# Interpretable Machine Learning

# SHAP (SHapley Additive exPlanation) Values

-		-	 -	-	 	 -	-	
-	1							
Pull-117	e 1979 - 5						1994 - 1994	

#### Learning goals

- Get an intuition of additive feature attributions
- Understand the concept of Kernel SHAP
- Ability to interpret SHAP plots
- Global SHAP methods



**Definition:** A kernel-based, model-agnostic method to compute Shapley values via local surrogate models (e.g. linear model)

- Sample coalitions
- 2 Transfer coalitions into feature space & get predictions by applying ML model
- Occupie weights through kernel
- Fit a weighted linear model
- 8 Return Shapley values



Step 1: Sample coalitions

• Sample K coalitions from the simplified feature space

$${f z}'^{(k)} \in \{0,1\}^p, \ \ k \in \{1,\ldots,K\}$$

• For our simple example, we have in total  $2^p = 2^3 = 8$  coalitions (without sampling)



Coalition	<b>z</b> ′ <sup>(k)</sup>	hum	temp	ws	
Ø	<b>z</b> ′ <sup>(1)</sup>	0	0	0	
hum	<b>z</b> ′ <sup>(2)</sup>	1	0	0	
temp	<b>z</b> ′ <sup>(3)</sup>	0	1	0	
WS	<b>z</b> ′ <sup>(4)</sup>	0	0	1	
hum, temp	<b>z</b> ′ <sup>(5)</sup>	1	1	0	
temp, ws	<b>z</b> ′ <sup>(6)</sup>	0	1	1	
hum, ws	<b>z</b> ′ <sup>(7)</sup>	1	0	1	
hum, temp, ws	<b>z</b> ′ <sup>(8)</sup>	1	1	1	

Step 2: Transfer Coalitions into feature space & get predictions by applying ML model

- $\mathbf{z}^{\prime(k)}$  is 1 if features are part of the *k*-th coalition, 0 if they are absent
- To calculate predictions for these coalitions, we need to define a function which maps the binary feature space back to the original feature space

	_					$\rightarrow$		
Coalition	<b>z</b> ′ <sup>(k)</sup>	hum	temp	ws	<b>x</b> <sup>coalition</sup>	hum	temp	ws
Ø	<b>z</b> ′ <sup>(1)</sup>	0	0	0	<b>x</b> {Ø}	Ø	Ø	Ø
hum	<b>z</b> ′ <sup>(2)</sup>	1	0	0	<b>x</b> { <i>hum</i> }	51.6	Ø	Ø
temp	<b>z</b> ′ <sup>(3)</sup>	0	1	0	<b>x</b> {temp}	ø	5.1	Ø
WS	<b>z</b> ′ <sup>(4)</sup>	0	0	1	<b>x</b> <sup>{<i>ws</i>}</sup>	ø	Ø	17.0
hum, temp	<b>z</b> ′ <sup>(5)</sup>	1	1	0	x <sup>{hum,tem</sup> }	<sup>o}</sup> 51.6	5.1	Ø
temp, ws	<b>z</b> ′ <sup>(6)</sup>	0	1	1	x <sup>{temp,ws</sup> ]	+ ø	5.1	17.0
hum, ws	<b>z</b> ′ <sup>(7)</sup>	1	0	1	<b>x</b> <sup>{hum,ws</sup> }	51.6	Ø	17.0
hum, temp, ws	<b>z</b> ′ <sup>(8)</sup>	1	1	1	<b>x</b> <sup>{hum,temµ</sup>	<sup>o,ws}</sup> 51.6	5.1	17.0



Step 2: Transfer Coalitions into feature space & get predictions by applying ML model

- Define h<sub>x</sub> (z'<sup>(k)</sup>) = z<sup>(k)</sup> where h<sub>x</sub> : {0, 1}<sup>p</sup> → ℝ<sup>p</sup> maps 1's to feature values of observation x for features part of the k-th coalition and 0's to feature values of a randomly sampled observation for features absent in the k-th coalition (feature values are permuted multiple times)
- Predict with ML model on this dataset  $\hat{f} : \hat{f} (h_x (\mathbf{z}'^{(k)}))$

	-				$h_x(\mathbf{z}^{\prime(k)})$				*	
Coalition	$\mathbf{z}^{\prime(k)}$	hum	temp	ws	_	<b>z</b> <sup>(k)</sup>	hum	temp	WS	$\hat{f}\left(h_{x}\left(\mathbf{z}^{\prime\left(k\right)}\right)\right)$
Ø	<b>z</b> ′ <sup>(1)</sup>	0	0	0		<b>z</b> <sup>(1)</sup>	64.3	28.0	14.5	6211
hum	<b>z</b> ′ <sup>(2)</sup>	1	0	0		<b>z</b> <sup>(2)</sup>	51.6	28.0	14.5	5586
temp	<b>z</b> ′ <sup>(3)</sup>	0	1	0		<b>z</b> <sup>(3)</sup>	64.3	5.1	14.5	3295
ws	<b>z</b> ′ <sup>(4)</sup>	0	0	1		<b>z</b> <sup>(4)</sup>	64.3	28.0	17.0	5762
hum, temp	<b>z</b> ′ <sup>(5)</sup>	1	1	0		<b>z</b> <sup>(5)</sup>	51.6	5.1	14.5	2616
temp, ws	<b>z</b> ′ <sup>(6)</sup>	0	1	1		<b>z</b> <sup>(6)</sup>	64.3	5.1	17.0	2900
hum, ws	<b>z</b> ′ <sup>(7)</sup>	1	0	1		<b>z</b> <sup>(7)</sup>	51.6	28.0	17.0	5411
hum, temp, ws	<b>z</b> ′ <sup>(8)</sup>	1	1	1		<b>z</b> <sup>(8)</sup>	51.6	5.1	17.0	2573



Step 3: Compute weights through Kernel

**Intuition**: We learn most about individual features if we can study their effects in isolation or at maximal interaction: Small coalitions (few 1's) and large coalitions (i.e. many 1's) get the largest weights





Step 3: Compute weights through Kernel 
see shapley\_kernel\_proof.pdf

**Intuition**: We learn most about individual features if we can study their effects in isolation or at maximal interaction: Small coalitions (few 1's) and large coalitions (i.e. many 1's) get the largest weights





### Step 3: Compute weights through Kernel

**Purpose**: to include this knowledge in the local surrogate model (linear regression), we calculate weights for each coalition which are the observations of the linear regression

$$\pi_{x}(\mathbf{z}') = \frac{(p-1)}{\binom{p}{|\mathbf{z}'|}} \rightsquigarrow \pi_{x}(\mathbf{z}' = (1,0,0)) = \frac{(3-1)}{\binom{3}{1} 1 (3-1)} = \frac{1}{3}$$

Coalition	$\mathbf{z}^{\prime(k)}$	hum	temp	ws	weight
Ø	<b>z</b> ′ <sup>(1)</sup>	0	0	0	$\infty$
hum	<b>z</b> ′ <sup>(2)</sup>	1	0	0	0.33
temp	<b>z</b> ′ <sup>(3)</sup>	0	1	0	0.33
ws	$z'^{(4)}$	0	0	1	0.33
hum, temp	<b>z</b> ′ <sup>(5)</sup>	1	1	0	0.33
temp, ws	<b>z</b> ′ <sup>(6)</sup>	0	1	1	0.33
hum, ws	$z'^{(7)}$	1	0	1	0.33
hum, temp, ws	<b>z</b> ′ <sup>(8)</sup>	1	1	1	$\infty$



### Step 3: Compute weights through Kernel

**Purpose**: to include this knowledge in the local surrogate model (linear regression), we calculate weights for each coalition which are the observations of the linear regression

Coalition	$\mathbf{z}^{\prime(k)}$	hum	temp	ws	weight
Ø	<b>z</b> ′ <sup>(1)</sup>	0	0	0	$\infty$
hum	<b>z</b> ′ <sup>(2)</sup>	1	0	0	0.33
temp	<b>z</b> ′ <sup>(3)</sup>	0	1	0	0.33
WS	<b>z</b> ′ <sup>(4)</sup>	0	0	1	0.33
hum, temp	<b>z</b> ′ <sup>(5)</sup>	1	1	0	0.33
temp, ws	<b>z</b> ′ <sup>(6)</sup>	0	1	1	0.33
hum, ws	<b>z</b> ′ <sup>(7)</sup>	1	0	1	0.33
hum, temp, ws	<b>z</b> ′ <sup>(8)</sup>	1	1	1	$\infty$

 $\leadsto$  weights for empty and full set are infinity and not used as observations for the linear regression

 $\leadsto$  instead constraints are used such that properties (local accuracy and missingness) are satisfied



### Step 4: Fit a weighted linear model

**Aim**: Estimate a weighted linear model with Shapley values being the coefficients  $\phi_i$ 

$$g\left(\mathbf{z}^{\prime(k)}\right) = \phi_0 + \sum_{j=1}^{p} \phi_j z_j^{\prime(k)}$$

and minimize by WLS using the weights  $\pi_x$  of step 3

$$L\left(\hat{f}, g, \pi_{x}\right) = \sum_{k=1}^{K} \left[\hat{f}\left(h_{x}\left(\mathbf{z}^{\prime(k)}\right)\right) - g\left(\mathbf{z}^{\prime(k)}\right)\right]^{2} \pi_{x}\left(\mathbf{z}^{\prime(k)}\right)$$

with  $\phi_0 = \mathbb{E}(\hat{f})$  and  $\phi_p = \hat{f}(x) - \sum_{j=0}^{p-1} \phi_j$  we receive a p-1 dimensional linear regression problem



### Step 4: Fit a weighted linear model

**Aim**: Estimate a weighted linear model with Shapley values being the coefficients  $\phi_i$ 

$$g\left(\mathbf{z}^{\prime(k)}\right) = \phi_0 + \sum_{j=1}^{p} \phi_j z_j^{\prime(k)} \rightsquigarrow g\left(\mathbf{z}^{\prime(k)}\right) = 4515 + 34 \cdot z_1^{\prime(k)} - 1654 \cdot z_2^{\prime(k)} - 323 \cdot z_3^{\prime(k)}$$

$\mathbf{z}^{\prime(k)}$	hum	temp	ws	weight	Î
<b>z</b> ′ <sup>(2)</sup>	1	0	0	0.33	4635
$z'^{(3)}$	0	1	0	0.33	3087
$\mathbf{z}^{\prime(4)}$	0	0	1	0.33	4359
$z'^{(5)}$	1	1	0	0.33	3060
$\mathbf{z}^{\prime(6)}$	0	1	1	0.33	2623
<b>z</b> ′ <sup>(7)</sup>	1	0	1	0.33	4450
	output				



### Step 5: Return SHAP values

Intuition: Estimated Kernel SHAP values are equivalent to Shapley values

$$g(\mathbf{z}^{\prime(8)}) = \hat{f}(h_x(\mathbf{z}^{\prime(8)})) = 4515 + 34 \cdot 1 - 1654 \cdot 1 - 323 \cdot 1$$
$$= \underbrace{\mathbb{E}(\hat{f})}_{\phi_0} + \phi_{hum} + \phi_{temp} + \phi_{ws} = \hat{f}(\mathbf{x}) = 2573$$



