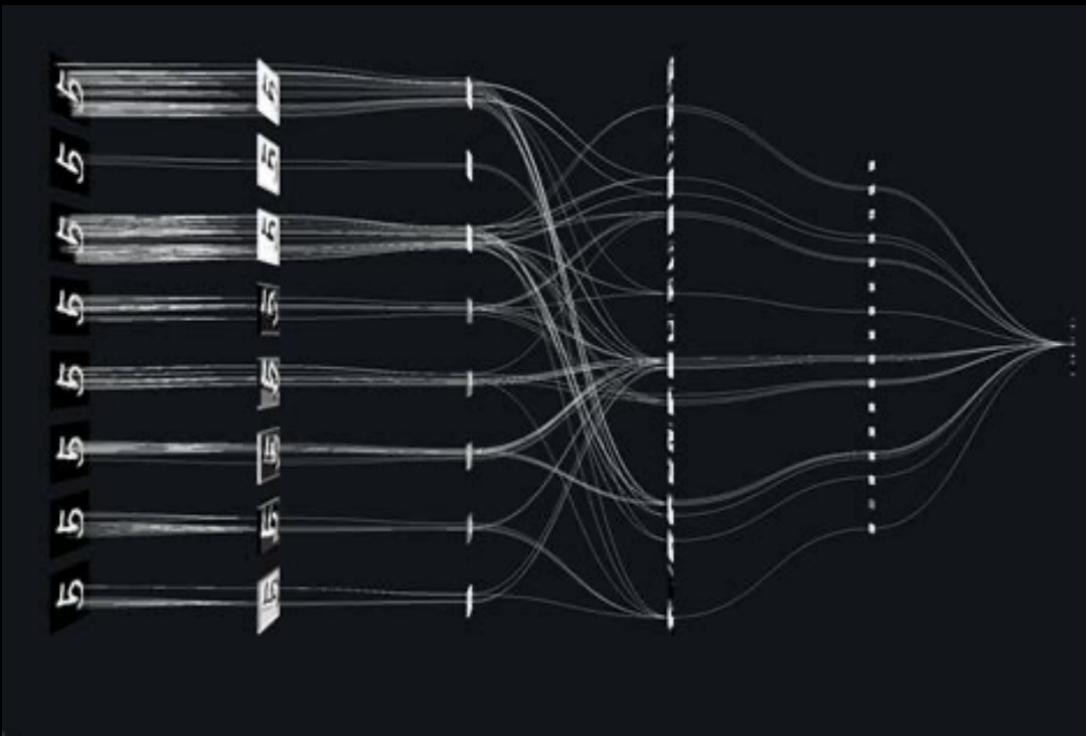


Explainable ML - XML

Pierre Gentine – Columbia University

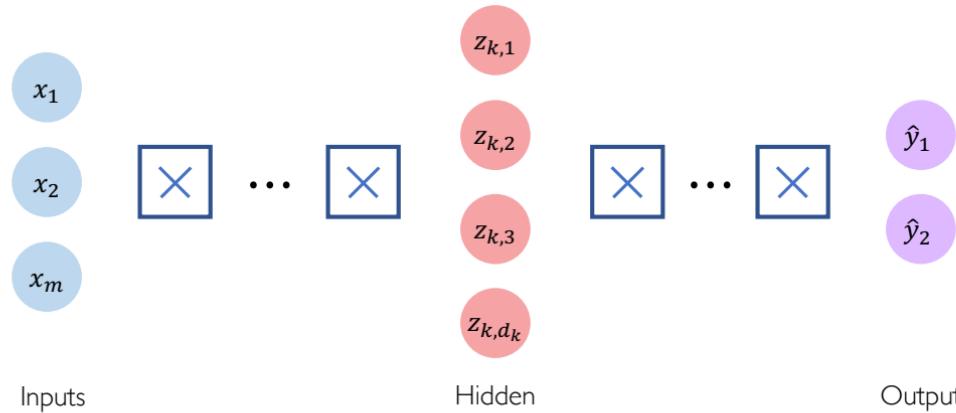


TRANSCENDING DISCIPLINES, TRANSFORMING LIVES

The problem

Major advances with the deep learning

Deep Neural Network



Many (many) weights → black box

How to open up the black box?

Interpreting what is in the gut of deep NN/CNN

Explainable ML

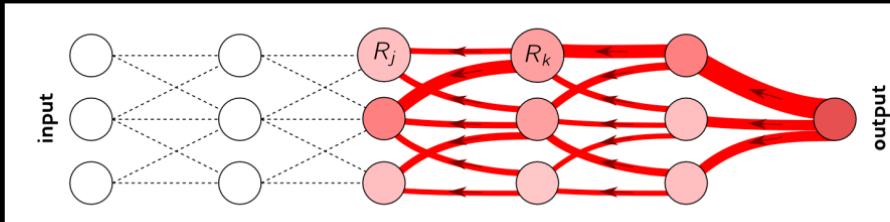
Also helps build trustworthiness in models
→ Right answer for the right reason

A few examples:

1. Layerwise Relevance Propagation (LRP)
2. Adjoint/gradient – Saliency maps
3. SHAPley values

1. Layerwise Relevance Propagation (LRP)

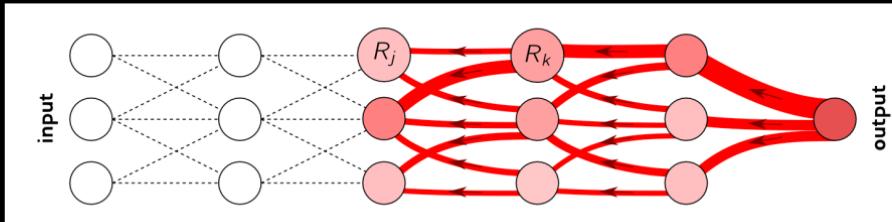
Trying to assess where the information is coming from.
Backward – from output to input that most explain the output



Pass information/relevance backward – a la Kirchoff law (conservation of current):

1. Layerwise Relevance Propagation (LRP)

Trying to assess where the information is coming from
Backward – from output to input that most explain the output



Pass information/relevance backward – a la Kirchoff law (conservation of current):

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

Layer j before layer k
 a_j – output from layer j (after activation function)
 w_{jk} - weights

1. Layerwise Relevance Propagation (LRP)

Different LRP rules

1. LRP-0

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

Equivalent to gradient x input

2. LRP-epsilon

Limit noisy information – filter by small noise epsilon

$$R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$$

3. LRP-gamma

Favor positive contributions compared to negative contributions

$$R_j = \sum_k \frac{a_j \cdot (w_{jk} + \gamma w_{jk}^+)}{\sum_{0,j} a_j \cdot (w_{jk} + \gamma w_{jk}^+)} R_k$$

1. Layerwise Relevance Propagation (LRP)

Different LRP rules → generic rule

$$R_j = \sum_k \frac{a_j \cdot \rho(w_{jk})}{\epsilon + \sum_{0,j} a_j \cdot \rho(w_{jk})} R_k$$

Can split computation into 4 substeps

$$\forall_k : z_k = \epsilon + \sum_{0,j} a_j \cdot \rho(w_{jk}) \quad (\text{forward pass})$$

$$\forall_k : s_k = R_k / z_k \quad (\text{element-wise division})$$

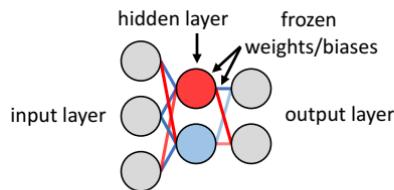
$$\forall_j : c_j = \sum_k \rho(w_{jk}) \cdot s_k \quad (\text{backward pass})$$

$$\forall_j : R_j = a_j c_j \quad (\text{element-wise product})$$

What we will do in the notebook – stay tuned

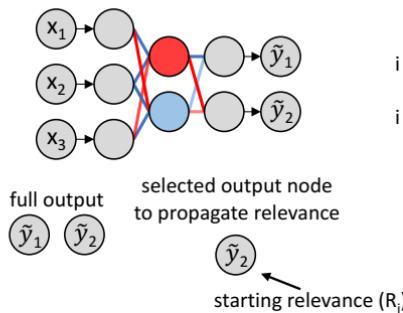
1. Layerwise Relevance Propagation (LRP)

- 1) Train network & freeze weights/biases

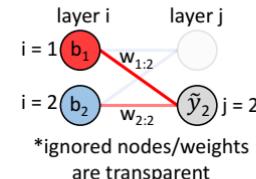


Information learned during training:
positive weights/biases
negative weights/biases

- 2) Input sample into frozen network and retain output



- 3) Propagate relevance from output node to previous layer



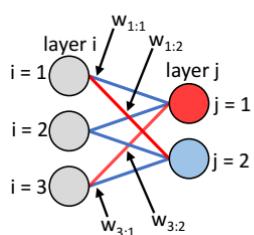
$$\text{propagation rule} \\ R_i = \sum_j \frac{a_i w_{ij}^+ + \max(0, b_j)}{\sum_i a_i w_{ij}^+ + \max(0, b_j)} R_j$$

example relevance calculations

$$R_{i=1} = \left(\frac{a_1 w_{1:2}}{a_1 w_{1:2} + a_2 w_{2:2}} \right) \tilde{y}_2$$

$$R_{i=2} = \left(\frac{a_2 w_{2:1}}{a_1 w_{1:2} + a_2 w_{2:2}} \right) \tilde{y}_2$$

- 4) Propagate relevance from hidden layer to input layer



propagation rule

$$R_i = \sum_j \left(\frac{w_{ij}^2}{\sum_i w_{ij}^2} \right) R_j$$

*all biases are ignored for this rule

relevance at input layer

$$(R_1) (R_2) (R_3)$$

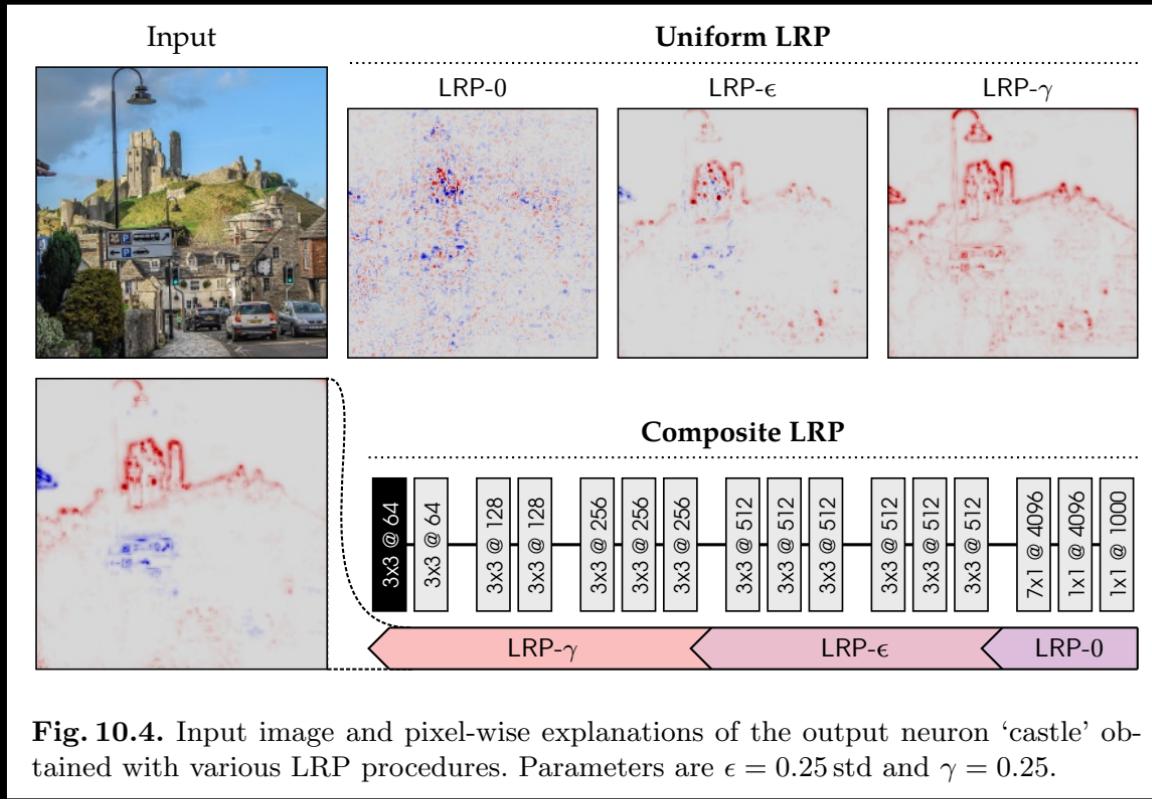
example relevance calculation

$$R_{i=1} = \left(\frac{w_{1:1}^2}{w_{1:1}^2 + w_{2:1}^2 + w_{3:1}^2} \right) R_{j=1} + \left(\frac{w_{1:2}^2}{w_{1:2}^2 + w_{2:2}^2 + w_{3:2}^2} \right) R_{j=2}$$

*relevance calculations are similar for other input nodes

- 5) Repeat for each sample of interest...

1. Layerwise Relevance Propagation (LRP)



2. Gradient/adjoint

Gradient/adjoint of output \mathbf{y} to inputs \mathbf{x}

$$\nabla_{\mathbf{x}} \mathbf{y}$$

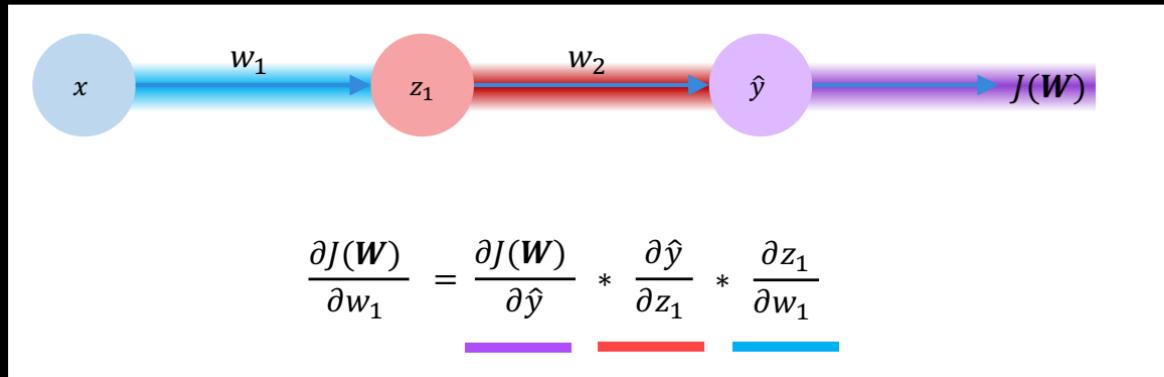
$$J = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial x_1} & \dots & \frac{\partial \mathbf{y}}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

2. Gradient/adjoint

Gradient/adjoint of output y to inputs x

$$J = \begin{pmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

Once upon a time: backpropagation
Chain's rule



Same thing!

2. Gradient/adjoint

Personal experience:

Non-linear so average over different samples

Use non-dimensional/normalized inputs (and ideally outputs)

Tends to work well for first cut – but can be noisy

2. Gradient/adjoint

Personal experience:

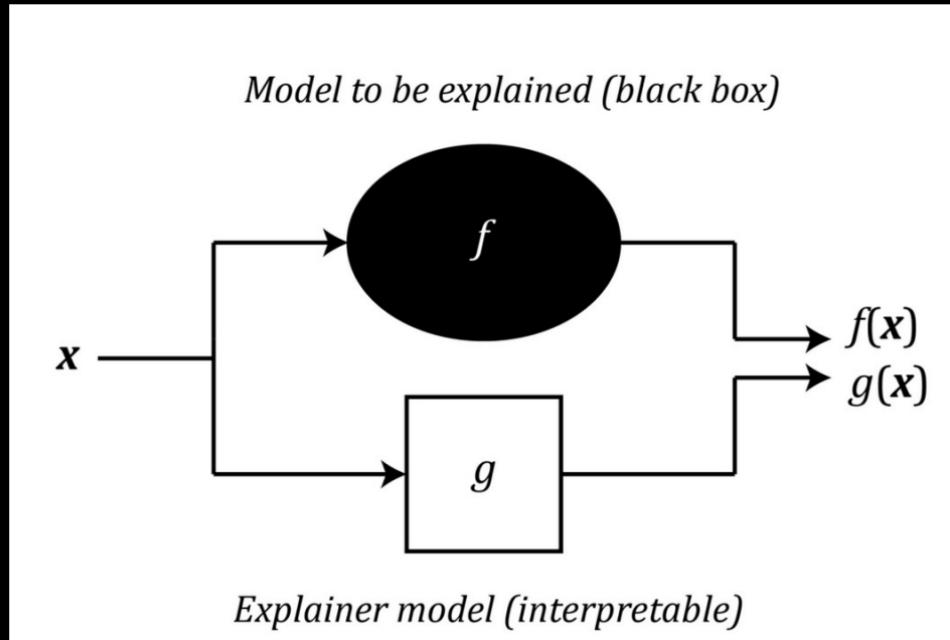
Non-linear so average over different samples

Use non-dimensional/normalized inputs (and ideally outputs)

Tends to work well for first cut – but can be noisy

3. Shapley Additive Explanations (SHAP) values

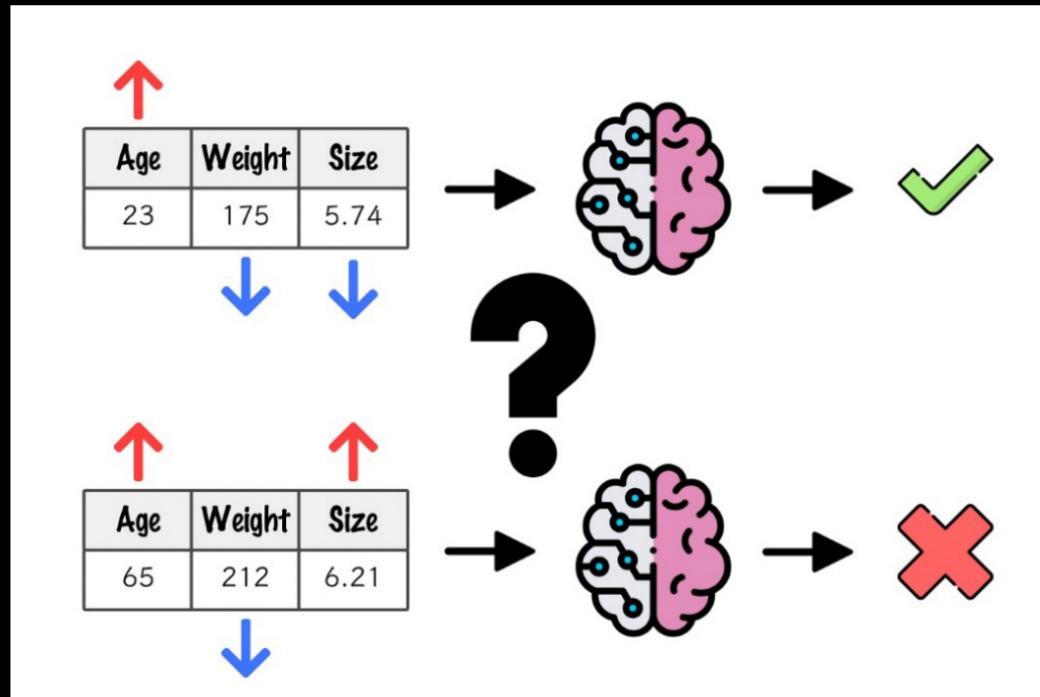
Inspired by Game theory.



g mimics f for a single prediction

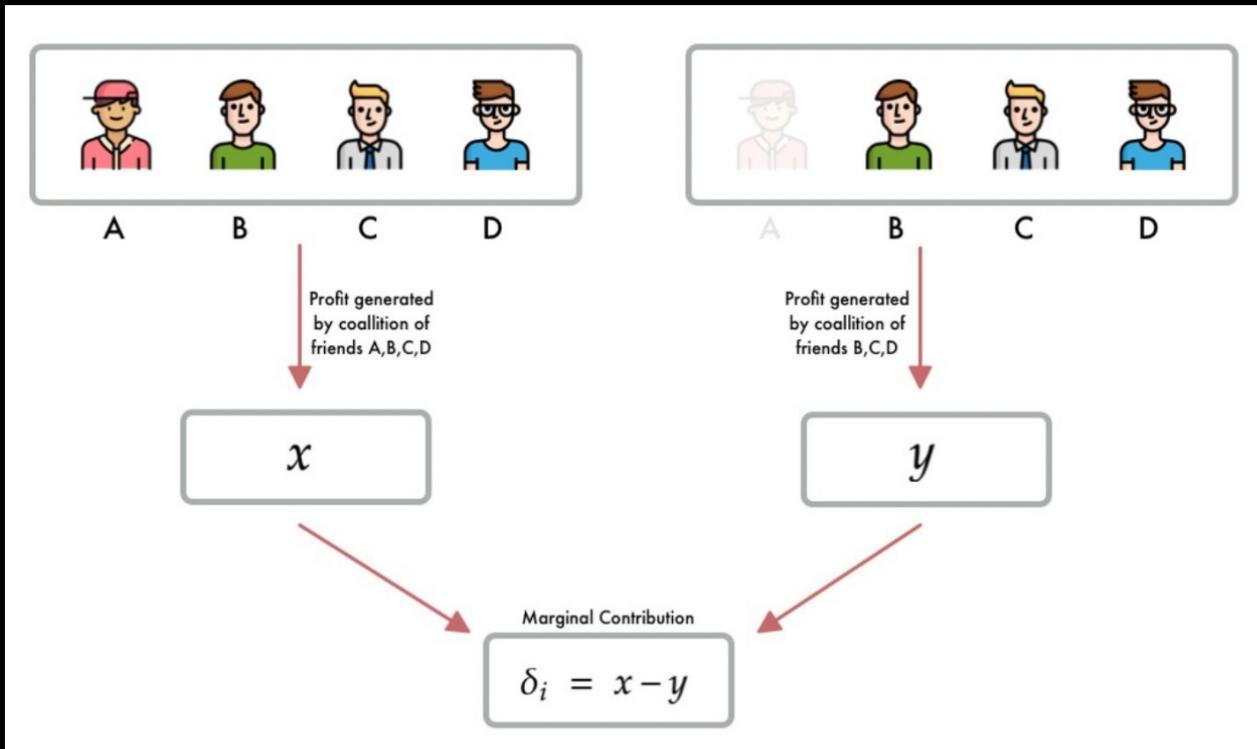
3. SHAP values

Inspired by Game theory.



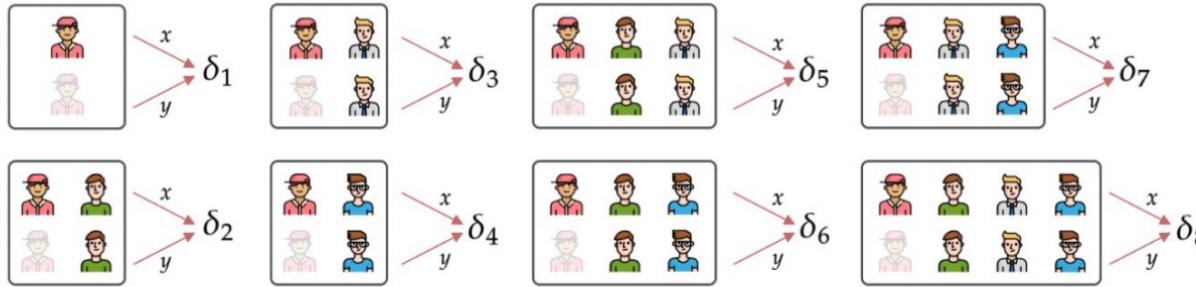
3. SHAP values

Marginal contribution of member A



3. SHAP values

Marginal contribution of member A



The Shapley value for member

is given by:

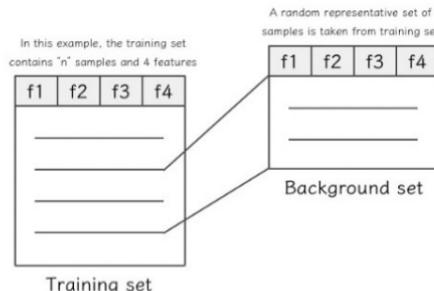
$$\phi_i = \frac{\delta_1 + \delta_3 + \delta_4 + \delta_5 + \delta_6 + \delta_7 + \delta_8}{8}$$

3. SHAP values

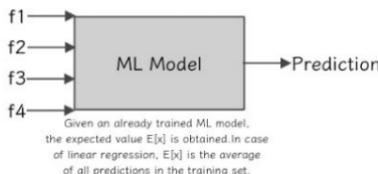
Say we want to explain feature "x":

f1	f2	f3	f4

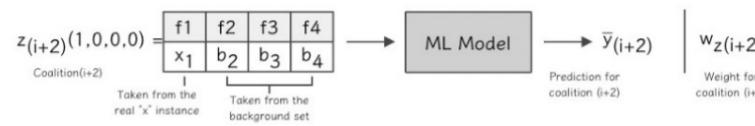
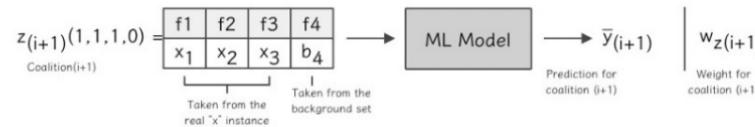
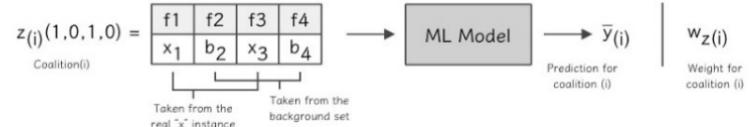
We have training set and background set:



And of course, an already trained ML model:



Coalitions, predictions and weights are calculated:



Weights are calculated with Equation 1!

The weighted linear model is fit

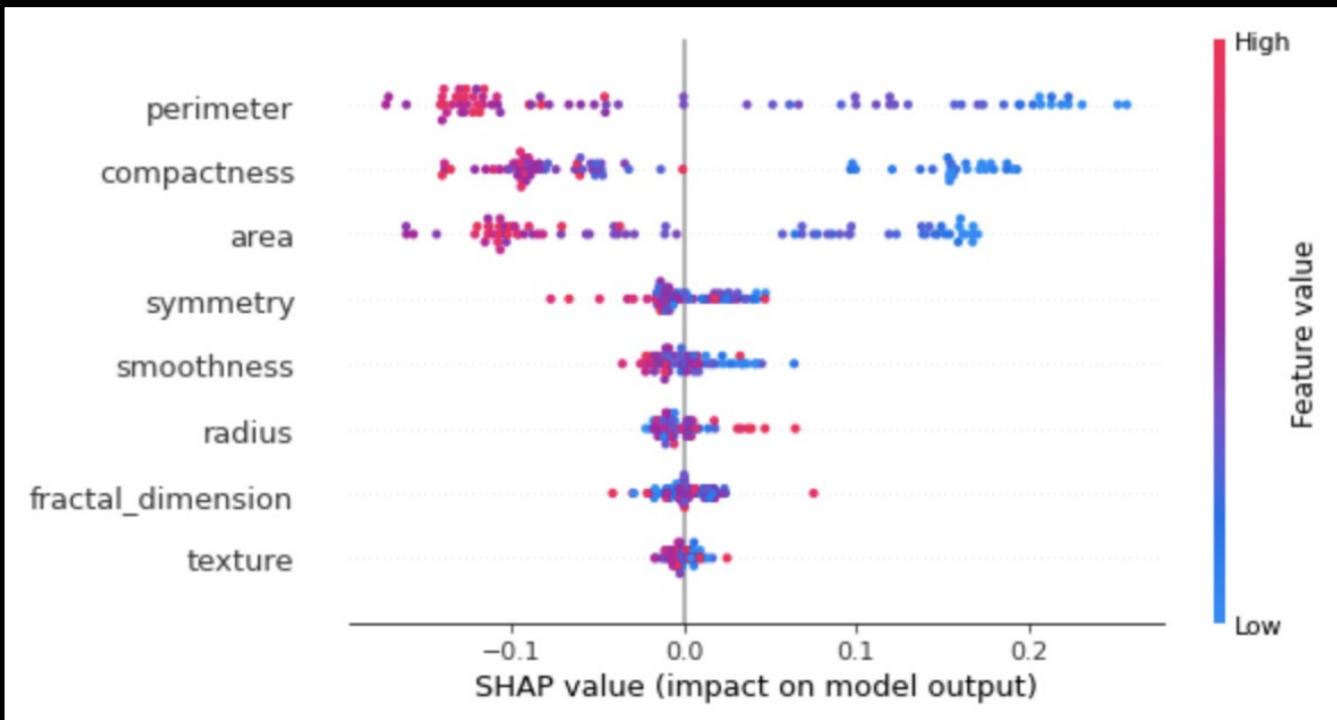
With coalitions, predictions and weights, the weighted linear model is built.

Shapley values are obtained!

Once optimized the weighted linear model, the coefficients are the Shapley values!

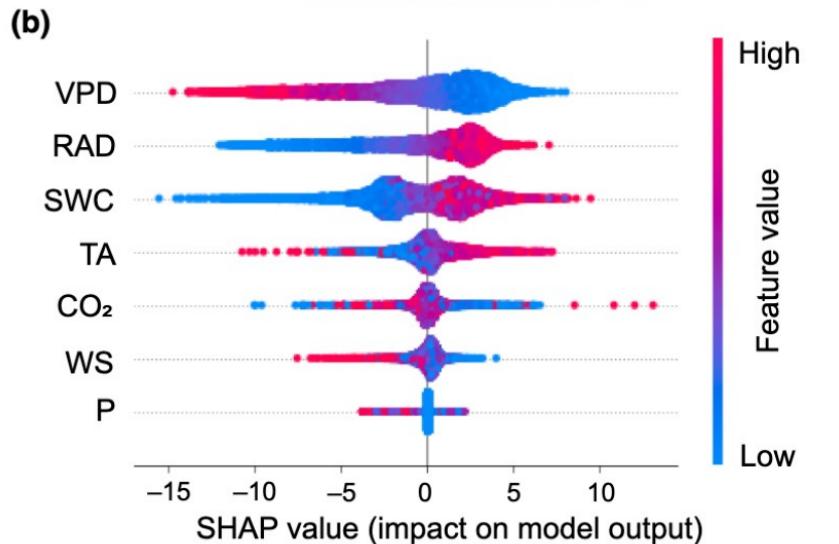
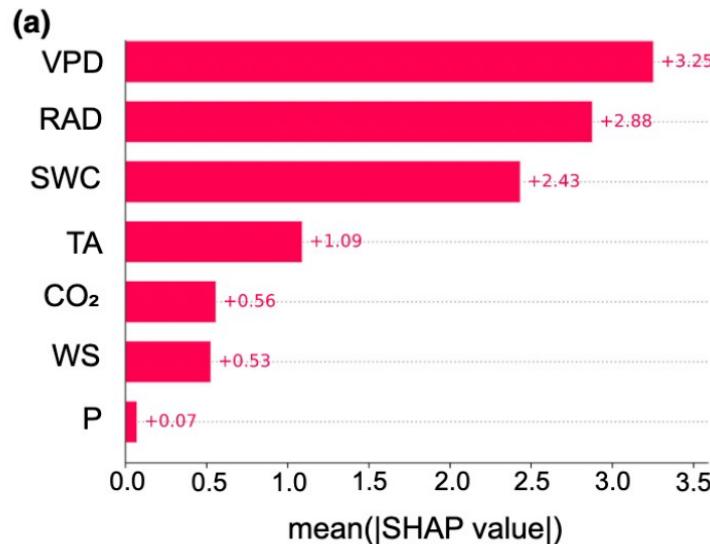
3. SHAP values

Generated results



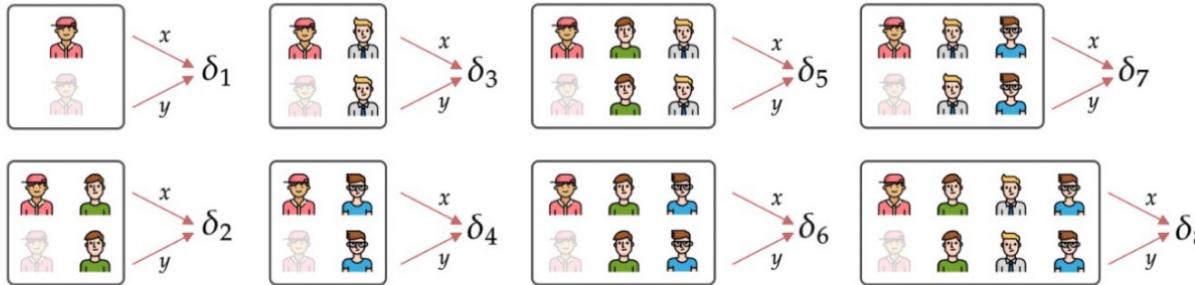
3. SHAP values

Example for continuous variables



3. SHAP values

Marginal contribution of member A



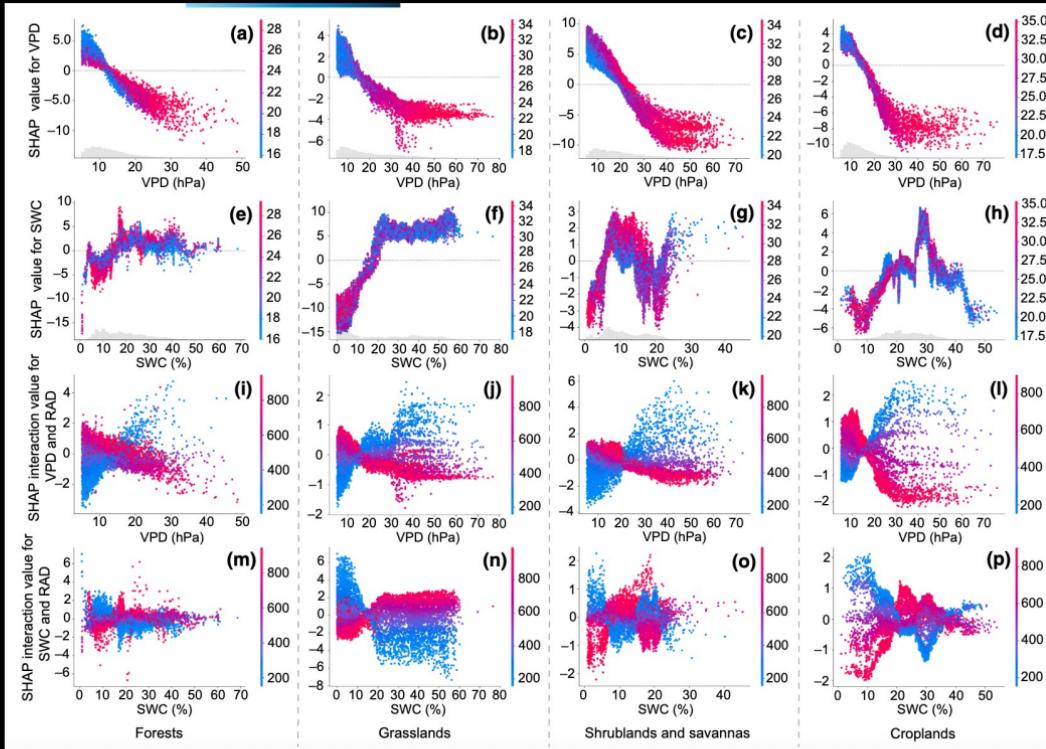
The Shapley value for member

is given by:

$$\phi_i = \frac{\delta_1 + \delta_3 + \delta_4 + \delta_5 + \delta_6 + \delta_7 + \delta_8}{8}$$

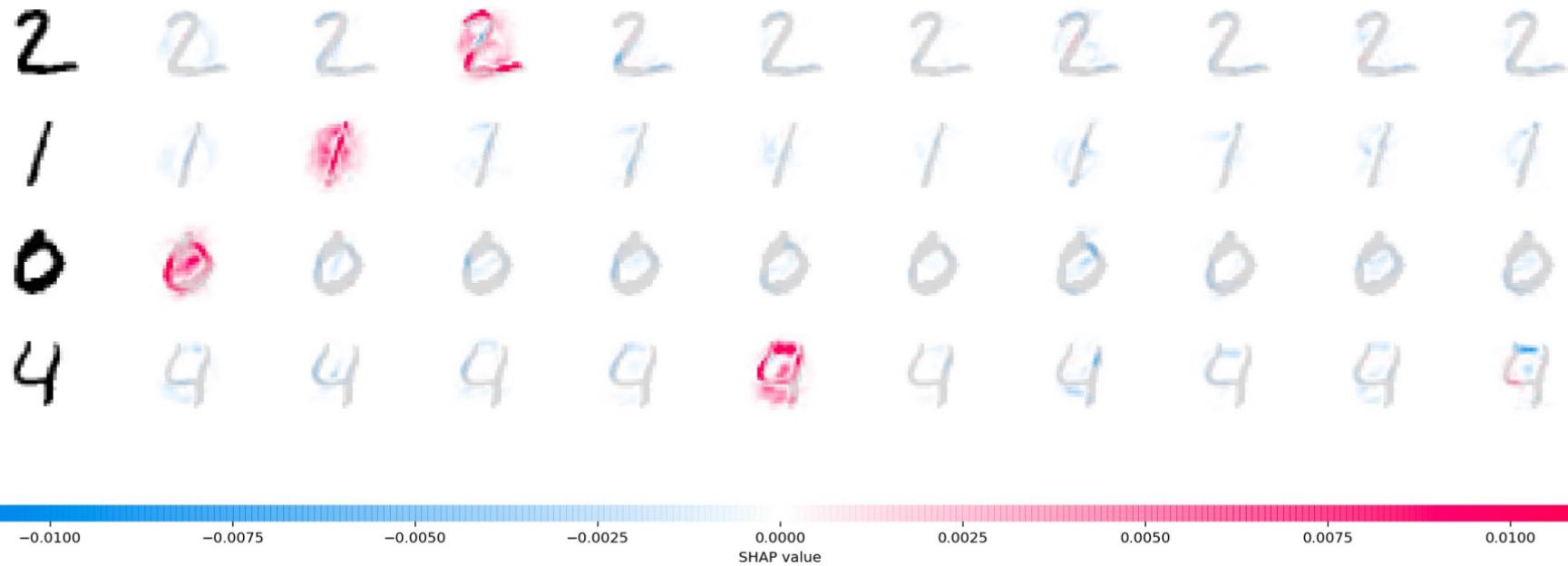
3. SHAP values

Example for continuous variables



3. SHAP values

Example for images



3. SHAP values

Example for images

