

Yapay Öğrenmede Yorumlanabilirlik

İlker Birbil

<https://sibirbil.github.io>



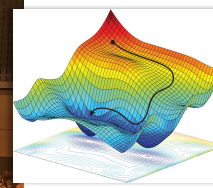
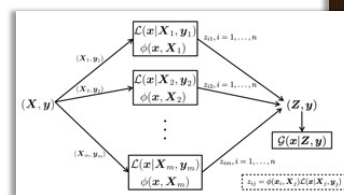
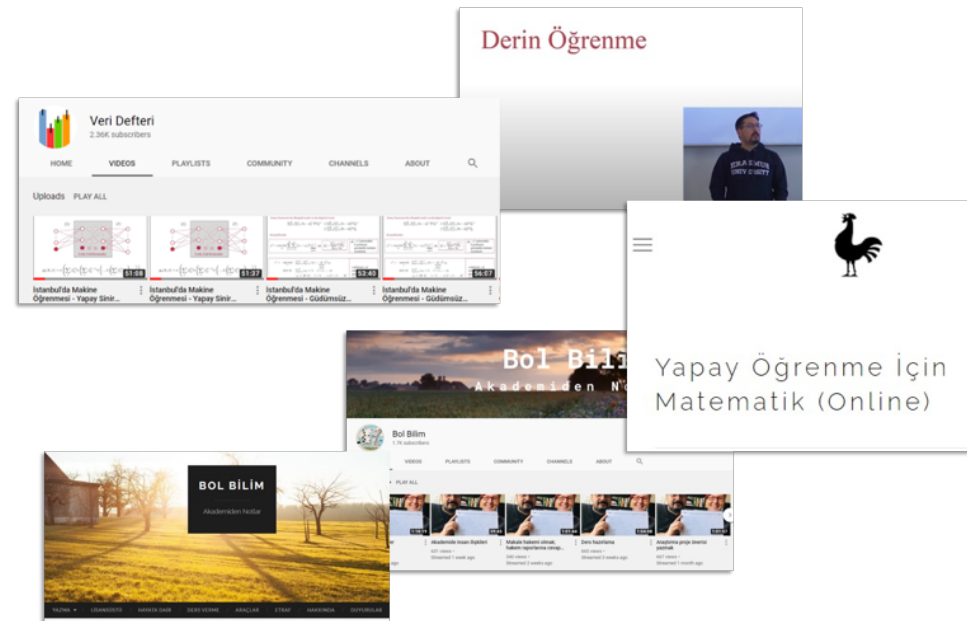
UNIVERSITEIT VAN AMSTERDAM



UvA
 d4c D4C
 Bol Bilim
 Veri Defteri
 @sibirbil



Veri Bilimi ve Optimizasyon



Robust and Fast Stochastic Gradient Decent with Model Building

```

Algorithm 1: SMB: Stochastic Model Building
1 Input:  $x_1 \in \mathbb{R}^n, f, g: \mathbb{R} \rightarrow \mathbb{R}$ , stepsize  $\{\alpha_k\}_{k=1}^N, c > 0$ 
2 for  $k = 1, \dots, N$  do
3    $s_k^i = -\alpha_k g_i$ 
4    $x_k^i = x_k + s_k^i, f_k^i = f(x_k^i, \xi_k), g_k^i = g(x_k^i, \xi_k)$ 
5   if  $f_k^i \leq f_k - c \cdot \alpha_k \|g_k^i\|^2$  then
6      $x_{k+1}^i = x_k^i, f_{k+1}^i = f_k^i, g_{k+1}^i = g_k^i$ 
7   else
8     for each parameter group  $p$  do
9        $y_{k,p} = g_{k,p} - g_{k,p}$ 
10       $s_{k,p} = c_{k,p} g(y_{k,p}) + c_{k,p} g(y_{k,p}) + c_{k,p} g(y_{k,p})$ , as defined in (4)
11       $x_{k+1} = x_k + s_k$ , where  $s_k = (s_{k,p_1}, \dots, s_{k,p_n})$  and  $n$  is the number of parameter groups;
12       $f_{k+1} = f(x_{k+1}, \xi_k), g_{k+1} = g(x_{k+1}, \xi_k)$ 
  
```

Data Privacy in Bid-Price Control for Network Revenue Management

Discovering Classification Rules for Interpretable Learning with Linear Programming

This work has been published in Computational Optimization and Applications. You can acc
 @hamsi-mf: HAMS! for matrix factorization



The Washington Post
Democracy Dies in Darkness

‘Creative ... motivating’ and fired



Sarah Wysocki was out of work for only a few days after she was fired by DCPS last year. She is now teaching at Hybla Valley Elementary School in Fairfax County. (Jahi Chikwendiu/The Washington Post)

By **Bill Turque**
 March 6, 2012

HDSR

Why Are We Using Black Box Models in AI When We Don't Need To? A Lesson From An Explainable AI Competition

by *Cynthia Rudin and Joanna Radin*

Published on Nov 22, 2019

 **CGAP**

BLOG 05 SEPTEMBER 2019

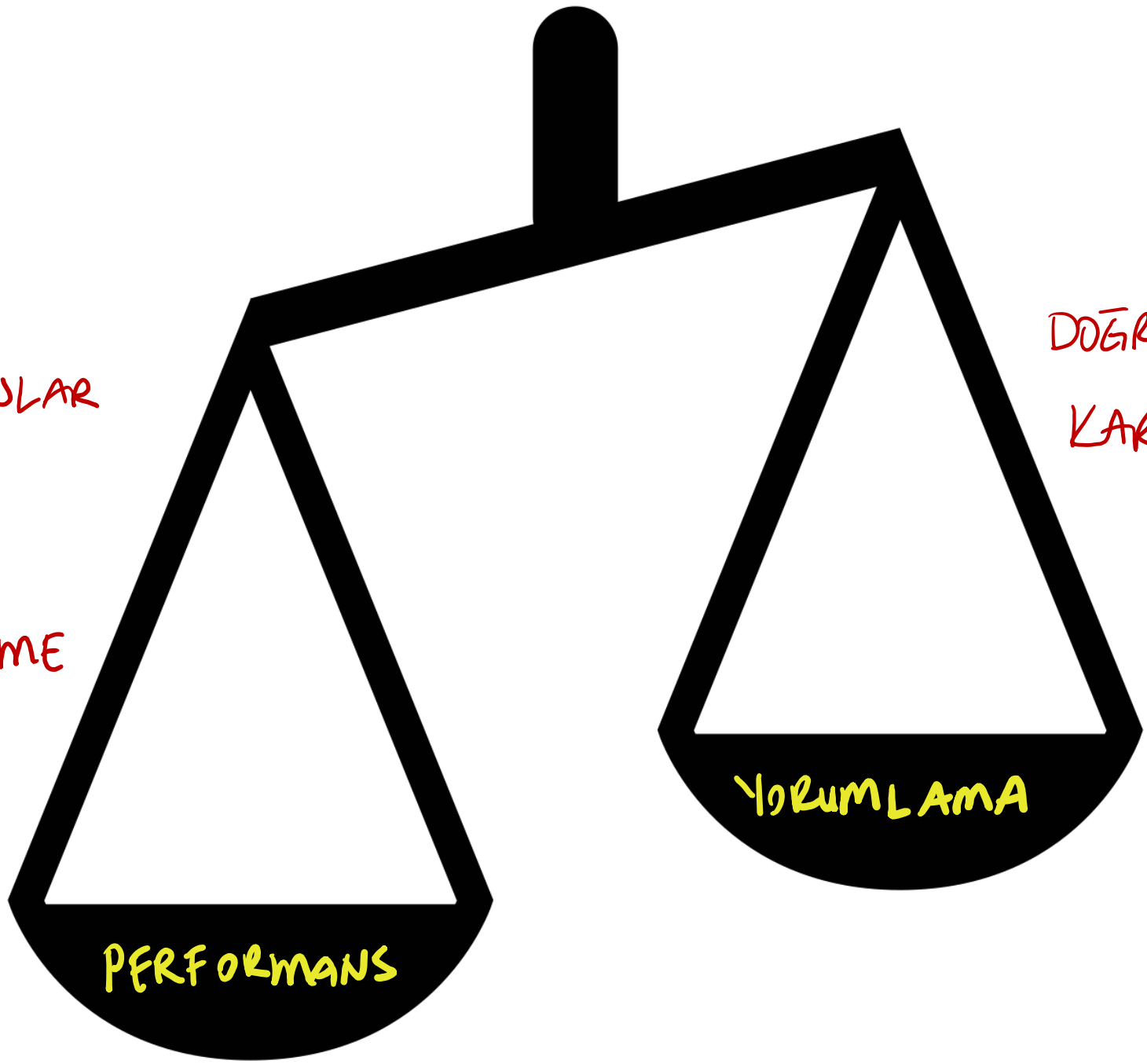
Algorithm Bias in Credit Scoring: What's Inside the Black Box?

By *Maria Fernandez Vidal, Jacobo Menajovsky*

”

The responsible use of algorithms requires providers to know which variables are being considered in their credit scoring models and how they are affecting people's scores.

RASSAL ORMANLAR
XGBOOST
DUM
DERİN ÖĞRENME



DOĞRUSAL BAĞLANIM
KARAR AĞAÇLARI



PERFORMANS

YORUMLAMA

YORUMLAMA YAKLAŞIMLARI

LOKAL ●
GLOBAL ●

ÖZEL ●
GENEL* ●

KENDİLİĞİNDEN ●
DIŞARIDAN ●

LIME ● ● ●

SHAP ● ● ●

SLIM ● ● ●

○
○
○

OCT ● ● ●

EBM ● ● ●

OSDT ● ● ●

○
○
○

BDR ● ● ●

RUX ● ● ●

RUG ● ● ●

○
○
○

● KARAR AĞACI TABANLI
● KURAL TABANLI

* AGNOSTİK

LIME ● ● ●
(RIBERIO VD., 2016)

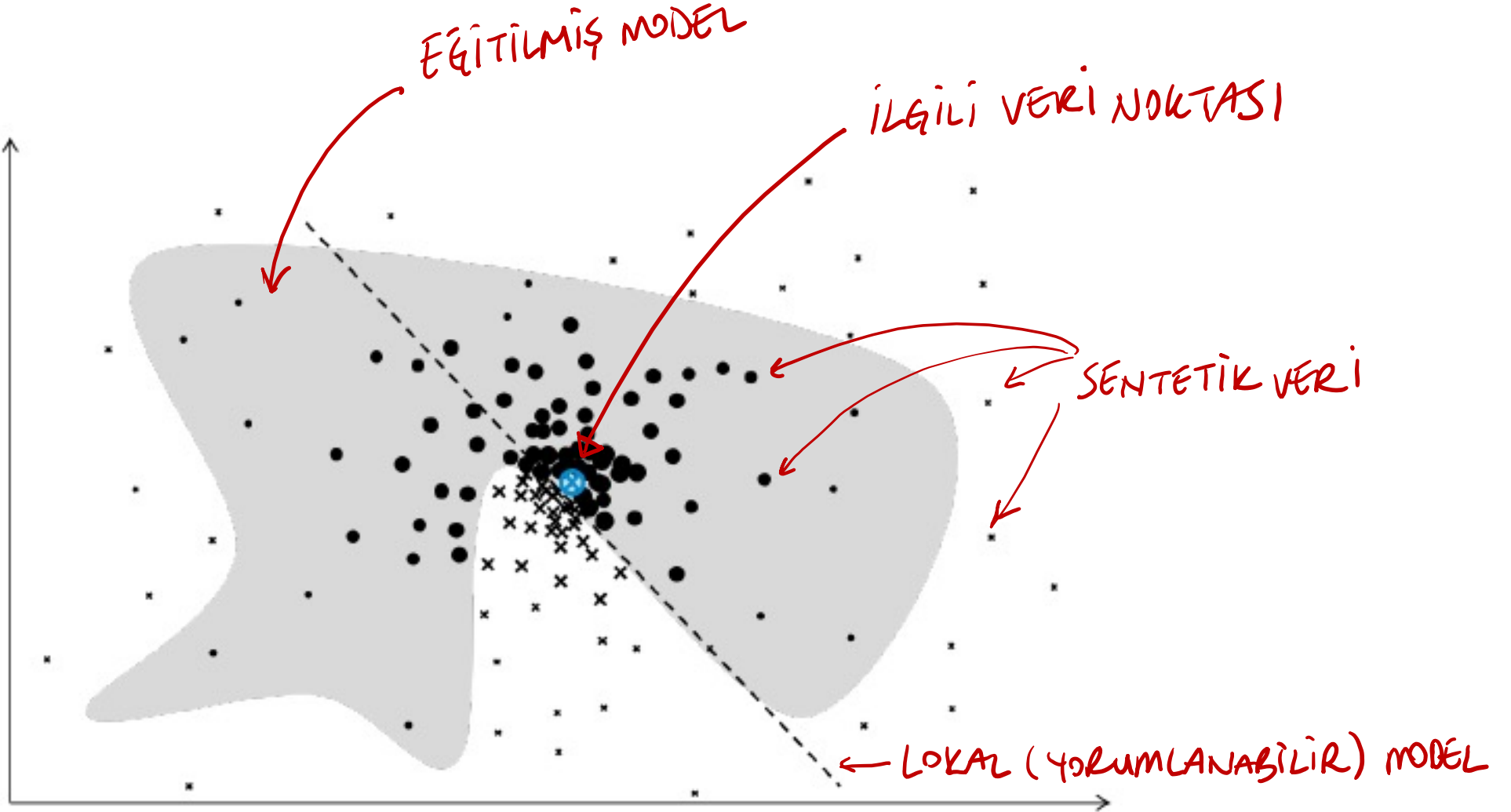
LOCAL INTERPRETABLE MODEL-AGNOSTIC EXPLANATIONS



MODEL
ETİKETLERİ
↓

● × } AĞIRLIKLAR
· x

↑
MESAFE İLE
TERS ORANTILI



LIME

AÇIKLAMA

$$\xi(x) = \arg \min_{g \in G} L(f, g, \pi_x) + \Omega(g)$$

İLGİLİ
VERİ
NOKTASI

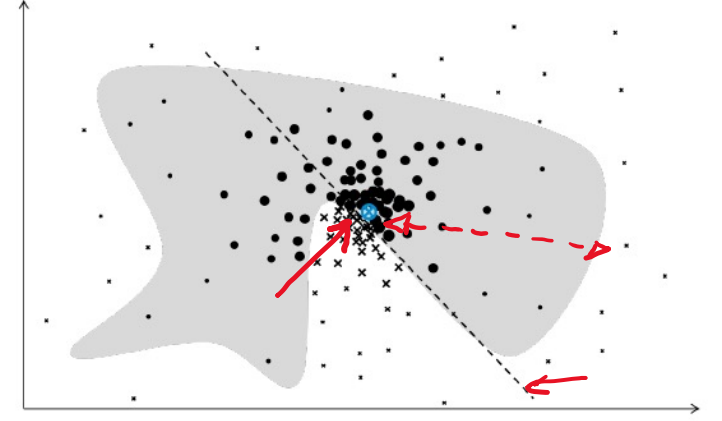
ÖĞRENME
HATASI

UZAKLIK

LOKAL
MODEL

LOKAL MODEL
KARMAŞIKLIĞI

- ÖZGİTİM SAYISI
- AĞAÇ DERİNLİĞİ
- ...



XGBoost model predicted an acceptance probability of 77.0%

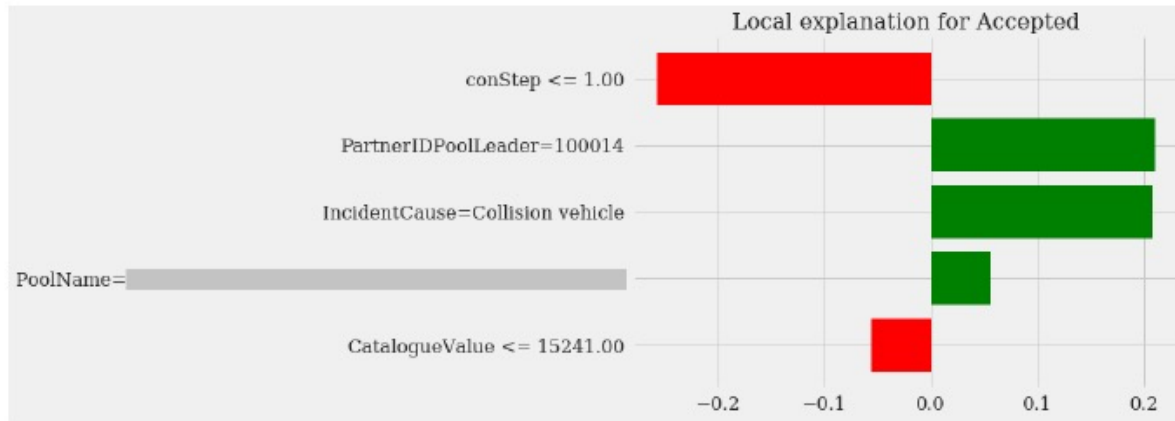


Figure 23: FP - LIME by Default (5 Features)

[Local Pred. = 0.58, Intercept = 0.42, $R^2 = 0.28$, RMSE = 0.26, Time = 4.67 s]

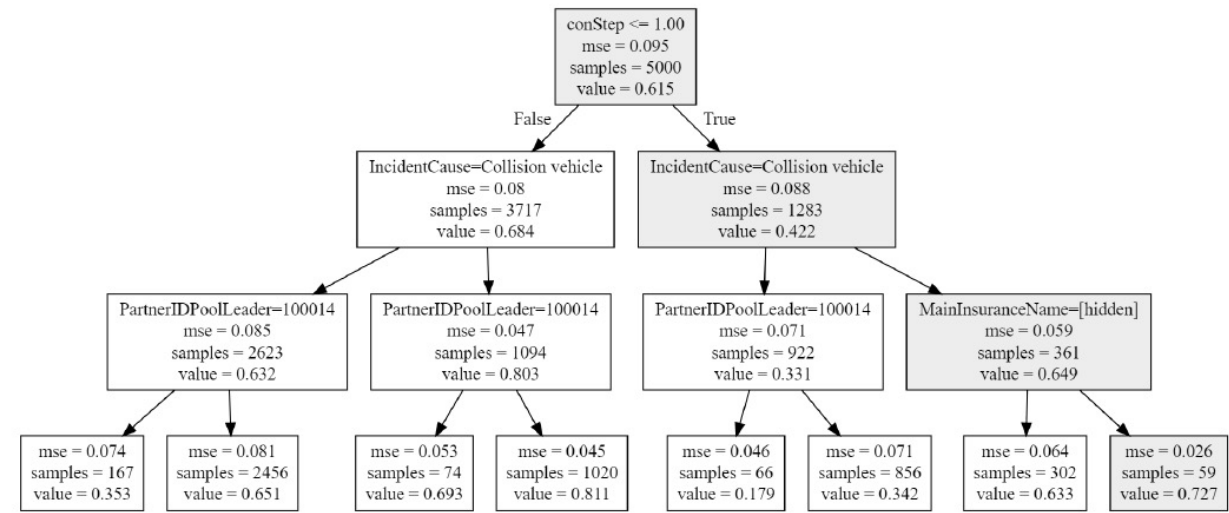


Figure 24: FP - LIME Decision Tree (3 Layers)

[Local Pred. = 0.73, $R^2 = 0.28$, RMSE = 0.26, Time = 2.15 s]

YORUMLAMA YAKLAŞIMLARI

LOKAL ●
GLOBAL ●

ÖZEL ●
GENEL* ●

KENDİLİĞİNDEN ●
DIŞARIDAN ●

LIME ● ● ●
SHAP ● ● ●
SLIM ● ● ●
○
○
○

OCT ● ● ●
EBM ● ● ●
OSDT ● ● ●
○
○
○

BDR ● ● ●
RUX ● ● ●
RUG ● ● ●
○
○
○

● KARAR AĞACI TABANLI
● KURAL TABANLI

* AGNOSTİK



SHAP ● ● ●

(LUNDBERG VE LEE, 2017)

SHAPLEY ADDITIVE EXPLANATIONS

AÇIKLAMA
MODELİ

$$g_i(z') = \phi_0 + \sum_{j=1}^M \phi_{ij} z'_{ij}$$

↑
ÖRNEK

↑
ÖZNETELİK

ÖZNETELİK BU ÖRNEKTE
VAR (1), YOK (0)

SHAPLEY
DEĞERİ

$$\phi_{ij} = \sum_{S \subseteq X_i \setminus \{j\}} \frac{|S|!(M - |S| - 1)!}{M!} [f_i(S \cup \{j\}) - f_i(S)]$$

↑
j HARİÇ
ÖZNETELİKLER

↑
MODEL
TAHMİNİ

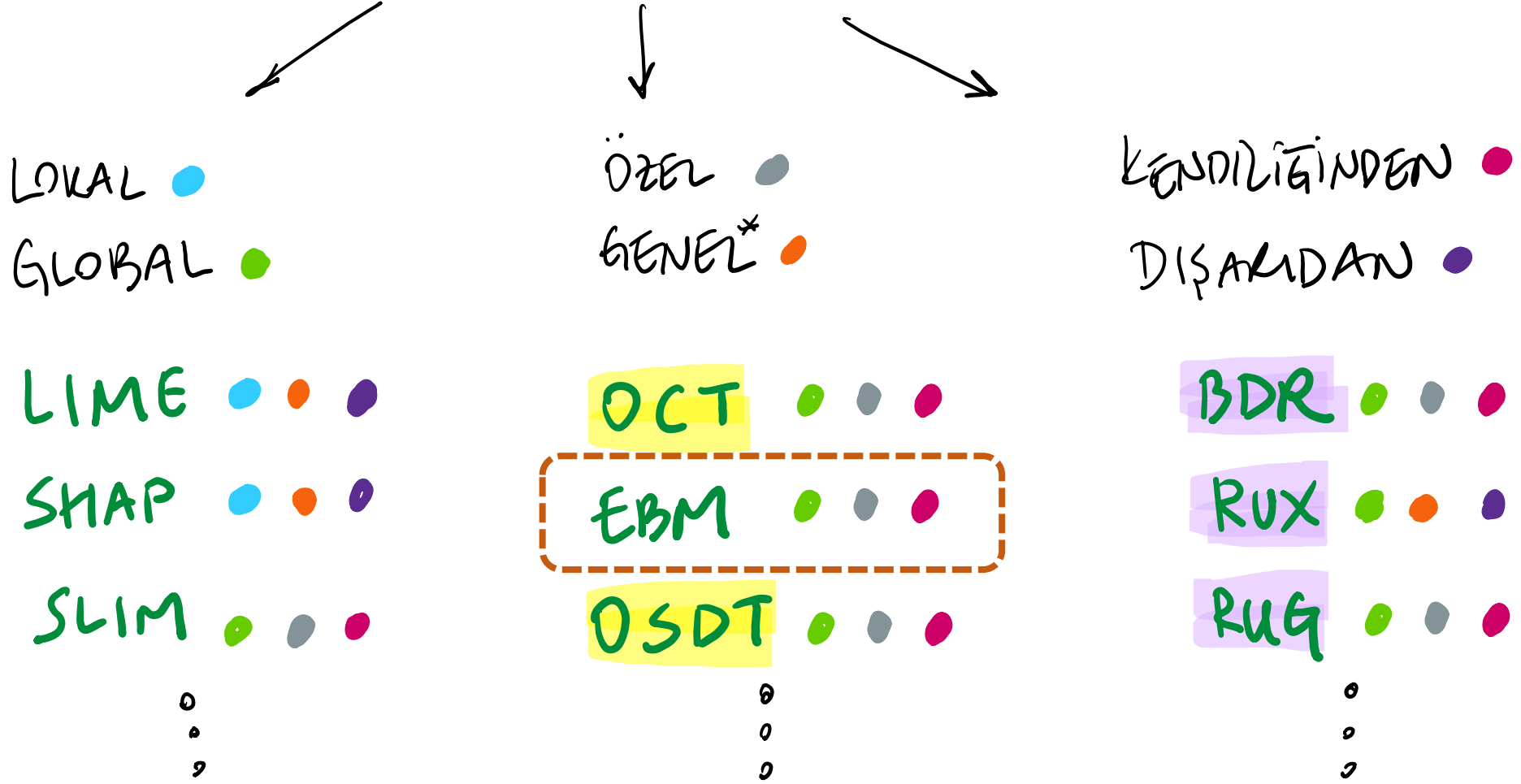
$$\phi_{ij} = \sum_{S \subseteq X_i \setminus \{j\}} \frac{|S|!(M - |S| - 1)!}{M!} [f_i(S \cup \{j\}) - f_i(S)]$$

Tüm
ALT KÜMELER
 2^M

YAKLAŞIKLAMA YÖNTEMLERİ

- KERNEL SHAP*
- LINEAR SHAP
- LOW-ORDER SHAP
- TREE SHAP
- DEEP SHAP

YORUMLAMA YAKLAŞIMLARI



* AGNOSTİK

KARAR AĞACI TABANLI
 KURAL TABANLI

EBM ● ● ●
(NDR1 v0.1 2019)

EXPLAINABLE BOOSTING MACHINE



“ÖRNEK İ İÇİN TAHMİN”

$$g(E[y_i]) = \beta_0 + \sum_{j=1}^M f_j(x_{ij}) + \sum_{j=1}^M \sum_{k \neq j} f_{jk}(x_{ij}, x_{ik})$$

ÖZNİTELİKLER ARASI İLİŞKİLER

GENELLEŞTİRİLMİŞ
DOĞRUSAL
MODEL

LINK
FONKSİYONU
(ÖRN. LOGIT)

ÖZNİTELİK
FONKSİYONLARI

TAKVİYE
(BOOSTING)
MODELLERİ

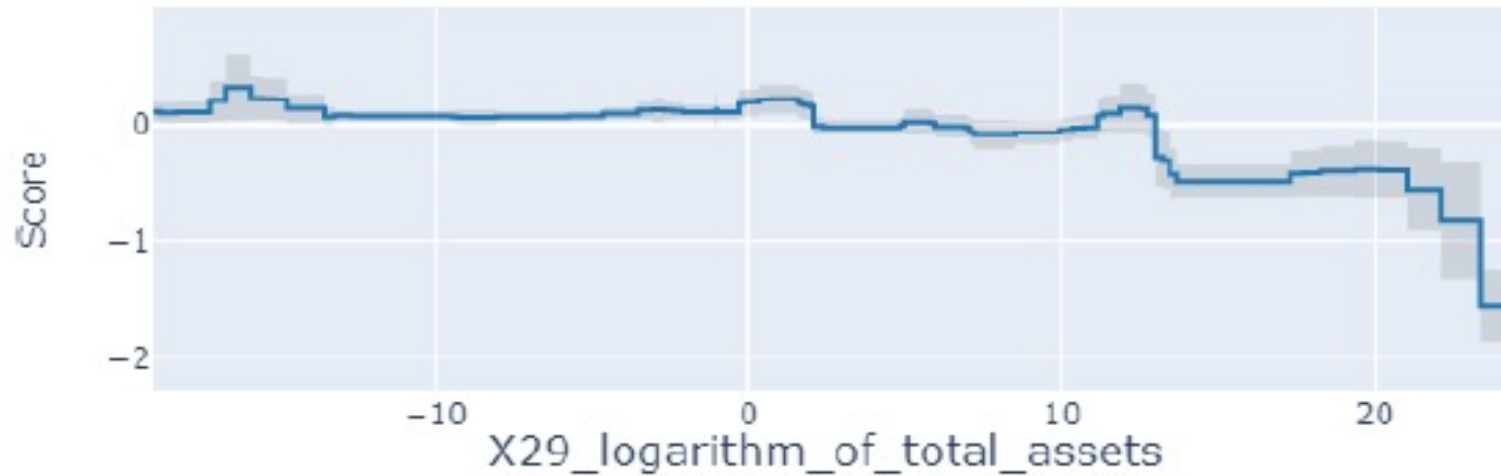


Figure 1: GAM plot of the effect of a feature $\log(\text{total assets})$ on the prediction of bankruptcy for a data set of Polish firms. There is quite some variance in the relation, but overall higher assets clearly decrease the bankruptcy risk as would be expected.

YORUMLAMA YAKLAŞIMLARI

↓

LOKAL ●
GLOBAL ●

↓

ÖZEL ●
GENEL* ●

↓

KENDİLİĞİNDEN ●
DIŞARIDAN ●

LIME ● ● ●
SHAP ● ● ●

SLIM ● ● ●

○
○
○

OCT ● ● ●
EBM ● ● ●
OSDT ● ● ●

○
○
○

BDR ● ● ●
RUX ● ● ●
RUG ● ● ●

○
○
○

● KARAR AĞACI TABANLI
● KURAL TABANLI

* AGNOSTİK



SLIM ● ● ●
 (üstün vd., 2013)

SUPERSPARSE LINEAR INTEGER MODEL

HATA

$$\min_{\lambda} \frac{1}{N} \sum_{i=1}^N \mathbf{1}[y_i x_i^T \lambda \leq 0] + C_0 \|\lambda\|_0 + C_1 \|\lambda\|_1,$$

s.t. $\lambda \in \mathcal{L}$

\uparrow
 KESİKLİ DEĞERLER (SKORLAR)
 $\{-1, 1\}$

$\underbrace{\hspace{10em}}$
 YORUMLANABİLİRLİK
 (AĞ SAHIDA "ÖZGÜTEK")

PREDICT MUSHROOM IS POISONOUS IF SCORE > 3

1.	<i>spore_print_color = green</i>	4 points	
2.	<i>stalk_surface_above_ring = grooves</i>	2 points	+
3.	<i>population = clustered</i>	2 points	+
4.	<i>gill_size = broad</i>	-2 points	+
5.	<i>odor ∈ {none, almond, anise}</i>	-4 points	+
ADD POINTS FROM ROWS 1-5		SCORE	=

Figure 2: The scoring system for mushroom edibility produced by SLIM as displayed in [Ustun and Rudin \(2015\)](#)

ANKET



14 VERİ ANALİZİ UZMANI
(11 KİŞİ YAPAY ÖĞRENME ALANINDA)

Table 9: Average scores of survey respondents on model specific questions as well as corresponding standard deviations (SD). Highest scores in bold.

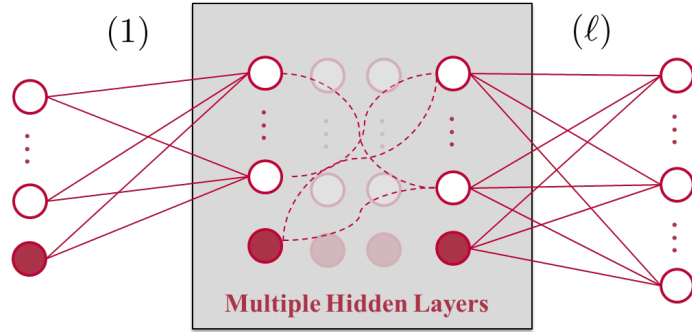
	Familiarity		Test	Understanding		Stakeholders		Practical Value	
	(1-3)	SD	% Correct	(1-7)	SD	(1-7)	SD	(1-7)	SD
Logit	3.00	0.00	78	5.21	1.25	4.57	1.50	5.36	0.93
SLIM	1.36	0.63	64	5.79	1.53	5.71	1.49	4.86	1.75
EBM	1.79	0.70	86	4.64	1.55	4.43	1.79	5.07	1.21
SHAP	2.36	0.63	71	5.43	1.09	5.00	1.24	5.86	0.77

XGBOOST +

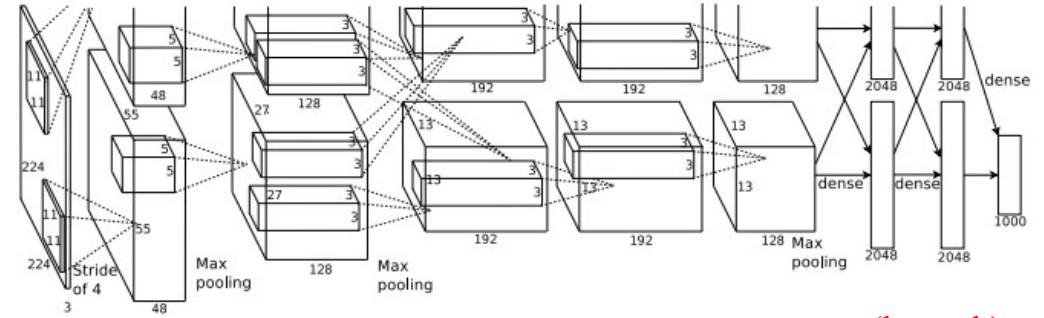
YÖNTEMLER İLGİLİ
DOĞRU-YANLIŞ SORULARI

- YÖNTEMLERİ ANLAMA (KİŞİSEL)
- YÖNTEMLERİ ANLATMA
- PRATİKTE KULLANIM ŞANSI

DERIN ÖĞRENME



$$\hat{y}_k(X, \beta) = \sigma \left(\sum_j \beta_{kj}^{(\ell)} h \left(\sum_s \beta_{js}^{(\ell-1)} h \left(\dots h \left(\sum_i \beta_{ji}^{(1)} X_i \right) \dots \right) \right) \right)$$



(kaynak)

LIME ● ● ●

SHAP ● ● ●

DİĞER ?



InterpretML: A Unified Framework for Machine Learning Interpretability

Harsha Nori
Samuel Jenkins
Paul Koch
Rich Caruana
Microsoft Corporation
1 Microsoft Way
Redmond, WA 98052, USA

HANORI@MICROSOFT.COM
SAJENKIN@MICROSOFT.COM
PAULKOCH@MICROSOFT.COM
RCARUANA@MICROSOFT.COM



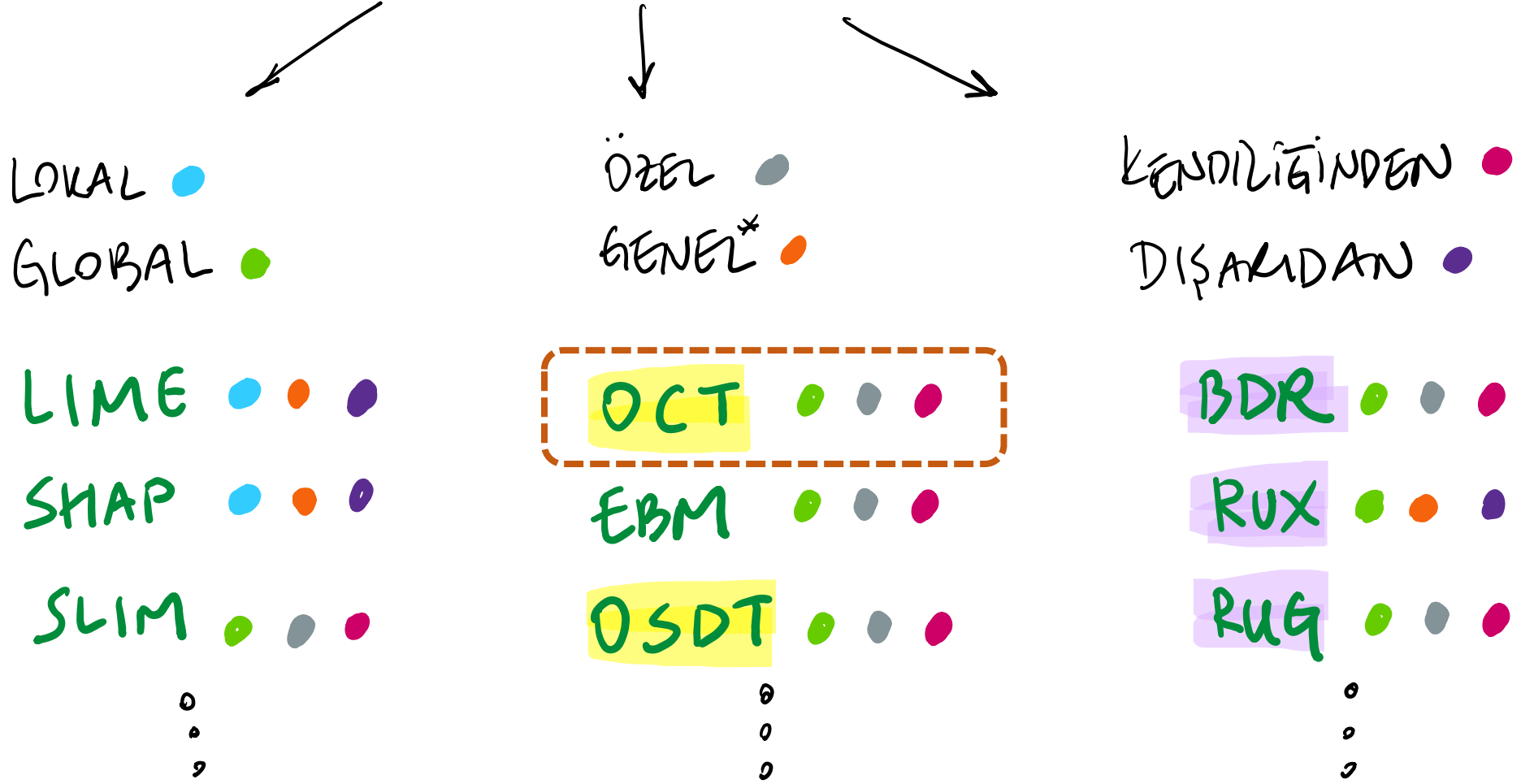
**Understand Models.
Build Responsibly.**

A toolkit to help understand models and enable responsible machine learning

Supported Techniques

Interpretability Technique	Type	Examples
Explainable Boosting	glassbox model	Notebooks
Decision Tree	glassbox model	Notebooks
Decision Rule List	glassbox model	Coming Soon
Linear/Logistic Regression	glassbox model	Notebooks
SHAP Kernel Explainer	blackbox explainer	Notebooks
SHAP Tree Explainer	blackbox explainer	Coming Soon
LIME	blackbox explainer	Notebooks
Morris Sensitivity Analysis	blackbox explainer	Notebooks
Partial Dependence	blackbox explainer	Notebooks

YORUMLAMA YAKLAŞIMLARI



* AGNOSTİK

● KARAR AĞACI TABANLI
● KURAL TABANLI

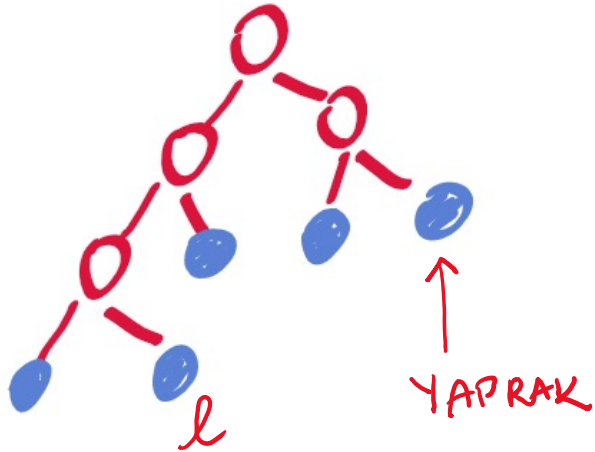
OCT ● ● ●

(BERTSIMAS VE DUNN, 2017)

OPTIMAL CLASSIFICATION TREES



(x_i, y_i)
↑ ↑
GİRİŞ ETİKET



SINIFLANDIRMA
HATASI

AĞAC BÜYÜKLÜĞÜ
(DÜĞÜM SAYISI)

$$\min R_{xy}(T) + \alpha |T|$$

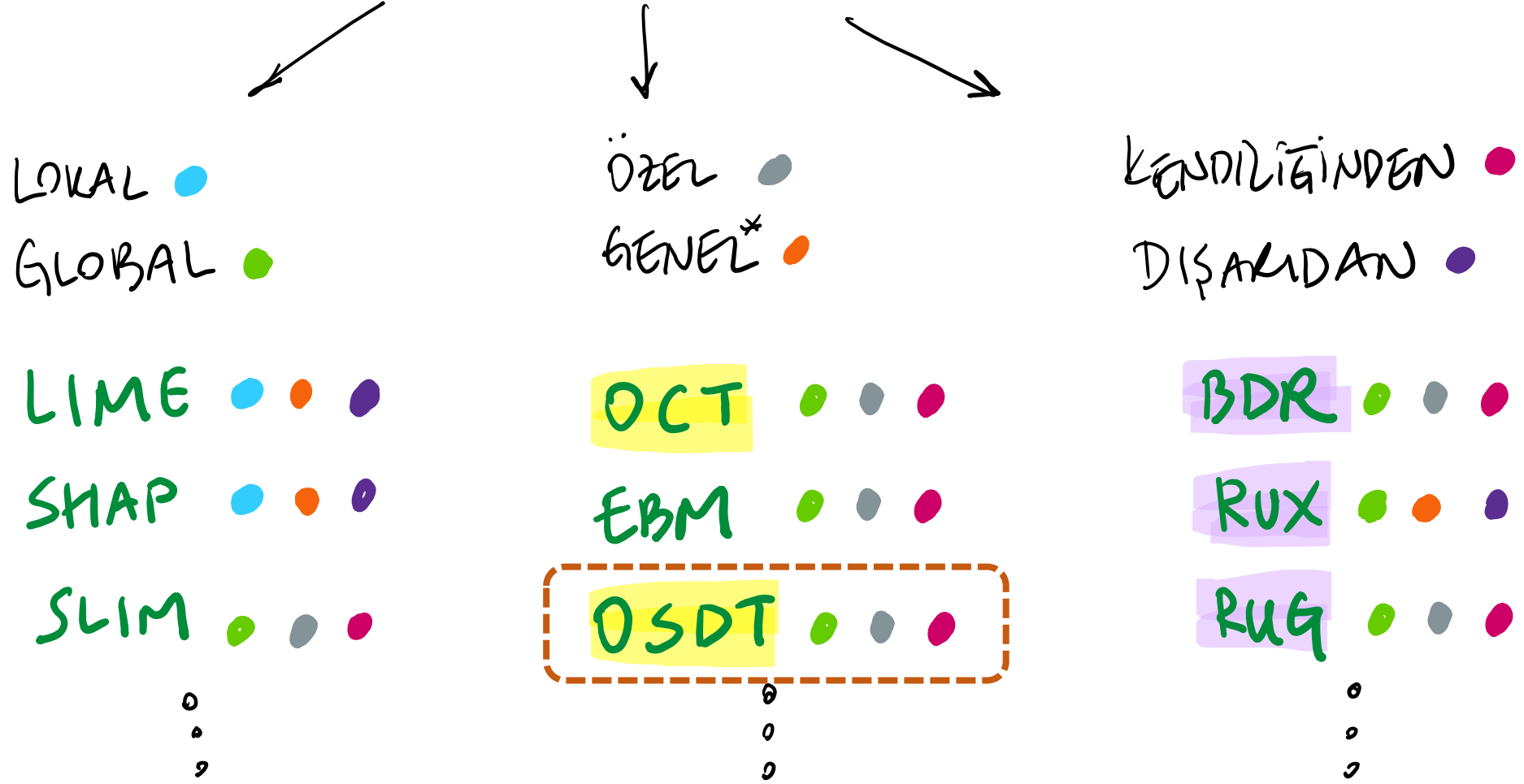
$$\text{s.t. } N_x(l) \geq N_{\min}$$



l YAPRAĞINDAKİ
ÖRNEK SAYISI

Table 12 Comparison of CART and OCT-H across a range of depths, showing the number of datasets for which each method had the highest out-of-sample accuracy, and the mean improvement in out-of-sample accuracy when using OCT-H across all datasets along with the p value indicating the statistical significance of this difference

Max. depth	CART wins	OCT-H wins	Ties	Accuracy improvement (%)	p value
1	3	36	14	5.12	$\sim 10^{-16}$
2	10	32	11	4.88	$\sim 10^{-14}$
3	13	31	9	3.59	$\sim 10^{-12}$
4	13	29	11	3.12	$\sim 10^{-11}$

YORUMLAMA YAKLAŞIMLARI



 KARAR AĞACI TABANLI
 KURAL TABANLI

* AGNOSTİK

OSDT ● ● ●

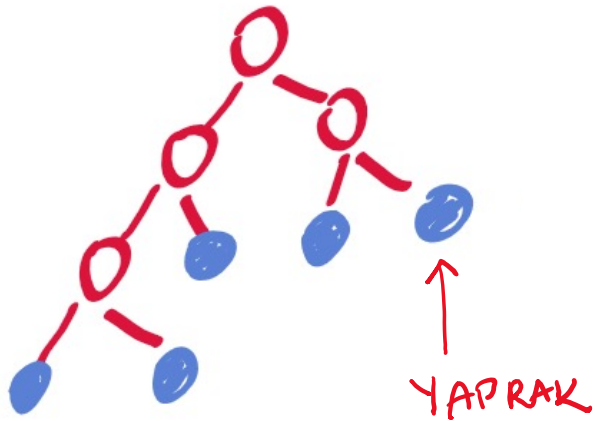
(HU VD., 2020)

OPTIMAL SPARSE DECISION TREES



(x_i, y_i)

↑
 $\{0,1\}$



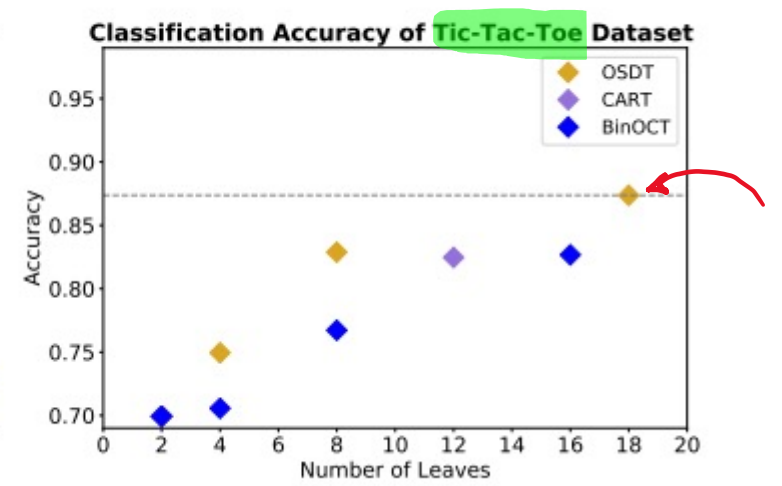
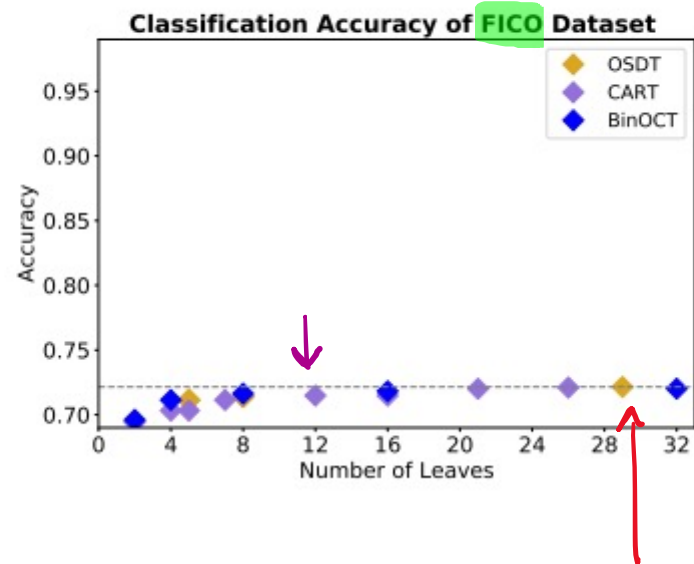
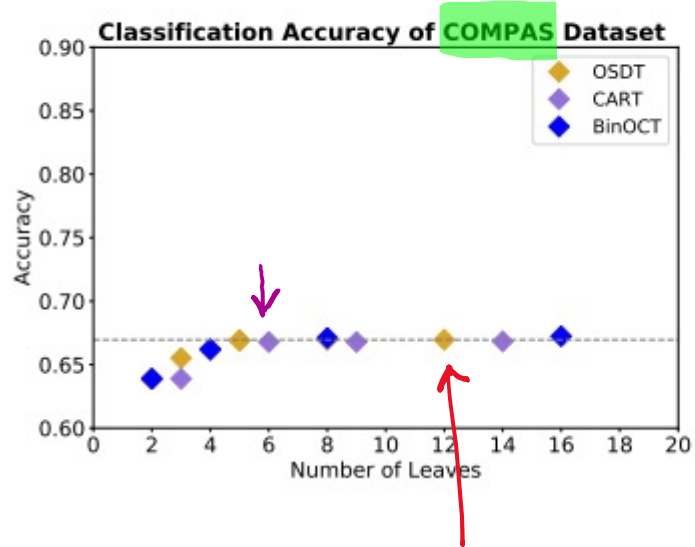
$$R(d, \mathbf{x}, \mathbf{y}) = \ell(d, \mathbf{x}, \mathbf{y}) + \lambda H_d$$

↑
AĞAÇ

↑
SINIFLANDIRMA
HATASI

↑
AĞAÇTAKİ
YAPRAK SAYISI

OSDT ● ● ●



SÜRE LİMİTİ 30 DK.

OSDT ● ● ●

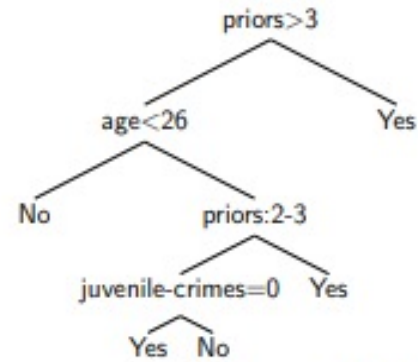


Figure 4: An optimal decision tree generated by OSDT on the COMPAS dataset. ($\lambda = 0.005$, accuracy: 66.90%)

YORUMLAMA YAKLAŞIMLARI



LOKAL ●
GLOBAL ●

LIME ● ● ●

SHAP ● ● ●

SLIM ● ● ●

○
○
○



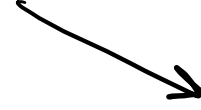
ÖZEL ●
GENEL* ●

OCT ● ● ●

EBM ● ● ●

OSDT ● ● ●

○
○
○



KENDİLİĞİNDEN ●
DIŞARIDAN ●

BDR ● ● ●

RUX ● ● ●

RUG ● ● ●

○
○
○

● KARAR AĞACI TABANLI
● KURAL TABANLI

* AGNOSTİK

BDR ● ● ●
(DASH VD., 2018)

BOOLEAN DECISION RULES



(x_i, y_i)
 \uparrow
 $\{0,1\}$

minimize $\sum_{i \in \mathcal{P}} \xi_i + \sum_{i \in \mathcal{Z}} \sum_{k \in \mathcal{K}_i} w_k$

TOPLAM SEÇILEN KURALLAR SAHISI

subject to $\xi_i + \sum_{k \in \mathcal{K}_i} w_k \geq 1, \quad \xi_i \geq 0, \quad i \in \mathcal{P}$

SINIFLANDIRMA KURALLARI

$$\sum_{k \in \mathcal{K}} c_k w_k \leq C$$

$$w_k \in \{0, 1\}, \quad k \in \mathcal{K}.$$

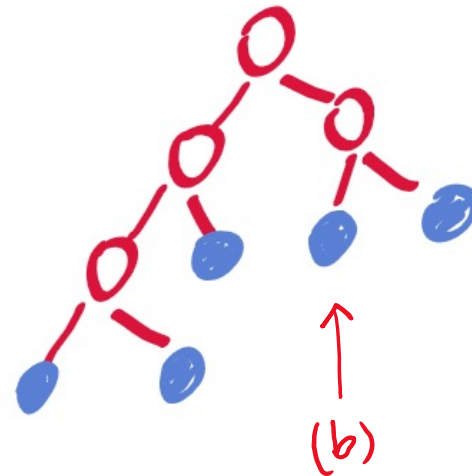
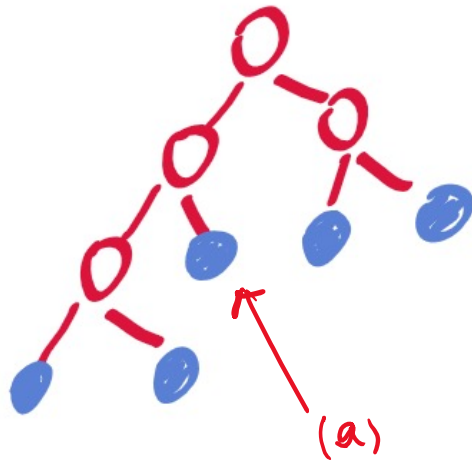
BDR ● ● ●



$$(\text{NumSatTrades} \geq 23) \wedge (\text{ExtRiskEstimate} \geq 70) \wedge (\text{NetFracRevolvBurden} \leq 63) \quad (a)$$

OR

$$(\text{NumSatTrades} \leq 22) \wedge (\text{ExtRiskEstimate} \geq 76) \wedge (\text{NetFracRevolvBurden} \leq 78) \quad (b)$$



YORUMLAMA YAKLAŞIMLARI



LOKAL ●
GLOBAL ●

LIME ● ● ●

SHAP ● ● ●

SLIM ● ● ●

○
○
○



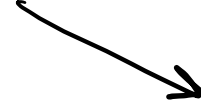
ÖZEL ●
GENEL* ●

OCT ● ● ●

EBM ● ● ●

OSDT ● ● ●

○
○
○



KENDİLİĞİNDEN ●
DIŞARIDAN ●

BDR ● ● ●

RUX ● ● ●

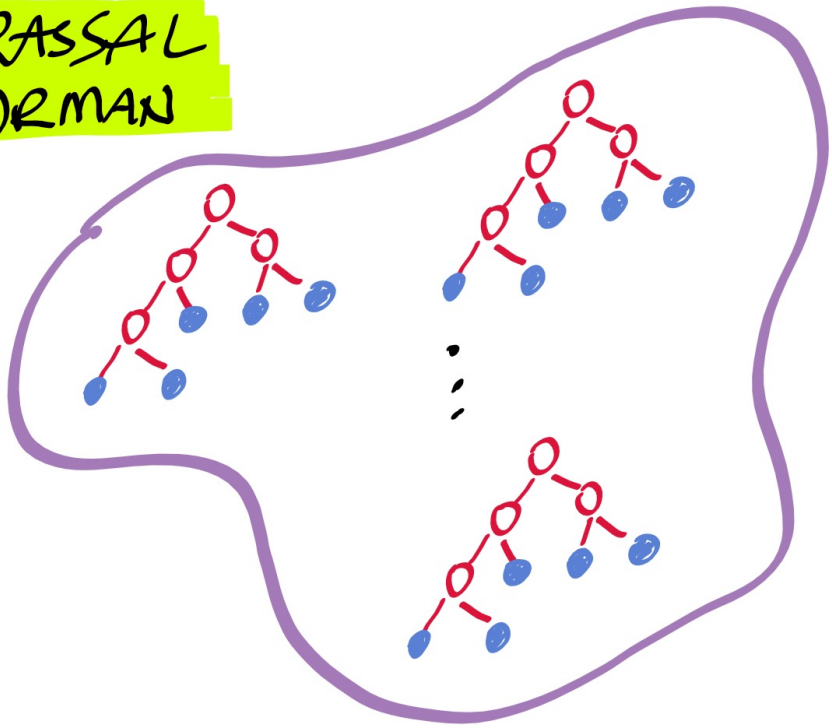
RUG ● ● ●

○
○
○

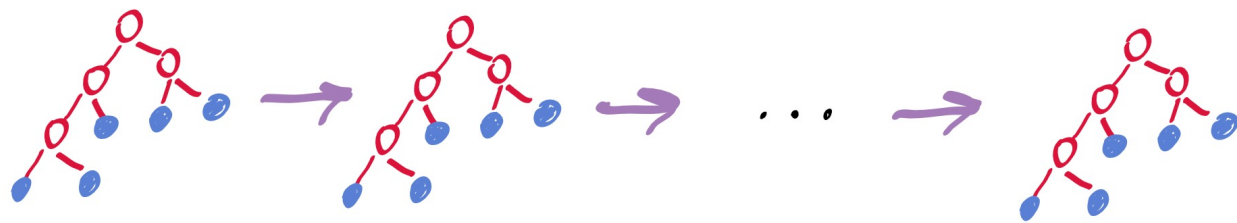
● KARAR AĞACI TABANLI
● KURAL TABANLI

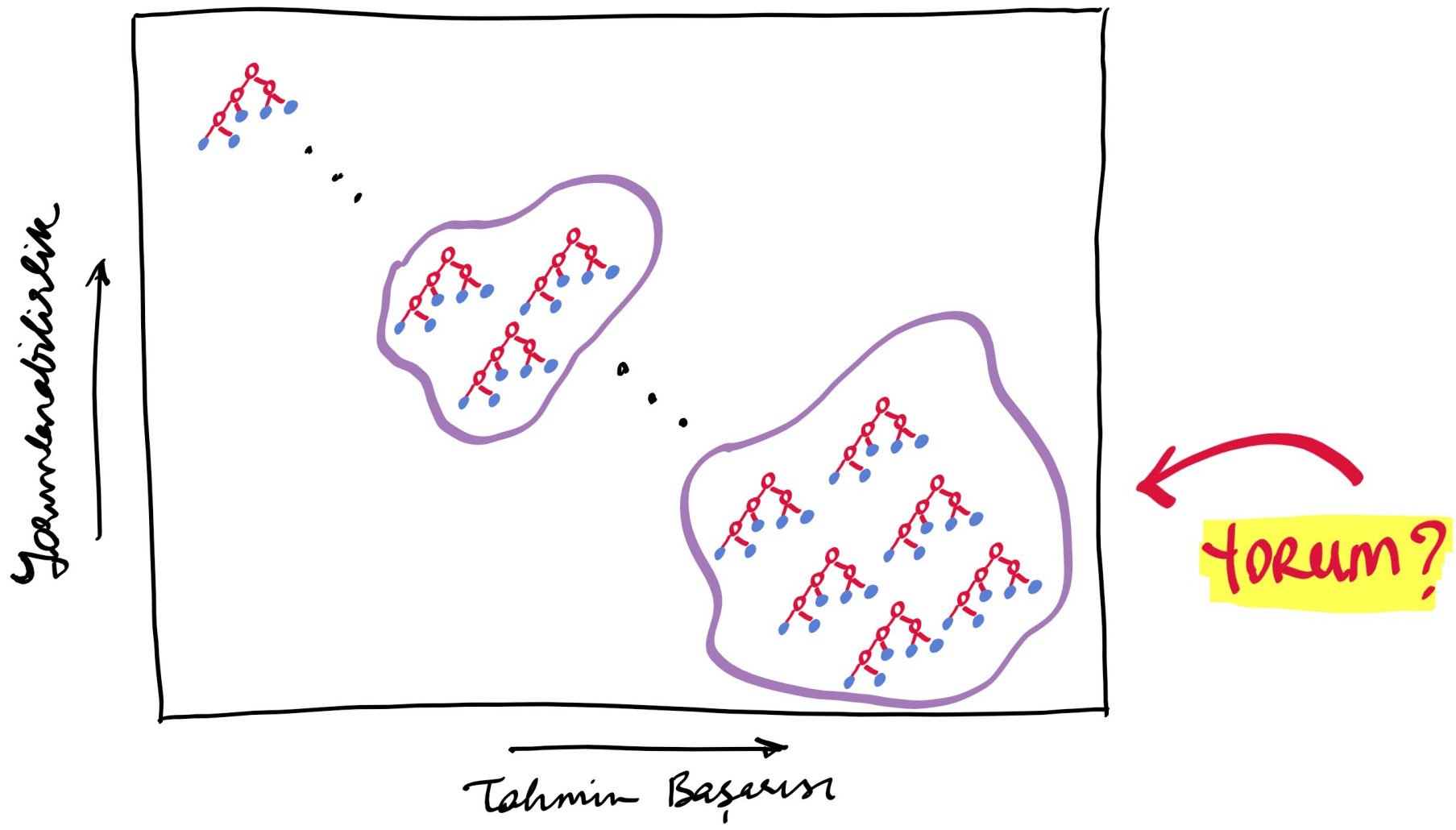
* AGNOSTİK

RASSAL
ORMAN



TAKVIYE





RUX ● ● ●

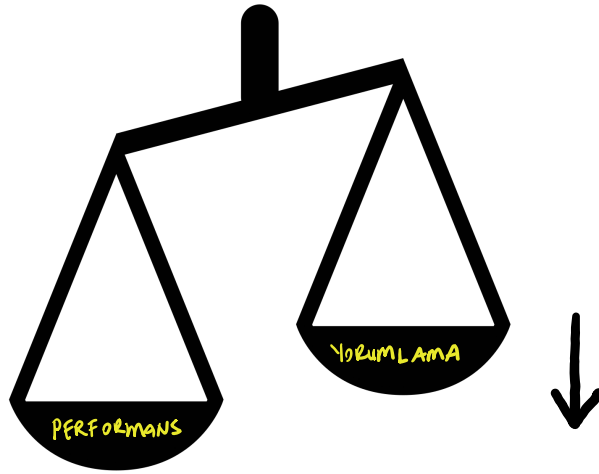
RUG ●



Discovering Classification Rules for Interpretable Learning with Linear Programming

Hakan Akyüz and Ş. İlker Birbil

<https://github.com/sibirbil/RuleDiscovery>



BASIT

HIZLI

KODLAMASI KOLAY

LP

minimize

subject to

$$\sum_{i \in \mathcal{I}} v_i + \sum_{j \in \mathcal{J}} c_j w_j$$

$$\sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i \geq 1, \quad i \in \mathcal{I};$$

$$\sum_{j \in \mathcal{J}} a_{ij} w_j \geq 1, \quad i \in \mathcal{I};$$

$$v_i \geq 0, \quad i \in \mathcal{I};$$

$$w_j \geq 0, \quad j \in \mathcal{J}$$

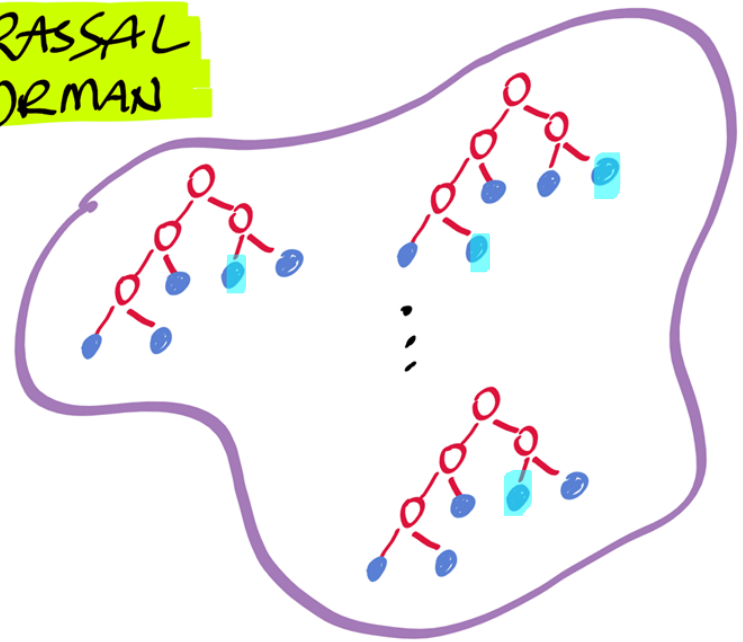
↓
TOPLAM
KAYIP

SEYREKLİK ✓
KURAL UZUNLUĞU ✓
SINIFLANDIRICI AĞIRLIKLARI ✓
YANLIŞ POZİTİF ETİKET SAYISI ✓
...
!

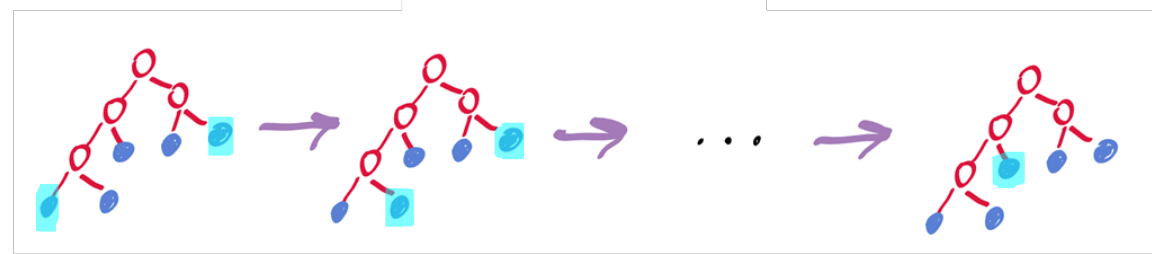
← KAPSAMA
KISITLARI

P(J)

RASSAL
ORMAN



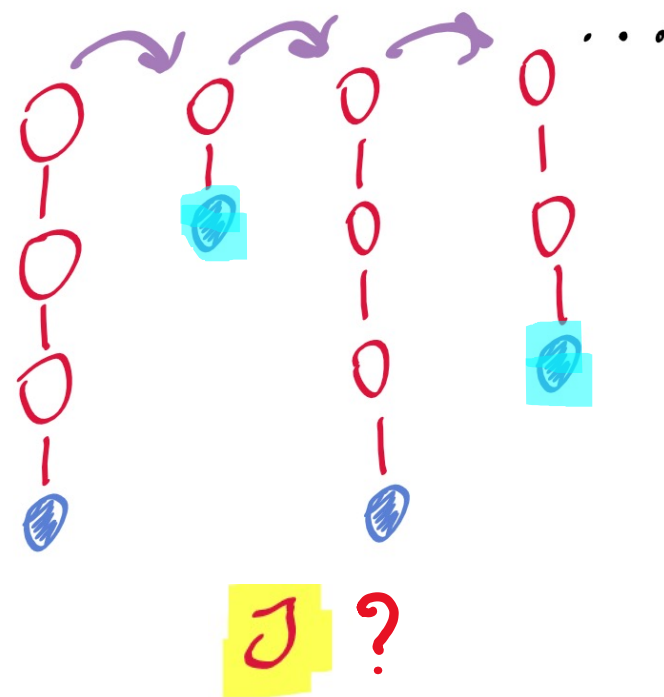
TAKVIYE



J ?

RULE EXTRACTION

RUX



RULE GENERATION

RUG

Tahmin Başarıları

Random Forest: 0.84

AdaBoost: 0.95

RUXRF: 0.90

RUXADA: 0.92

RUG: 0.92

Kural Sayıları

Random Forest: 614

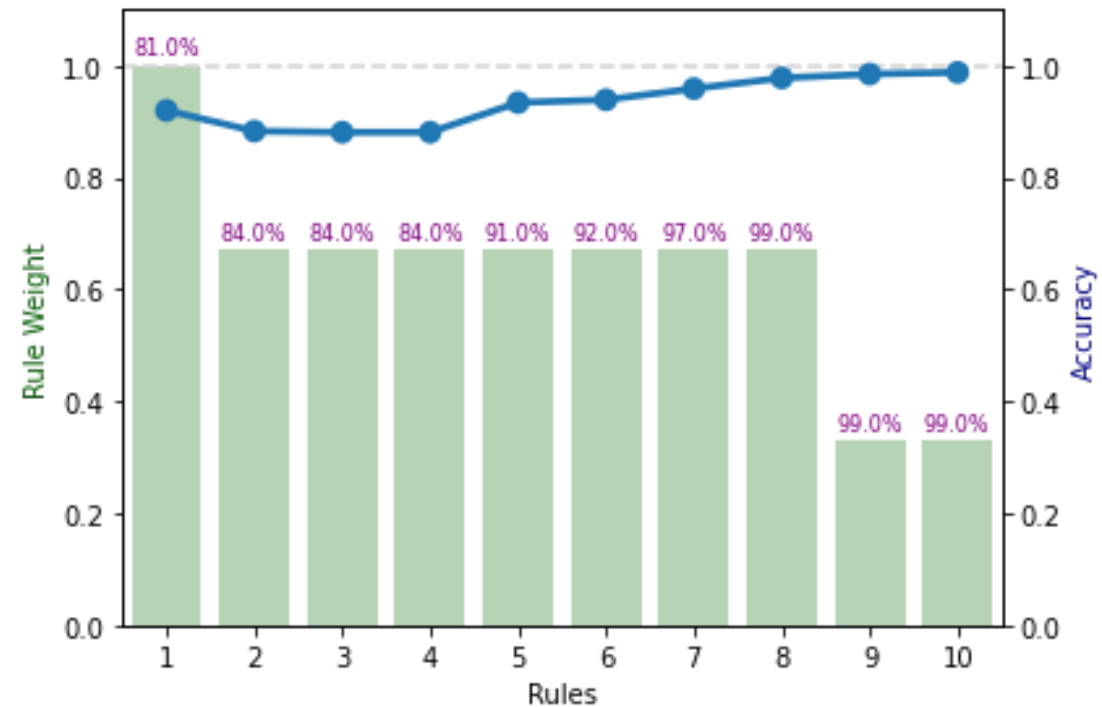
AdaBoost: 742

RUXRF: 19

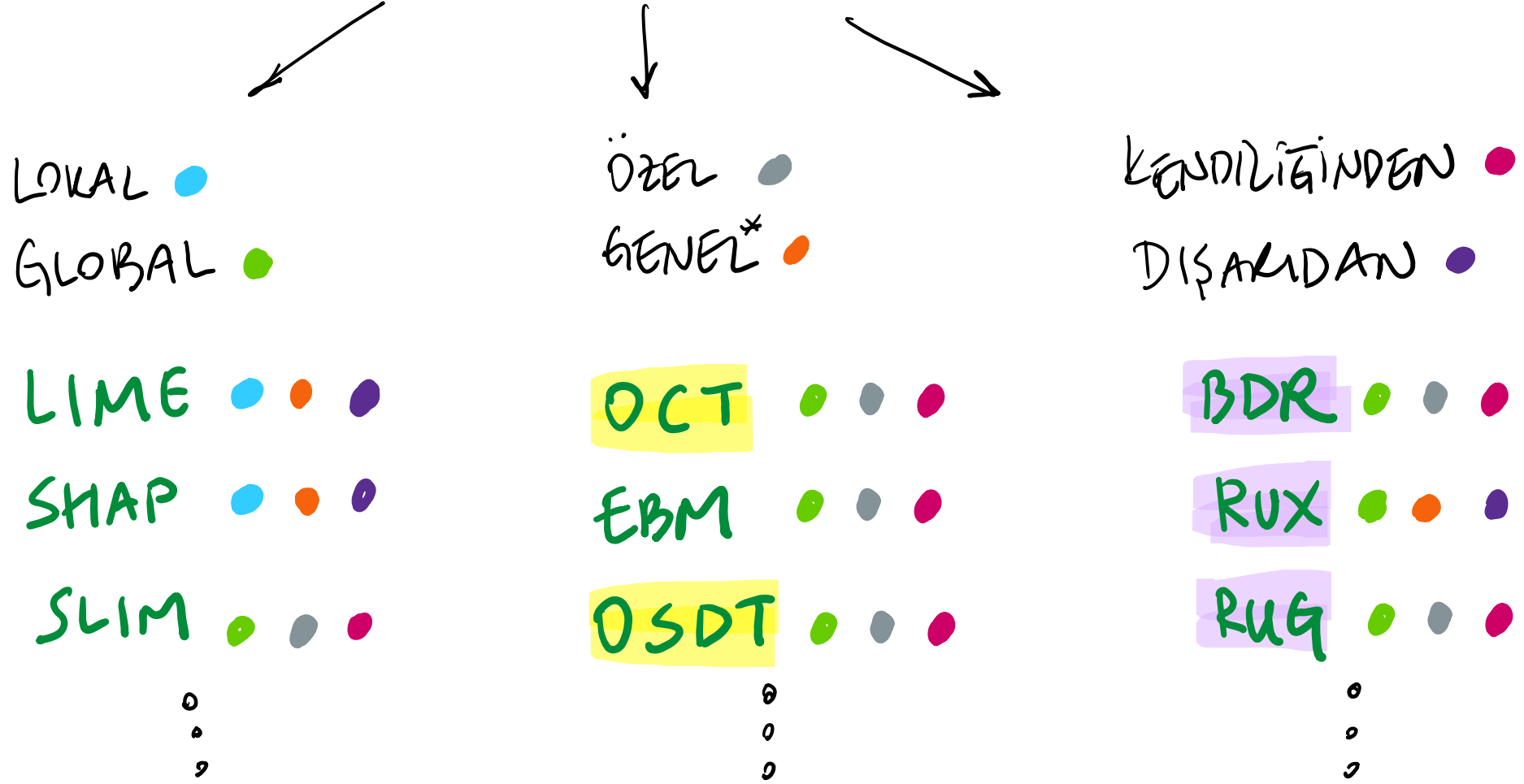
RUXADA: 20

RUG: 24

$$\begin{aligned} &\text{minimize} && \sum_{i \in \mathcal{I}} v_i + \sum_{j \in \mathcal{J}} c_j w_j \\ &\text{subject to} && \sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i > 1 \quad i \in \mathcal{T}. \end{aligned}$$



YORUMLAMA YAKLAŞIMLARI



 KARAR AĞACI TABANLI
 KURAL TABANLI

* AGNOSTİK

Teşekkürler



-  [UvA](#)
-  [D4C](#)
-  [Bol Bilim](#)
-  [Veri Defteri](#)
-  [@sibirbil](#)