

S = a subset of the M features

$v(S) = E[f(x)|x_S = x_S^*]$

$w(S)$ = weight function



groupSHAP

Efficient Shapley value explanation through feature groups

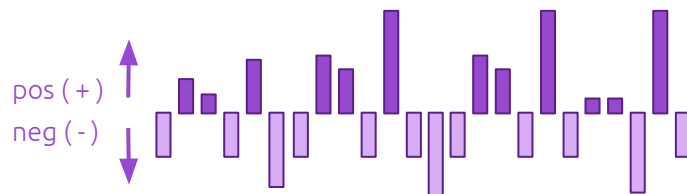
SHAP formula for feature j

$$\phi_j = \sum_{\substack{\text{all } S \text{ without } j \\ 2^M \text{ terms}}} w(S) \underbrace{(v(S \cup \{j\}) - v(S))}_{\text{marginal contribution}}$$

computationally intractable for many features



Shapley values contribution ϕ_j per feature



loong list for many features hard to interpret

Divide the M features into G feature groups

SHAP formula for feature group g

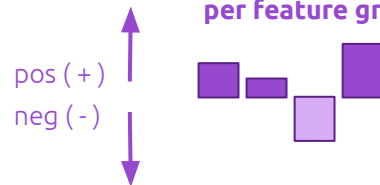
$$\psi_g = \sum_{\text{all } T \text{ without } g} w(T) (v(T \cup \{g\}) - v(T))$$

$2^G \text{ terms} \ll 2^M$



computationally tractable unless G too large

Shapley value contribution ψ_g per feature group



short list easy to interpret





groupSHAP

**Efficient Shapley value
explanation through
feature groups**

SHAP (SHapley Additive exPlanations)

The practical issue

Notation

S = a subset of the M features

$v(S) = E[f(x)|x_S = x_S^*]$

$w(S)$ = weight function

- Problem: Want to explain predictions from a model $f(x)$ with M features
- Shapley value formula for feature j

$$\phi_j = \sum_{\substack{\text{all } S \text{ without } j \\ 2^M \text{ terms}}} w(S) \underbrace{(v(S \cup \{j\}) - v(S))}_{\text{marginal contribution}}$$

- Computational complexity

$$M = 5 \Rightarrow 2^M = 32$$

$$M = 10 \Rightarrow 2^M = 1024$$

$$M = 20 \Rightarrow 2^M = 1048676$$

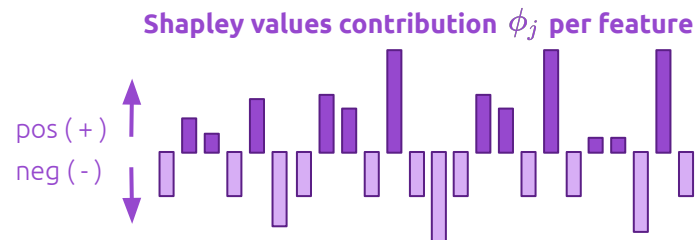
$$M = 40 \Rightarrow 2^M > 10^{12}$$

$$M = 100 \Rightarrow 2^M > 10^{30}$$

$$M = 1000 \Rightarrow 2^M > 10^{301}$$

Generally nonsatisfactory approximation methods exists

- KernelSHAP*
 - Inaccurate for large M
- TreeSHAP/TreeExplainer*
 - Limited to tree-based models
- DASP (Deep Approximate Shapley Propagation)
 - Limited to neural Networks



- In any case: Potentially too long list of contributions from dependent features

*The popular python library *shap* uses KernelSHAP/TreeSHAP

groupSHAP

The idea

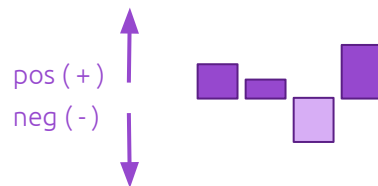
- A fundamentally simple solution:
 - Divide the M features into G feature groups
 - Replace the feature subsets S in the Shapley formula by feature group subsets T

- Shapley formula for feature group g

$$\psi_g = \sum_{\text{all } T \text{ without } g} w(T)(v(T \cup \{g\}) - v(T))$$

- 2^G terms $\ll 2^M \Rightarrow$ Computationally tractable
- Still a Shapley value, so all properties remains
- Potential grouping criteria
 - Type of feature
 - Origin/source of feature
 - Feature dependence (high-dependence features in same group)

Shapley value contribution ψ_j per feature group



Example groups, car insurance:

- Car info
- main driver info
- other driver info
- previous incident info

groupSHAP

Theoretical result

Is group-wise Shapley values (using groupSHAP) ever the same as summing feature-wise Shapley values?

I.e. do we ever have? $\psi_g = \sum_{j \in g} \phi_j$

YES!

Any partially additively separable function with between-group feature independence

$$(f(\mathbf{x}) = \sum_{g=1}^G f_g(\mathbf{x}_g)) \\ (\mathbf{x}_g \perp\!\!\!\perp \mathbf{x}_{g'})$$

has $\psi_g = \sum_{j \in g} \phi_j$

groupSHAP

Practical use (through R-package *shapr*)

groupSHAP branch
on GitHub
(under development)

Code example

```
1 remotes::install_github("NorskRegnesentral/shapr",ref = "groupSHAP")
2
3 # Loading the Boston housing data set
4 data("Boston", package = "MASS")
5 x_var <- c("lstat", "rm", "dis", "indus", "nox", "tax")
6 y_var <- "medv"
7
8 x_train <- as.matrix(Boston[1:50, x_var])
9 y_train <- Boston[1:500, y_var]
10 x_test <- as.matrix(Boston[501:504, x_var])
11
12 # Fitting a basic xgboost model
13 model <- xgboost::xgboost(
14   data = x_train,
15   label = y_train,
16   nround = 20,
17   verbose = FALSE
18 )
19
20 # Define the feature groups
21 group <- list(c("lstat", "rm", "dis"),
22             c("indus", "nox"),
23             "tax")
24
25 # Prepare the data for explanation
26 explainer <- shapr::shapr(x_train, model, group = group)
27
28 # Run the explainer
29 explanation <- shapr::explain(x_test,
30                             explainer,
31                             approach = "empirical",
32                             prediction_zero = mean(y_train))
33
34 # Plot the group-wise shapley values
35 plot(explanation, plot_phi0 = F)
```



Resulting explanation figure

