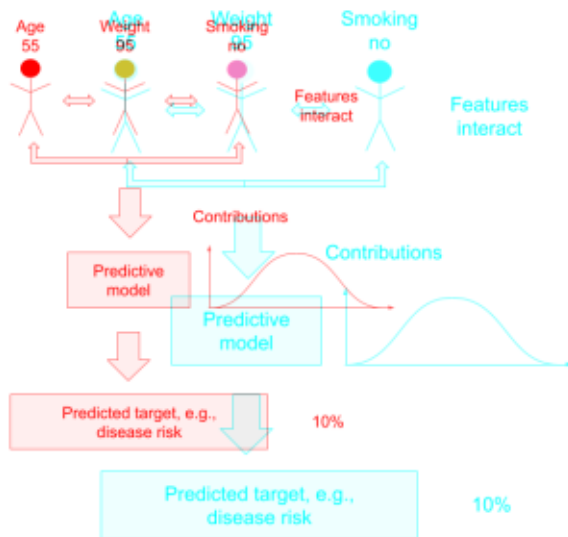




# FROM GAME THEORY TO MACHINE LEARNING



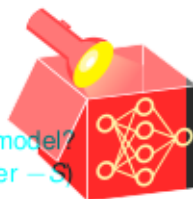
# FROM GAME THEORY TO MACHINE LEARNING

- Game: Make prediction of  $\hat{y}(x_1, x_2, \dots, x_p)$  for a single observation  $x$



# FROM GAME THEORY TO MACHINE LEARNING

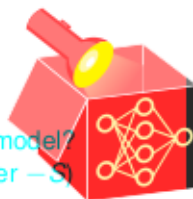
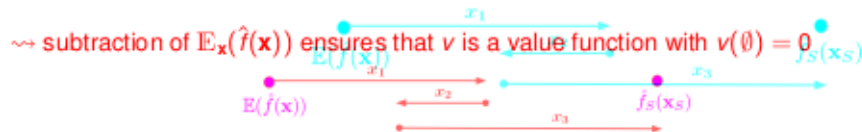
- Game: Make prediction  $\hat{f}(x_1, x_2, \dots, x_p)$  for a single observation  $\mathbf{x}$
- Players: Features  $x_j, j \in \{1, \dots, p\}$ , which cooperate to produce a prediction  
~> How can we make a prediction with a subset of features without changing the model?  
PD function:  $\hat{f}_S(\mathbf{x}_S) := \int_{X_{-S}} \hat{f}(\mathbf{x}_S, X_{-S}) d\mathbb{P}_{X_{-S}}$  ("removing" by marginalizing over  $-S$ )  
~> PD function:  $\hat{f}_S(\mathbf{x}_S) := \int_{X_{-S}} \hat{f}(\mathbf{x}_S, X_{-S}) d\mathbb{P}_{X_{-S}}$  ("removing" by marginalizing over  $-S$ )



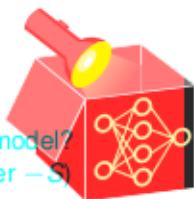
# FROM GAME THEORY TO MACHINE LEARNING

- Game: Make prediction of  $\hat{f}(x_1, x_2, \dots, x_p)$  for a single observation  $\mathbf{x}$
- Players: Features  $x_j, j \in \{1, \dots, p\}$ , which cooperate to produce a prediction  
 ~ How can we make a prediction with a subset of features without changing the model?  
 PD function:  $\hat{f}_S(\mathbf{x}_S) := \int_{\mathcal{X}_{-S}} \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}) d\mathbb{P}_{\mathcal{X}_{-S}}$  ("removing" by marginalizing over  $-S$ )  
 ~ PD function:  $\hat{f}_S(\mathbf{x}_S) := \int_{\mathcal{X}_{-S}} \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}) d\mathbb{P}_{\mathcal{X}_{-S}}$  ("removing" by marginalizing over  $-S$ )
- Value function / payout of coalition  $S \subseteq P$  for observation  $\mathbf{x}$ :  
 $v(S) = \hat{f}_S(\mathbf{x}_S) - \mathbb{E}_{\mathbf{x}}(\hat{f}(\mathbf{x}))$  where  $\hat{f}_S : \mathcal{X}_S \mapsto \mathcal{Y}$

~ subtraction of  $\mathbb{E}_{\mathbf{x}}(\hat{f}(\mathbf{x}))$  ensures that  $v$  is a value function with  $v(\emptyset) = 0$

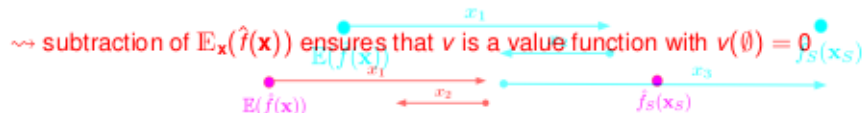


# FROM GAME THEORY TO MACHINE LEARNING



- Game: Make prediction of  $\hat{f}(x_1, x_2, \dots, x_p)$  for a single observation  $\mathbf{x}$
- Players: Features  $x_j, j \in \{1, \dots, p\}$ , which cooperate to produce a prediction
  - How can we make a prediction with a subset of features without changing the model?
- PD function:  $\hat{f}_S(\mathbf{x}_S) := \int_{\mathcal{X}_{-S}} \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}) d\mathbb{P}_{\mathbf{x}_{-S}}$  ("removing" by marginalizing over  $-S$ )
- Value function / payout of coalition  $S \subseteq P$  for observation  $\mathbf{x}$ :
  - Value function:  $\hat{f}_S(\mathbf{x}_S) := \int_{\mathcal{X}_{-S}} \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}) d\mathbb{P}_{\mathbf{x}_{-S}}$  ("removing" by marginalizing over  $-S$ )
- Value function / payout of coalition  $S \subseteq P$  for observation  $\mathbf{x}$ :
  - Value function:  $\hat{f}_S(\mathbf{x}_S) := \int_{\mathcal{X}_{-S}} \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}) d\mathbb{P}_{\mathbf{x}_{-S}}$  ("removing" by marginalizing over  $-S$ )

subtraction  $v(S) = \hat{f}_S(\mathbf{x}_S) - \mathbb{E}_{\mathbf{x}}(\hat{f}(\mathbf{x}))$ , where  $\hat{f}_S$  is a value function with  $v(\emptyset) = 0$



- Marginal contribution:  $v(S \cup \{j\}) - v(S) = \hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) - \hat{f}_S(\mathbf{x}_S)$
- Marginal contribution:  $v(S \cup \{j\}) - v(S) = \hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) - \hat{f}_S(\mathbf{x}_S)$ 
  - $\mathbb{E}_{\mathbf{x}}(\hat{f}(\mathbf{x}))$  cancels out due to the subtraction of value functions

# SHAPLEY VALUE - DEFINITION

Shapley (1953)3

Stumbelj et al. (2014)4



Shapley value  $\phi_j$  of feature  $j$  for observation  $\mathbf{x}$  via **order definition**:

$$\phi_j(\mathbf{x}) = \frac{1}{|P|!} \sum_{\tau \in \Pi} \underbrace{\left( \hat{f}_{S^\tau \cup \{j\}}(\mathbf{x}_{S^\tau \cup \{j\}}) - \hat{f}_{S^\tau}(\mathbf{x}_{S^\tau}) \right)}_{\text{marginal contribution of feature } j}$$

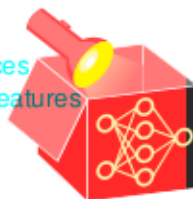
- Interpretation: Feature  $x_j$  contributed  $\phi_j$  to difference between  $\hat{f}(\mathbf{x})$  and average prediction
- Note: Marginal contributions and Shapley values can be negative
- Note: Marginal contributions and Shapley values can be negative
- For exact computation of  $\phi_j(\mathbf{x})$ , the PD function  $f_S(\mathbf{x}_S) = \frac{1}{n} \sum_{i=1}^n \hat{f}(\mathbf{x}_S, \mathbf{x}_{-S}^{(i)})$  for any set of features  $S$  can be used which yields

$$\phi_j(\mathbf{x}) = \frac{\phi_j(\mathbf{x})}{|P|! \cdot n} \sum_{\tau \in \Pi} \sum_{i=1}^n \left( \hat{f}^{(i)}(\mathbf{x}_{S^\tau \cup \{j\}}) - \hat{f}^{(i)}(\mathbf{x}_{S^\tau}, \mathbf{x}_{-S^\tau}^{(i)}) \right)$$

- Note:  $\hat{f}_S$  marginalizes over all other features  $-S$  using all observations  $i = 1, \dots, n$
- Note:  $f_S$  marginalizes over all other features  $-S$  using all observations  $i = 1, \dots, n$

## ESTIMATION: A PRACTICAL PROBLEM

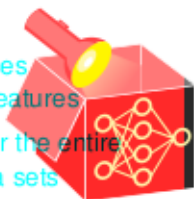
- Exact Shapley value computation is problematic for high-dimensional feature spaces
- For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features
- $\leadsto$  For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features





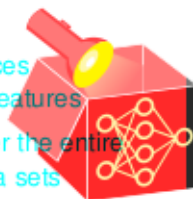
## ESTIMATION: A PRACTICAL PROBLEM

- Exact Shapley value computation is problematic for high-dimensional feature spaces  
For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets

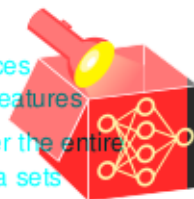


# ESTIMATION: A PRACTICAL PROBLEM

- Exact Shapley value computation is problematic for high-dimensional feature spaces
  - For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features
  - For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets
- Solution to both problems is sampling: Instead of averaging over  $|P|! \cdot n$  terms, we approximate it using a limited amount of  $M$  random samples of  $\tau$  to build coalitions  $S_j^\tau$



# ESTIMATION: A PRACTICAL PROBLEM



- Exact Shapley value computation is problematic for high-dimensional feature spaces  
For 10 features, there are already  $|P|! = 10! \approx 3.6$  million possible orders of features
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets
- Additional problem due to estimation of the marginal prediction  $\hat{f}_{S_j^\tau}$ : Averaging over the entire data set for each coalition  $S_j^\tau$  introduced by  $\tau$  can be very expensive for large data sets
- Solution to both problems is sampling: Instead of averaging over  $|P|! : n$  terms, we approximate it using a limited amount of  $M$  random samples of  $\tau$  to build coalitions  $S_j^\tau$
- $M$  is a tradeoff between accuracy of the Shapley value and computational costs
- Solution to both problems is sampling: Instead of averaging over  $|P|! : n$  terms, we approximate it using a limited amount of  $M$  random samples of  $\tau$  to build coalitions  $S_j^\tau$
- The higher  $M$ , the closer to the exact Shapley values, but the more costly the computation
- $M$  is a tradeoff between accuracy of the Shapley value and computational costs
- The higher  $M$ , the closer to the exact Shapley values, but the more costly the computation

# APPROXIMATION ALGORITHM → Strumbelj et al. (2014)

Estimation of  $\phi_j$  for observation  $x$  of model  $f$  fitted on data  $D$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

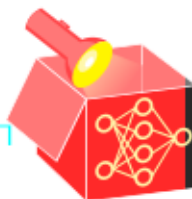


# APPROXIMATION ALGORITHM Strumbelj et al. (2014)

Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $\hat{f}$  fitted on data  $D$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(p)}) \in \Pi$

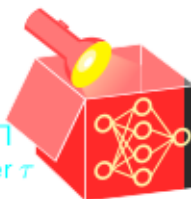


# APPROXIMATION ALGORITHM Strumbelj et al. (2014)

Estimation of  $\phi_j$  for observation  $x$  of model  $f$  fitted on data  $D$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

- 1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(p)}) \in \Pi$
- 2 Determine coalition  $S_m = S_j^\tau$ , i.e. the set of feat. before feat.  $j$  in order  $\tau$



# APPROXIMATION ALGORITHM

Strumbelj et al. (2014)

Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $\hat{f}$  fitted on data  $\mathcal{D}$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

- 1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(p)}) \in \Pi$
- 2 Determine coalition  $S_m = S_j^\tau$ , i.e., the set of feat. before feat.  $j$  in order  $\tau$
- 3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$



Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $\hat{f}$  fitted on data  $\mathcal{D}$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

- 1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(p)}) \in \Pi$
- 2 Determine coalition  $S_m = S_j^\tau$ , i.e., the set of feat. before feat.  $j$  in order  $\tau$
- 3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$
- 4 Construct two artificial obs. by replacing feature values from  $\mathbf{x}$  with  $\mathbf{z}^{(m)}$ :





Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $f$  fitted on data  $\mathcal{D}$  using sample size  $M$ :



1 For  $m=1, \dots, M$  do:

1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \dots, \tau^{(p)}) \in \Pi$

2 Determine coalition  $S_m = S_j^\tau$ , i.e., the set of feat. before feat.  $j$  in order  $\tau$

3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$

4 Construct two artificial obs. by replacing feature values from  $\mathbf{x}$  with  $\mathbf{z}^{(m)}$ :

$$\bullet \mathbf{x}_{+j}^{(m)} = \underbrace{(X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}}, X_j, Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}})}_{\substack{\mathbf{x}_{S_m \cup \{j\}} \\ \mathbf{x}_{S_m \cup \{j\}}}} \quad \underbrace{(Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}})}_{\substack{\mathbf{z}_{- \{S_m \cup \{j\}\}}^{(m)} \\ \mathbf{z}_{- \{S_m \cup \{j\}\}}^{(m)}}} \text{ takes features } S_m \cup \{j\} \text{ from } \mathbf{x}$$

Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $f$  fitted on data  $\mathcal{D}$  using sample size  $M$ :



1 For  $m=1, \dots, M$  do:

1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \dots, \tau^{(p)}) \in \Pi$

2 Determine coalition  $S_m = S_j^\tau$ , i.e., the set of feat. before feat.  $j$  in order  $\tau$

3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$

4 Construct two artificial obs. by replacing feature values from  $\mathbf{x}$  with  $\mathbf{z}^{(m)}$ :

$$\mathbf{x}_{+j}^{(m)} = \underbrace{(X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}})}_{\mathbf{x}_{S_m \cup \{j\}}} \underbrace{, X_j, Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}}}_{\mathbf{z}_{-S_m \cup \{j\}}^{(m)}} \text{ takes features } S_m \cup \{j\} \text{ from } \mathbf{x}$$

$$\mathbf{x}_{-j}^{(m)} = \underbrace{(X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}})}_{\mathbf{x}_{S_m}} \underbrace{, Z_j, Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}}}_{\mathbf{z}_{-S_m}^{(m)}} \text{ takes features } S_m \text{ from } \mathbf{x}$$

$S_m$  from  $\mathbf{x}$



Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $\hat{f}$  fitted on data  $\mathcal{D}$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \dots, \tau^{(p)}) \in \Pi$

2 Determine coalition  $S_m = S_j^\tau$ , i.e. the set of feat. before feat.  $j$  in order  $\tau$

3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$

4 Construct two artificial obs. by replacing feature values from  $\mathbf{x}$  with  $\mathbf{z}^{(m)}$ :

$$\mathbf{x}_{+j}^{(m)} = \underbrace{(X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}}, X_j, \dots, X_{\tau^{(p)}})}_{\mathbf{x}_{S_m \cup \{j\}}} \underbrace{(Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}})}_{\mathbf{z}_{-S_m \cup \{j\}}^{(m)}} \text{ takes features } S_m \cup \{j\} \text{ from } \mathbf{x}$$

$$\mathbf{x}_{-j}^{(m)} = \underbrace{(X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}})}_{\mathbf{x}_{S_m}} \underbrace{(Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}})}_{\mathbf{z}_{-S_m}^{(m)}} \text{ takes features } S_m \text{ from } \mathbf{x}$$

5 Compute difference  $\phi_j^m = \hat{f}(\mathbf{x}_{+j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)})$

6 Compute difference  $\phi_j^m = \hat{f}(\mathbf{x}_{+j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)})$   
 $\leadsto \hat{f}_{S_m}(\mathbf{x}_{S_m})$  is approximated by  $\hat{f}(\mathbf{x}_{-j}^{(m)})$  and  $\hat{f}_{S_m \cup \{j\}}(\mathbf{x}_{S_m \cup \{j\}})$  by  $\hat{f}(\mathbf{x}_{+j}^{(m)})$  over  $M$  iters  
 $\leadsto \hat{f}_{S_m}(\mathbf{x}_{S_m})$  is approximated by  $\hat{f}(\mathbf{x}_{-j}^{(m)})$  and  $\hat{f}_{S_m \cup \{j\}}(\mathbf{x}_{S_m \cup \{j\}})$  by  $\hat{f}(\mathbf{x}_{+j}^{(m)})$  over  $M$  iters.



Estimation of  $\phi_j$  for observation  $\mathbf{x}$  of model  $\hat{f}$  fitted on data  $\mathcal{D}$  using sample size  $M$ :

1 For  $m=1, \dots, M$  do:

1 Select random order / perm. of feature indices  $\tau = (\tau^{(1)}, \dots, \tau^{(p)}) \in \Pi$

2 Determine coalition  $S_m = S_j^\tau$ , i.e. the set of feat. before feat.  $j$  in order  $\tau$

3 Select random data point  $\mathbf{z}^{(m)} \in \mathcal{D}$

4 Construct two artificial obs. by replacing feature values from  $\mathbf{x}$  with  $\mathbf{z}^{(m)}$ :

•  $\mathbf{x}_{+j}^{(m)} = (\underbrace{X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}}}_{\mathbf{x}_{S_m \cup \{j\}}}, \underbrace{X_j, Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}}}_{\mathbf{z}_{-S_m \cup \{j\}}^{(m)}})$  takes features  $S_m \cup \{j\}$  from  $\mathbf{x}$

•  $\mathbf{x}_{-j}^{(m)} = (\underbrace{X_{\tau^{(1)}}, \dots, X_{\tau^{(|S_m|-1)}}}_{\mathbf{x}_{S_m}}, \underbrace{Z_j, Z_{\tau^{(|S_m|+1)}}, \dots, Z_{\tau^{(p)}}}_{\mathbf{z}_{-S_m}^{(m)}})$  takes features  $S_m$  from  $\mathbf{x}$

5 Compute difference  $\phi_j^m = \hat{f}(\mathbf{x}_{+j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)})$

6 Compute difference  $\phi_j^m = \hat{f}(\mathbf{x}_{+j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)})$   
 $\leadsto \hat{f}_{S_m}(\mathbf{x}_{S_m})$  is approximated by  $\hat{f}(\mathbf{x}_{-j}^{(m)})$  and  $\hat{f}_{S_m \cup \{j\}}(\mathbf{x}_{S_m \cup \{j\}})$  by  $\hat{f}(\mathbf{x}_{+j}^{(m)})$  over  $M$  iters  
 $\leadsto \hat{f}_{S_m}(\mathbf{x}_{S_m})$  is approximated by  $\hat{f}(\mathbf{x}_{-j}^{(m)})$  and  $\hat{f}_{S_m \cup \{j\}}(\mathbf{x}_{S_m \cup \{j\}})$  by  $\hat{f}(\mathbf{x}_{+j}^{(m)})$  over  $M$  iters.

2 Compute Shapley value  $\phi_j = \frac{1}{M} \sum_{m=1}^M \phi_j^m$

# SHAPLEY VALUE APPROXIMATION - ILLUSTRATION

## Definition

$\mathbf{x}$ : obs. of interest

obs. of interest

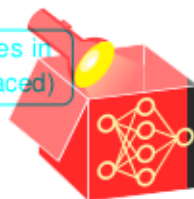
$\mathbf{x}$  with feature values in  $S_m$  (other are replaced)

$\mathbf{x}$  with feature values in  $S_m$  (other are replaced)

$$\phi_j(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \left[ \hat{f}(\mathbf{x}_{-j}^{(m)}) - \hat{f}(\mathbf{x}_{+j}^{(m)}) \right]$$

$\mathbf{x}$  with feature values in  $S_m \cup \{j\}$

$\mathbf{x}$  with feature values in  $S_m \cup \{j\}$



|                   | Temperature | Humidity | Windspeed  | Year |
|-------------------|-------------|----------|------------|------|
| $\mathbf{x}$      | 10.66       | 56       | 11         | 2012 |
| $\mathbf{x}_{+j}$ | 10.66       | 56       | random : 2 | 2012 |
| $\mathbf{x}_{-j}$ | 10.66       | 56       | random : 2 | 2012 |

$j$ 
 $j$

# SHAPLEY VALUE APPROXIMATION - ILLUSTRATION

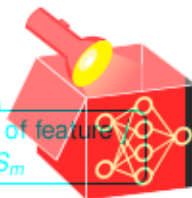
## Definition

Contribution of feature  $j$   
to coalition  $S_m$

$$\phi_j(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \underbrace{\left[ \hat{f}(\mathbf{x}_{-j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)}) \right]}_{:= \Delta(j, S_m)} - \hat{f}(\mathbf{x}_{-j}^{(m)})$$

$:= \Delta(j, S_m) \quad := \Delta(j, S_m)$

Contribution of feature  $j$   
to coalition  $S_m$



- $\Delta(j, S_m) = \hat{f}(\mathbf{x}_{-j}^{(m)}) - \hat{f}(\mathbf{x}_{-j}^{(m)})$  is the marginal contribution of feature  $j$  to coalition  $S_m$
- Here: Feature **year** contributes +700 bike rentals if it joins coalition  $S_m = \{\text{temp, hum}\}$
- Here: Feature **year** contributes +700 bike rentals if it joins coalition

$S_m = \{\text{temp, hum}\}$

|                  | Temperature | Humidity | Windspeed     | Year | Count |
|------------------|-------------|----------|---------------|------|-------|
| $\mathbf{x}$     | 10.66       | 56       | 11            | 2012 |       |
| $\mathbf{x} + j$ | 10.66       | 56       | random : 2012 | 2012 | 5600  |
| $\mathbf{x} - j$ | 10.66       | 56       | random : 2012 | 2012 | 4900  |
| $\mathbf{x} + j$ | 10.66       | 56       | random : 2012 | 2012 | 5600  |
| $\mathbf{x} - j$ | 10.66       | 56       | random : 2012 | 2012 | 4900  |

$\Delta(j, S_m)$   
marginal contribution  
700  
 $\Delta(j, S_m)$   
marginal contribution  
700

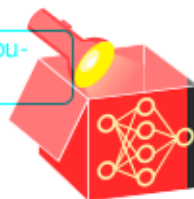
# SHAPLEY VALUE APPROXIMATION - ILLUSTRATION

## Definition

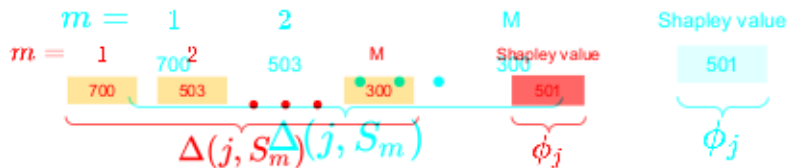
$$\phi_j(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \left[ \hat{f}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left[ \hat{f}(\mathbf{x}_{-j}^{(m)}) \right] - \hat{f}(\mathbf{x}_{-j}^{(m)}) \right]$$

average the contributions of feature  $j$

average the contributions of feature  $j$



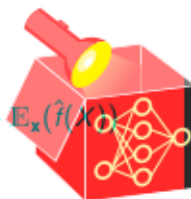
- Compute marginal contribution of feature  $j$  towards the prediction across all randomly drawn feature coalitions  $S_1, \dots, S_m$
- Average all  $M$  marginal contributions of feature  $j$
- Average all  $M$  marginal contributions of feature  $j$
- Shapley value  $\phi_j$  is the payout of feature  $j$ , i.e., how much feature  $j$  contributed to the overall prediction in bicycle counts of a specific observation  $\mathbf{x}$
- Shapley value  $\phi_j$  is the payout of feature  $j$ , i.e., how much feature  $j$  contributed to the overall prediction in bicycle counts of a specific observation  $\mathbf{x}$



# REVISITED: AXIOMS FOR FAIR ATTRIBUTIONS

We take the general axioms for Shapley Values and apply it to predictionists:

- **Efficiency:** Shapley values add up to the (centered) prediction:  $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$

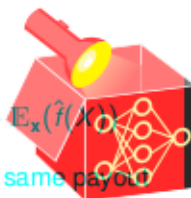




# REVISITED: AXIOMS FOR FAIR ATTRIBUTIONS

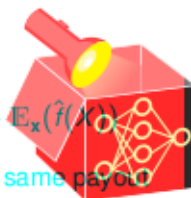
We take the general axioms for Shapley Values and apply it to predictionists:

- **Efficiency:** Shapley values add up to the (centered) prediction:  $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$   
 $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$
- **Symmetry:** Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout  
Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout =  $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$   
 $\leadsto$  interaction effects between features are fairly divided  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$



# REVISITED: AXIOMS FOR FAIR ATTRIBUTIONS

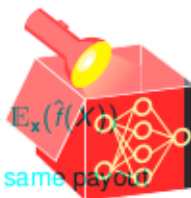
We take the general axioms for Shapley Values and apply it to predictionists:



- Efficiency:** Shapley values add up to the (centered) prediction:  $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$   
 $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$
- Symmetry:** Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout  
 interaction effects between features are fairly divided
- Symmetry:** Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$   
 $\leadsto$  interaction effects between features are fairly divided
- Dummy / Null Player:** Shapley value of a feature that does not influence the prediction is zero  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$   
 $\leadsto$  if a feature was not selected by the model (e.g., tree or LASSO), its Shapley value is zero
- Dummy / Null Player:** Shapley value of a feature that does not influence the prediction is zero  $\leadsto$  if a feature was not selected by the model (e.g., tree or LASSO), its Shapley value is zero  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_S(\mathbf{x}_S)$  for all  $S \subseteq P$  then  $\phi_j = 0$

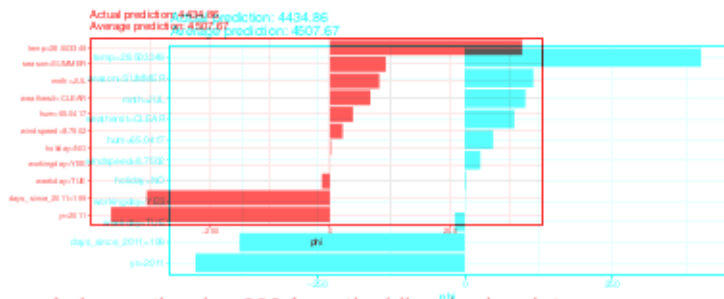
# REVISITED: AXIOMS FOR FAIR ATTRIBUTIONS

We take the general axioms for Shapley Values and apply it to predictionists:



- Efficiency:** Shapley values add up to the (centered) prediction:  $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$   
 $\sum_{j=1}^p \phi_j = \hat{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}(\hat{f}(X))$
- Symmetry:** Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout  
 interaction effects between features are fairly divided
- Symmetry:** Two features  $j$  and  $k$  that contribute the same to the prediction get the same payout  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$   
 $\leadsto$  interaction effects between features are fairly divided
- Dummy / Null Player:** Shapley value of a feature that does not influence the prediction is zero  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_{S \cup \{k\}}(\mathbf{x}_{S \cup \{k\}})$  for all  $S \subseteq P \setminus \{j, k\}$  then  $\phi_j = \phi_k$   
 $\leadsto$  if a feature was not selected by the model (e.g., tree or LASSO), its Shapley value is zero
- Dummy / Null Player:** Shapley value of a feature that does not influence the prediction is zero  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_S(\mathbf{x}_S)$  for all  $S \subseteq P$  then  $\phi_j = 0$   
 $\leadsto$  if a feature was not selected by the model (e.g., tree or LASSO), its Shapley value is zero
- Additivity:** For a prediction with combined payouts, the payout is the sum of payouts:  
 $\hat{f}_{S \cup \{j\}}(\mathbf{x}_{S \cup \{j\}}) = \hat{f}_S(\mathbf{x}_S) + \phi_j$  for all  $S \subseteq P$  then  $\phi_j = 0$   
 $\leadsto$  Shapley values for model ensembles can be combined
- Additivity:** For a prediction with combined payouts, the payout is the sum of payouts:  $\phi_j(v_1) + \phi_j(v_2) \leadsto$  Shapley values for model ensembles can be combined

# BIKE SHARING DATASET



- Shapley values of observation  $i = 200$  from the bike sharing data
- Difference between model prediction of this observation and the average prediction of the data is fairly distributed among the features (i.e.,  $4434 - 4507 \approx -73$ )
- Difference between model prediction of this observation and the average prediction of the data is fairly distributed among the features (i.e.,  $4434 - 4507 \approx -73$ )
- Feature value  $\text{temp} = 28.5$  has the most positive effect, with a contribution (increase of prediction) of about +400

# ADVANTAGES AND DISADVANTAGES

## Advantages:

- **Solid theoretical foundation** in game theory
- Prediction is **fairly distributed** among the feature values → easy to interpret for a user
  - **Contrastive explanations** that compare the prediction with the average prediction
- **Contrastive explanations** that compare the prediction with the average prediction

## Disadvantages:

- Without sampling, Shapley values need a lot of computing time to inspect all possible coalitions

## Disadvantages:

- Like many other IML methods, Shapley values suffer from the inclusion of unrealistic data observations when features are correlated
- Without sampling, Shapley values need a lot of computing time to inspect all possible coalitions
- Like many other IML methods, Shapley values suffer from the inclusion of unrealistic data observations when features are correlated

