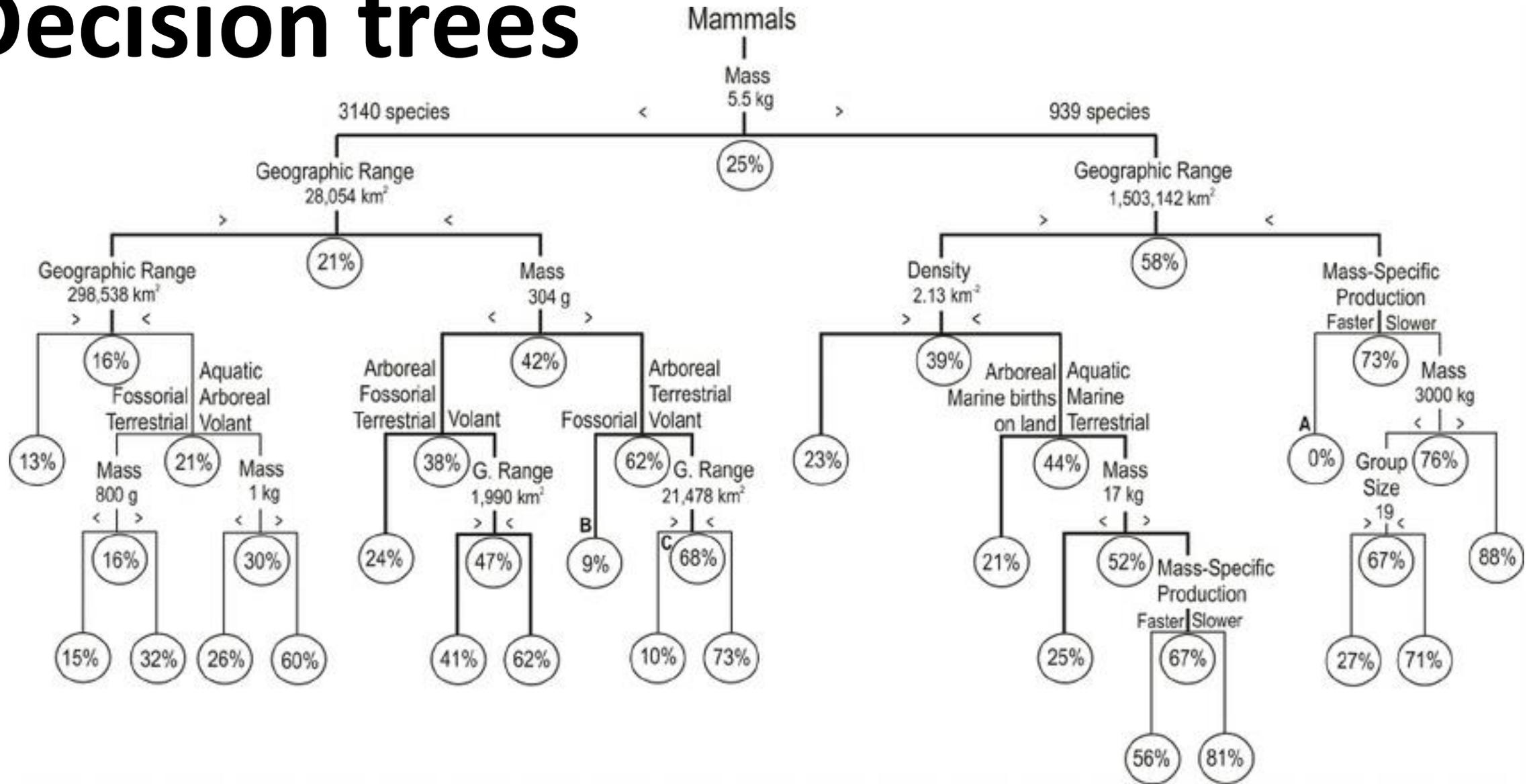| Obs | X | Y |
|-----|-----|-----|
| 3 | 5.3 | 2.8 |
| 1 | 4.3 | 2.4 |
| 3 | 5.3 | 2.8 |

$Z^{*1}$ $\longrightarrow$ $\hat{\alpha}^{*1}$

| Obs | X | Y |
|-----|-----|-----|
| 1 | 4.3 | 2.4 |
| 2 | 2.1 | 1.1 |
| 3 | 5.3 | 2.8 |

Original Data (Z)

$Z^{*2}$

| Obs | X | Y |
|-----|-----|-----|
| 2 | 2.1 | 1.1 |
| 3 | 5.3 | 2.8 |
| 1 | 4.3 | 2.4 |

$\longrightarrow$ $\hat{\alpha}^{*2}$

$Z^{*B}$

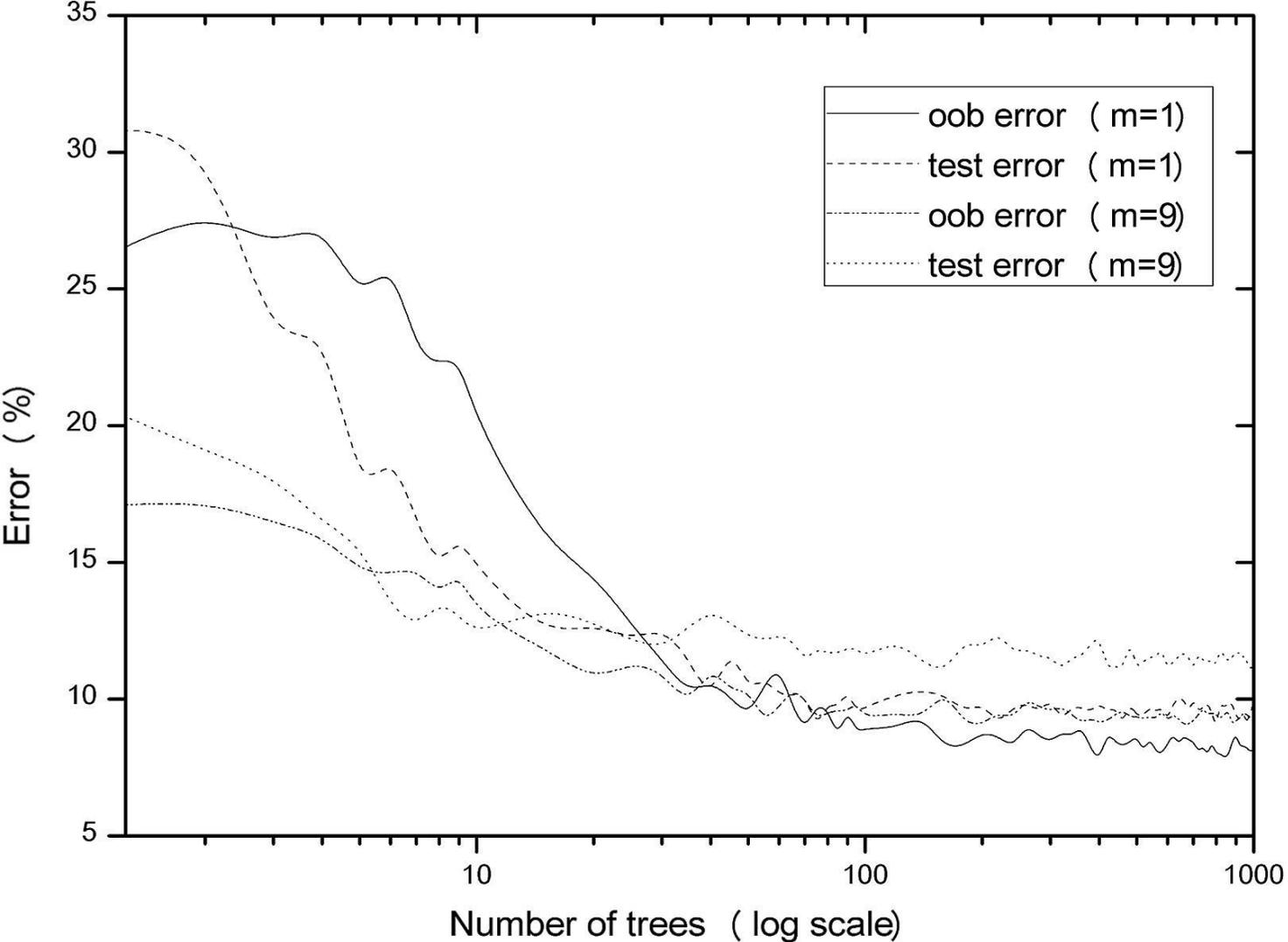| Obs | X | Y |
|-----|-----|-----|
| 2 | 2.1 | 1.1 |
| 2 | 2.1 | 1.1 |
| 1 | 4.3 | 2.4 |

$\longrightarrow$ $\hat{\alpha}^{*B}$

# Decision trees



Decision tree showing extinction risk based on ecological traits (body mass, geographic range size, mass-specific production rate, population density, group size, home range, activity period, type of landmass, habitat mode, sociality, trophic group).

https://www.pnas.org/doi/10.1073/pnas.0901956106

**The more trees in the random forest, the more accurate the model (with diminishing returns)**

Error is estimated during bootstrapping – called out of bag (OOB) error

An assessment of the effectiveness of a random forest classifier for land-cover classification

V.F. Rodriguez-Galiano [a] ✉, B. Ghimire [b] ✉, J. Rogan [b] ✉, M. Chica-Olmo [a] ✉, J.P. Rigol-Sanchez [c] ✉

# A Random Forest Machine Learning Approach for the Retrieval of Leaf Chlorophyll Content in Wheat

by Syed Haleem Shah [1,*] ✉ ⓘ, Yoseline Angel [1] ⓘ, Rasmus Houborg [2] ⓘ, Shawkat Ali [3] ⓘ and Matthew F. McCabe [1] ⓘ

A Random Forest Machine Learning Approach for the Retrieval of Leaf Chlorophyll Content in Wheat

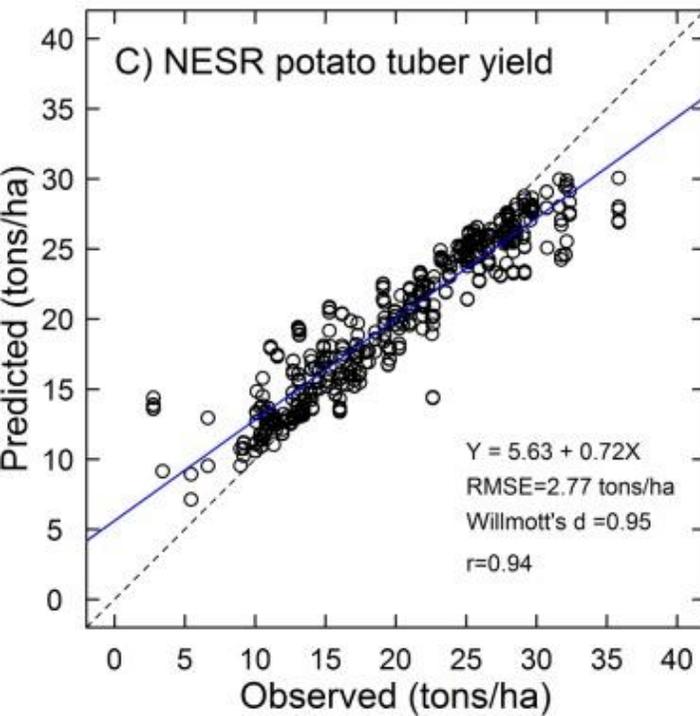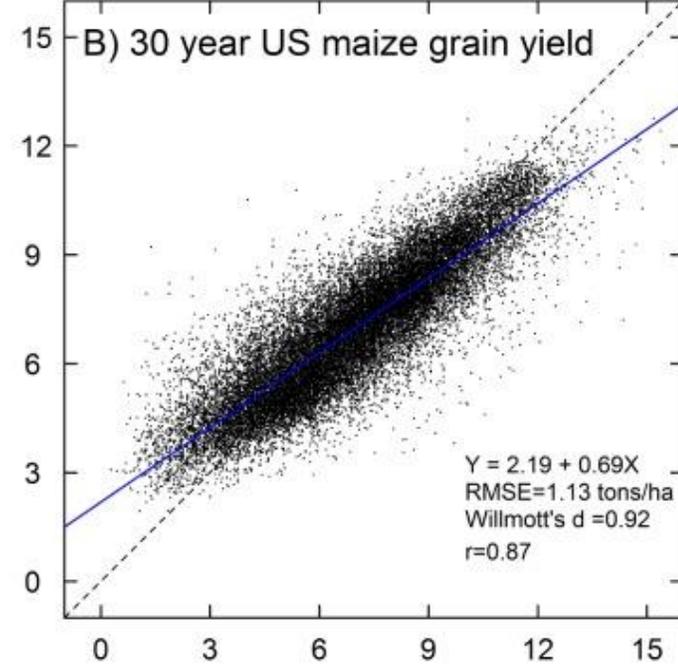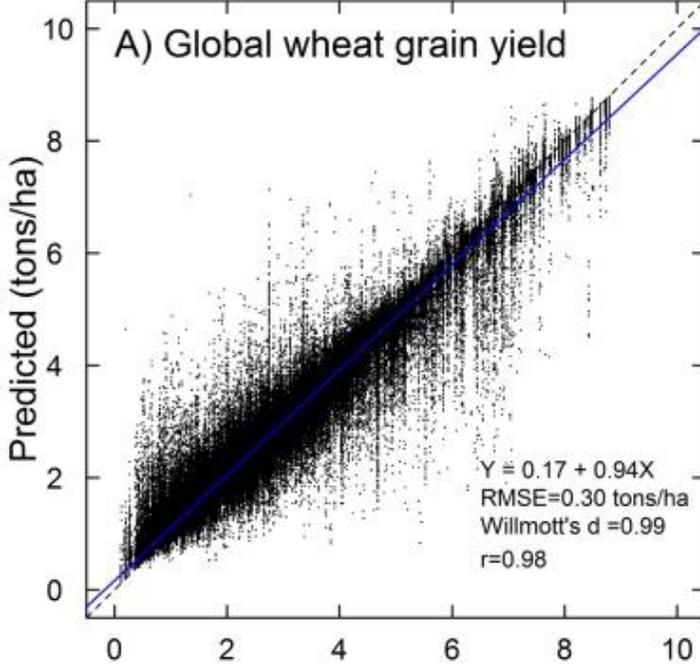by Syed Haleem Shah [1,*], Yoseline Angel [1], Rasmus Houborg [2], Shawkat Ali [3] and Matthew F. McCabe [1]

A) Global wheat grain yield
Y = 0.17 + 0.94X
RMSE=0.30 tons/ha
Willmott's d =0.99
r=0.98

B) 30 year US maize grain yield
Y = 2.19 + 0.69X
RMSE=1.13 tons/ha
Willmott's d =0.92
r=0.87

C) NESR potato tuber yield
Y = 5.63 + 0.72X
RMSE=2.77 tons/ha
Willmott's d =0.95
r=0.94

D) NESR maize silage yield
Y = 6.10 + 0.81X
RMSE=1.9 tons/ha
Willmott's d =0.97
r=0.95
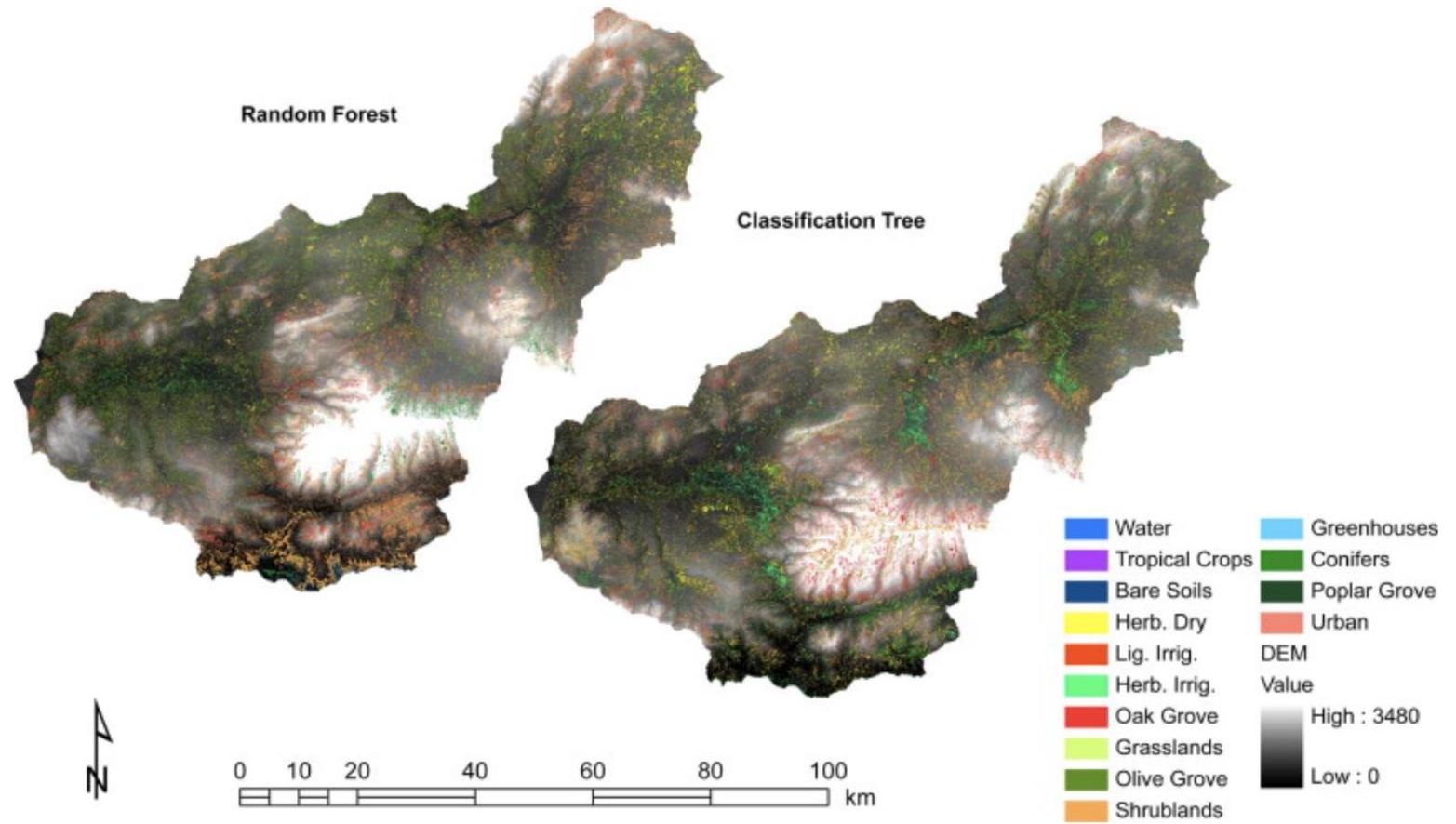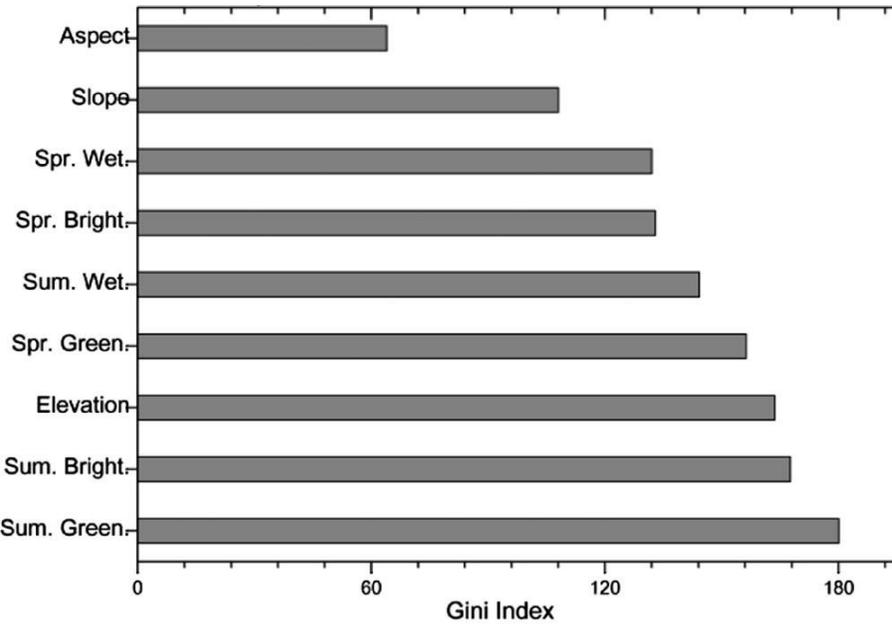
Axis labels: Observed (tons/ha), Predicted (tons/ha)

# Random Forests for Global and Regional Crop Yield Predictions

Jig Han Jeong,[1] Jonathan P. Resop,[2,3] Nathaniel D. Mueller,[4,5] David H. Fleisher,[3] Kyungdahm Yun,[1] Ethan E. Butler,[6] Dennis J. Timlin,[3] Kyo-Moon Shim,[7] James S. Gerber,[8] Vangimalla R. Reddy,[3] and Soo-Hyung Kim[1,*]

| Global wheat and US maize grain yields | | | Importance Rank | |
| --- | --- | --- | --- | --- |
| Variable | Abbreviation | Unit | wheat | maize |
| Averaged monthly temperature | AVT | °C | 8 | 9 |
| Annual evapotranspiration | EVA | mm | 2 | 6 |
| Summer solstice day length | DAYL | hour | 4 | 7 |
| Maximum monthly temperature | MAX | °C | 7 | 3 |
| Mean coldest quarter Temperature | MCQ | °C | 9 | 11 |
| Minimum monthly temperature | MIN | °C | 10 | 8 |
| Mean warmest quarter temperature | MWQ | °C | 6 | 10 |
| Nitrogen fertilizer application rate | NFERT | kg/ha | 1 | 2 |
| Growing season precipitation | PRE49 | mm | 3 | 4 |
| Annual precipitation | PRECI | mm | 5 | 5 |
| Year (US maize only) | YR | | - | 1 |

# Gini score helps interpretability
# Provides measure of variable importance



Random Forest

Classification Tree

Legend:
- Water
- Tropical Crops
- Bare Soils
- Herb. Dry
- Lig. Irrig.
- Herb. Irrig.
- Oak Grove
- Grasslands
- Olive Grove
- Shrublands
- Greenhouses
- Conifers
- Poplar Grove
- Urban
- DEM Value — High : 3480, Low : 0

0 10 20 40 60 80 100 km

Gini Index variables (top to bottom): Sum. Green, Sum. Bright, Elevation, Spr. Green, Sum. Wet, Spr. Bright, Spr. Wet, Slope, Aspect

Download : Download high-res image (1MB)    Download : Download full-size image

An assessment of the effectiveness of a random forest classifier for land-cover classification

V.F. Rodriguez-Galiano [a], B. Ghimire [b], J. Rogan [b], M. Chica-Olmo [a], J.P. Rigol-Sanchez [c]

## Variable importance for *Verbascum thapsus*
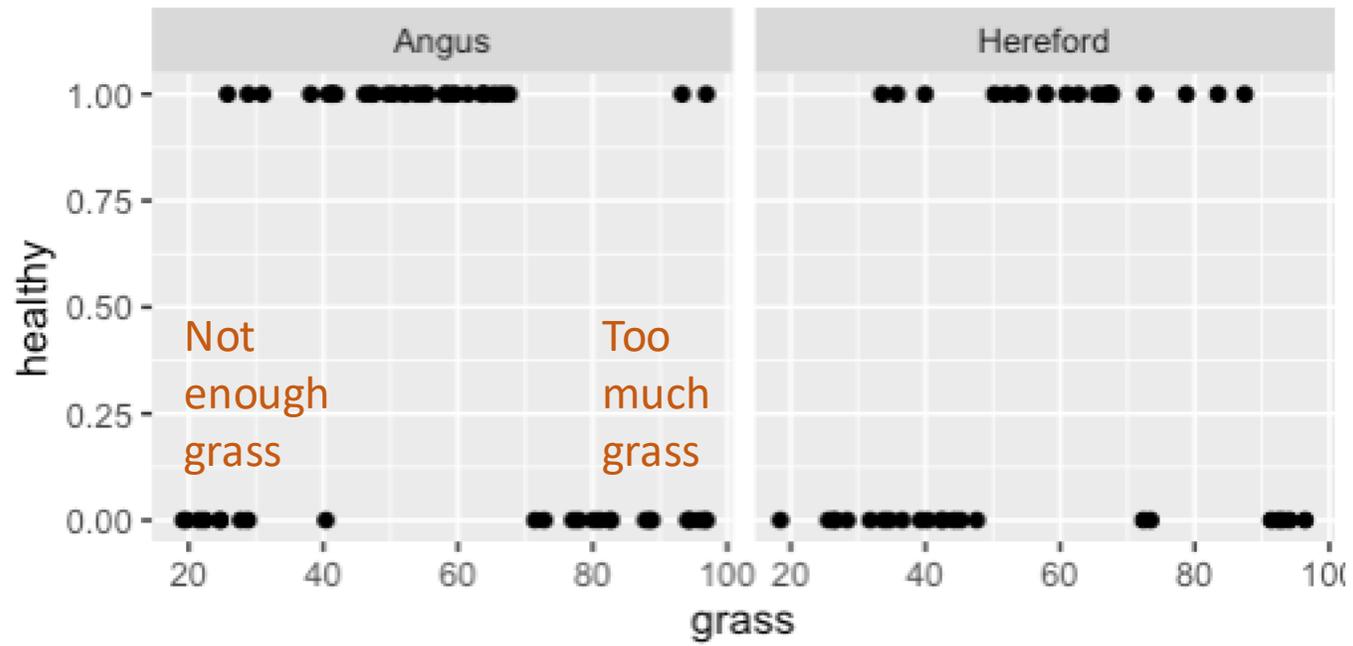
D. Richard Cutler,[1,7] Thomas C. Edwards, Jr.,[2] Karen H. Beard,[3] Adele Cutler,[4] Kyle T. Hess,[4] Jacob Gibson,[5] and Joshua J. Lawler[6]
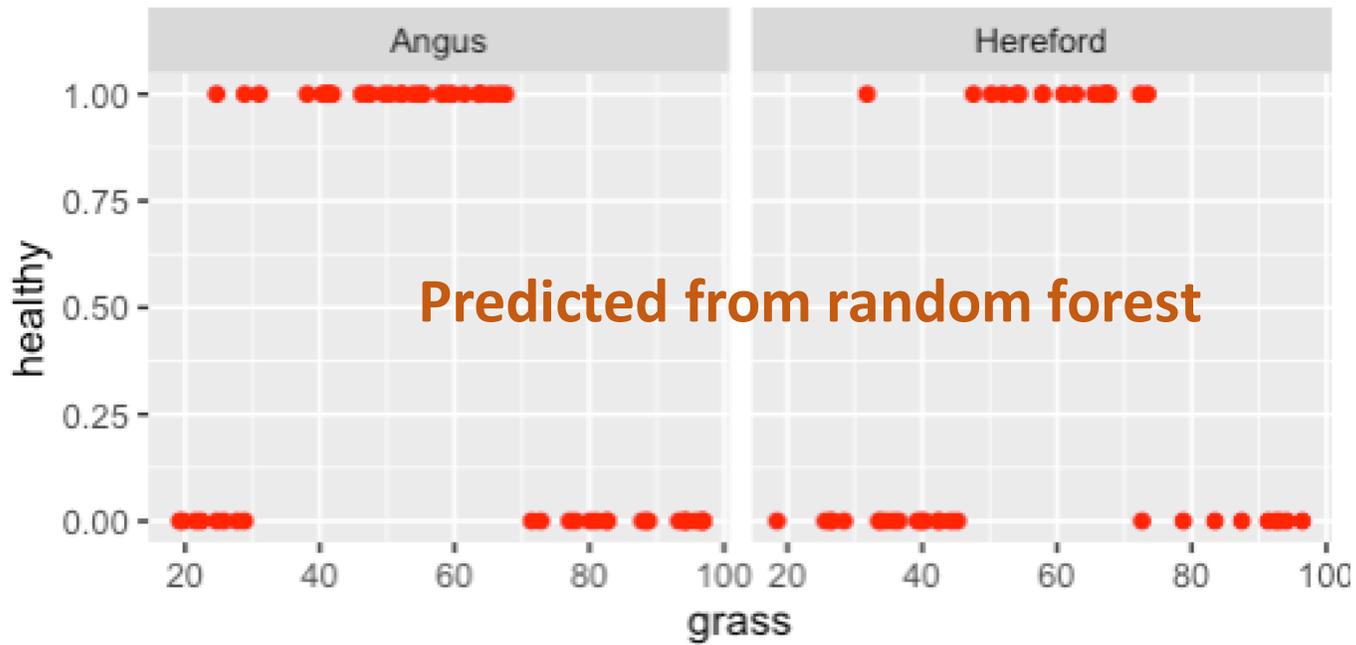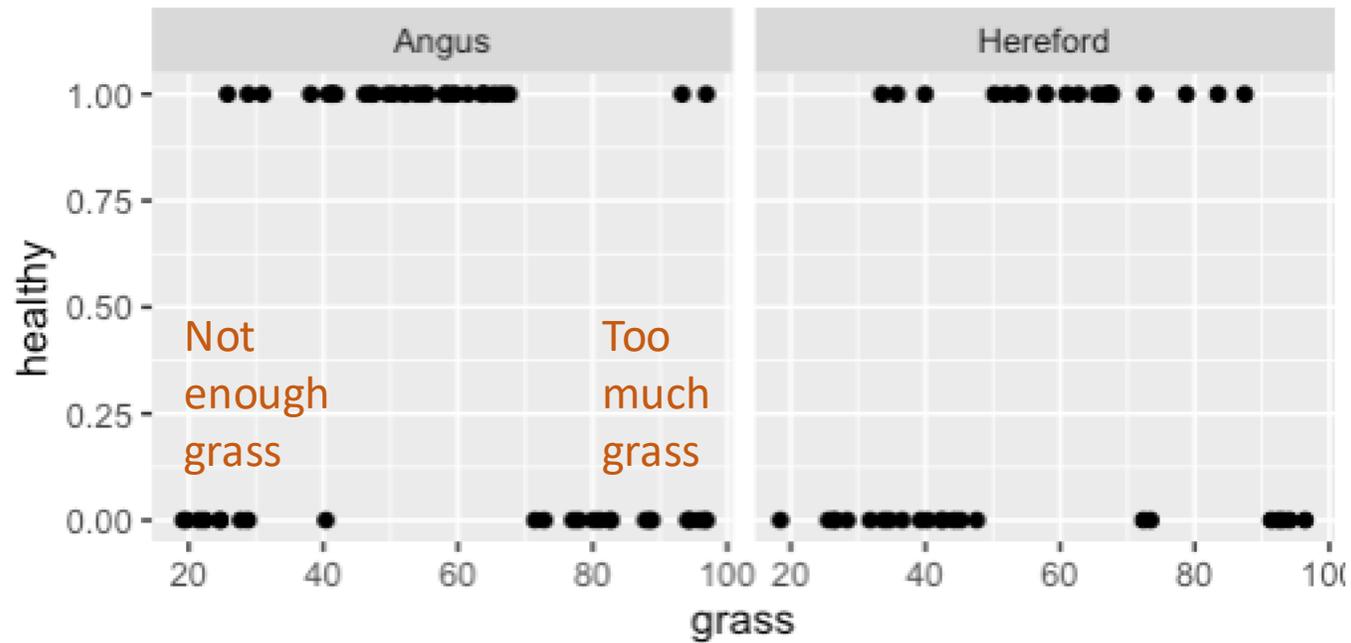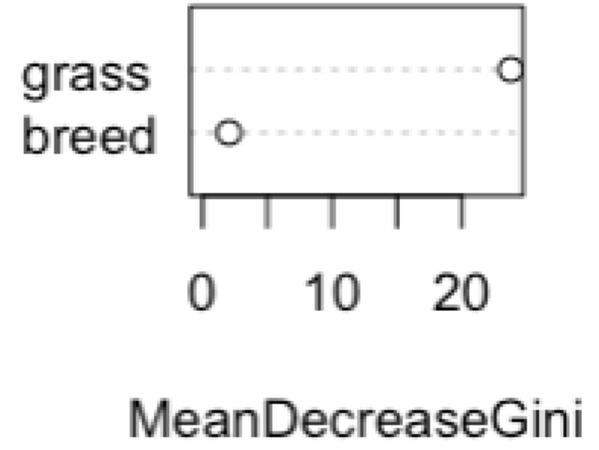
Common Mullen

**Non-linear relationships are difficult to model with standard linear models**

```
model<-randomForest(factor(healthy)~grass+breed, data=data, ntree=100)
```

Not enough grass

Too much grass

Predicted from random forest

|  | MeanDecreaseGini |
| --- | --- |
| grass | 23.778664 |
| breed | 2.030991 |

Grass is important predictor

# Random Forest

## Advantages

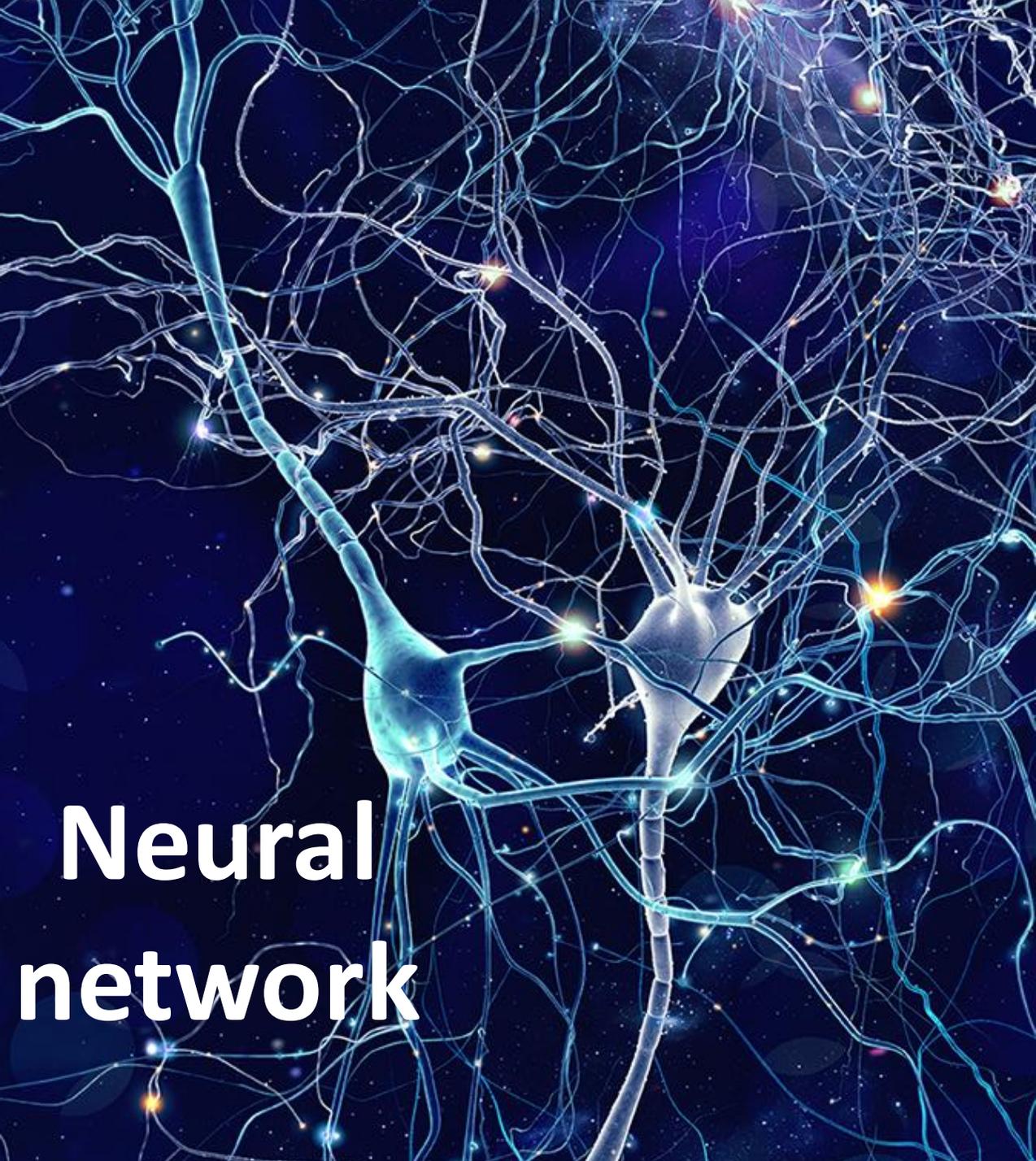**1.Accuracy under complexity:** Random Forest is known for its high accuracy when relationships between variables are not straightforward (interactions and non-linearities).

**2.Variable Importance:** It provides insights into which variables are most important in prediction, aiding in hypothesis testing and feature selection (though multicollinearity is an issue)

**3.Versatility:** Suitable for both classification and regression tasks, making it a versatile tool for various types of data.

## Disadvantages

**1.Model Interpretability:** Less interpretable compared to simpler models like linear regression.

**2.Computationally Intensive:** It can be computationally expensive, especially with large datasets and a large number of trees.

**3.Bias in Multiclass Problems:** Tends to be biased towards classes with more instances in classification problems (address by sub-setting data for balanced classes to train model)
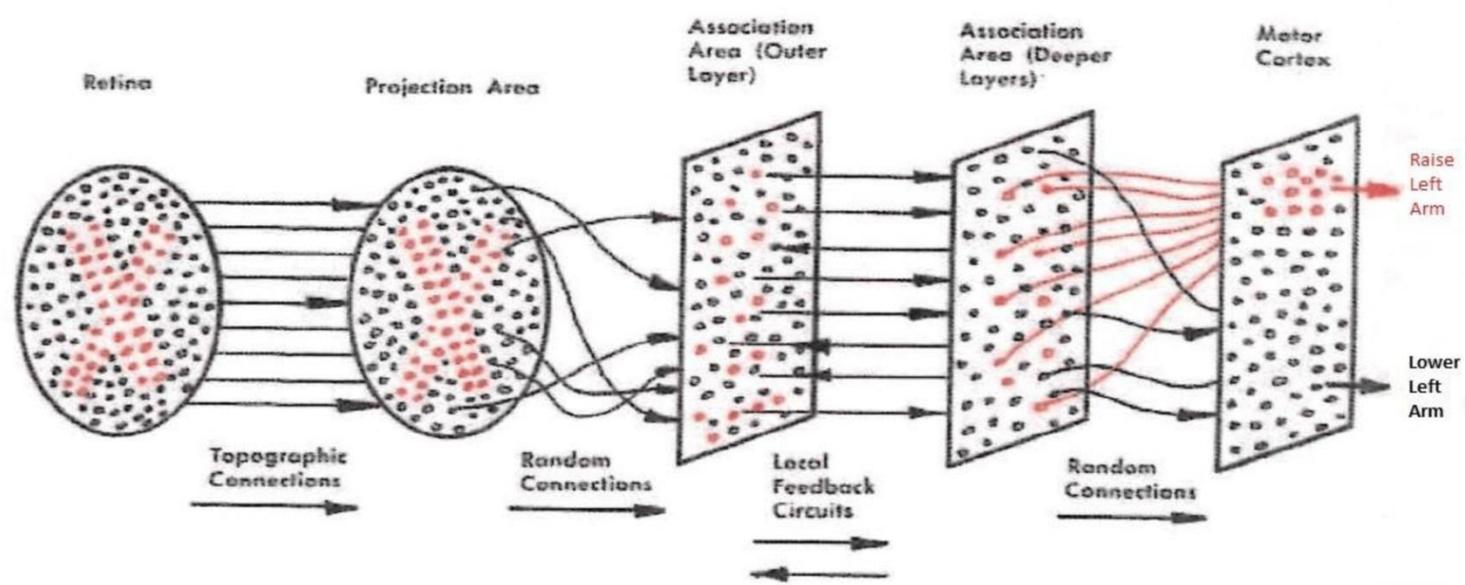
Random forest

Neural network

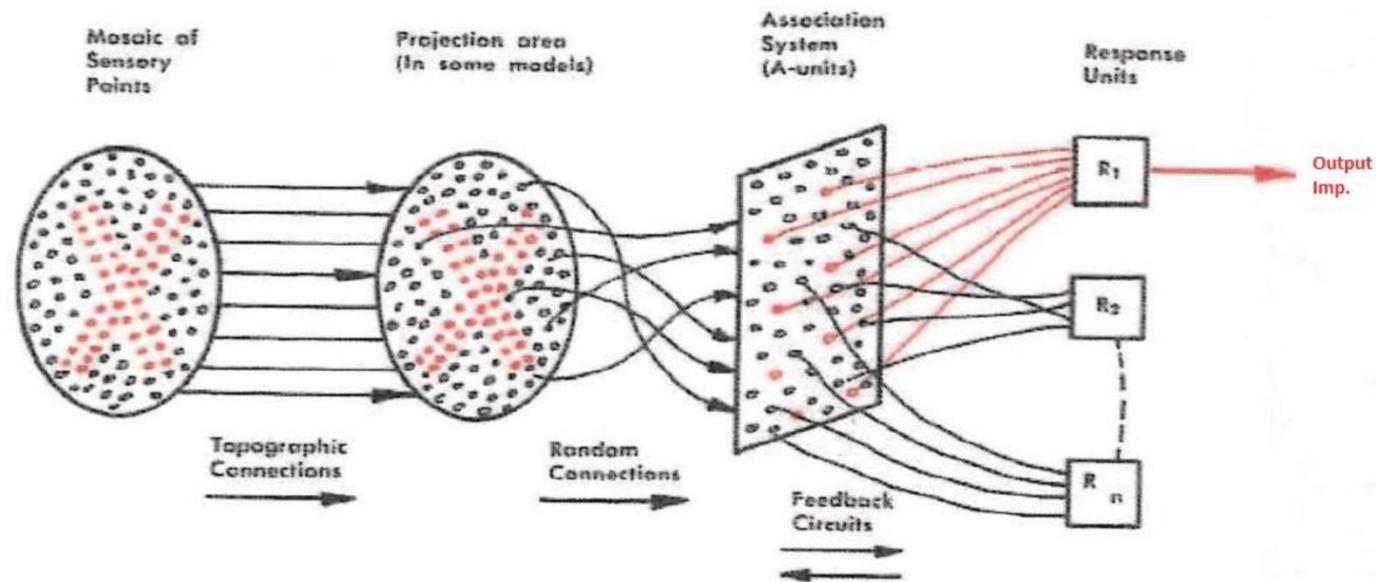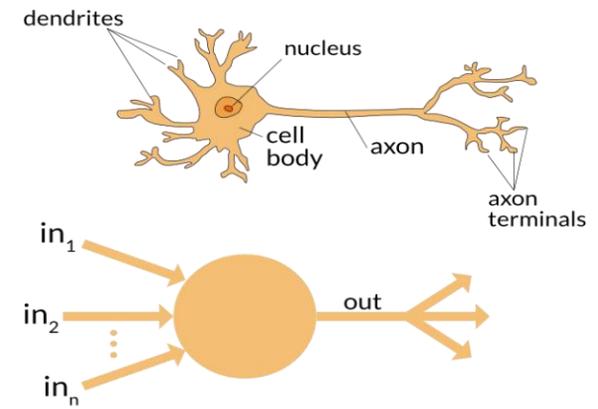FIG. 1 — Organization of a biological brain. (Red areas indicate active cells, responding to the letter X.)
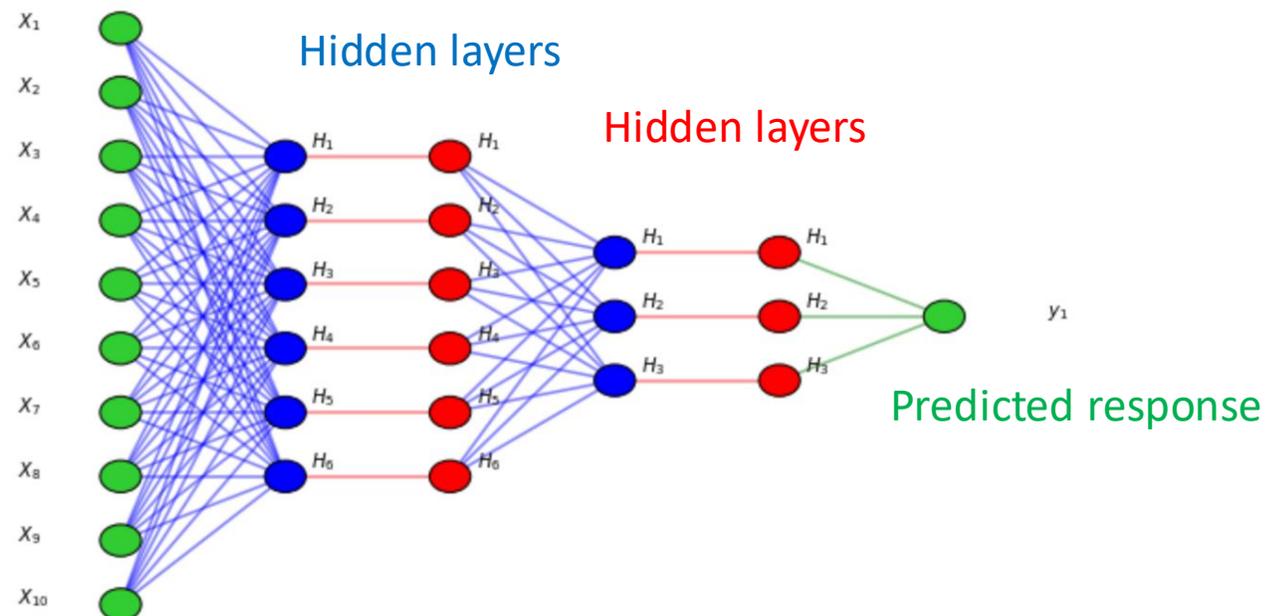


FIG. 2 — Organization of a perceptron.
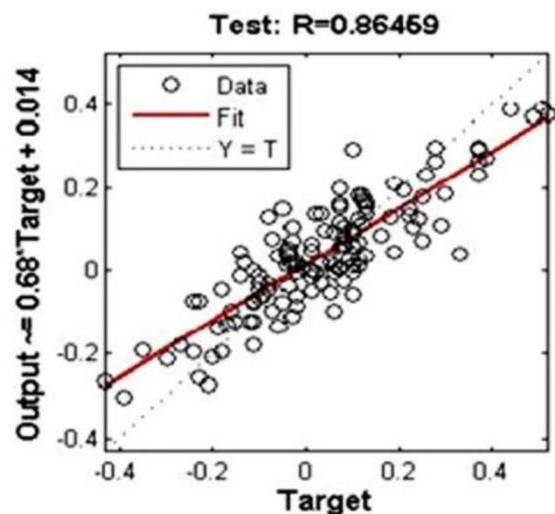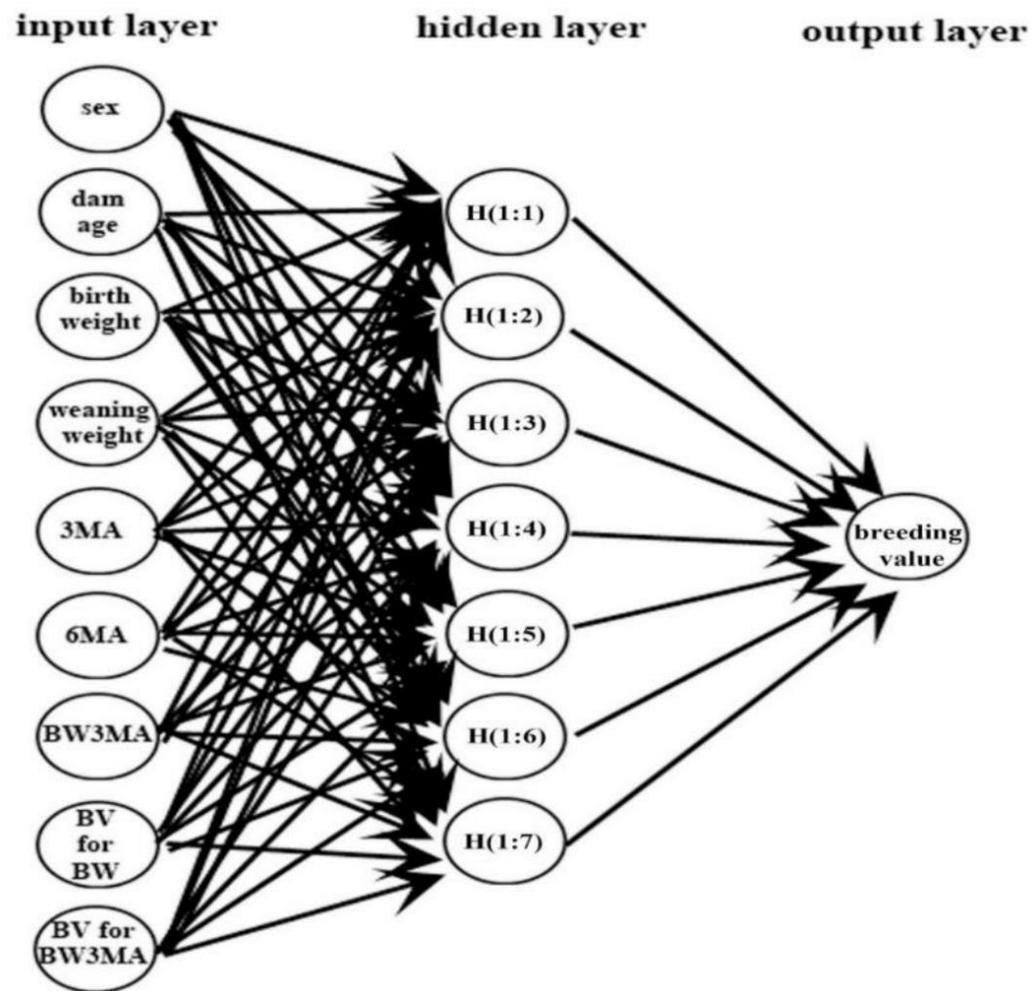
# Neural network

- Internal nodes are functions

- Edges between nodes influence outcome

- Outcome is product of effect of all input and transformations in internal nodes

- Values of weights and edges are those that fit the data best



Predictor variables

Hidden layers

Hidden layers

Predicted response

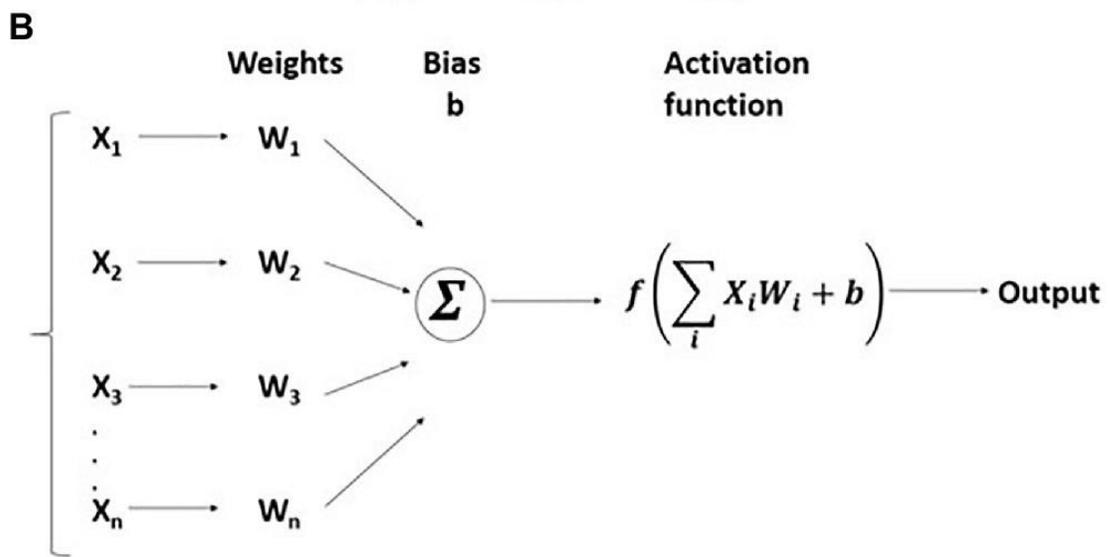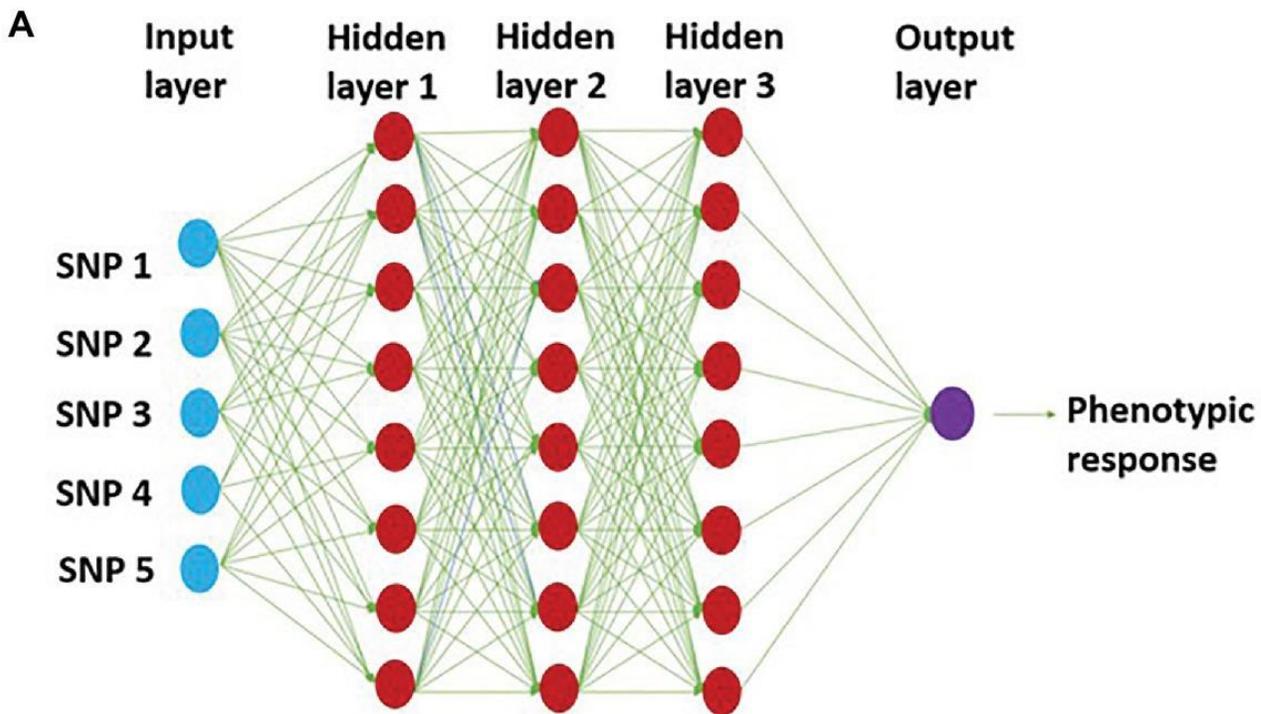"Deep learning" refers to having multiple hidden layers

# Predicting breeding value of body weight at 6-month age using Artificial Neural Networks in Kermani sheep breed

Hamidreza Ghotbaldini[1], Mohammadreza Mohammadabadi[1*] (iD), Hossein Nezamabadi-pour[1], Olena Ivanivna Babenko[2], Maryna Vitaliivna Bushtruk[1] and Serhii Vasyliovych Tkachenko[1]

**A**

Input layer — Hidden layer 1 — Hidden layer 2 — Hidden layer 3 — Output layer

SNP 1
SNP 2
SNP 3
SNP 4
SNP 5

→ Phenotypic response

**B**

Weights — Bias b — Activation function

$X_1 \longrightarrow W_1$
$X_2 \longrightarrow W_2$
$X_3 \longrightarrow W_3$
.
.
$X_n \longrightarrow W_n$

$\Sigma$ → $f\left(\sum_i X_i W_i + b\right)$ → Output

# Deep Learning for Predicting Complex Traits in Spring Wheat Breeding Program

Karansher S. Sandhu[1], Dennis N. Lozada[2], Zhiwu Zhang[1], Michael O. Pumphrey[1] and Arron H. Carter[1*]

| Model | Grain yield | Grain protein content | Test weight | Plant height | Heading date |
|---|---|---|---|---|---|
| rrBLUP | 0.39 | 0.48 | 0.45 | 0.52 | 0.46 |
| MLP | **0.44** | **0.53** | **0.48** | **0.57** | **0.51** |
| CNN | 0.39 | 0.48 | 0.47 | 0.55 | 0.49 |

*The highest prediction accuracy is bolded for each trait under each model scenario.*
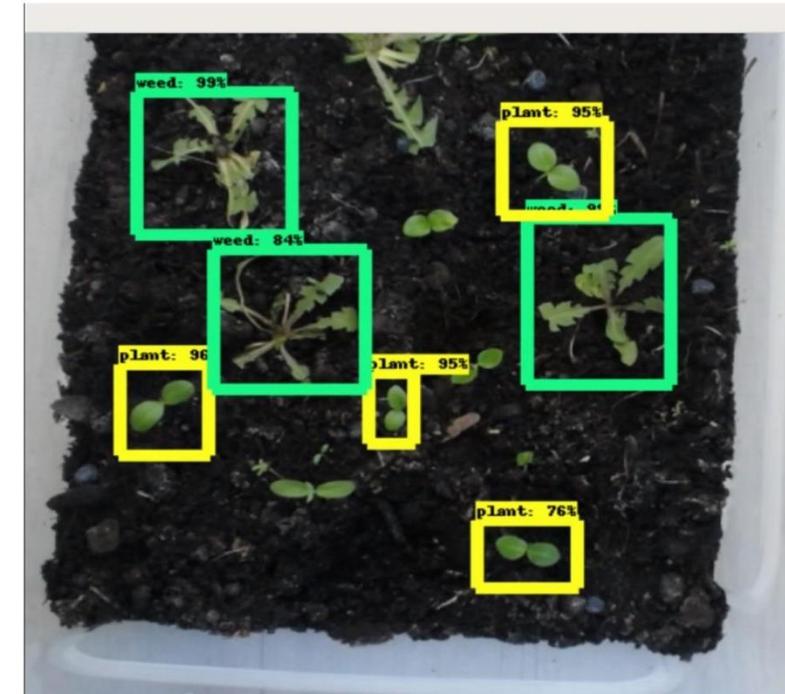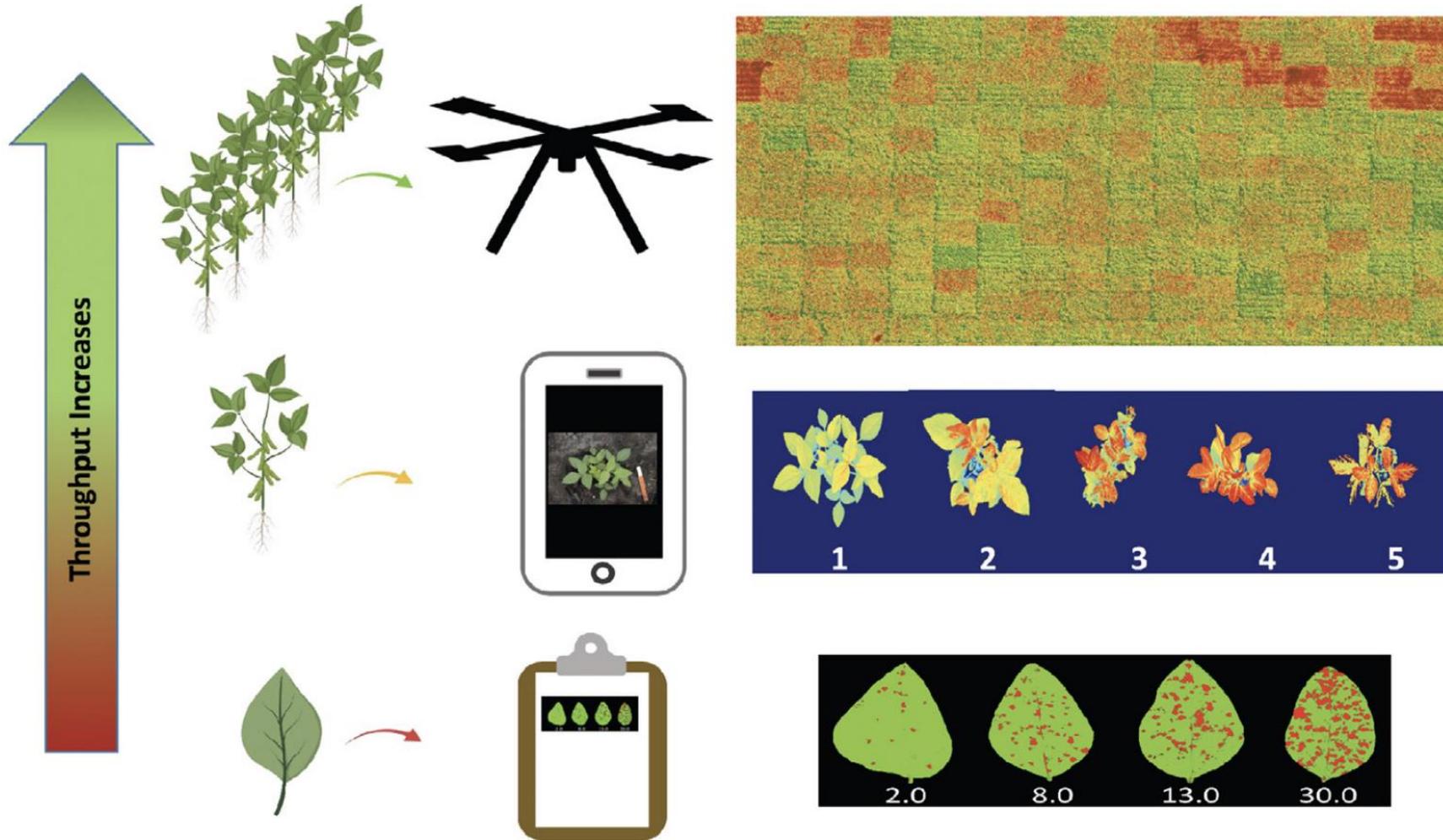
# Plant and Weed Identifier Robot as an Agroecological Tool Using Artificial Neural Networks for Image Identification

by 👤 **Tavseef Mairaj Shah** *,† ✉ 🆔, 👤 **Durga Prasad Babu Nasika** † ✉ 🆔 and 👤 **Ralf Otterpohl** ✉

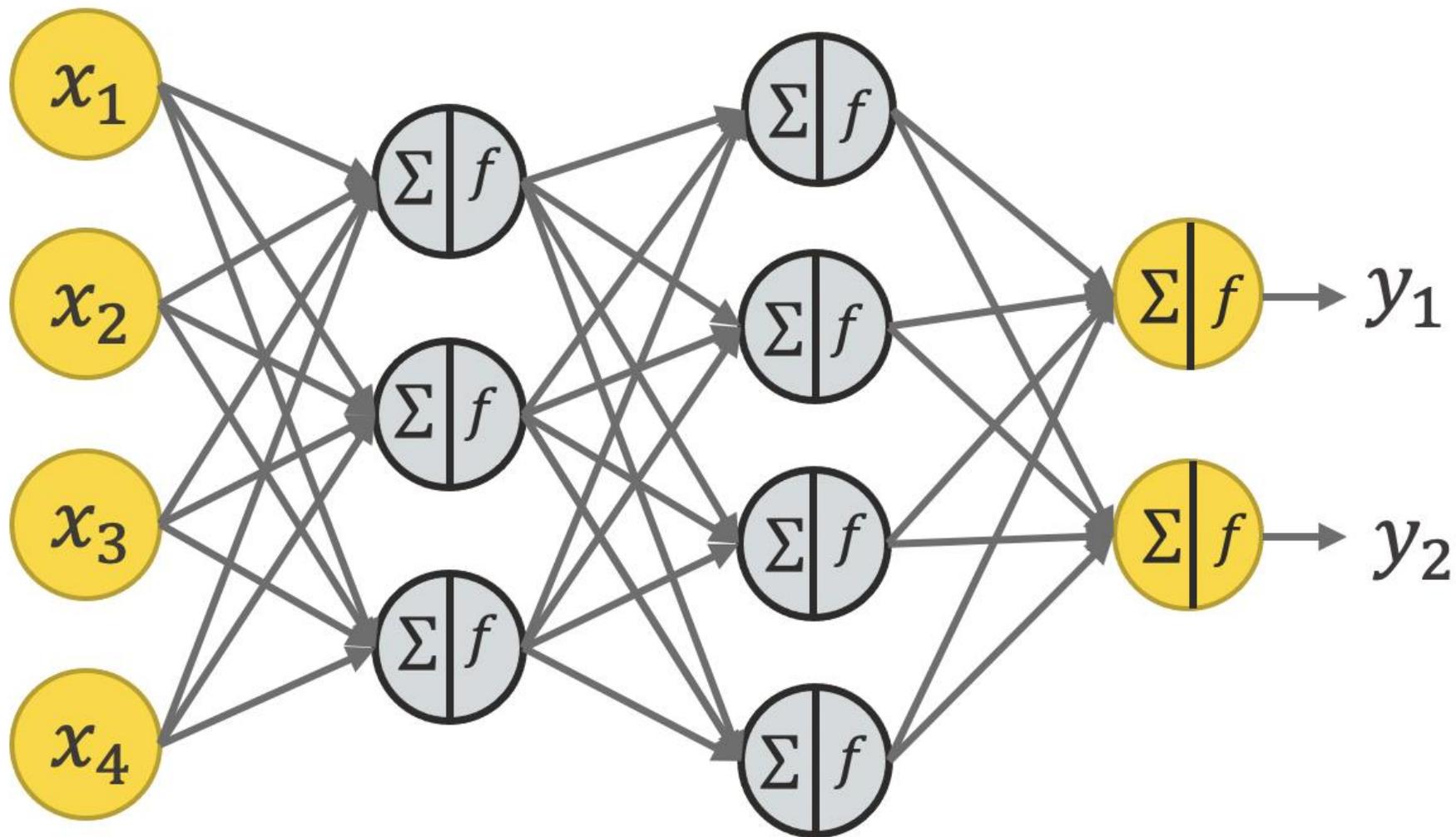# Challenges and Opportunities in Machine-Augmented Plant Stress Phenotyping

Arti Singh,[1,*] Sarah Jones,[1] Baskar Ganapathysubramanian,[2] Soumik Sarkar,[2] Daren Mueller,[3] Kulbir Sandhu,[1] and Koushik Nagasubramanian[4]
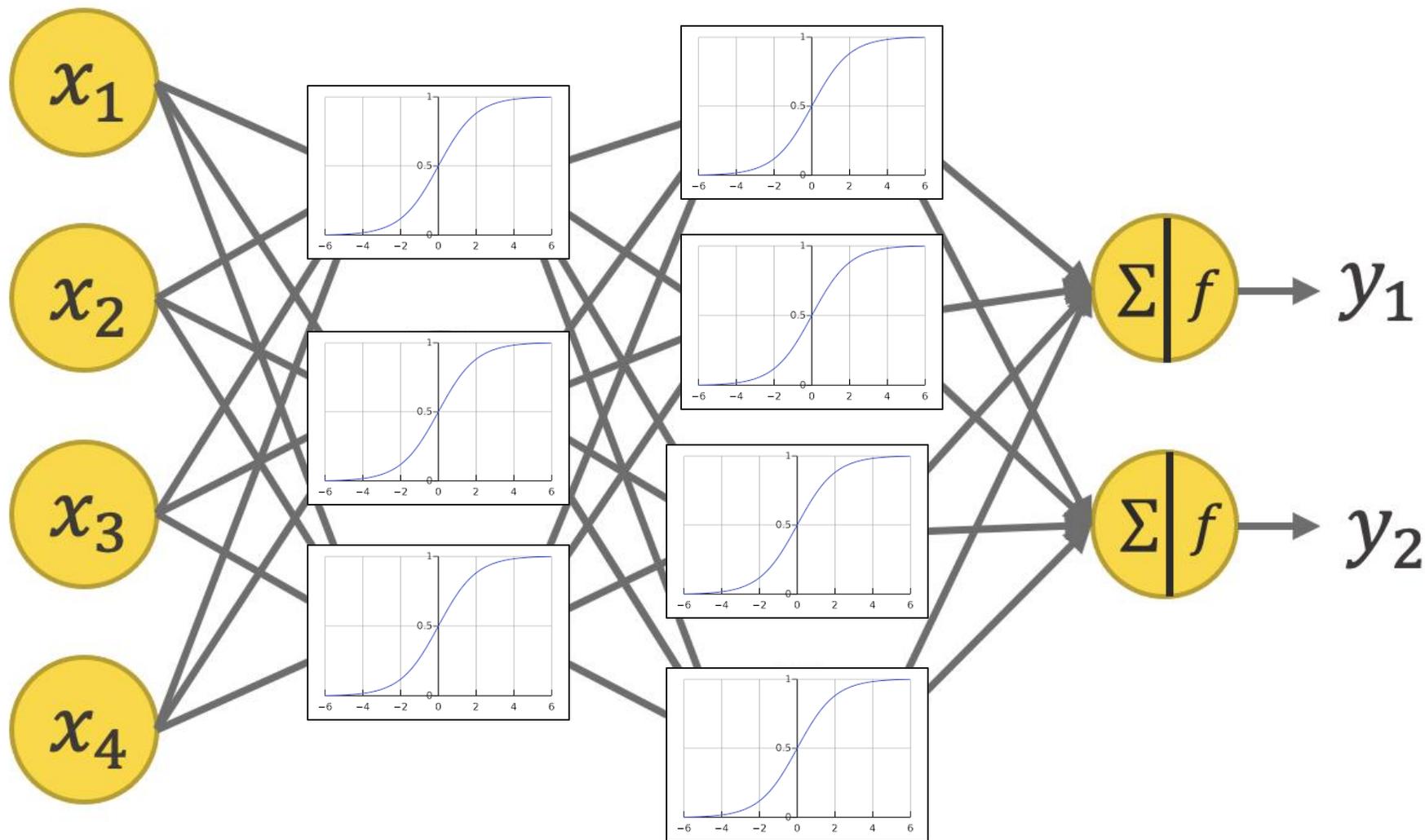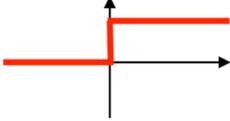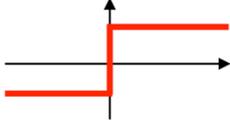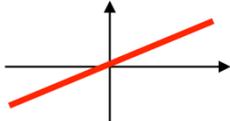
Input layer

Hidden layers
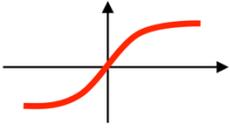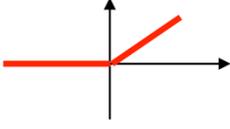
Output layer

Input layer

Hidden layers

Output layer

$x_1$

$x_2$

$x_3$

$x_4$

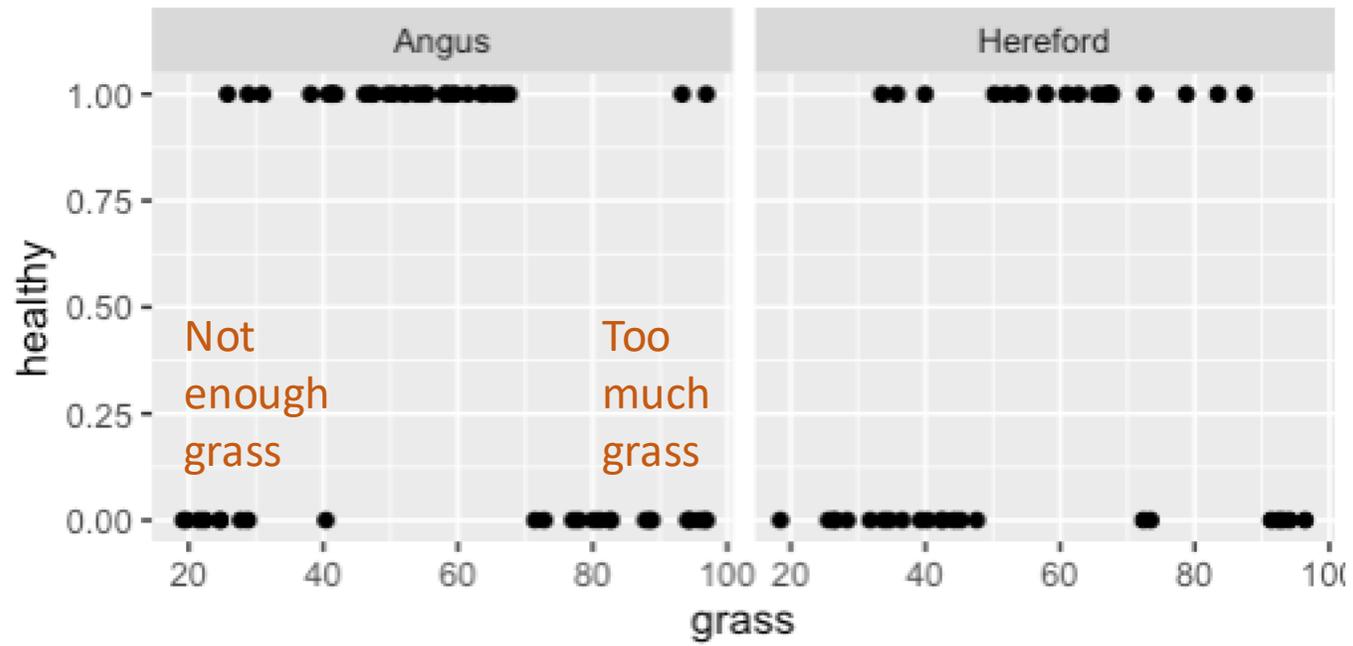$\Sigma \mid f$ → $y_1$

$\Sigma \mid f$ → $y_2$

# Activation functions

| Activation function | Equation | Example | 1D Graph |
|---|---|---|---|
| Unit step (Heaviside) | $\phi(z) = \begin{cases} 0, & z < 0, \\ 0.5, & z = 0, \\ 1, & z > 0, \end{cases}$ | Perceptron variant | |
| Sign (Signum) | $\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$ | Perceptron variant | |
| Linear | $\phi(z) = z$ | Adaline, linear regression | |
| Piece-wise linear | $\phi(z) = \begin{cases} 1, & z \geq \frac{1}{2}, \\ z + \frac{1}{2}, & -\frac{1}{2} < z < \frac{1}{2}, \\ 0, & z \leq -\frac{1}{2}, \end{cases}$ | Support vector machine | |
| Logistic (sigmoid) | $\phi(z) = \dfrac{1}{1 + e^{-z}}$ | Logistic regression, Multi-layer NN | |
| Hyperbolic tangent | $\phi(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | Multi-layer Neural Networks | |
| Rectifier, ReLU (Rectified Linear Unit) | $\phi(z) = max(0, z)$ | Multi-layer Neural Networks | |
| Rectifier, softplus | $\phi(z) = \ln(1 + e^z)$ | Multi-layer Neural Networks | |

```
model<-neuralnet(factor(healthy)~grass+breedHereford, data = mod_mat, hidden = c(2), rep=100,  threshold=.2, act.fct="logistic")
```