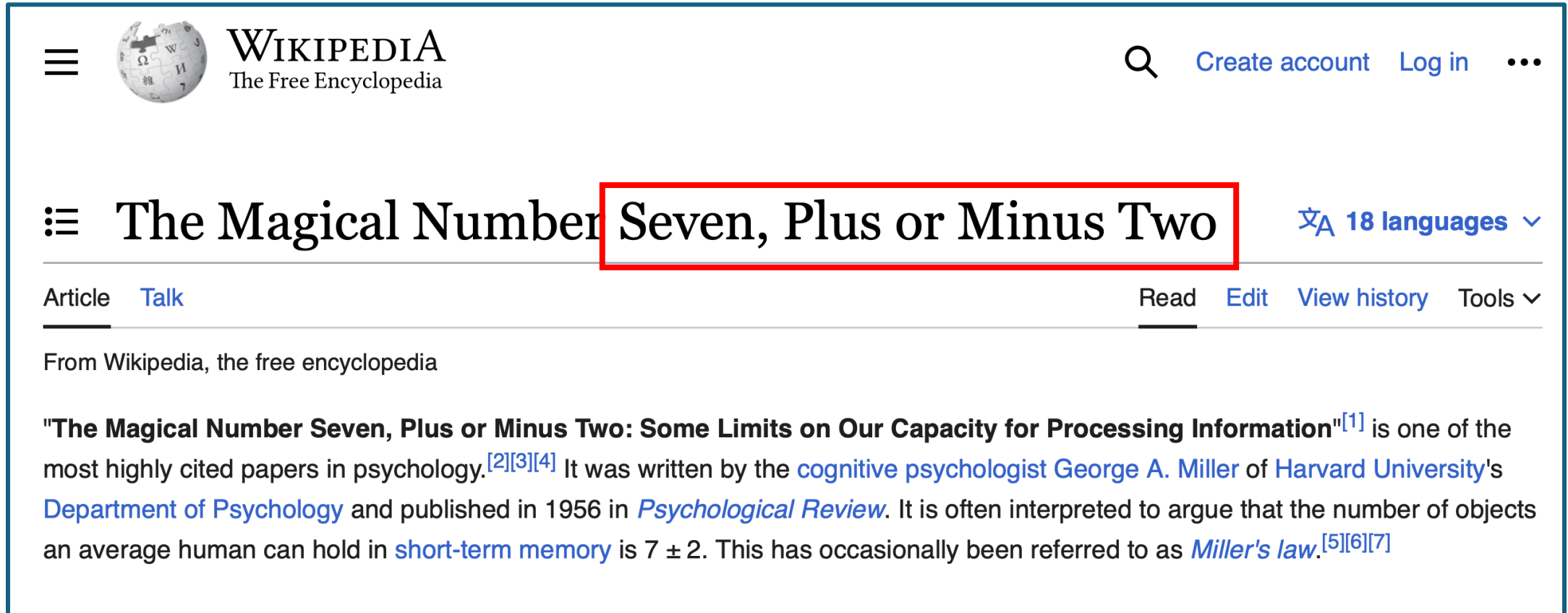


Improving Decision Sparsity

Yiyang Sun, Tong Wang, and Cynthia Rudin

Sparsity is important for interpretability



The screenshot shows the Wikipedia interface for the article "The Magical Number Seven, Plus or Minus Two". The title is highlighted with a red box. The article text discusses George A. Miller's 1956 paper on short-term memory capacity, mentioning "Miller's law".

WIKIPEDIA
The Free Encyclopedia

Search [Create account](#) [Log in](#) ...

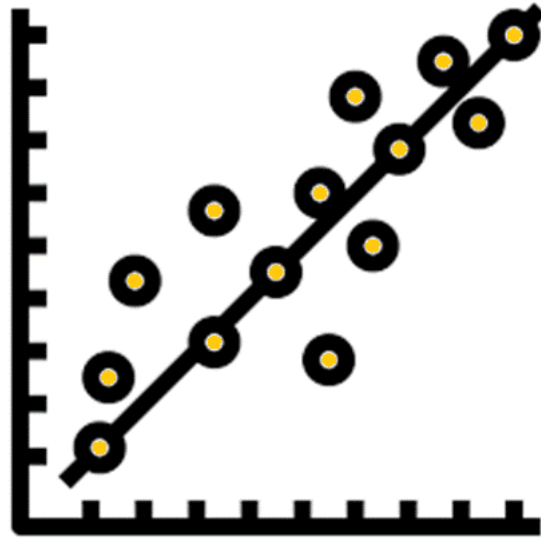
[☰](#) **The Magical Number Seven, Plus or Minus Two** [🌐 18 languages](#) ▾

[Article](#) [Talk](#) [Read](#) [Edit](#) [View history](#) [Tools](#) ▾

From Wikipedia, the free encyclopedia

"**The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information**"^[1] is one of the most highly cited papers in psychology.^{[2][3][4]} It was written by the [cognitive psychologist George A. Miller](#) of [Harvard University's Department of Psychology](#) and published in 1956 in *Psychological Review*. It is often interpreted to argue that the number of objects an average human can hold in [short-term memory](#) is 7 ± 2 . This has occasionally been referred to as *Miller's law*.^{[5][6][7]}

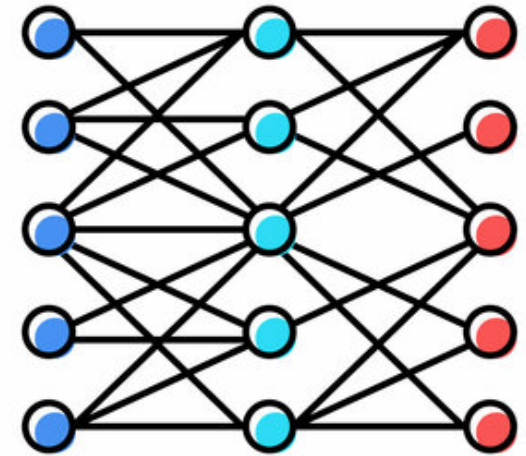
Common metrics for evaluating sparsity of the model



of terms
in linear regressions



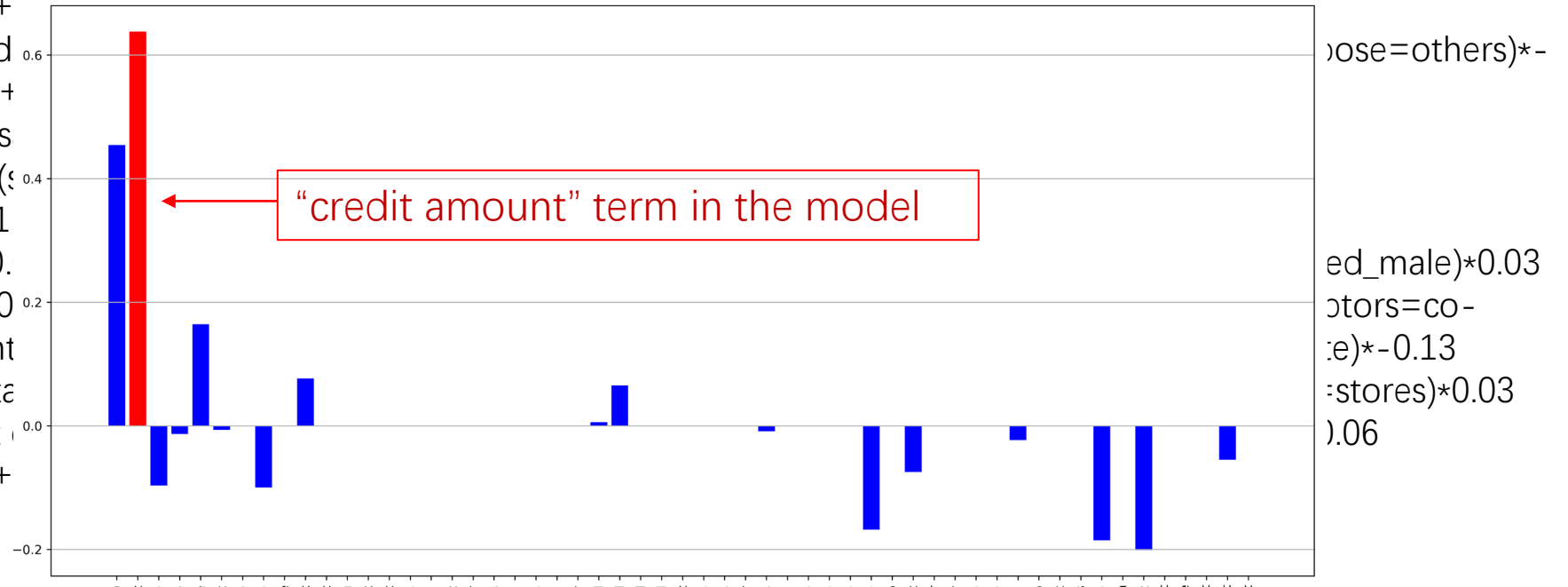
of leaf nodes
in decision trees



of parameters
in neural networks

This is not a sparse model.

$(\text{duration}(\text{month})) * 0.21 + (\text{credit_amount}) * 0.33 + (\text{installment_rate}) * 0.13 + (\text{present_residence_since}) * -0.01 + (\text{age}) * -0.11 + (\text{number_of_credit}) * 0.03 + (\text{number_of_people_being_liable}) * -0.01 + (\text{provide_maintenance}) * -0.09 + (\text{telephone}) * -0.07 + (\text{check_account_existence} = 0 - 200\text{DM}) * 0.06 + (\text{checking_account} = \text{none}) * 0.35 + (\text{credit_history} = \text{all_credit}) * 0.08 + (\text{credit_history} = \text{critical_account}) * -0.19 + (\text{credit_history} = \text{delay_in_paying}) * 0.01 + (\text{purpose} = \text{business}) * 0.02 + (\text{purpose} = \text{education}) * 0.01 + (\text{purpose} = \text{radio/television}) * -0.12 + (\text{saving_account} = 100 - 500\text{DM}) * 0.02 + (\text{saving_account} \geq 1000\text{DM}) * -0.04 + (\text{present_employment} = 4 - 7\text{years}) * -0.11 + (\text{present_employment} = \text{unemployed}) * 0.01 + (\text{personal_status_sex} = \text{married_male}) * -0.01 + (\text{applicant} = \text{other_debtors}) * 0.02 + (\text{other_debtors} = \text{guarantor}) * -0.01 + (\text{property} = \text{unknown}) * 0.04 + (\text{other_installment_stores}) * 0.03 + (\text{housing} = \text{free}) * 0.0 + (\text{housing} = \text{doesn't own}) * -0.06 + (\text{job} = \text{unemployed_non-resident}) * -0.01 + (\text{job} = \text{unemployed_resident}) * -0.01 + (\text{job} = \text{employed_male}) * 0.03 + (\text{job} = \text{employed_female}) * -0.13 + (\text{job} = \text{retired}) * 0.06$



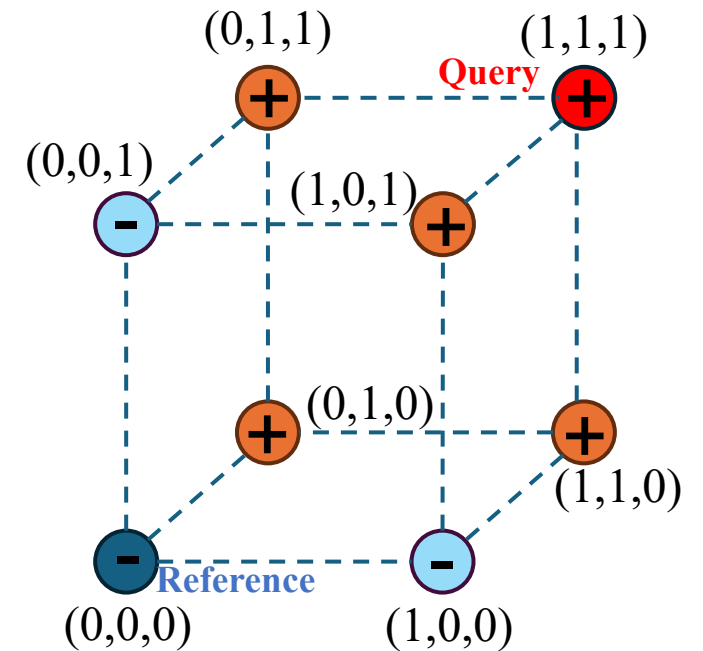
Why was Arun's loan application denied?

The prediction is determined by *one feature!*

He asked for over \$10K!

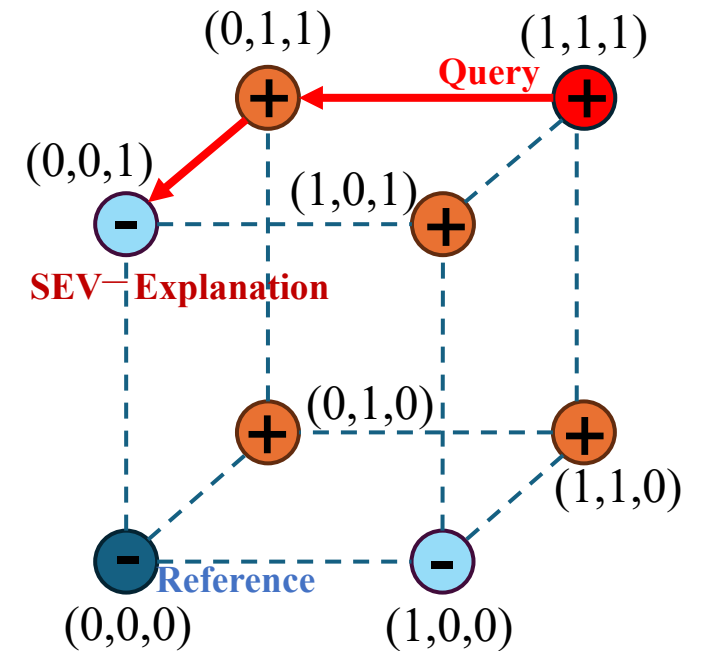
Recap of Decision Sparsity: SEV⁻

- Step 1: Define reference (usually either **0** or average of negative class)
- Step 2: Define Boolean hypercube for query **x**:
 - coordinate_{*j*} is 1 if feature *j* is **x_j**
 - coordinate_{*j*} is 0 if feature *j* is at reference value **r_j**

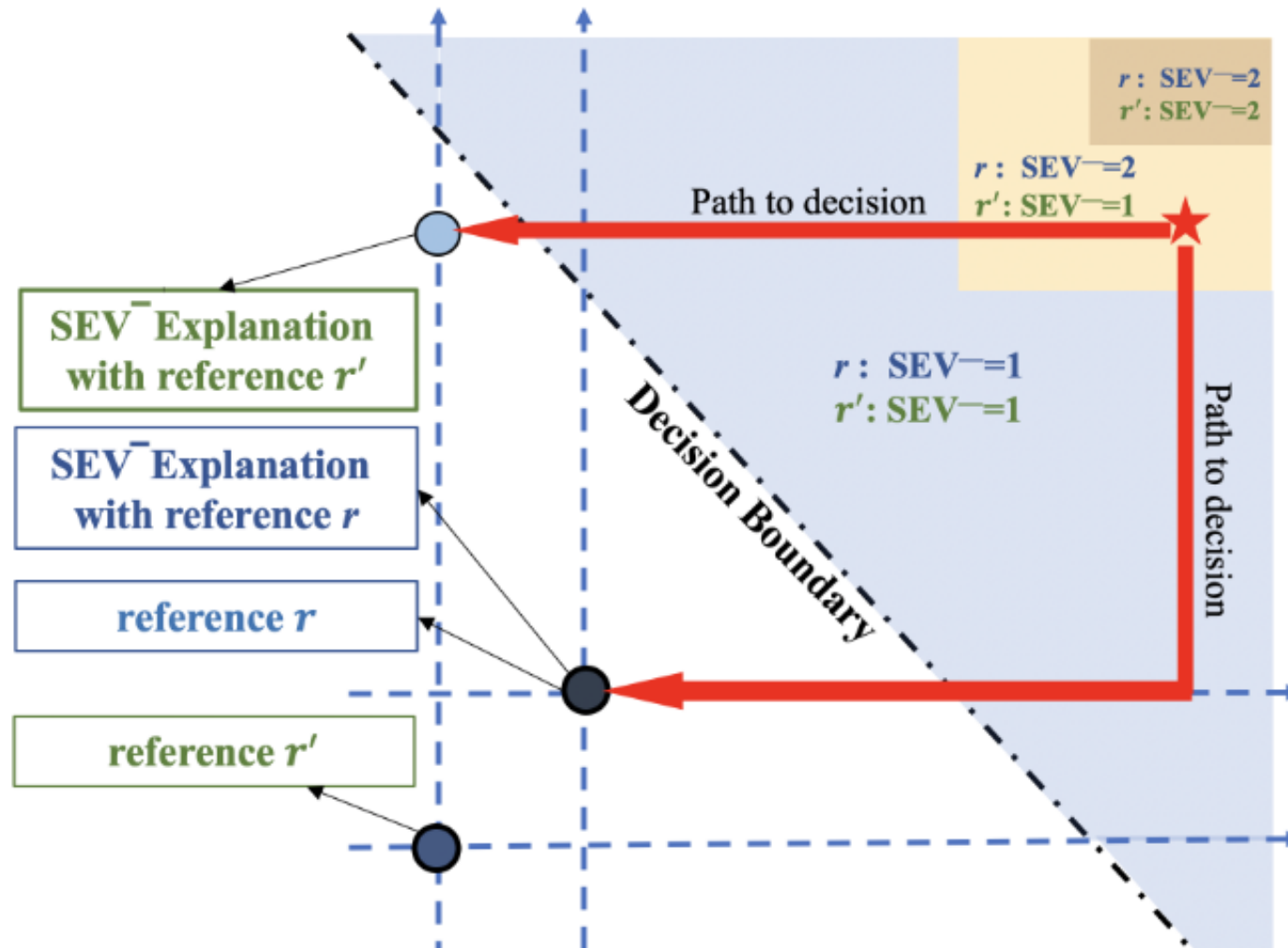


Recap of Explanation Sparsity: SEV^-

- Step 1: Define reference (usually either $\mathbf{0}$ or average of negative class)
- Step 2: Define Boolean hypercube for query \mathbf{x} :
 - coordinate j is 1 if feature j is x_j
 - coordinate j is 0 if feature j is at reference value r_j
- Step 3: Define SEV^- for \mathbf{x} . Moving from \mathbf{x} towards reference, SEV^- is the minimum l_0 distance to a negative label.



SEV⁻ is sensitive to the reference selection



Selection criteria for references

Minimum number of features used for explaining decisions

Reference and Query are under the same subpopulations

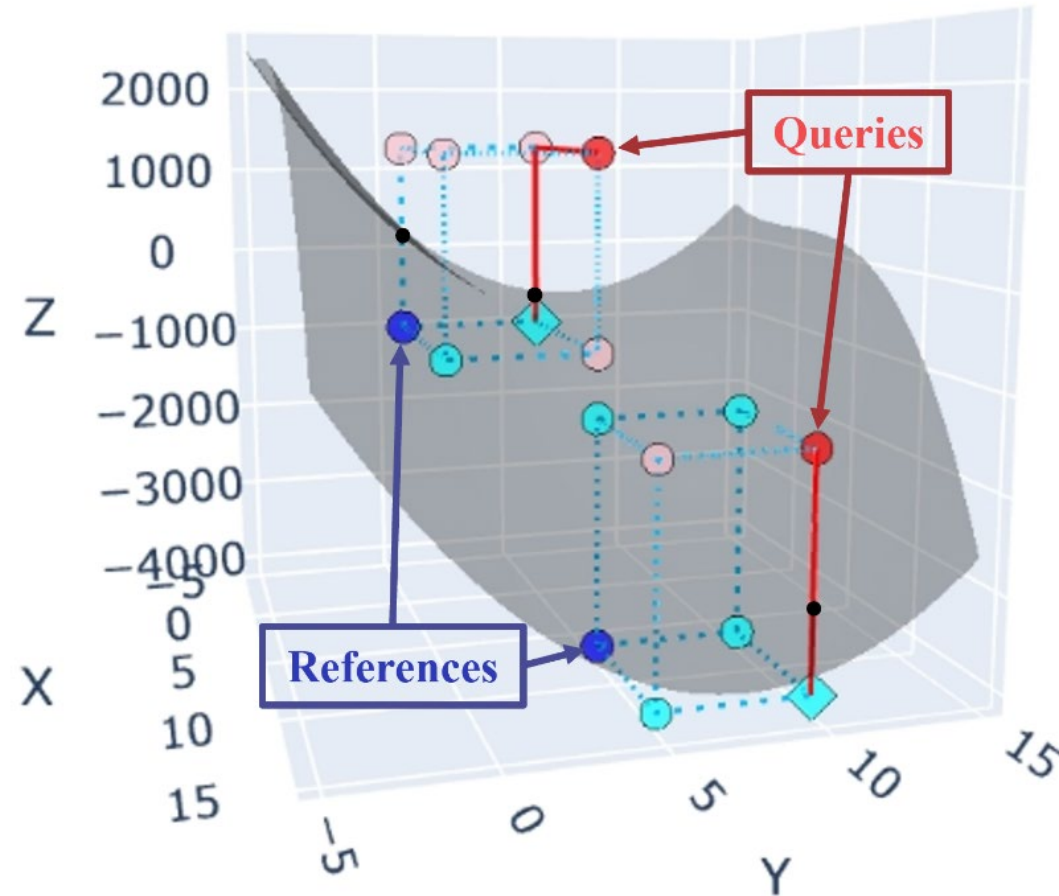
$$\|\mathbf{x}_i - \mathbf{r}_i\|, \mathbf{r}_i \in \mathcal{R} \quad (\text{Closeness})$$

$$\text{SEV}^-(f, \mathbf{x}_i, \mathbf{r}_i), \mathbf{r}_i \in \mathcal{R} \quad (\text{Sparsity})$$

$$-P(\mathbf{x}_i^{\text{expl}, \mathbf{r}_i} | X^-) \quad (\text{Negated Credibility})$$

Sparse explanations are under the original data distribution

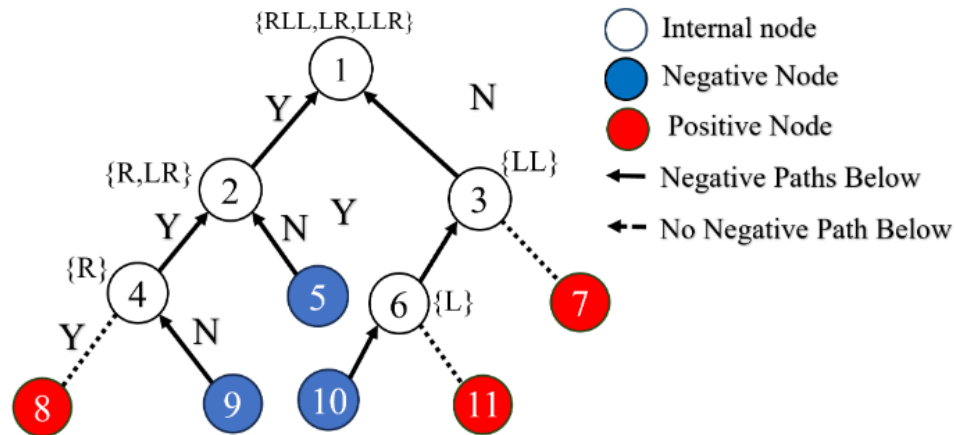
Improving Closeness: Cluster-based SEV (SEV-C)



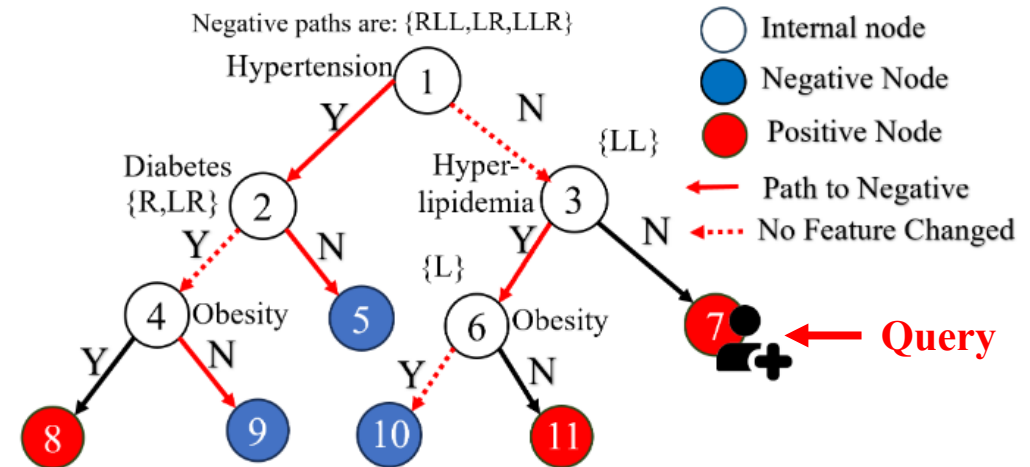
The process for calculating the SEV-C:

- **Step 1:** Selects the references by clustering the negative populations, and regarding the cluster centroid points as the references.
- **Step 2:** Assign each query their closest cluster centroid as their reference points.
- **Step 3:** Go over the original SEV Calculation

Special case for SEV-C: Tree-based SEV



Step 1: SEV-T Preprocessing: Collect negative leaf node information for each internal nodes



Step 2: Efficient SEV-T Calculation: Go over internal node of decision path, and check all negative paths

It have many useful properties and computational benefits!

Optimizing SEV Variants for Models

Gradient-based Optimization (AllOpt)

- Maximize the fraction of points with $SEV^- = 1$

Search-based Optimization (TOpt)

- Find a model with the lowest SEV with
in a set of classification models with
the best performance

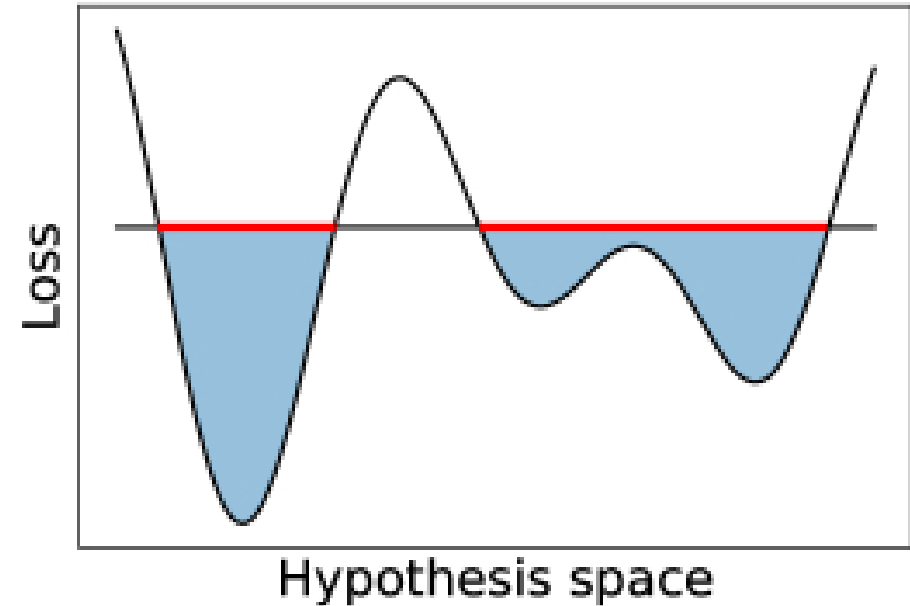
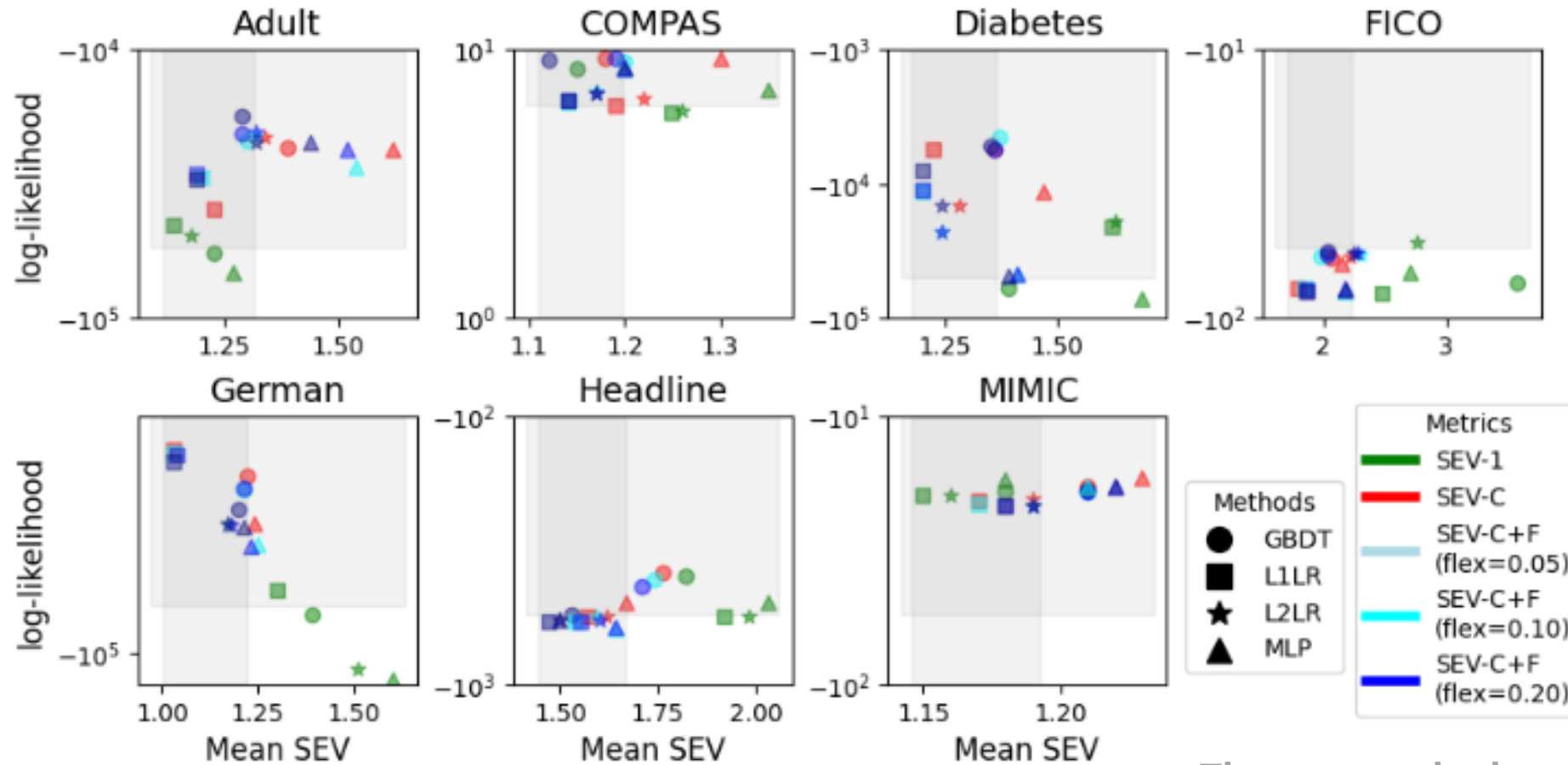


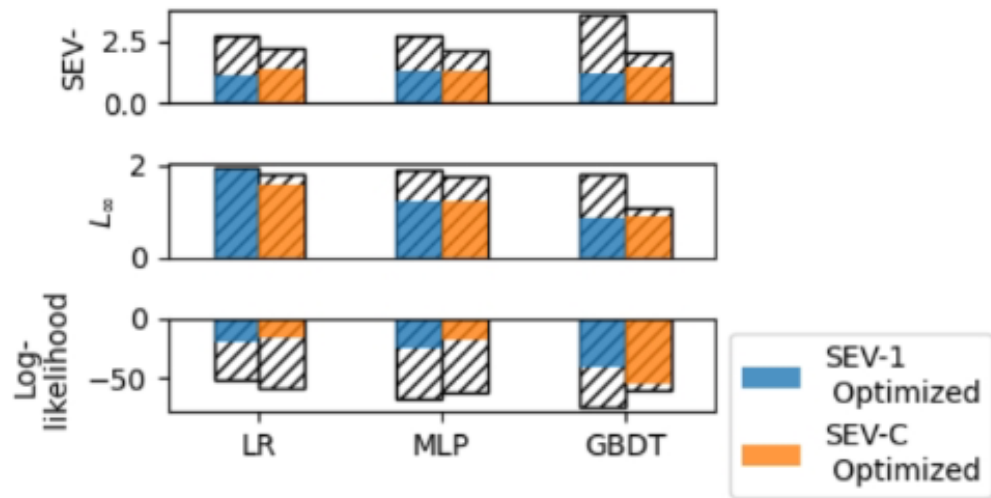
Figure 1: Rashomon Sets

SEV-C improves in credibility and closeness



The grayer, the better.

Optimization improves explanations while preserving model performance



(a) All-Opt⁻ Performance

	TRAIN ACC	TEST ACC	MEAN SEV ^T
CART	0.71 ± 0.01	0.71 ± 0.01	1.10 ± 0.03
C4.5	0.71 ± 0.01	0.71 ± 0.01	1.13 ± 0.05
GOSDT	0.70 ± 0.01	0.70 ± 0.01	1.08 ± 0.02
TOpt	0.70 ± 0.01	0.70 ± 0.01	1.00 ± 0.00

(b) SEV^T performance on different tree-based models

More in the paper

- SEV Variants for further improving the sparsity and credibility of the explanations.
- Sparsity and Credibility comparison with counterfactual explanation methods.
- Score-based soft K-Means for avoiding positive predicted references
- Detailed Algorithms for tree-based SEV
- Timing experiments
- ...

SEV Paper



NeurIPS Paper

