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



Question: How much does a feat. *j* contribute to the prediction of a single obs. **Idea:** Use Shapley values from cooperative game theory



Question: How much does a feat. *j* contribute to the prediction of a single obs. **Idea:** Use Shapley values from cooperative game theory **Procedure:**

- Compare "reduced prediction function" of feature coalition S with $S \cup \{j\}$
- Iterate over possible coalitions to calculate marginal contribution of feature *j* to sample x

$$\phi_{j} = \frac{1}{p!} \sum_{\tau \in \Pi} \underbrace{\hat{f}_{S_{j}^{\tau} \cup \{j\}}(\mathbf{x}_{S_{j}^{\tau} \cup \{j\}}) - \hat{f}_{S_{j}^{\tau}}(\mathbf{x}_{S_{j}^{\tau}})}_{\text{marginal contribution of feature } j}$$



Question: How much does a feat. *j* contribute to the prediction of a single obs. **Idea:** Use Shapley values from cooperative game theory **Procedure:**

- Compare "reduced prediction function" of feature coalition S with $S \cup \{j\}$
- Iterate over possible coalitions to calculate marginal contribution of feature *j* to sample x

$$\phi_{j} = \frac{1}{p!} \sum_{\tau \in \Pi} \hat{f}_{S_{j}^{\tau} \cup \{j\}} (\mathbf{x}_{S_{j}^{\tau} \cup \{j\}}) - \hat{f}_{S_{j}^{\tau}} (\mathbf{x}_{S_{j}^{\tau}})$$
marginal contribution of feature *j*

Remember:

- \hat{f} is the prediction function, *p* denotes the number of features
- Non-existent feat. in a coalition are replaced by values of random feat. values
- Recall S_i^{τ} defines the coalition as the set of players before player *j* in order

$$\tau = (\tau^{(1)}, \dots, \tau^{(p)})$$

$$\tau^{(1)} \qquad \tau^{(|S|)} \qquad \tau^{(|S|+1)} \qquad \tau^{(|S|+2)} \qquad \cdots \qquad \tau^{(p)}$$

 S_j^{τ} : Players before player j player j Players after player j



Example:

- Train a random forest on bike sharing data only using features humidity (hum), temperature (temp) and windspeed (ws)
- Calculate Shapley value for an observation **x** with $\hat{f}(\mathbf{x}) = 2573$
- Mean prediction is $\mathbb{E}(\hat{f}) = 4515$



Example:

- Train a random forest on bike sharing data only using features humidity (hum), temperature (temp) and windspeed (ws)
- Calculate Shapley value for an observation **x** with $\hat{f}(\mathbf{x}) = 2573$
- Mean prediction is $\mathbb{E}(\hat{f}) = 4515$

Exact Shapley calculation for humidity:

S	$\boldsymbol{S}\cup\{j\}$	\hat{f}_S	$\hat{f}_{S\cup\{j\}}$	weight
Ø	hum	4515	4635	2/6
temp	temp, hum	3087	3060	1/6
ws	ws, hum	4359	4450	1/6
temp, ws	hum, temp, ws	2623	2573	2/6

$$\phi_{hum} = \frac{2}{6}(4635 - 4515) + \frac{1}{6}(3060 - 3087) + \frac{1}{6}(4450 - 4359) + \frac{2}{6}(2573 - 2623) = 34$$



FROM SHAPLEY TO SHAP

Example continued: Same calculation can be done for temperature and windspeed:

• $\phi_{temp} = ... = -1654$

•
$$\phi_{ws} = \ldots = -323$$

Remember: Shapley values explain difference between actual and average pred.:

$$\begin{array}{ll} 2573 - 4515 & = 34 - 1654 - 323 & = -1942 \\ \hat{f}(\mathbf{x}) - \mathbb{E}(\hat{f}) & = \phi_{hum} + \phi_{temp} + \phi_{ws} \end{array}$$



 \rightsquigarrow can be rewritten to

$$\hat{f}(\mathbf{x}) = \underbrace{\mathbb{E}(\hat{f})}_{\phi_0} + \phi_{hum} + \phi_{temp} + \phi_{ws}$$



Actual prediction: 2572.67;

SHAP DEFINITION • Lundberg et al. 2017

Aim: Find an additive combination that explains the prediction of an observation \mathbf{x} by computing the contribution of each feature to the prediction using a (more efficient) estimation procedure.

Definition

- Simplified (binary) coalition feat. space $\mathbf{Z}' \in \{0,1\}^{K \times p}$ with K rows and p cols.
- Rows are referred to as $\mathbf{z}'^{(k)} = \{z'^{(k)}_1, \dots, z'^{(k)}_p\}$ with $k \in \{1, \dots, K\}$ (indexes *k*-th coalition)
- Cols are referred to as \mathbf{z}_j with $j \in \{1, \dots, p\}$ being the index of the original feat.

Example:

Coalition	$\mathbf{z}^{\prime(k)}$	hum	temp	WS	
Ø	z ′ ⁽¹⁾	0	0	0	
hum	z ′ ⁽²⁾	1	0	0	
temp	$z'^{(3)}$	0	1	0	
ws	$z'^{(4)}$	0	0	1	
hum, temp	z ′ ⁽⁵⁾	1	1	0	
temp, ws	z ′ ⁽⁶⁾	0	1	1	
hum, ws	$z'^{(7)}$	1	0	1	
hum, temp, ws	z ′ ⁽⁸⁾	1	1	1	



SHAP DEFINITION • Lundberg et al. 2017

Aim: Find an additive combination that explains the prediction of an observation \mathbf{x} by computing the contribution of each feature to the prediction using a (more efficient) estimation procedure.

Definition

- Simplified (binary) coalition feat. space $\mathbf{Z}' \in \{0,1\}^{K \times p}$ with K rows and p cols.
- Rows are referred to as $\mathbf{z}'^{(k)} = \{z'^{(k)}_1, \dots, z'^{(k)}_p\}$ with $k \in \{1, \dots, K\}$ (indexes *k*-th coalition)
- Cols are referred to as z_j with $j \in \{1, ..., p\}$ being the index of the original feat.





SHAP DEFINITION Lundberg et al. 2017

Aim: Find an additive combination that explains the prediction of an observation \mathbf{x} by computing the contribution of each feature to the prediction using a (more efficient) estimation procedure.

Definition

- Simplified (binary) coalition feat. space $\mathbf{Z}' \in \{0,1\}^{K \times p}$ with K rows and p cols.
- Rows are referred to as $\mathbf{z}'^{(k)} = \{z'^{(k)}_1, \dots, z'^{(k)}_p\}$ with $k \in \{1, \dots, K\}$ (indexes *k*-th coalition)
- Cols are referred to as z_j with $j \in \{1, ..., p\}$ being the index of the original feat.





SHAP DEFINITION Lundberg et al. 2017

Aim: Find an additive combination that explains the prediction of an observation **x** by computing the contribution of each feature to the prediction using a (more efficient) estimation procedure.

Definition

- Simplified (binary) coalition feat. space $\mathbf{Z}' \in \{0,1\}^{K \times p}$ with K rows and p cols.
- Rows are referred to as $\mathbf{z}'^{(k)} = \{z'^{(k)}_1, \dots, z'^{(k)}_p\}$ with $k \in \{1, \dots, K\}$ (indexes *k*-th coalition)
- Cols are referred to as z_j with $j \in \{1, ..., p\}$ being the index of the original feat.





SHAP DEFINITION • Lundberg et al. 2017

Aim: Find an additive combination that explains the prediction of an observation \mathbf{x} by computing the contribution of each feature to the prediction using a (more efficient) estimation procedure.





Problem

How do we estimate the Shapley values ϕ_j ?

Local Accuracy

$$f(\mathbf{x}) = g(\mathbf{x}') = \phi_0 + \sum_{j=1}^{p} \phi_j x'_j$$

Intuition: If the coalition includes all features $(\mathbf{x}' \in \{1\}^p)$, the attributions ϕ_j and the null output ϕ_0 sum up to the original model output $f(\mathbf{x})$

Local accuracy corresponds to the axiom of efficiency in Shapley game theory



Local Accuracy

$$f(\mathbf{x}) = g(\mathbf{x}') = \phi_0 + \sum_{j=1}^{p} \phi_j x'_j$$



$$x'_j = \mathbf{0} \Longrightarrow \phi_j = \mathbf{0}$$

Intution: A missing feature gets an attribution of zero



Local Accuracy

$$f(\mathbf{x}) = g(\mathbf{x}') = \phi_0 + \sum_{j=1}^{p} \phi_j x'_j$$

Missingness

 $\mathbf{x}'_j = \mathbf{0} \Longrightarrow \phi_j = \mathbf{0}$

Consistency

$$\hat{f}_x\left(\mathbf{z}^{\prime(k)}
ight) = \hat{f}\left(h_x\left(\mathbf{z}^{\prime(k)}
ight)
ight)$$
 and $\mathbf{z}_{-j}^{\prime(k)}$ denote setting $z_j^{\prime(k)} = 0$. For any two models \hat{f} and \hat{f}' , if
 $\hat{f}'_x\left(\mathbf{z}^{\prime(k)}
ight) - \hat{f}'_x\left(\mathbf{z}_{-j}^{\prime(k)}
ight) \ge \hat{f}_x\left(\mathbf{z}^{\prime(k)}
ight) - \hat{f}_x\left(\mathbf{z}_{-j}^{\prime(k)}
ight)$

for all inputs $\mathbf{z}^{\prime(k)} \in \{0, 1\}^p$, then

$$\phi_j\left(\hat{f}',\mathbf{x}\right) \geq \phi_j(\hat{f},\mathbf{x})$$



Local Accuracy

$$f(\mathbf{x}) = g(\mathbf{x}') = \phi_0 + \sum_{j=1}^{p} \phi_j x'_j$$



 $x'_j = \mathbf{0} \Longrightarrow \phi_j = \mathbf{0}$

Consistency

$$\hat{f}'_{x}\left(\mathbf{z}^{\prime(k)}\right) - \hat{f}'_{x}\left(\mathbf{z}^{\prime(k)}_{-j}\right) \geq \hat{f}_{x}\left(\mathbf{z}^{\prime(k)}_{-j}\right) - \hat{f}_{x}\left(\mathbf{z}^{\prime(k)}_{-j}\right) \Longrightarrow \phi_{j}\left(\hat{f}', \mathbf{x}\right) \geq \phi_{j}(\hat{f}, \mathbf{x})$$

Intution: If a model changes so that the marginal contribution of a feature value increases or stays the same, the Shapley value also increases or stays the same

From consistency the Shapley axioms of additivity, dummy and symmetry follow

