

Prediction explanation with Shapley values

Martin Jullum (jullum@nr.no)



Imperial College London, 03.02.2022



Prediction explanation

- ▶ Assume a model $f(\mathbf{x}) \in \mathbb{R}$ that predicts some unknown outcome based on a set of features $\mathbf{x} = (x_1, \dots, x_M)$
- ▶ We apply the predictive model for a specific input $\mathbf{x} = \mathbf{x}^*$, reaching a certain prediction $f(\mathbf{x}^*)$
- ▶ Individual prediction explanation
 - Want to understand how the different **features**, or **types of features** affect this specific prediction value $f(\mathbf{x}^*)$
 - I.e. **explain the predicted outcome** in terms of the input $\mathbf{x} = \mathbf{x}^*$ (**local explanation**)
- ▶ Frameworks...
 - LIME
 - Anchors
 - Counterfactual explanations
 - Explanation Vectors
 - PredDiff
 - **Shapley values**

Prediction explanation – by example

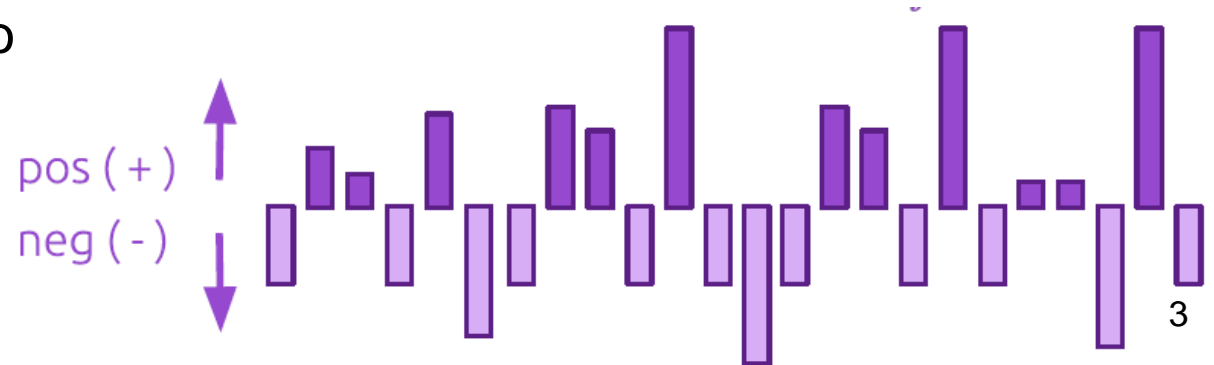
► Car insurance

- Response y : Insured crashed or not
- Features $\mathbf{x} = (x_1, \dots, x_M)$: Data about the insured, his/her car and crashing history
- Predictive model f : Model trained to predict probability of crash: $f(\mathbf{x}) \approx \Pr(y = \text{yes}|\mathbf{x})$



► Prediction explanation

- Why did a guy with features \mathbf{x}^* get a predicted probability of crashing equal to $f(\mathbf{x}^*) = 0.3$?



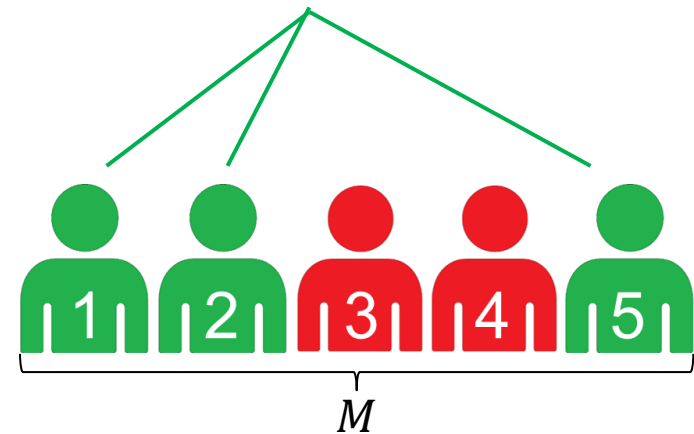
Shapley values

- ▶ Concept from (cooperative) game theory in the 1950s
- ▶ Used to distribute the total payoff to the players
- ▶ Explicit formula for the “fair” payment to every player j :

$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

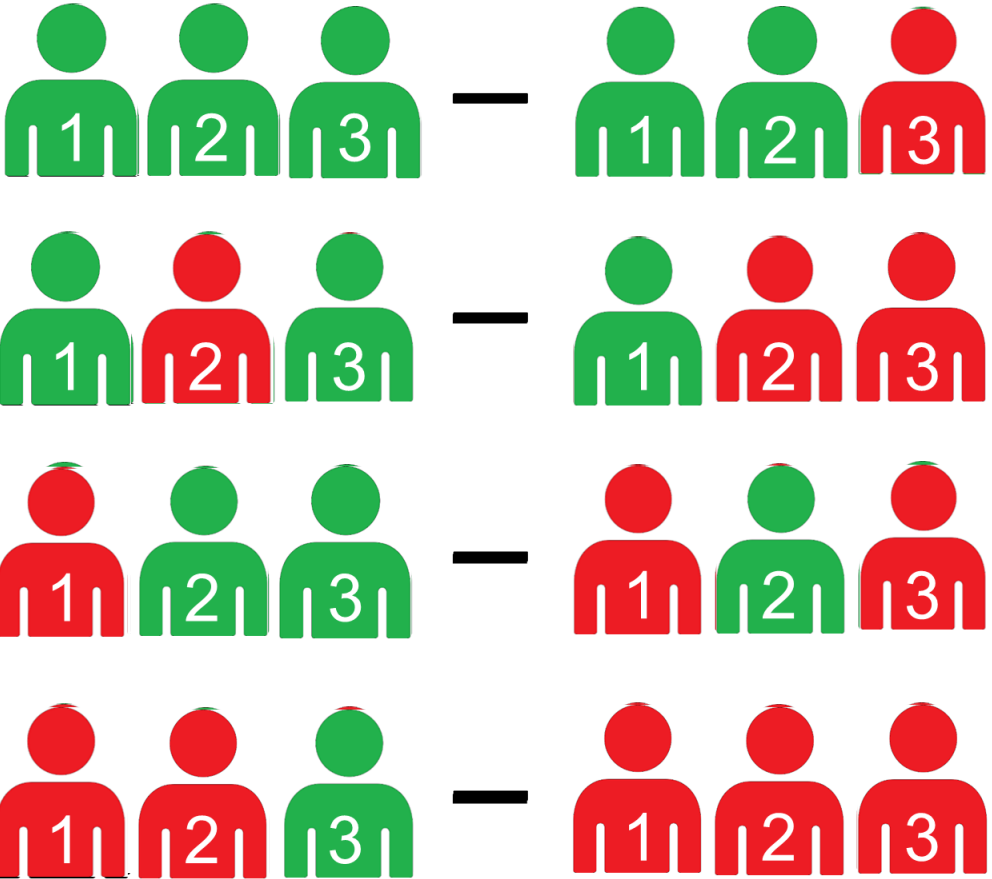
$v(S)$ is the payoff with only players in subset S

- ▶ Several mathematical optimality properties



Intuition behind the Shapley formula

Game with 3 players



Shapley values for taxi sharing

Costs: \$3/mi

$$v(\{R, B, G\}) = (4 + 6 + 2)mi * \$3 = \$36$$

$$v(\{\}) = \$0$$

$$v(\{R\}) = 4mi * \$3 = \$12$$

$$v(\{B\}) = (5 + 2)mi * \$3 = \$21$$

$$v(\{G\}) = 5mi * \$3 = \$15$$

$$v(\{R, B\}) = (4 + 6)mi * \$3 = \$30$$

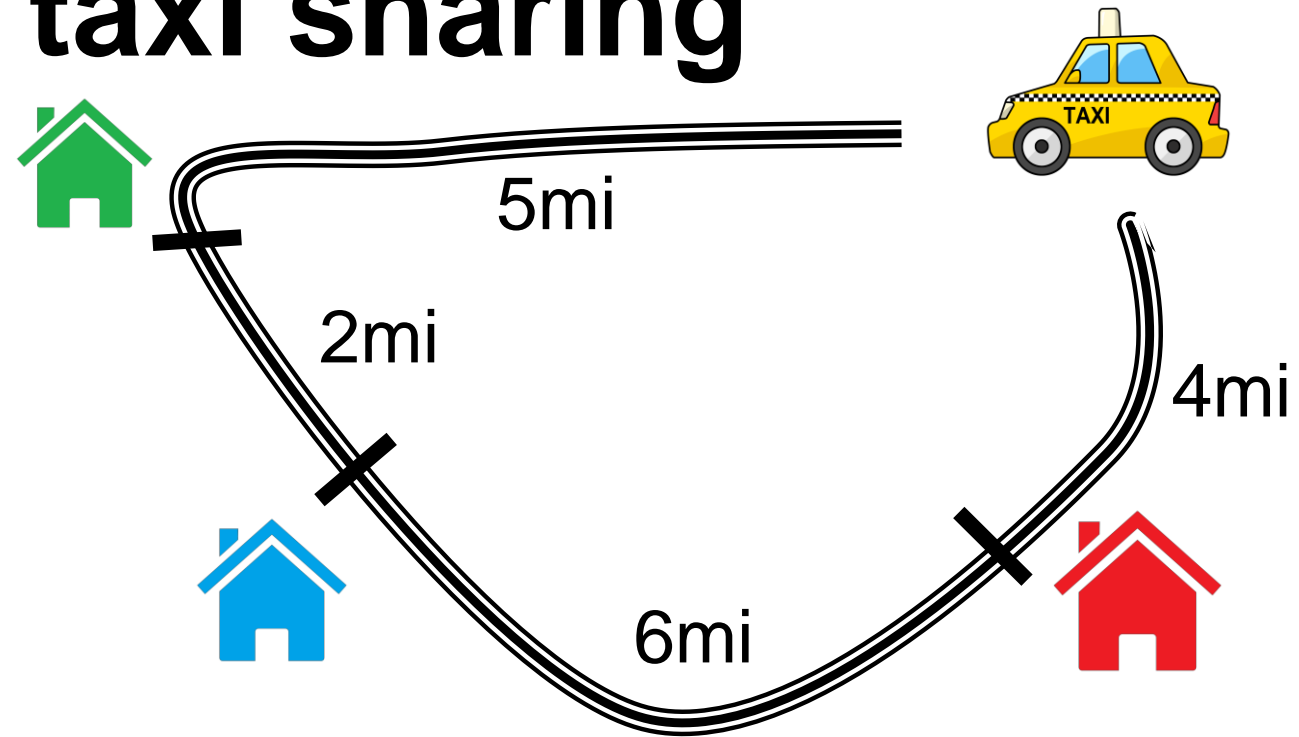
$$v(\{R, G\}) = (4 + 6 + 2)mi * \$3 = \$36$$

$$v(\{B, G\}) = (5 + 2)mi * \$3 = \$21$$

$$\phi_R = \frac{1}{3}(v(\{R, B, G\}) - v(\{B, G\})) + \frac{1}{6}(v(\{R, B\}) - v(\{B\})) + \frac{1}{6}(v(\{R, G\}) - v(\{G\})) + \frac{1}{3}(v(\{R\}) - v(\{\})) = \$14$$

$$\phi_B = \frac{1}{3}(v(\{R, B, G\}) - v(\{R, G\})) + \frac{1}{6}(v(\{R, B\}) - v(\{R\})) + \frac{1}{6}(v(\{B, G\}) - v(\{G\})) + \frac{1}{3}(v(\{B\}) - v(\{\})) = \$11$$

$$\phi_G = \frac{1}{3}(v(\{R, B, G\}) - v(\{R, B\})) + \frac{1}{6}(v(\{R, G\}) - v(\{R\})) + \frac{1}{6}(v(\{B, G\}) - v(\{B\})) + \frac{1}{3}(v(\{G\}) - v(\{\})) = \$11$$



Shapley values for prediction explanation

► Approach popularised by Lundberg & Lee (2017)

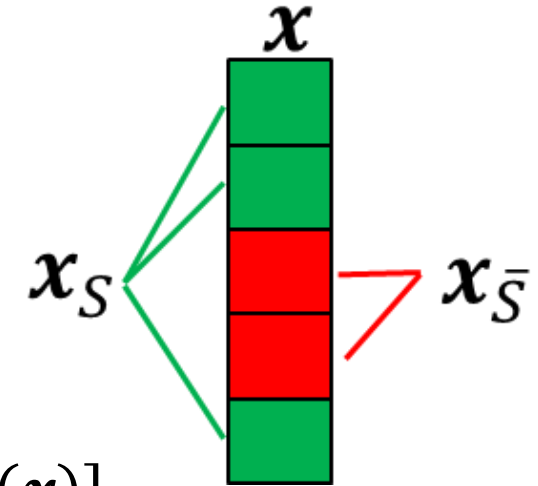
- Players = features (x_1, \dots, x_M)
- Payoff = prediction ($f(\mathbf{x}^*)$)
- Contribution function: $v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$
- Properties

$$\phi_0 + \sum_{j=1}^M \phi_j = f(\mathbf{x}^*)$$

$$\phi_0 = E[f(\mathbf{x})]$$

$$f(\mathbf{x}) \perp\!\!\!\perp x_j \\ \text{implies } \phi_j = 0$$

$$x_i, x_j \text{ same contribution} \\ \text{implies } \phi_i = \phi_j$$



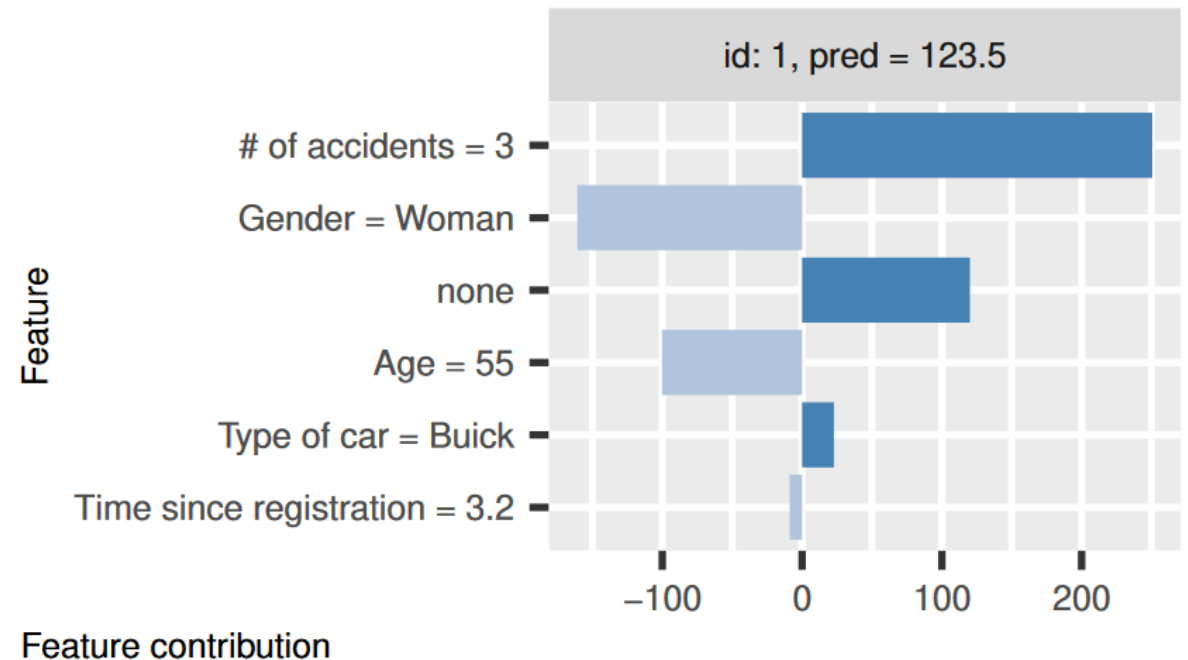
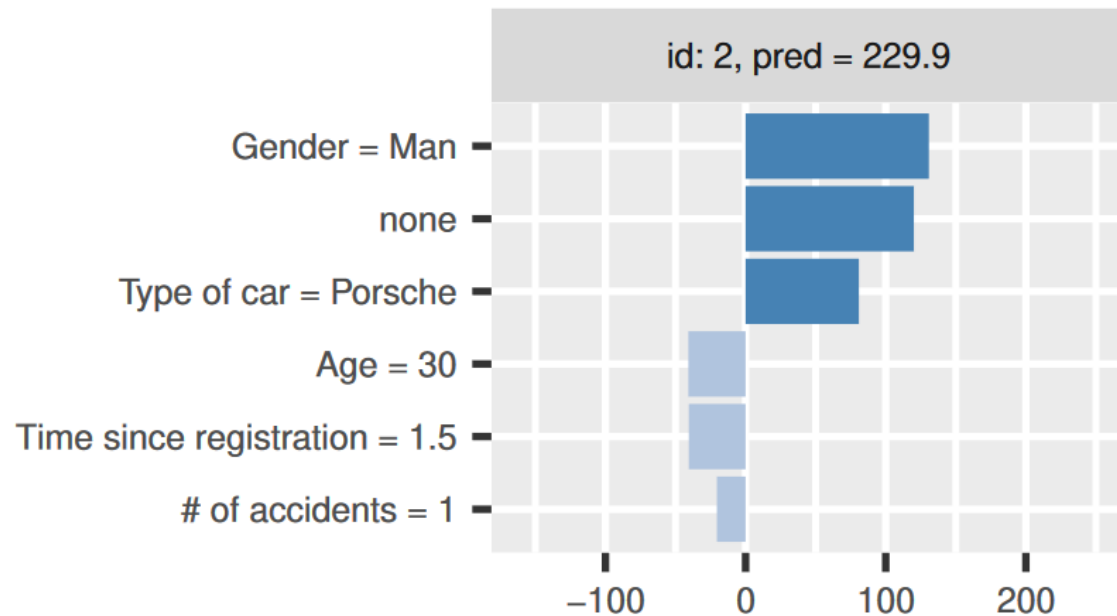
- Interpretation of ϕ_j : **The prediction change caused by observing the value of x_j – averaged over whether the other features were observed or not**

Example of Shapley value explanation

- ▶ Consider a model $f(x)$ trained to predict a fair price of a car insurance based on the following features x :
 - Owner's age, owner's gender, type of car, time since the car was registered, number of accidents the last 5 years



Shapley value prediction explanation



Two main challenges

1. The computational complexity in the Shapley formula is of size 2^M

$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

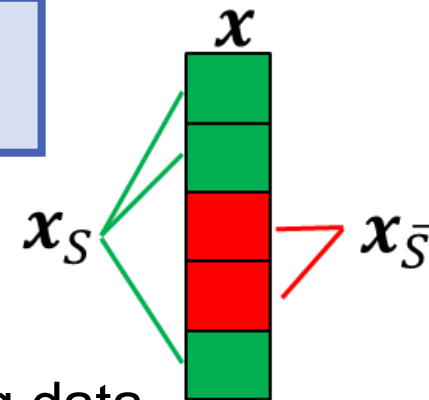
- Approximate solutions may be obtained by using a finite sample of subsets S (KernelSHAP; Lundberg & Lee, 2017)

2. Estimating the contribution function

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

- Lundberg & Lee (2017)

- Approximates $v(S) \approx \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}}) d\mathbf{x}_{\bar{S}}$,
- Estimates $p(\mathbf{x}_{\bar{S}})$ using the empirical distribution of the training data
- Monte Carlo integration to solve the integral



This assumes the features are independent!

Two main challenges

1. The computational complexity in the Shapley formula is of size 2^M

$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

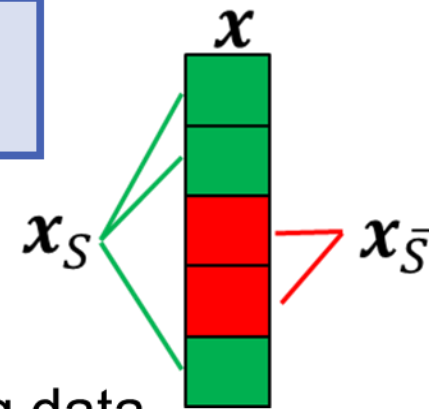
- Approximate solutions may be obtained by using a finite sample of subsets S (KernelSHAP; Lundberg & Lee, 2017)

2. Estimating the contribution function

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

- Lundberg & Lee (2017)

- Approximates $v(S) \approx \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}}) d\mathbf{x}_{\bar{S}}$,
- Estimates $p(\mathbf{x}_{\bar{S}})$ using the empirical distribution of the training data
- Monte Carlo integration to solve the integral



This assumes the features are independent!

Consequences of the independence assumption

- ▶ Requires evaluating $f(x_{\bar{S}}, x_S)$ at potentially unlikely or illegal combinations of $x_{\bar{S}}$ and x_S

- ▶ Example 1

- Number of transactions to Switzerland: 0
- Average transaction amount to Switzerland: 100 €



- ▶ Example 2

- Age: 17
- Marital status: Widow
- Profession: Professor



Estimating $v(S)$ properly



Artificial Intelligence
Volume 298, September 2021, 103502

Explaining individual predictions when features are dependent: More accurate approximations to Shapley values[☆]

Kjersti Aas*, Martin Jullum, Anders Løland

Depend. Model. 2021; 9:62–81

DE GRUYTER

Kjersti Aas*, Thomas Nagler, Martin Jullum, and Anders Løland

Explaining predictive models using Shapley values and non-parametric vine copulas

CD-MAKE 2020: Machine Learning and Knowledge Extraction pp 117-137

Explaining Predictive Models with Mixed Features Using Shapley Values and Conditional Inference Trees

Authors

Annabelle Redelmeier✉, Martin Jullum✉, Kjersti Aas

arXiv:2111.13507v1 [stat.ML] 26 Nov 2021

Using Shapley Values and Variational Autoencoders to Explain Predictive Models with Dependent Mixed Features

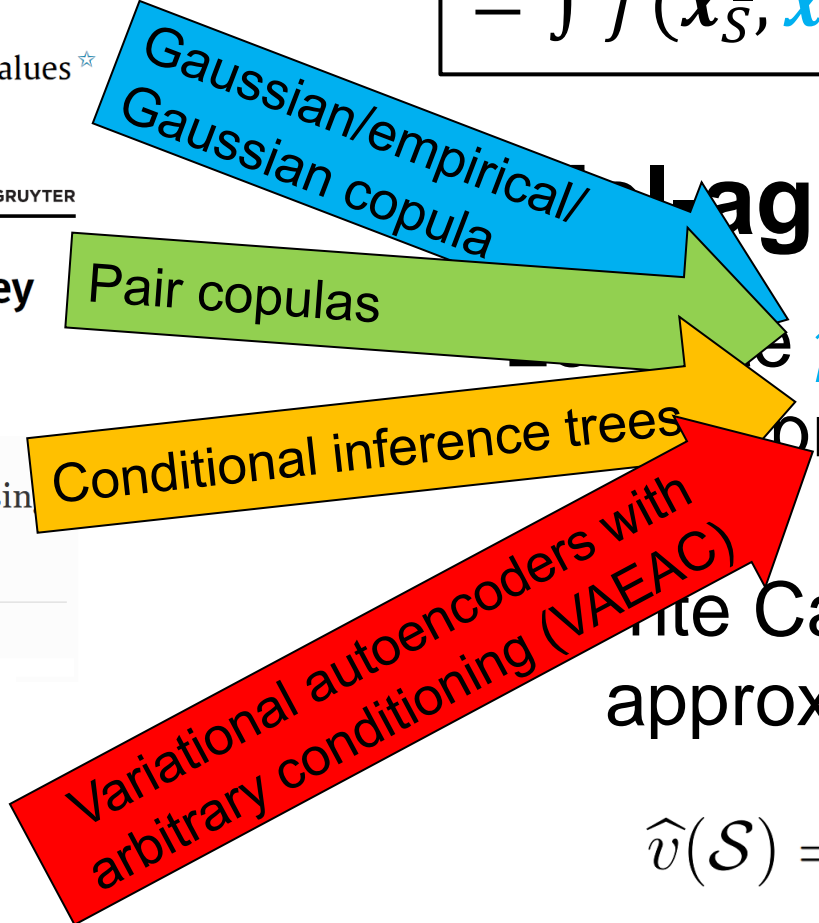
Lars Henry Berge Olsen*¹, Ingrid Kristine Glad¹,
Martin Jullum², and Kjersti Aas²

¹Department of Mathematics, University of Oslo

²Norwegian Computing Center

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$$

$$= \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$



diagnostic idea:

estimate $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ in a proper way

+ Monte Carlo integration to approximate integral:

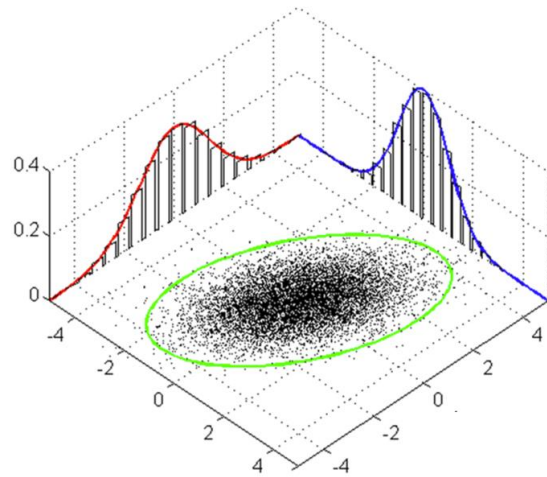
$$\hat{v}(S) = \frac{1}{K} \sum_{k=1}^K f(\mathbf{x}_{\bar{S}}^k, \mathbf{x}_S^*)$$

*Frye et. al (2020) briefly outline a VAEAC approach

Approaches to estimate and sample from $p(x_{\bar{S}} | x_S = x_S^*)$

1. Continuous features: Assume $p(x)$ is Gaussian $N(\mu, \Sigma)$

Estimate μ, Σ



Gaussian distribution

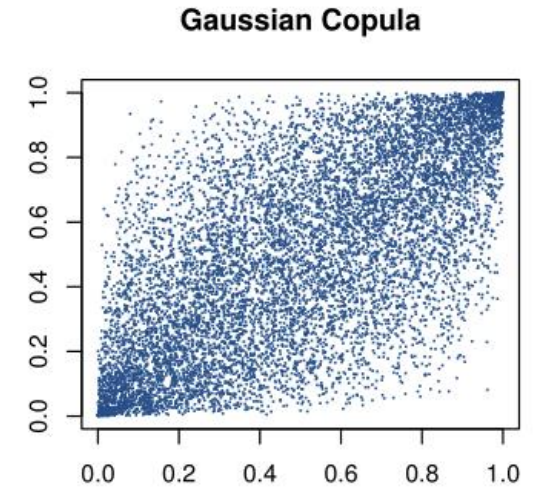
Compute conditional means and covariances to obtain analytical expression for

$$p(x_{\bar{S}} | x_S = x_S^*)$$

Sample from that Gaussian distribution

Approaches to estimate and sample from $p(x_{\bar{S}} | x_S = x_S^*)$

2. Continuous features: Assume $p(x)$ is a Gaussian copula



Transform each x_i to $u_i \sim U[0,1]$
with inverse empirical CDF

Transform each u_i to $v_i \sim N(0,1)$

Estimate the correlation Σ^v of (v_1, \dots, v_M)

Obtain analytical expression for
the Gaussian distribution

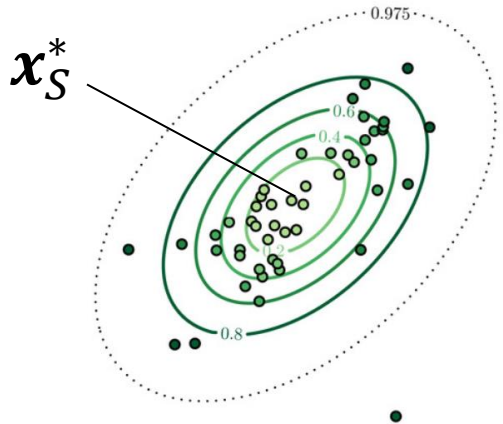
$$p_v(v_{\bar{S}} | v_S = v_S^*)$$

Sample from $p_v(v_{\bar{S}} | v_S = v_S^*)$
+ transform back to original scale

Approaches to estimate and sample from $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$

3. Continuous features: Use an empirical (conditional) distribution which weights the training observations ($\mathbf{x}_{\bar{S}}^i$) by their proximity to \mathbf{x}_S^* :

Mahalanobis distance to \mathbf{x}_S^*



$$D_S(\mathbf{x}^*, \mathbf{x}^i) = \sqrt{\frac{(\mathbf{x}_S^* - \mathbf{x}_S^i)^T \Sigma_S^{-1} (\mathbf{x}_S^* - \mathbf{x}_S^i)}{|S|}}$$

Gaussian kernel on distance as weight

$$w_S(\mathbf{x}_S^*, \mathbf{x}_{\bar{S}}^i) = \exp\left(-\frac{D_S(\mathbf{x}_S^*, \mathbf{x}^i)^2}{2\sigma^2}\right)$$

Approx $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ by the empirical distribution with weight on each training observation $\mathbf{x}_{\bar{S}}^i$

$$\frac{w_S(\mathbf{x}_S^*, \mathbf{x}_{\bar{S}}^i)}{\sum_{i=1}^{n_{\text{train}}} w_S(\mathbf{x}_S^*, \mathbf{x}_{\bar{S}}^i)}$$

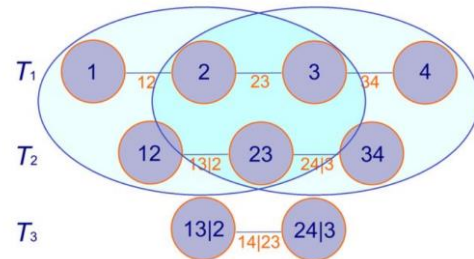
Approaches to estimate and sample from $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$

4. Continuous features: Estimate the dependence structure with a pair copula and weight the training observations using this construction

Transform each \mathbf{x}_i to $u_i \sim U[0,1]$ with inverse empirical CDF

Define $w_S(\mathbf{x}_S, \mathbf{x}_{\bar{S}}) = \hat{c}(u_{\bar{S}}, u_S) / \hat{c}(u_{\bar{S}})$ as the dependence structure when conditioning on \mathbf{x}_S

For each S , estimate the dependence structure $\hat{c}(u_{\bar{S}}, u_S)$ with a non-parametric pair copula



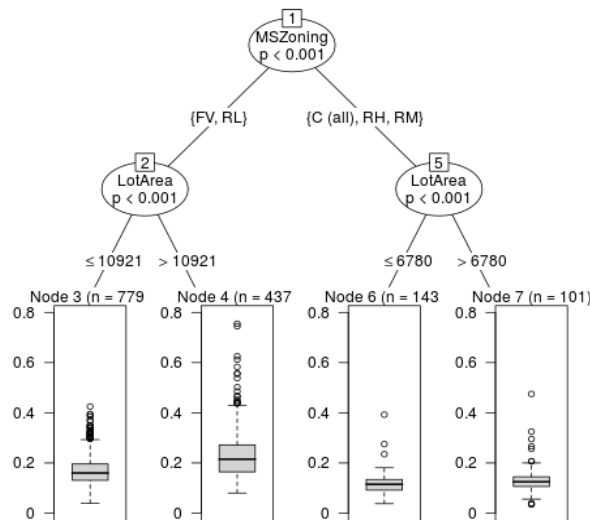
Approx $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ by the empirical distribution with weight on each training observation $\mathbf{x}_{\bar{S}}^i$

$$\frac{w_S(\mathbf{x}_S^*, \mathbf{x}_{\bar{S}}^i)}{\sum_{i=1}^{n_{\text{train}}} w_S(\mathbf{x}_S^*, \mathbf{x}_{\bar{S}}^i)}$$

Approaches to estimate and sample from $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$

5. Mixed data: Use a multivariate decision trees as empirical distribution

For each S , fit a multivariate decision tree* to response $\mathbf{y} = \mathbf{x}_{\bar{S}}$ based on \mathbf{x}_S



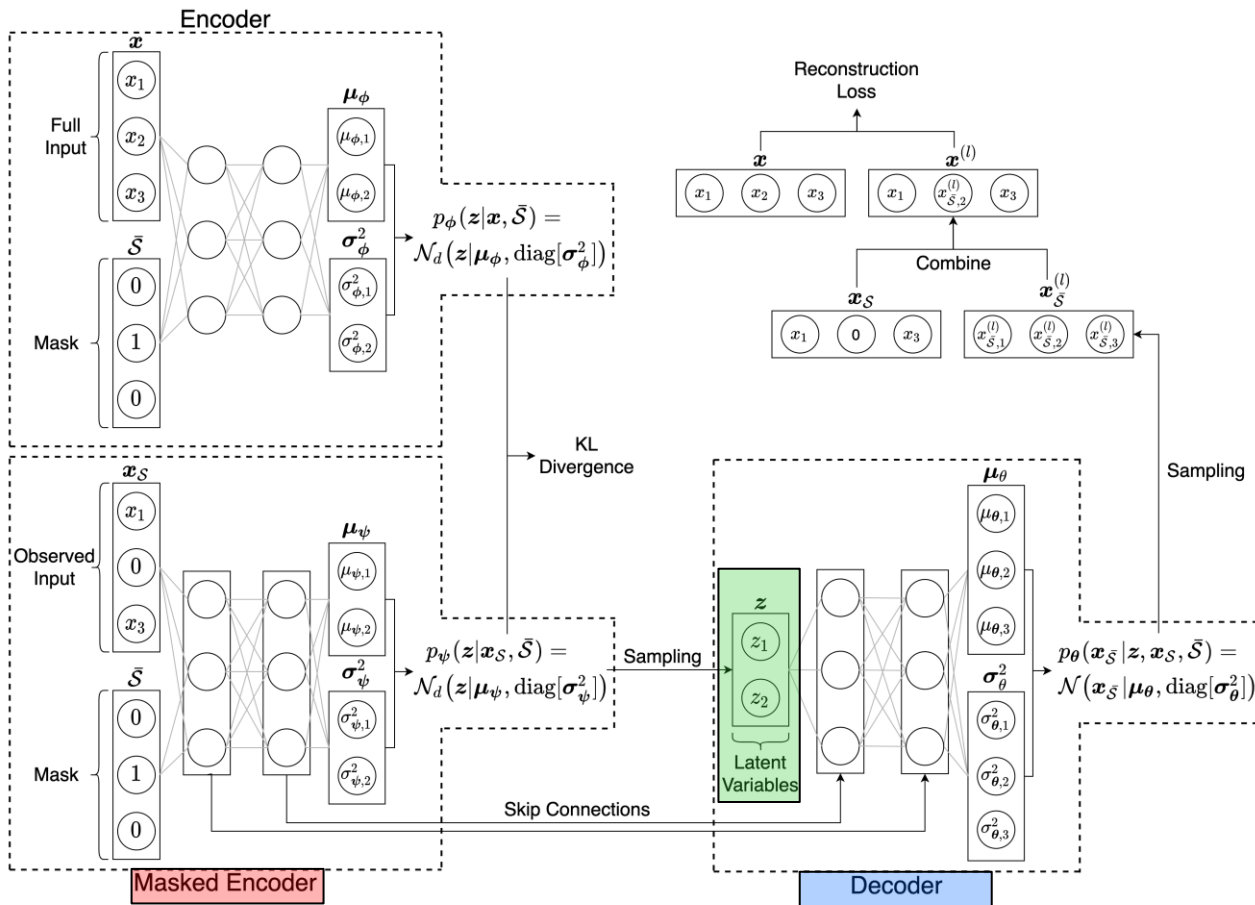
Approx $p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*)$ by the empirical distribution of the training observations $\mathbf{x}_{\bar{S}}^i$ in the terminal node of $\mathbf{x}_S = \mathbf{x}_S^*$

*We use conditional inference trees

Approaches to estimate and sample from $p(x_{\bar{S}} | x_S = x_S^*)$

6. Mixed data: Use an variational autoencoder with arbitrary conditioning (VAEAC)

Fit a VAEAC to all conditional distributions



For each S :

1. Feed the **masked encoder** with x_S^* (masking \bar{S})
2. Generate **latent variables z**
3. Feed z to the **decoder** and use it to generate samples x_S^k
4. Use x_S^k as samples from $p(x_{\bar{S}} | x_S = x_S^*)$

When to use the different approaches

- ▶ Trade-offs between speed and accuracy
- ▶ Performance depends on data type and dependence structure
- ▶ General advice
 - Continuous data: Empirical approach
 - Gaussian if large n_{train} or M
 - Pair-copula if very heavy tails
 - Categorical/mixed data: Ctree
 - Future: Maybe VAEAC when properly implemented
- ▶ TreeSHAP and KernelSHAP available from the shap Python library do NOT give useable estimates of $v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$ unless all features are close to independent

Different prediction explanation games

- ▶ Observational/conditional Shapley values uses

- $v_C(S) = E[f(x)|x_S = x_S^*] = \int f(x_{\bar{S}}, x_S^*) p(x_{\bar{S}}|x_S = x_S^*) dx_{\bar{S}}$

Janzing et al. (2019)

- ▶ Interventional Shapley values uses

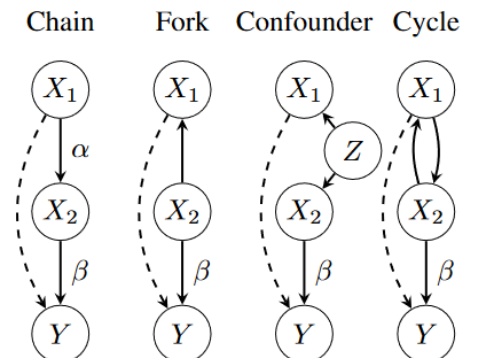
- $v_{do}(S) = E[f(x)|do(x_S = x_S^*)] = \int f(x_{\bar{S}}, x_S^*) p(x_{\bar{S}}|do(x_S = x_S^*)) dx_{\bar{S}} = \int f(x_{\bar{S}}, x_S^*) p(x_{\bar{S}}) dx_{\bar{S}} = v_I(S)$

- ▶ Chen et al. (2020) states that whether $v_C(S)$ or $v_I(S)$ is most appropriate depends on the application

- $v_C(S)$ is most appropriate if you want to learn about the actual relationship between features and modelled response
 - $v_I(S)$ is appropriate if you are debugging your model

- ▶ Heskes et al. (2020): Causal Shapley values

- $v_{do}(S) = E[f(x)|do(x_S = x_S^*)] = \int f(x_{\bar{S}}, x_S^*) p(x_{\bar{S}}|do(x_S = x_S^*)) dx_{\bar{S}}$, but $p(x_{\bar{S}}|do(x_S = x_S^*)) \neq p(x_{\bar{S}})$, so $v_{do}(S) \neq v_I(S)$!
 - Rather *model* $p(x_{\bar{S}}|do(x_S = x_S^*))$ with assumed causal ordering
 - Requires estimate conditional distributions, but not all combinations



Different prediction explanation games

My viewpoint

- ▶ Independence/Interventional Shapley values ($v_I(S)$) is only appropriate if
 - all features close to independent
 - or all dependence between features are due to a common confounder
 - You are debugging/testing robustness of your model
- ▶ Use Causal Shapley values
 - when you have confident knowledge about causal dependence between features
- ▶ All other cases: Use observational/conditional Shapley values
 - Observational/conditional Shapley values \Leftrightarrow Causal Shapley values if no features are assumed to causally affect other features

Two main challenges

1. The computational complexity in the Shapley formula is of size 2^M

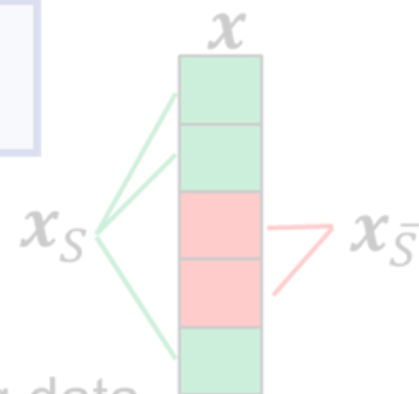
$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

- Approximate solutions may be obtained by using a finite sample of subsets S (KernelSHAP; Lundberg & Lee, 2017)

2. Estimating the contribution function

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

- Lundberg & Lee (2017)
 - Approximates $v(S) \approx \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}}) d\mathbf{x}_{\bar{S}}$,
 - Estimates $p(\mathbf{x}_{\bar{S}})$ using the empirical distribution of the training data
 - Monte Carlo integration to solve the integral



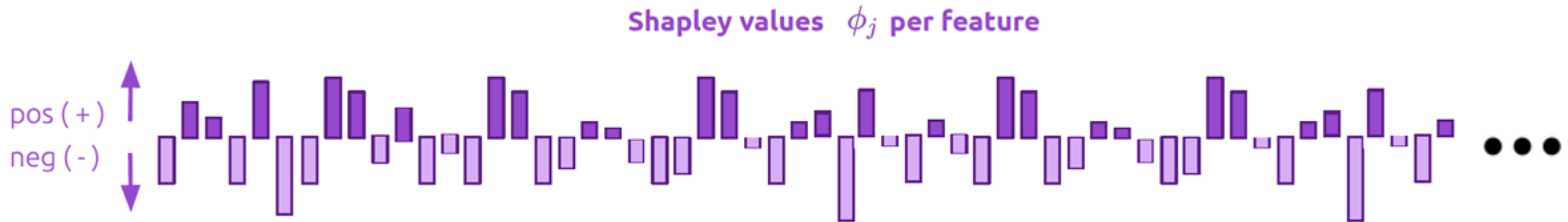
This assumes the features are independent!

Computational bottlenecks

1. The sum in the Shapley value formula is of size 2^M , growing exponentially in the number of features

$$\begin{aligned}M = 5 &\Rightarrow 2^M = 32 \\M = 10 &\Rightarrow 2^M = 1024 \\M = 20 &\Rightarrow 2^M = 1048676 \\M = 40 &\Rightarrow 2^M > 10^{12} \\M = 100 &\Rightarrow 2^M > 10^{30} \\M = 1000 &\Rightarrow 2^M > 10^{301}\end{aligned}$$

2. How can we visualize, interpret and extract knowledge from 100s or 1000s of Shapley values?



- Typically: the sum of many small ϕ_j > sum of the few large ones
- Many highly dependent features complicates the interpretation

groupShapley

Efficient and simple prediction explanations with groupShapley: A practical perspective

Martin Jullum¹, Annabelle Redelmeier¹ and Kjersti Aas¹

¹Norwegian Computing Center, P.O. Box 114, Blindern, N-0314 Oslo, Norway

► Fundamentally very simple approach

- Divide the M features into a small number of G disjoint groups $\{G_1, \dots, G_G\}$.
- Replace the feature subsets S in the Shapley formula by group subsets T :

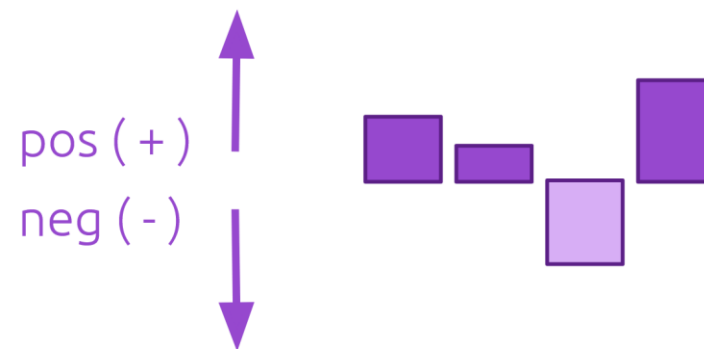
$$\phi_{G_i} = \sum_{T \subseteq G \setminus \{G_i\}} w(|T|) (v(T \cup G_i) - v(T))$$

- The scores are still Shapley values, so all mathematical properties are kept (on group level)

► What about the bottlenecks?

- $2^G \ll 2^M \Rightarrow$ computationally tractable
- G small \Rightarrow easy to visualize

Shapley value contribution ϕ_{G_i} per feature group



How to group the features?

- ▶ Crucial to group features based on the desired explanation
- ▶ Grouping based on feature dependence
 - Highly dependent features grouped together, using e.g. a **clustering method**.
 - Easier to study theoretically
 - Often **difficult to extract knowledge** from in practice
- ▶ Grouping based on application/feature knowledge
 - Group features of **similar type** or general category
 - Gives directly **meaningful interpretations** of computed groupShapley values
 - May perform multiple explanations with different groupings for increased understanding
- ▶ **We advocate grouping based on feature knowledge in practical applications**

Practical example 1: Car insurance

► US Car insurance dataset

- 10 302 customers with records of **crash/no crash** + 21 features
- Fit a **random forest** model with 500 trees to **predict crash** based on the 21 features

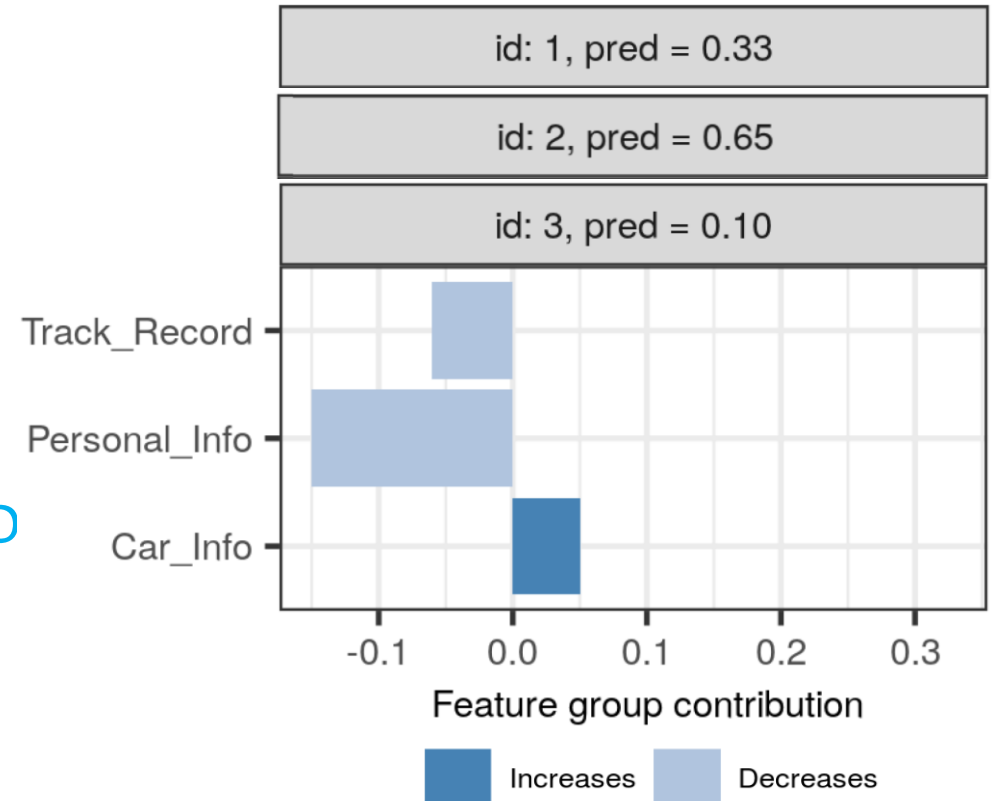


- 3 feature groups based on type
 - **Track record** (4 features): # claims last 5 years, # licence record points, previous licence revokes, time as customer
 - **Personal information** (13 features): age of driver, education level, # children, job type, # driving children, marital status, gender, distance to work +++
 - **Car information** (4 features) value of car, age of car, type of car, whether car is red

Practical example 1: Car insurance

► Explain predictions for 3 individuals

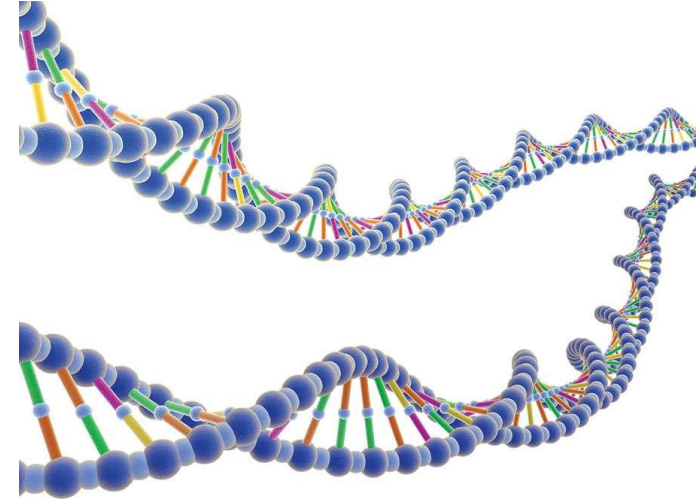
1. 1 claim last 5 years, 3 licence record points.
Single mother of 4 (2 driving).
Driving a SUV, 27 miles to work.
2. Got licence revoked and 10 licence record points.
37 year old father of 2 (1 driving).
3. 3 claims last 5 years, no licence record points
60 year old married doctor with no children, with a PhD
Red sports car.



Practical example 2: Gene data

► Disease classification with high dimensional gene data

- 127 patients where 85 are diseased with either Crohn's disease (CD) or Ulcerative colitis (UC) + 42 healthy controls.
- 4 834 genes (after pre-processing)
- Using 100 random individuals, we fit a Lasso penalized linear regression model to predict $P(\text{diseased with either CD or UC})$ based on the patient's genes

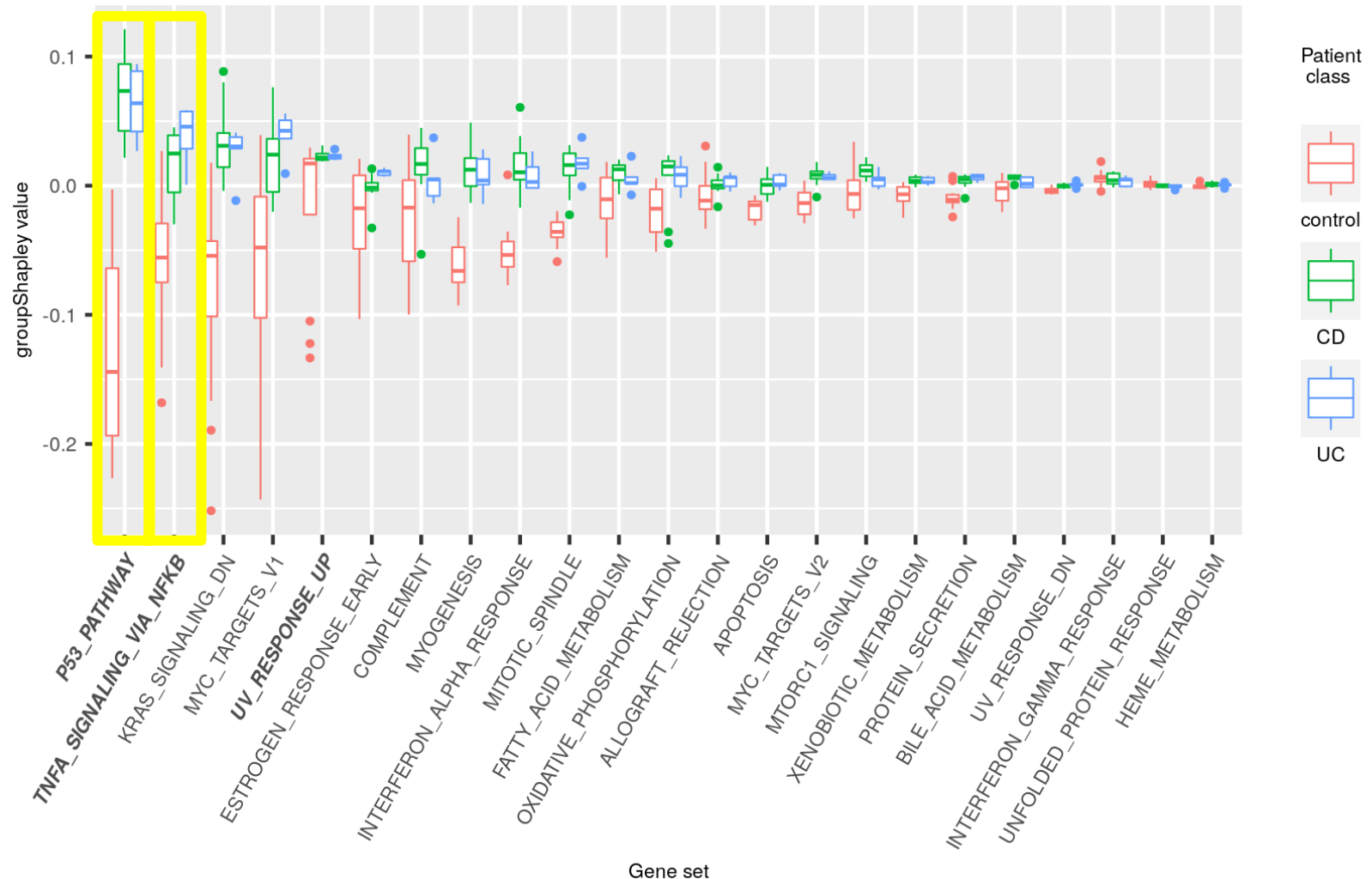


► Feature groups

- Use the so-called Hallmark gene set to group the features (genes) into 23 different groups commonly used in gene set enrichment analysis
- The Hallmark gene set “conveys a specific biological state or process” (Liberzon et al., 2015)

Practical example 2: Gene data

- ▶ Compute **groupShapley values** for the remaining 27 patients
- ▶ Make separate groupShapley **boxplots** for UC, CD and controls
- ▶ Can we identify **genetical similarities** and **differences** for UC and CD?
- ▶ Note: **Model not trained to separate UC and CD**



Software

shapr: An R-package for explaining machine learning models with dependence-aware Shapley values

Nikolai Sellereite¹ and Martin Jullum¹

- ▶ R-package shapr
 - Computes Shapley values for any model $f(x)$ with different dependence-aware methods for estimating $v(S)$
 - All functionality works for both feature-wise and group-wise Shapley values
 - Currently undergoing heavy restructuring to allow
 - Parallellization
 - Reduce memory usage
 - Causal Shapley values
 - Improved user experience +++
 - Python wrapper



Take home points

- ▶ The Shapley value framework from game theory can be used to explain predictions from any ML model
- ▶ Shapley value measures the value of observing each feature
- ▶ There are two main challenges with such explanations
 - Computational complexity -> Approximate by sample of subsets S , or explain feature groups instead
 - Estimating contribution function $v(S)$ -> Several different methods for different settings
- ▶ There exists other prediction explanation games:
 - Interventional Shapley values can be used for “debugging”
 - Causal Shapley values is promising if you have prior causal knowledge
- ▶ You can do most (hopefully all quite soon) types of prediction explanations efficiently with the shapr R-package
 - See package vignette at <https://norskregnesentral.github.io/shapr/> for an intro

References

1. Aas, K., Jullum, M., Løland, A.: Explaining individual predictions when features are dependent: More accurate approximations to shapley values. *Artificial Intelligence* **298**, 103502 (2021)
2. Aas, K., Nagler, T., Jullum, M., Løland, A.: Explaining predictive models using shapley values and non-parametric vine copulas. *Dependence Modeling* **9**, 62–81 (2021)
3. Chen, H., Janizek, J.D., Lundberg, S., Lee, S.I.: True to the model or true to the data? arXiv preprint arXiv:2006.16234 (2020)
4. Heskes, T., Sijben, E., Bucur, I.G., Claassen, T.: Causal shapley values: Exploiting causal knowledge to explain individual predictions of complex models. *Advances in Neural Information Processing Systems* **33** (2020)
5. Janzing, D., Minorics, L., Blöbaum, P.: Feature relevance quantification in explainable ai: A causal problem. In: *International Conference on Artificial Intelligence and Statistics*. pp. 2907–2916. PMLR (2020)
6. Jullum, M., Redelmeier, A., Aas, K.: Efficient and simple prediction explanations with groupshapley: A practical perspective. In: *Proceedings of the 2nd Italian Workshop on Explainable Artificial Intelligence*. pp. 28–43. CEUR Workshop Proceedings (2021)
7. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: *Proceedings of the 31st international conference on neural information processing systems*. pp. 4768–4777 (2017)
8. Olsen, L.H.B., Glad, I.K., Jullum, M., Aas, K.: Using shapley values and variational autoencoders to explain predictive models with dependent mixed features. arXiv preprint arXiv:2111.13507 (2021)
9. Redelmeier, A., Jullum, M., Aas, K.: Explaining predictive models with mixed features using shapley values and conditional inference trees. In: *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*. pp. 117–137. Springer (2020)
10. Sellereite, N., Jullum, M.: shapr: An r-package for explaining machine learning models with dependence-aware shapley values. *Journal of Open Source Software* **5**(46), 2027 (2020)
11. Shapley, L.: A value for n -person games. *Ann. Math. Study*28, *Contributions to the Theory of Games*, ed. by HW Kuhn, and AW Tucker pp. 307–317 (1953)

Contact

jullum@nr.no

<https://martinjullum.netlify.app>