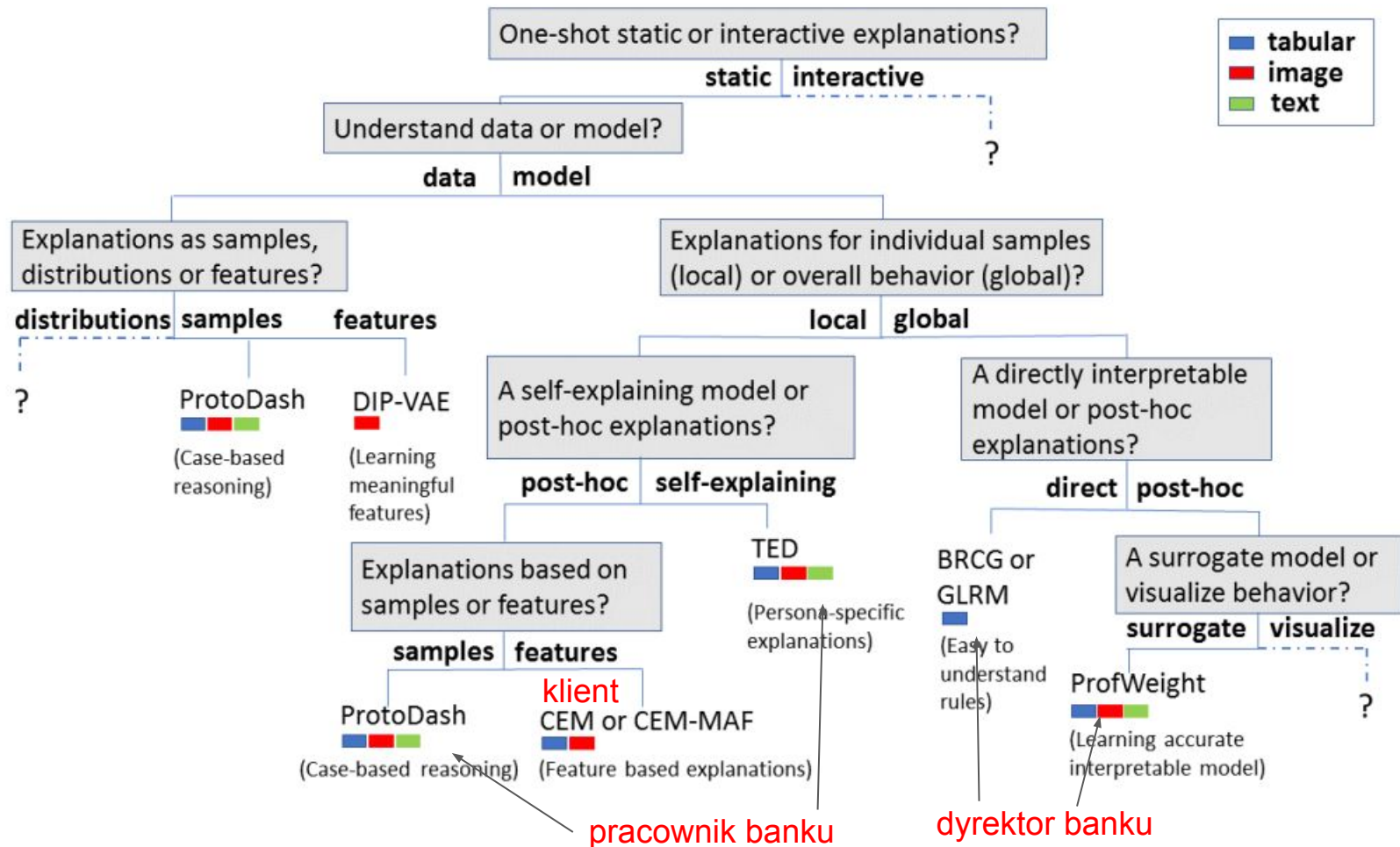


One Explanation Does Not Fit All: A Toolkit and Taxonomy of AI Explainability Techniques



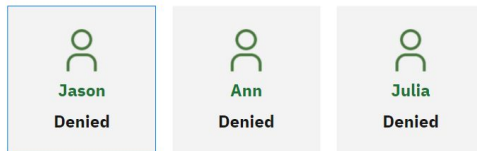
CEM

- SHAP, LIME

-

$$\theta_j^{pp} = 1 - \exp \frac{-|x_{[j]}^{pp}|}{\sigma_j}$$

$$\theta_j^{pn} = 1 - \exp \frac{-|x_{[j]} - x_{[j]}^{pn}|}{\sigma_j}$$



Several features in Jason's application fall outside the acceptable range. All would need to improve before acceptance was recommended.

Factors contributing to Jason's application denial

1. The value of **Consolidated risk markers** is **65**. It needs to be around **72** for the application to be approved.
2. The value of **Average age of accounts in months** is **52**. It needs to be around **68** for the application to be approved.
3. The value of **Months since most recent credit inquiry not within the last 7 days** is **2**. It needs to be around **3** for the application to be approved.

ProtoDash i TED (Teaching Explanations for Decisions)

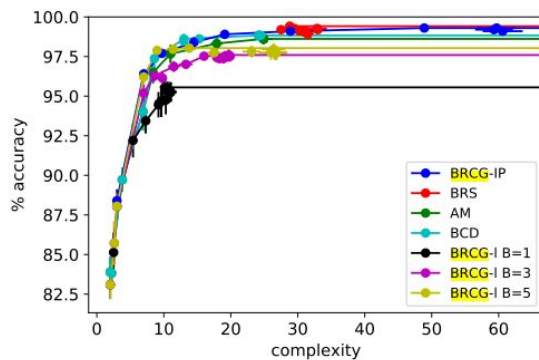
	Robert	James	Danielle	Franklin
Outcome	-	Defaulted	Defaulted	Defaulted
Similarity to Robert (from 0 to 1)	-	0.690	0.114	0.108
ExternalRiskEstimate	78	71	72	69
MSinceOldestTradeOpen	82	95	166	193
MSinceMostRecentTradeOpen	5	1	12	12
AverageMInFile	54	43	74	167
NumSatisfactoryTrades	33	33	37	36
NumTrades60Ever2DerogPubRec	0	0	1	0
NumTrades90Ever2DerogPubRec	0	0	1	0
PercentTradesNeverDelq	100	100	95	100
MSinceMostRecentDelq	0	0	7	0
MaxDelq2PublicRecLast12M	7	7	4	7
MaxDelqEver	8	8	4	8
NumTotalTrades	41	41	41	8
NumTradesOpeninLast12M	2	4	0	0
PercentInstallTrades	15	17	15	6
MSinceMostRecentInqexcl7days	0	0	0	0
NumInqLast6M	3	4	1	0
NumInqLast6				
NetFractionRc				

TED:

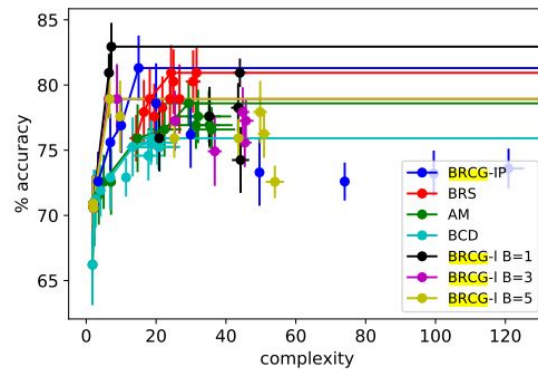
- tworzenie reguł wyjaśnienia przez specjalistów,
- model zwraca predykcję i wyjaśnienie swojej predykcji.

$$\theta_j = \exp \frac{-|x_{[j]} - x'_{[j]}|}{\sigma_j}$$

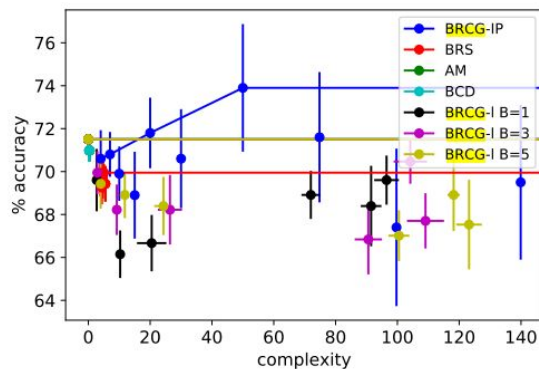
ProfWeight i BRCG



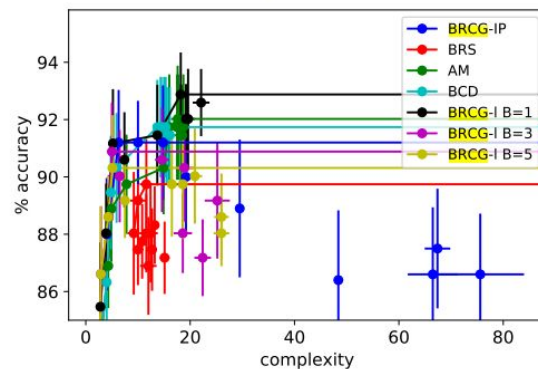
(a) banknote



(b) heart



(c) ILPD



(d) ionosphere

Co nam daje GLRM? (Generalized Linear Rule Models)

Odpowiada na pytania:

- jaka jest ogólna logika modelu przy podejmowaniu decyzji?
- czy ta logika jest rozsądna abyśmy z pewnością mogli ją wdrożyć?

AIX360 - metryki do oceniania wyjaśnień

Faithfulness

Monotonicity

$$\phi = -\rho(\theta, \mathbf{p})$$