

# SHAPLEY'S NALUE: FROMGAMENTHEORY TO EXPLAINABLE A

#### **Colin Rowat**



#### UNIVERSITY<sup>OF</sup> BIRMINGHAM

# 





IVAL

### von Neumann matches gloves



**Colin Rowat** 

Z has a left glove K & S each have a right glove a coalition only scores if it has a pair: v(S,Z) = v(K,Z) = v(S,K,Z) = 1a von Neumann-Morgenstern solution: Z has a veto K & S together have a veto any outcome where S,K share equally yuck!

- 1.
- 2.
- 3. 4.





# BIRMINGHAM

many outcomes in the solution many asymmetric solutions solutions might not exist solutions might be really ugly **#ESRCFestival** 



S, Z

TIVAL

 $\frac{1}{3}$ 

# enter the Shapley value, I

K, S, Z

**Colin Rowat** 

 $\frac{1}{3}$ 

what should each player receive? average over value it adds to coalitions  $p = \frac{1}{3}$ : K joins empty coalition; adds 0  $p = \frac{1}{3} \times \frac{1}{2}$ : K joins S; adds 0  $p = \frac{1}{3} \times \frac{1}{2}$ : K joins Z; adds 1  $p = \frac{1}{3}$ : K joins S, Z; adds O • thus, K's Shapley value is 1/6

#### yum!

- 2. 3.
- 4.



# BIRMINGHAM

a single, unique result, (1/6, 1/6, 2/3) always exists easy to calculate symmetric





FESTIVAL

### exit von Neumann, Shapley

**OF/SOCIAL SCIENCE** GAULISH VILLAGE BELGICA COMPENDIUM LAUDANUM AQUARIUM LUTETIA OTORUM ARMORICA GAUL ROMAN CONQUEST) 50 B.C. CELTICA PROVINCIA AQUITANIA



# 1. 2.

**Colin Rowat** 

asterix.fandom.com



#### UNIVERSITYOF BIRMINGHAM

Why does Nash's equilibrium take over? handles information better constructive: how do we get an outcome?



20

WI

 $w_2$ 

W3

Dr.

 $x_1$ 

 $(x_2)$ 

 $(x_n)$ 

#### rise of the machines FESTIVAL OF/SOCIAL SCIENCE

 $w_i x_i$ 

When an neural network/AI system makes a prediction, how do we explain it?

 $\sigma$ 

 $\Sigma | \sigma$ 

 $w_0$ 

**Colin Rowat** 



#### NIVERSITYOF BIRMINGHAM

#### Hidden layer

Input

layer

11

19

Output layer

 $()_1$ 

 $O_2$ 

#### github.com/PetarV-/TikZ



FESTIVAL

### enter the Shapley value, II

APPLIED STOCHASTIC MODELS IN BUSINESS AND INDUSTRY Appl. Stochastic Models Bus. Ind., 2001; **17**:319–330 (DOI: 10.1002/asmb.446)

Analysis of regression in game theory approach

Stan Lipovetsky<sup>\*,†</sup> and Michael Conklin

#### SUMMARY

Working with multiple regression analysis a researcher usually wants to know a comparative importance of predictors in the model. However, the analysis can be made difficult because of multicollinearity among regressors, which produces biased coefficients and negative inputs to multiple determination from presumably useful regressors. To solve this problem we apply a tool from the co-operative games theory, the Shapley Value imputation. We demonstrate the theoretical and practical advantages of the Shapley Value and show that it provides consistent results in the presence of multicollinearity. Copyright © 2001 John Wiley & Sons, Ltd.

Knowl Inf Syst (2014) 41:647–665 DOI 10.1007/s10115-013-0679-x Explaining prediction mod with feature contributions Erik Štrumbelj · Igor Kononenko

This leads to the following explicit definition (see [28] for proof):

 $\varphi_i(x) =$ 

Equation (9) is equivalent to the Shapley value [25], a concept from coalitional game theory. In a coalitional game, it is usually assumed that n players form a grand coalition **#ESRCFestiva** 

#### **Colin Rowat**



#### UNIVERSITY<sup>OF</sup> BIRMINGHAM

# Explaining prediction models and individual predictions with feature contributions

$$\frac{|Q|!(|S| - |Q| - 1)!}{|S|!} (\Delta Q \cup \{i\}(x) - \Delta Q(x)).$$
(9)



TIVAL

**K**, S

 $\frac{1}{3}$ 

**K**, **S**, **Z** 

# how did Shapley get here?

# K, S, Z

**K**, S, Z

K, S, Z

K, S, Z

players  $\Rightarrow$  features

in the coalition  $\Rightarrow$  specific value of feature out of the coalition  $\Rightarrow$  average value e.g.  $f(\overline{K},S_i,Z_i)$  prediction given average K, specific values of S and Z

'value added' ⇒ the prediction's change  $f(K,S,Z) - f(\overline{K},S,Z)$ 

#### $\frac{1}{3}$ **Colin Rowat**

**K**, S, **Z** 





# BIRMINGHAM

e.g.  $K \Rightarrow$  e.g. `kids allowed' ...



FESTIVAL

# Shapley redux?



#### **Colin Rowat**



UNIVERSITY<sup>OF</sup> BIRMINGHAM

The Shapley value might be the only method to deliver a full explanation. In situations where the law requires explainability like EU's "right to explanations" the Shapley value might be the only legally compliant method, because it is based on a solid theory and distributes the effects fairly. ... The Shapley value is the only explanation method with a solid theory. (Molnar, 03/11/19)



K, S

### what if we break symmetry?

K, S, Z

K, S, Z

K. S. 7

, S, Z

### is this just obscure theory? Probabilistic values for games

books.google.com

The study of methods for measuring the "value" of playing a particular role in an nperson game is motivated by several considerations. One is to determine an equitable distribution of the wealth available to the players through their participation in the game. Another is to help an individual assess his prospects from participation in the game. When a method of valuation is used to determine equitable distribu-tions, a natural defining property is "efficiency": The sum of the individual values should equal the total payoff achieved through ....

☆ Cited by 546 Related articles ≫

#### Variations on the Shapley value

D Monderer, D Samet - Handbook of game theory with economic ..., 2002 - Elsevier

This survey captures the main contributions in the area described by the title that were published up to 1997.(Unfortunately, it does not capture all of them.) The variations that are the subject of this chapter are those axiomatically characterized solutions which are ...

☆ Cited by 85 Related articles - **>>>** 

### **Colin Rowat**

<u>к, s, z</u>

C

р





#### VERSITYOF BIRMINGHAM

# don't assume equal arrival probabilities

RJ Weber - The Shapley Value. Essays in Honor of Lloyd S ..., 1988 -



TIVAL

# Correlation isn't causation

0

#### **Colin Rowat**



#### VERSITYOF BIRMINGHAM

#### e.g. hotel occupancy data price ∝ occupancy rates hike prices to increase profits?

#### use DAGs to set Shapley arrival order? let causes arrive before effects



FESTIVAL

# explaining high earnings





#### UNIVERSITY<sup>OF</sup> BIRMINGHAM

# causal ancestors age, sex, native country, race

causal descendants
marital-status, education...

some values seen as less importante.g. marital-status

others seen as more importante.g. sex (largest), race

Frye, Feige & Rowat (2019)



# conclusions

Wikipedia Commons

**Colin Rowat** 

... about explainable AI/ML more sensible explanations 1. 2. what are the right axioms? 3. single number search misguided? e.g. Wasserstein, Schirm, Lazar (2019) on p < 0.05 what about allowing joint arrival?

4.



#### IVERSITYOF BIRMINGHAM

... about ideas and tools sometimes make unexpected leaps