



# Introduction to explainable artificial intelligence

## SHAP and GNNExplainer

Agnieszka Wojtuch<sup>1</sup>

<sup>1</sup>Department of Informatics and Chemistry,  
Vysoká škola chemicko-technologická v Praze

June 2025



Introduction

Examples

SHAP

GNNEExplainer



Introduction

Examples

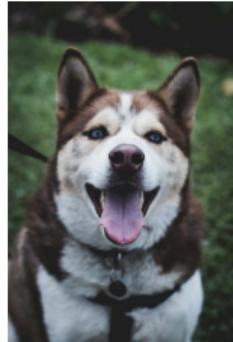
SHAP

GNNExplainer

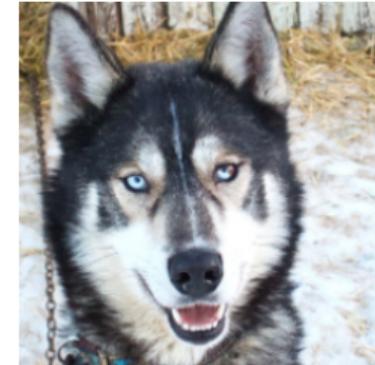


## training data

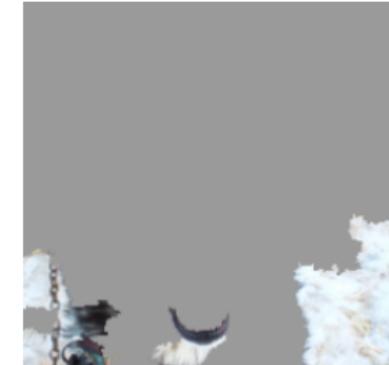
huskies:



input image



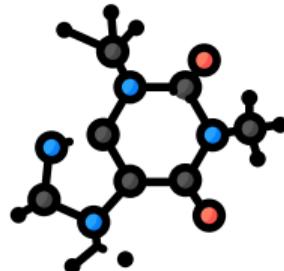
explanation



wolves:



correct: husky  
predicted: wolf

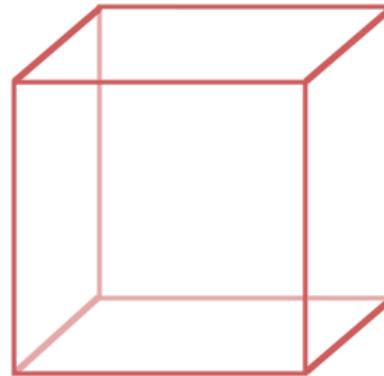


**Explainability** consists of methods that explain how a particular AI model works in a comprehensible way.

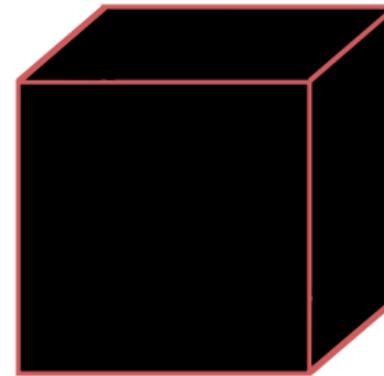
- ▶ why was this sample classified as negative? (right to explanation)
- ▶ which features does the model use? (fairness, scientific discovery)
- ▶ how this feature affects the prediction? (right to explanation, sample optimisation)



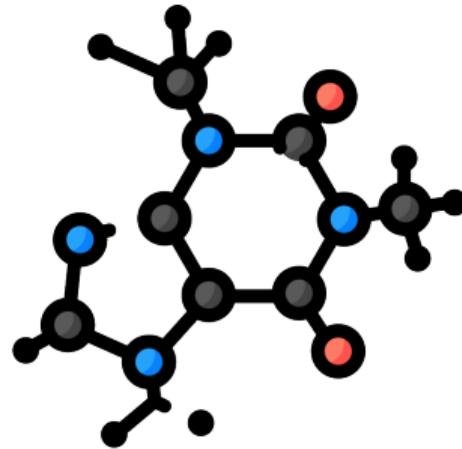
interpretability



explainability

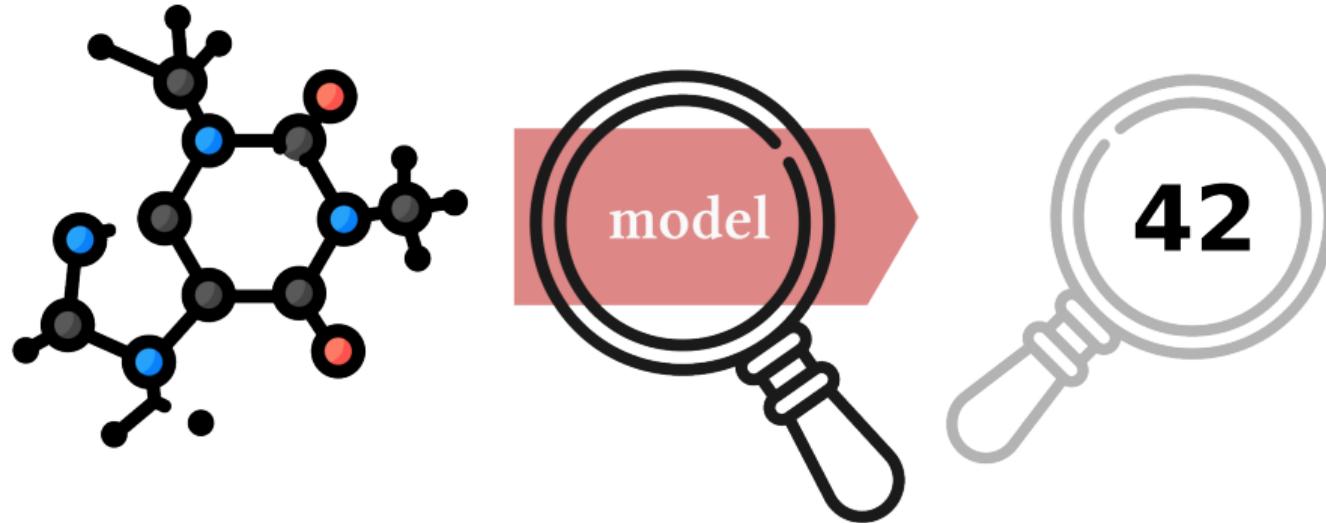


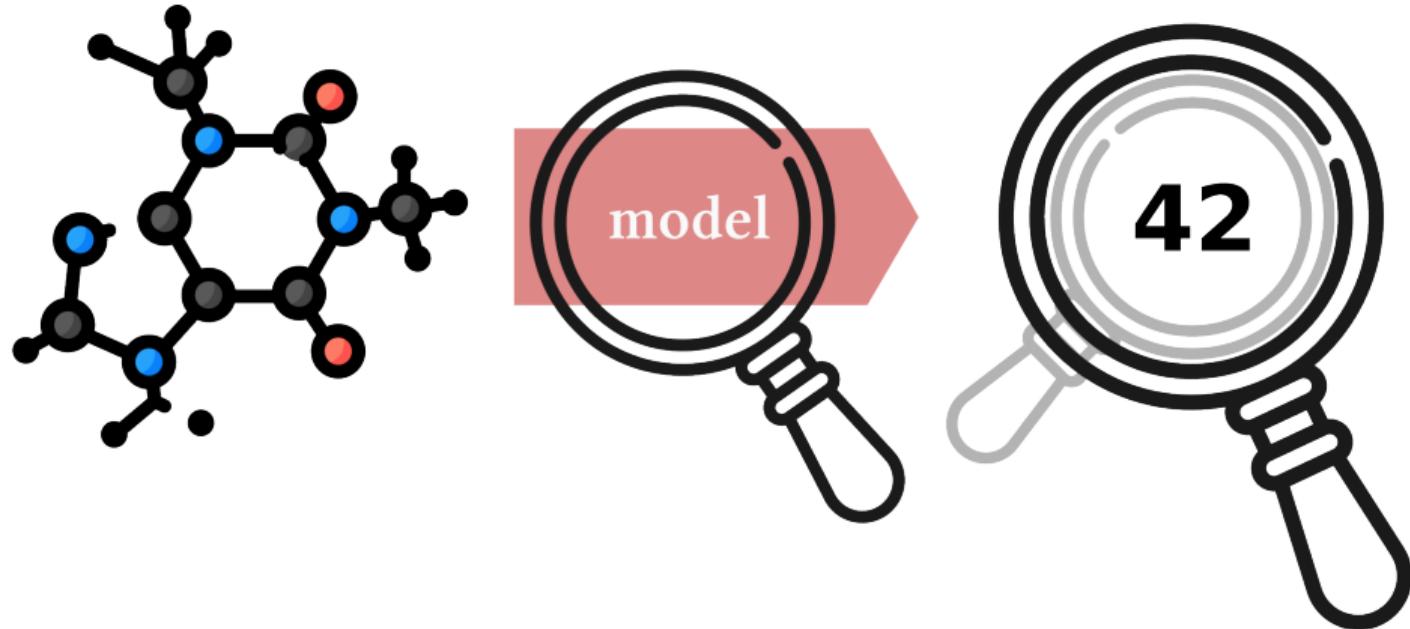
- ▶ **interpretability** – we can interpret what the model does e.g. a small decision tree
- ▶ **explainability** – someone has to explain the model's inner working to us  
e.g. a neural network



model









1. understanding the model
  - ▶ model diagnostics
  - ▶ model trust
2. understanding the modelled process
  - ▶ scientific discovery
  - ▶ molecule optimisation

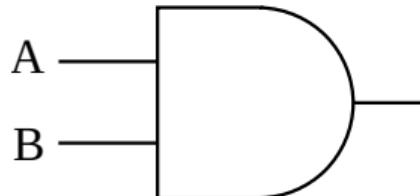


Introduction

Examples

SHAP

GNNEExplainer



model: logic AND gate

inputs: A, B

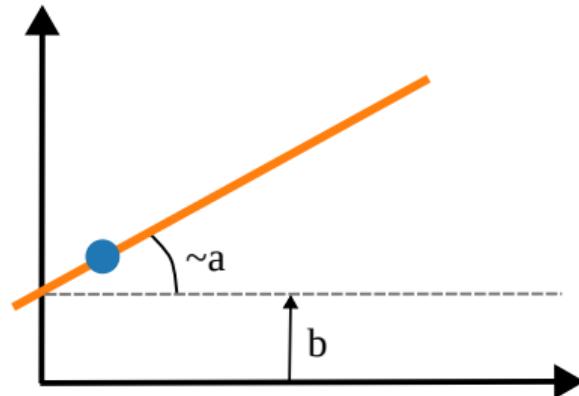
output: 1 if both A and B are 1; otherwise 0

What is the relative importance  
of these features?

$A = 1, B = 1$

And in this case?

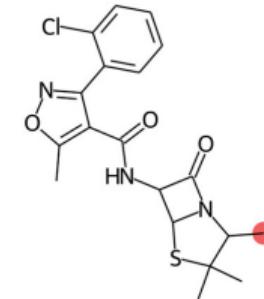
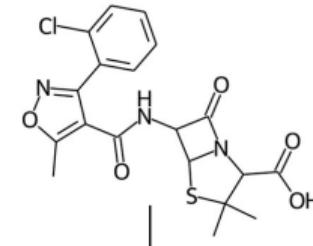
$A = 0, B = 1$



$$y = \underbrace{X \cdot a}_{\text{explanation}} + b$$



1	0	0	1	0	1	1	0
1.7	-0.3	0.45	0.97	-2.1	0.03	1.25	-1.4





Introduction

Examples

SHAP

GNNExplainer

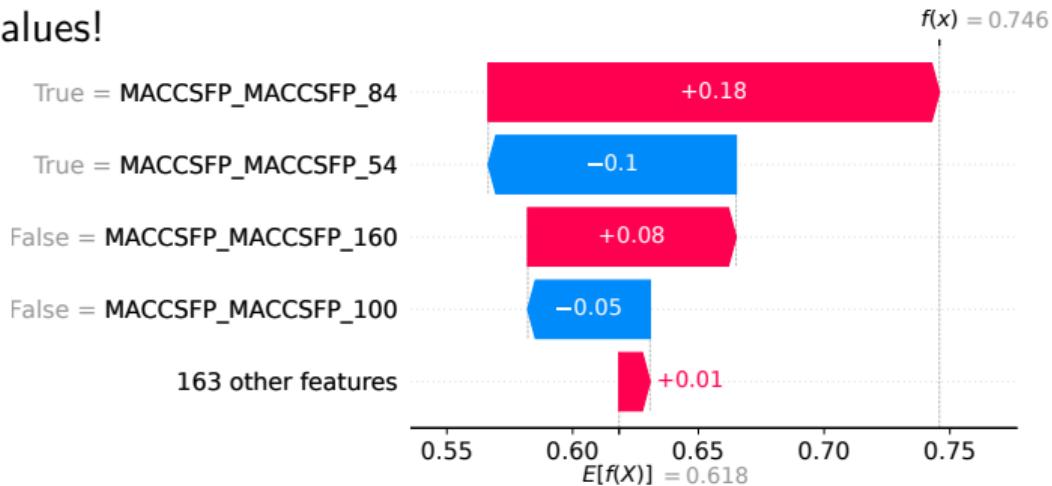


The linear regression-like explanations sum to the prediction so they're easy to interpret.

Could we use them for other models somehow?



We could use Shapley values!





- ▶ multi-player coalition game – the players work together to gain the biggest reward they can
- ▶ the reward is then divided between the players

**Problem:** How to split the total payout (reward) between the players in a **fair way**, i.e. taking into consideration their contribution?

Shapley showed that there exists a **unique solution** to this problem.



## Efficiency

The entire reward must be split.

## Symmetry

If two players contributed equally then their payout should be the same.

## Dummy

A player who has no impact on the reward value gets no reward.

## Additivity

If players play one game after another then calculating their rewards after each game should give the same result as calculating their total rewards once, after the last game.

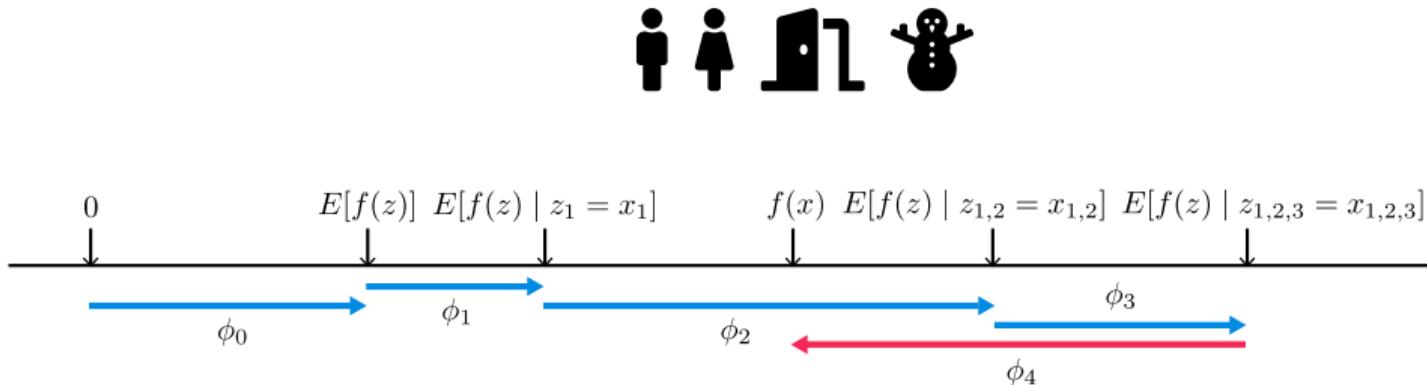


**Definition:** The Shapley value is the average marginal contribution of a feature value across all possible coalitions.

$$\phi_j(val) = \sum_{\substack{S \subseteq \\ \text{all possible coalitions}}} \underbrace{\frac{|S|!(p - |S| - 1)!}{p!}}_{\text{weighting factor}} \left( \underbrace{val(S \cup \{j\})}_{\text{reward with player } j \text{ present}} - \underbrace{val(S)}_{\text{reward without player } j} \right)$$

↑  
of  $p$  players  
without player  $j$

↑  
 $j$ 's contribution for coalition  $S$



When calculating Shapley values, all possible orderings are considered and an average contribution is calculated.

Order matters when the function is nonlinear or the variables are not independent.



**Nonlinearity:** let's go back to our AND example.  $A = 0, B = 1$

- ▶ if we learn about A first, then it gets the entire contribution – we already know what the output would be – and so B will get 0 contribution
- ▶ but if we learn first about B, then we get some information but we still don't know what the output would be – B gets non-zero contribution



### Dependent variables:

- ▶ we have four variables: A, B, C, D;  $B = C + D$
- ▶ we predict the sum of all the variables
- ▶ if we learn about the variables in order: A, B, C, D; then C and D have no contribution – once we know A and B, we know that the sum is  $A + 2B$
- ▶ if we learn about the variables in order: A, C, D, B, then C and D have nonzero contribution, but B has zero contribution – once we learn about A, C, D, we know that the sum is  $A + 2 * (C + D)$

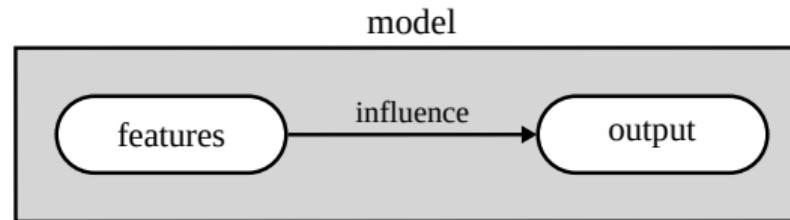
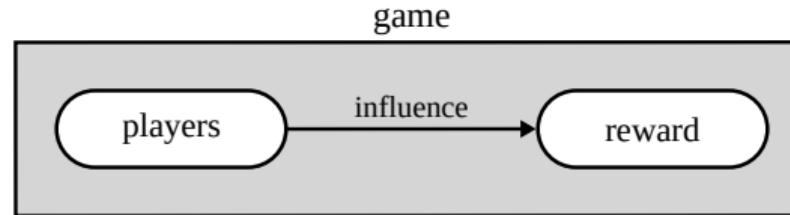


When calculating  $j$ 's contribution for some coalition  $S$  we need to know how many sequences correspond to this situation:

- ▶  $|S|$  players can be ordered in  $|S|!$  sequences
  - ▶ the remaining players can be ordered in  $(p - |S| - 1)!$  sequences
  - ▶ coalition  $S$  corresponds to  $|S|!(p - |S| - 1)!$  sequences

There are  $p!$  sequences of  $p$  players, so each sequence has weight  $\frac{1}{p!}$   
 Therefore, the weight for  $j$ 's contribution to coalition  $S$  is:

$$\frac{|S|!(p - |S| - 1)!}{p!}$$





$$\phi_j(val) = \sum_{S \subseteq \{x_1, \dots, x_p\} \setminus \{x_j\}} \frac{|S|!(p - |S| - 1)!}{p!} (val(S \cup \{j\}) - val(S))$$



- ▶ When we calculate Shapley values we play out an imaginary scenario in which some players are not taking part in the game.
- ▶ In AI, this corresponds to not knowing values of some features. But most models don't accept missing features!

**Solution:** Although one cannot calculate a prediction when a feature is missing, one can pretend not to know its value. The value of the 'missing' feature is replaced with any value from the dataset and prediction can be easily calculated.

SHAP values are not exactly the same as Shapley values.

SHAP values are the Shapley values of a **conditional expectation function of the original model**.

$$\hat{y} = f(x) \rightarrow \bar{y} = \mathbb{E}[f(x)]$$



$$\tilde{y} = f(\tilde{x})$$



$$\bar{y} = \mathbb{E}_{|x_{1:4}}[f(x)]$$

In simple English: SHAP values are the Shapley values of a **fancy version of our model**.



Let's predict apartment price based on four features:

- ▶ 🏅, 🥈, 🏆 – floor
- ▶ 1 2 3 4 – size
- ▶ 🌳, 🚨 – park nearby?
- ▶ 😺, 😵 – are cats allowed?

$$(\text{🏅}, \text{1 2 3 4} \text{50m}^2, \text{🌳}, \text{😺}) \rightarrow 300\text{K}$$

What is the contribution of the cat feature 😺?

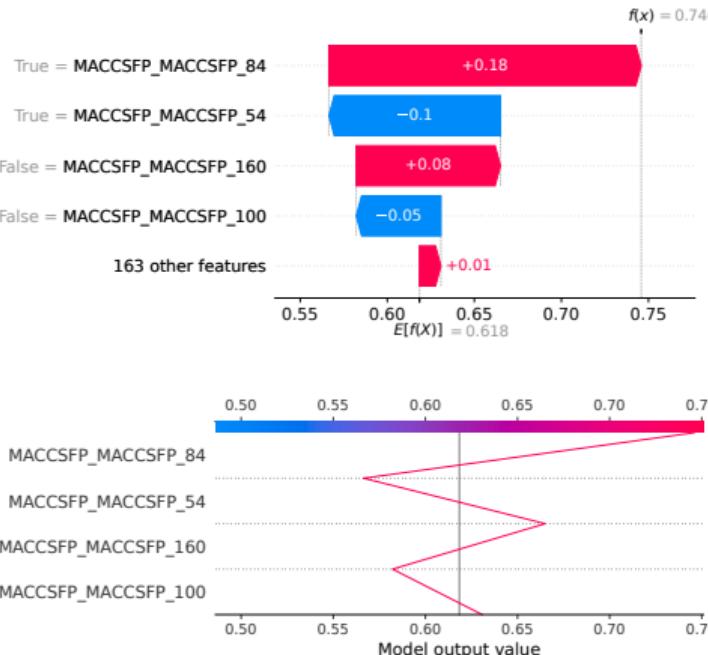
1. Sample a coalition of features: 1 2 3 4 50m<sup>2</sup>, 🌳
2. Sample the values for the missing features: 🏅
3. Prediction for the new sample: 🏅, 1 2 3 4 50m<sup>2</sup>, 🌳, 😵 → 310K
4. Replace the cat feature with a randomly drawn value: 😺 → 😵
5. Prediction for the new sample: 🏅, 1 2 3 4 50m<sup>2</sup>, 🌳, 😺 → 320K
6. The contribution of the cat feature is: 310K - 320K = -10K
7. This is a very rough estimation! Average over more samples → repeat steps 2-6.
8. Do this for every coalition → repeat steps 1-7.

$$\sum_{S \subseteq \{x_1, \dots, x_p\} \setminus \{x_j\}} \frac{|S|!(p-|S|-1)!}{p!} (\text{val}(S \cup \{j\}) - \text{val}(S))$$



Be careful to interpret the Shapley value correctly: the Shapley value is the **average contribution** of a feature value to the prediction **in different coalitions**.

The Shapley value is **NOT** the difference in prediction when we would remove the feature from the model.



- ▶ additive feature attribution method
- ▶ explains how to get from the base prediction (mean prediction of the model) to the actual prediction
- ▶ suitable for tabular data (e.g. fingerprints)
- ▶ both classification and regression
- ▶ model-agnostic – works with any model
- ▶ provides local explanations – explains predictions (not the entire model)
- ▶ perturbation-based – calculates explanations by perturbing (introducing small changes into), the input instance
- ▶ based on Shapley values from game theory – nice mathematical guarantees
- ▶ can be rather slow but there exist model-specific fast alternatives
- ▶ **Limitation:** prone to unrealistic data instances



Introduction

Examples

SHAP

GNNExplainer



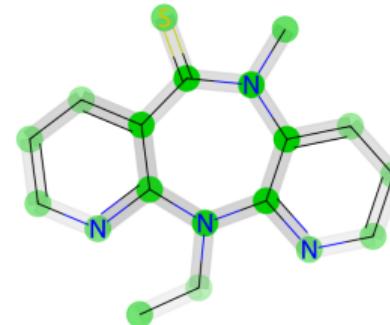
How to calculate explanation of this prediction?

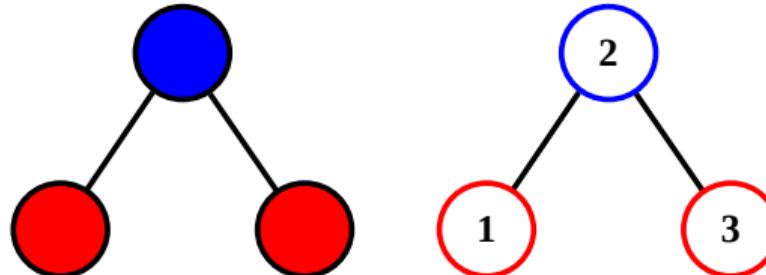


Let's do what we always do – train a neural network to solve the problem for us!



We want to see which nodes are important for the prediction. And edges.  
And features.





$$A = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

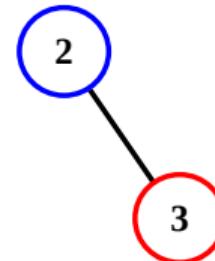
$$A = [[1, 2], [2, 1], [2, 3], [3, 2]]$$

feature order: B, C, N, O, F

$$X = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$M_v = [0, 1, 1]$$

$$M_e = [0, 0, 1, 1]$$





Bernoulli distribution:

$$P(X = 1) = p$$

$$P(X = 0) = 1 - p$$

$$P(x_1 = 1, x_2 = 0, x_3 = 1) =$$

$$P(X = 1) \cdot P(X = 0) \cdot P(X = 1) =$$

$$p \cdot (1 - p) \cdot p$$

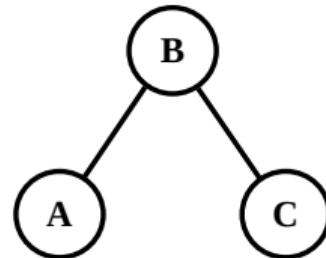


Multivariate Bernoulli distribution:

$$P(A = 1) = p_A$$

$$P(B = 1) = p_B$$

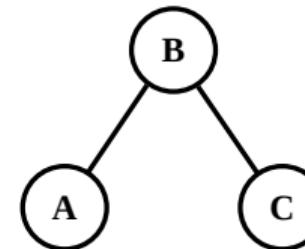
$$P(A = 1, B = 0) = p_A \cdot (1 - p_B)$$



$$P(A = 1) = p_A$$

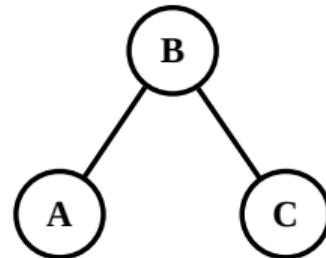
$$P(B = 1) = p_B$$

$$P(C = 1) = p_C$$



$$P(A = 1, B = 1, C = 1) = p_A \cdot p_B \cdot p_C$$

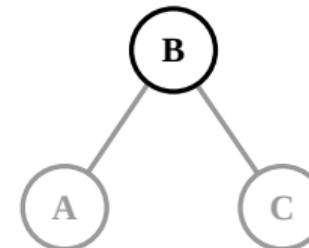




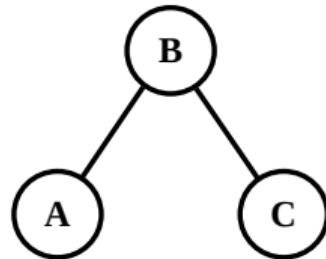
$$P(A = 1) = p_A$$

$$P(B = 1) = p_B$$

$$P(C = 1) = p_C$$



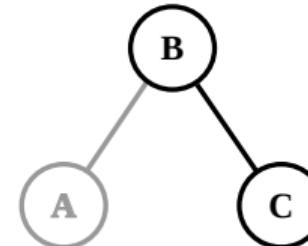
$$P(A = 0, B = 1, C = 0) = (1 - p_A) \cdot p_B \cdot (1 - p_C)$$



$$P(A = 1) = p_A$$

$$P(B = 1) = p_B$$

$$P(C = 1) = p_C$$



$$P(A = 0, B = 1, C = 1) = (1 - p_A) \cdot p_B \cdot p_C$$



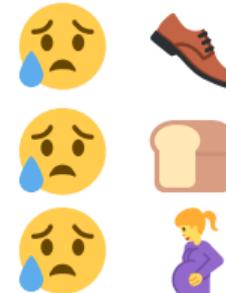
entropy:

$$H[X] = - \sum_{x \in X} p(x) \log p(x)$$



mutual information:

$$MI[X, Y] = H[X] - H[X|Y] = H[Y] - H[Y|X]$$





A proxy model is trained to generate such a mask  $M$  that defines a subgraph of the input graph which maximises the probability of class  $c$  predicted by black-box model  $\Phi$ .

In other words:

- ▶ we look for a smaller graph that would give the same prediction as the input graph,
- ▶ this smaller graph is defined by masking the input graph with mask  $M$ ,
- ▶ mask  $M$  is generated by a proxy model.



$$\max_{G_S} \underbrace{MI(Y, (G_S, X_S))}_{\substack{\text{random variable} \\ \text{representing labels} \\ \uparrow \\ \text{how much do we know about } Y \\ \text{knowing } (G_S, X_S)?}} \quad \begin{array}{c} \text{features of } G_S \\ \downarrow \\ X_S \end{array}$$

↑  
examine all subgraphs of  $G$



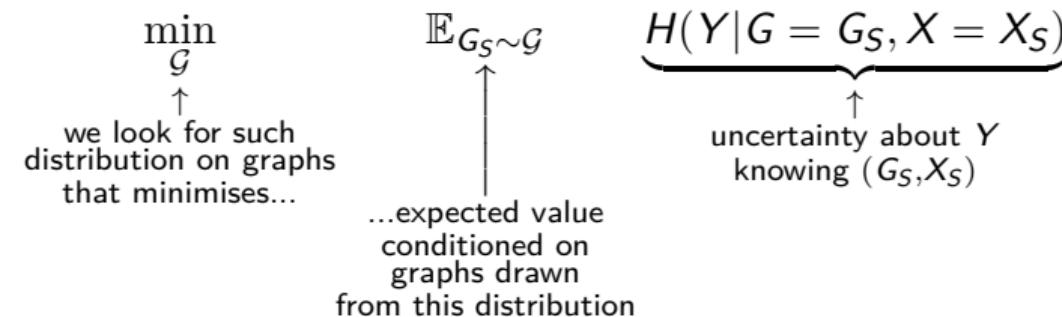
We express mutual information in terms of entropy:

$$\max_{G_S} MI(Y, (G_S, X_S)) = H(Y) - \underbrace{H(Y|G = G_S, X = X_S)}_{\substack{\text{uncertainty about } Y \\ \text{knowing } (G_S, X_S)}}$$

↑  
uncertainty about  $Y$   
(constant)



$$\min_{G_S} \longrightarrow \min_{\mathcal{G}} \mathbb{E}_{G_S \sim \mathcal{G}}$$





Using Jensen's inequality ( $\mathbb{E}[\varphi(X)] \leq \varphi(\mathbb{E}[X])$ ), we get the following upper bound:

$$\mathbb{E}_{G_S \sim \mathcal{G}} H(Y|G = G_S, X = X_S) \leq H(Y|G = \mathbb{E}_{\mathcal{G}}[G_S], X = X_S)$$

so we can do:

$$\min_{\mathcal{G}} H(Y|G = \mathbb{E}_{\mathcal{G}}[G_S], X = X_S)$$

## How to estimate $\mathbb{E}_{\mathcal{G}}$ ?

We will parametrise the distribution over graphs,  $\mathcal{G}$ , as a multivariate Bernoulli distribution:

$$P_{\mathcal{G}} \quad (G_S) = \prod_{(j,k) \in G_c} A_S[j, k]$$

graph sampled from  $\mathcal{G}$

$P_{\mathcal{G}}$

↑ probability defined by distribution  $\mathcal{G}$

$(G_S)$

↑

$A_S$

[ $j, k$ ]

↑ expectation on whether edge  $(v_j, v_k)$  exists

↑ fractional adjacency matrix of  $G_S$

↑

edge between  $j^{th}$  and  $k^{th}$  node

edges from **computational** input graph



Thus

$$H(Y|G = \mathbb{E}_{\mathcal{G}}[G_S], X = X_S) = H(Y|G = \underbrace{A_c \odot \sigma(M)}_{\text{expected graph}}, X = X_S)$$

- ▶  $\odot$  – element-wise multiplication
- ▶  $M$  – mask that we need to learn
- ▶  $\sigma(\cdot)$  – sigmoid



$$H[X] = - \sum_{x \in X} p(x) \log p(x)$$

$$H(Y|G=G_S, X=X_S) = -\mathbb{E}_{Y|G_S, X_S}$$

↑  
expected distribution of  $Y$   
conditioned on  $G_S, X_S$

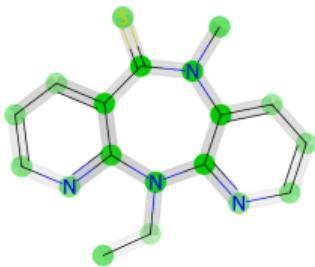
probability defined  
by model  $\Phi$   
 $\downarrow$   
 $P_\Phi$   
 $\underbrace{(Y|G=G_S, X=X_S)}_{\text{distribution of labels  
given only a subgraph  
and a subset of features}}$



$$\min_M - \underbrace{\sum_{c=1}^C \mathbb{1}[y = c] \log}_{\text{cross-entropy loss}} \underbrace{P_\Phi(Y = y | G = A_c \odot \sigma(M), X = X_c)}_{\Phi's \text{ prediction for the masked graph}}$$



- ▶ simplifying the optimisation criterion
- ▶ modelling graphs with a probabilistic distribution (helped with tractability)
- ▶ using Jensen's inequality to get an upper bound
- ▶ tricking neural networks into returning values from  $[0, 1]$  range



- ▶ suitable for graph representations
- ▶ model-agnostic – works with any graph neural-network
- ▶ uses a proxy-model – a small graph neural network is trained to predict the explanation
- ▶ explanation is a subgraph of the input graph and a subset of features that are important for the prediction
- ▶ both classification and regression
- ▶ provides local explanations
- ▶ produces graph-level explanations – graph-level tasks depend on the entire graph, e.g. activity prediction
- ▶ produces node-level explanations – node-level tasks concern a single node, e.g. predicting aromaticity of each atom
- ▶ **Limitation:** the authors assume (by using Jensen's inequality) that the function represented by GNN being explained is convex



### Attributions:

- ▶ Motivation
  - ▶ dataset photographs come from Unsplash.com
  - ▶ input and explanation photographs come from „*Why Should I Trust You?*“  
*Explaining the Predictions of Any Classifier* M. T. Ribeiro et al.
- ▶ Example 3 – different explanations – an example of counterfactual explanation was adapted from  
*A Perspective on Explanations of Molecular Prediction Models*, Geemi P. Wellawatte, et al.
- ▶ Shapley values, The unique solution, Missing players – missing features
  - ▶ image of Lloyd Shapley – CC BY-SA 2.0 DE copyright is with Mathematisches Institut Oberwolfach, background was removed from the image
  - ▶ cards – from Flaticon.com
- ▶ Order matters – *A Unified Approach to Interpreting Model Predictions*, Scott M. Lundberg and Su-In Lee
- ▶ twemoji icons – The Twitter emoji graphics provided by twemojis are licensed under CC-BY 4.0.  
Copyright ©2019, Twitter, Inc. and other contributors.