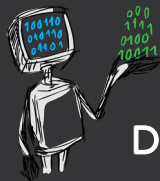


Interpretability in Data-Centric ML

Alexandra (Ola) Zyttek
Introduction to Data-Centric AI
IAP 2023



Massachusetts
Institute of
Technology



DATA TO AI

experimentation, we found out that security analysts highly value *easily interpretable features* when analyzing outputs of machine learning algorithms.

good accuracy as well as recall and that *use human-understandable features*. Our findings indicate that moderators would appreciate the ability to understand outputs based on such features. As

MGM has two core elements which perform *interpretable feature extraction* and selection. At the

of the process, even before an actual model is developed. For example, P11, referring to *feature engineering*, remarked: “... *this is the first step toward making interpretable models, even though we don't have any model yet.*”. In particular, we found several data scientists complement feature

Roadmap

- Introduction to interpretable ML
- **Why** do we care about interpretable features?
- **What** are interpretable features *really*?
- **How** do we get interpretable features?

- **Introduction to interpretable ML**

- **Why** do we care about interpretable features?

- **What** are interpretable features *really*?

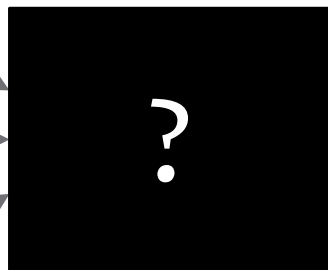
- **How** do we get interpretable features?



Location:
Denfield

Population:
120

Median Income:
\$32,000

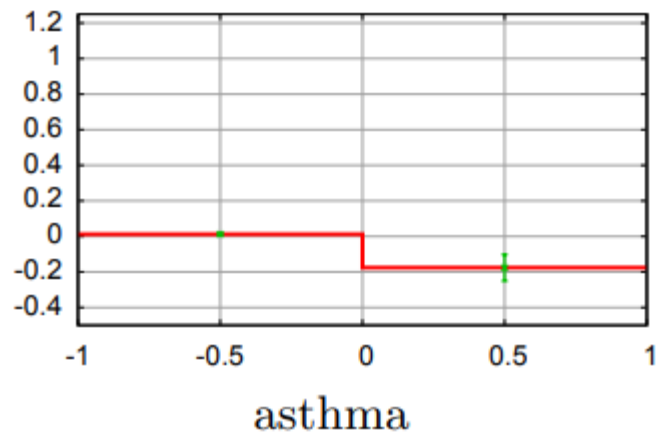
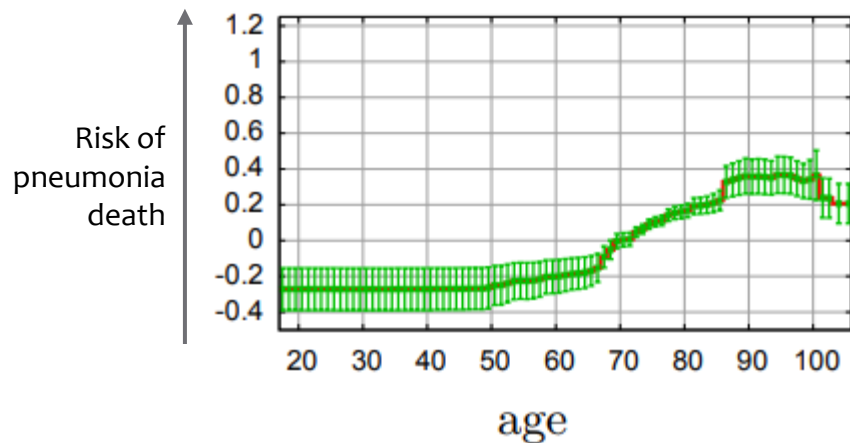


Median House Price:
\$198,000

Why do we need interpretable ML?

1. Debugging and validation
2. Reviewing decisions
3. Improving usability

Debugging and Validation



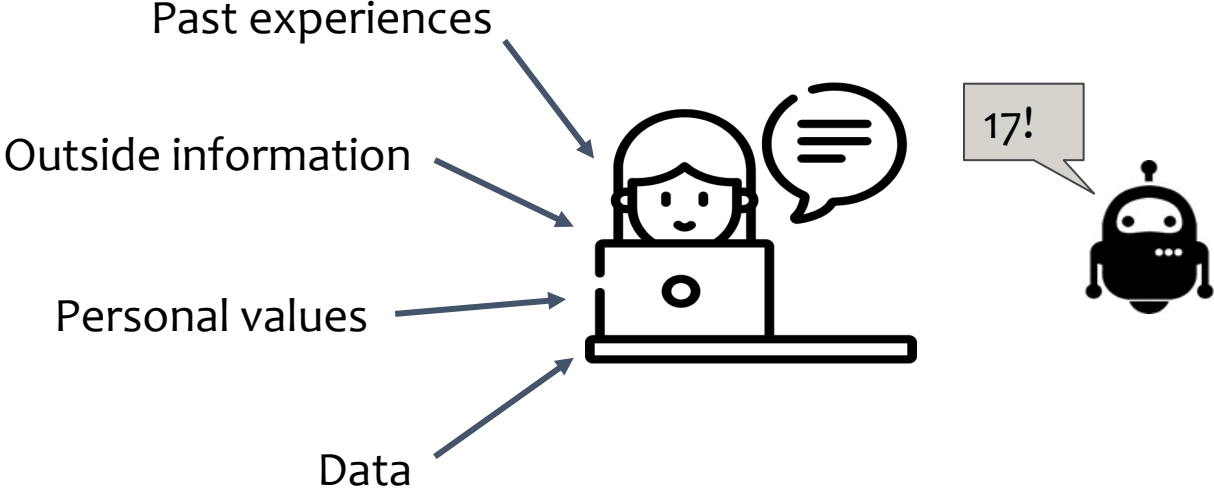
Reviewing Decisions

What self-driving cars can't recognize may be a matter of life and death

Engineers are racing to program artificial intelligence to recognize different scenarios that human drivers know inherently



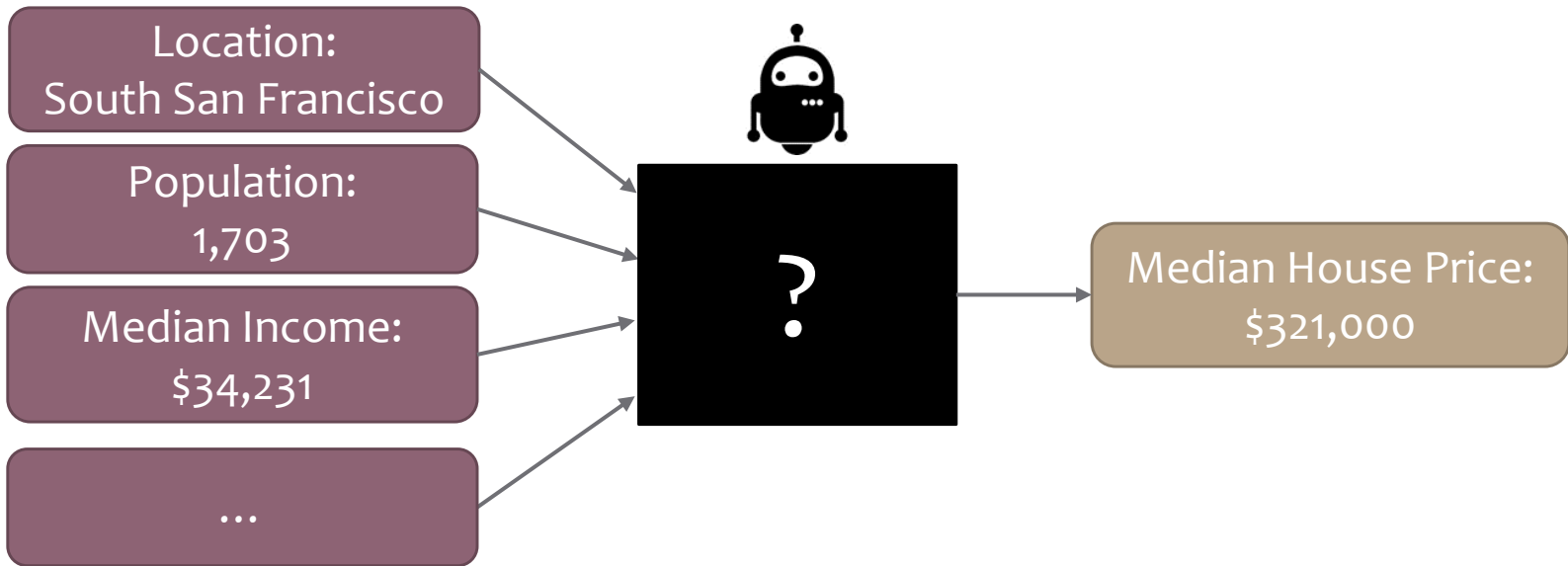
Improving Usability

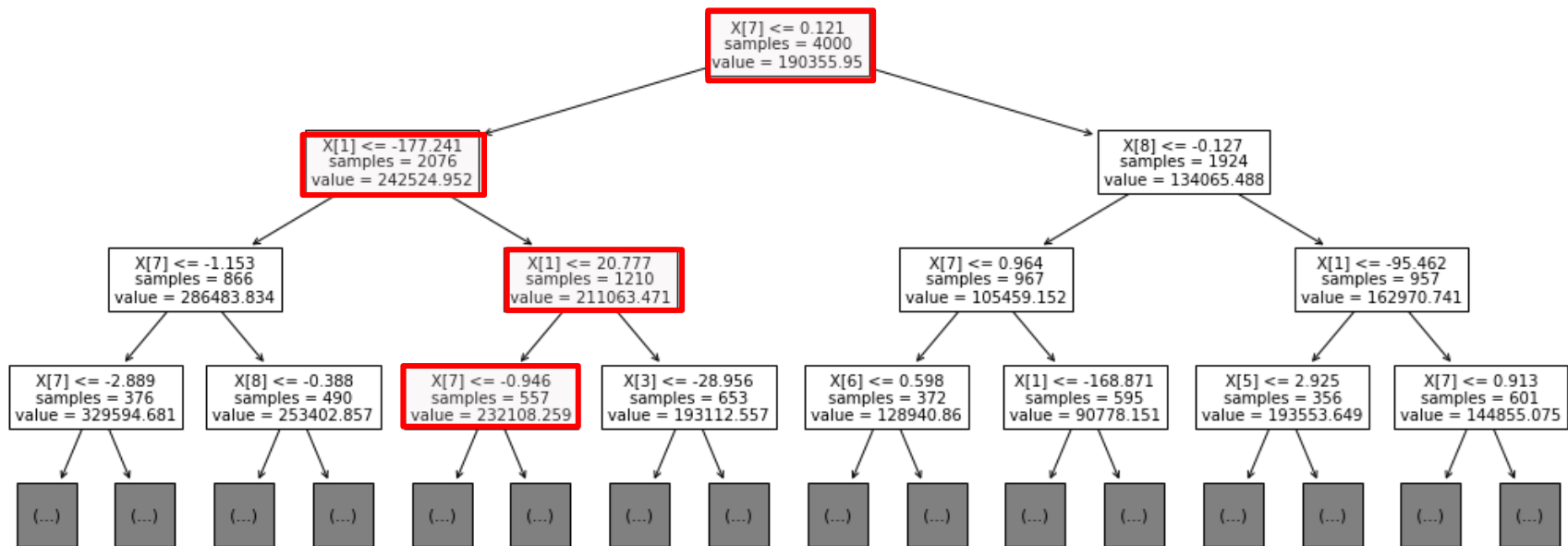


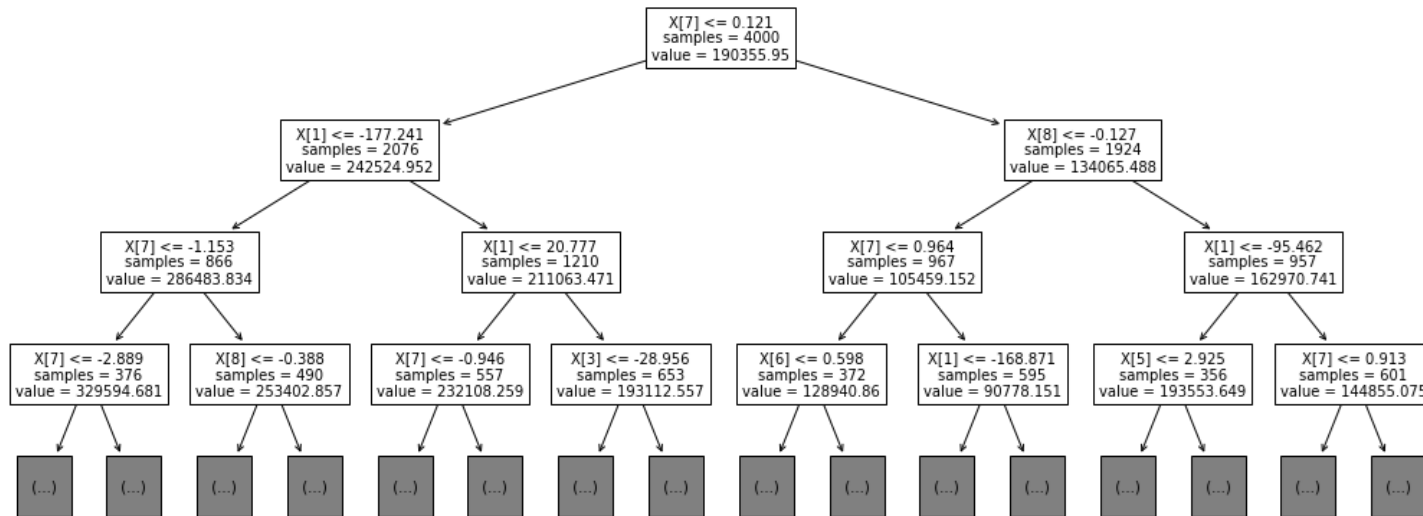
We need interpretable ML...

- When the problem formulation is **incomplete**
- When there is associated **risk**
- When **humans are involved** in decision-making

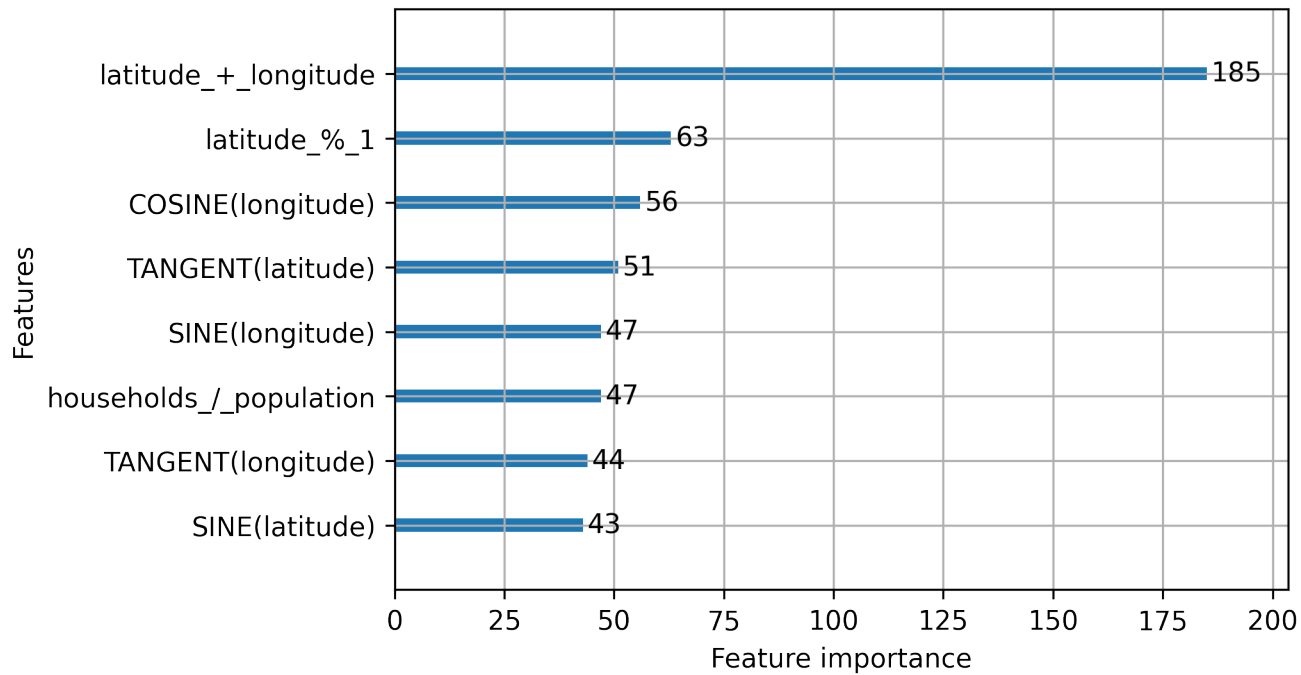
- Introduction to interpretable ML
- **Why do we care about interpretable features?**
- **What** are interpretable features *really*?
- **How** do we get interpretable features?



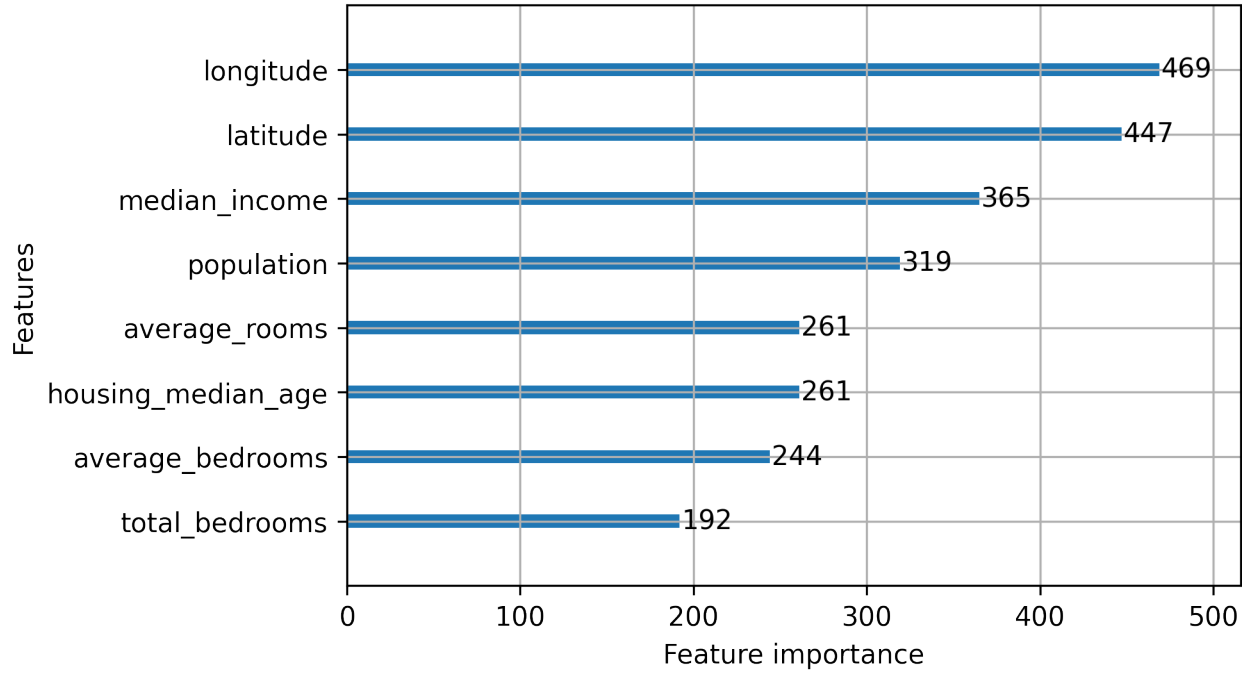




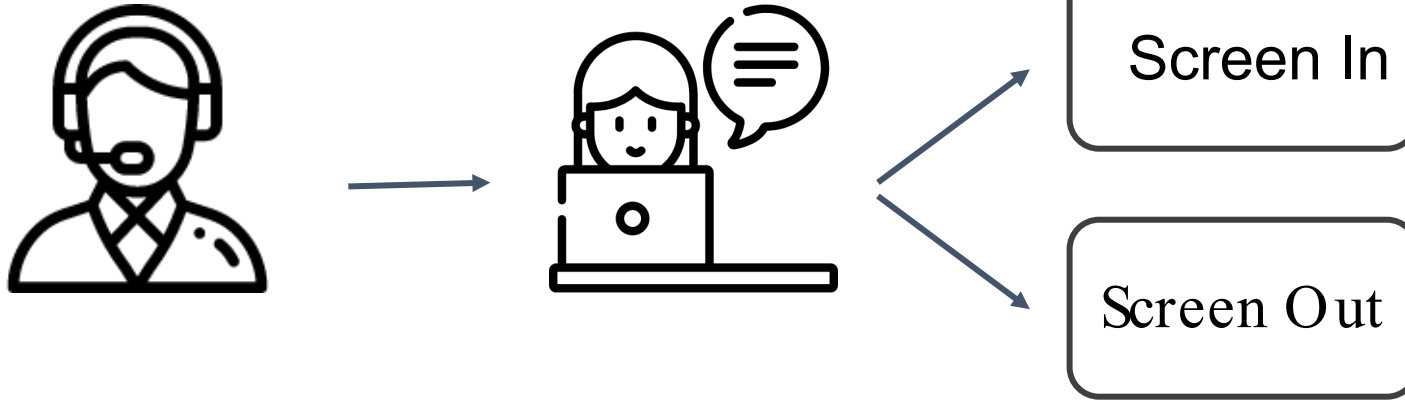
Feature importance



Feature importance



Case Study: Child Welfare



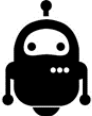
Case Study: Child Welfare






Screen In

Screen Out

1- 20



[Click "Show All Factors" to enable Search and Filter](#)

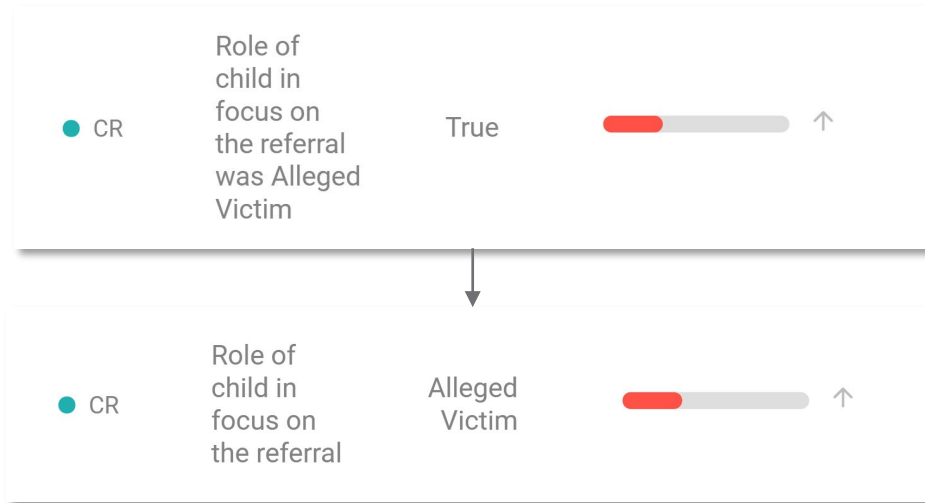
Category	Factor	Value	Contribution
● DG	Age range of child in focus	<1 year	↓  ↑
● DG	Age of the child in focus at time of referral	0	↓  ↑
● RO	Number of other children (non victims) on the referral	0	↓  ↑

Problem: Confusing Features

 PH	Count of days the child in focus was in a child welfare placement in the last 365 days	1	↓		↑
 PH	Count of days the child in focus was in a child welfare placement in the last 730 days	1	↓		↑

Problem: Confusing Language

“The ‘true’ and ‘false’ is hard to interpret... Would rather have a positive statement (e.g., no perpetrator named)” –Child Welfare Screener

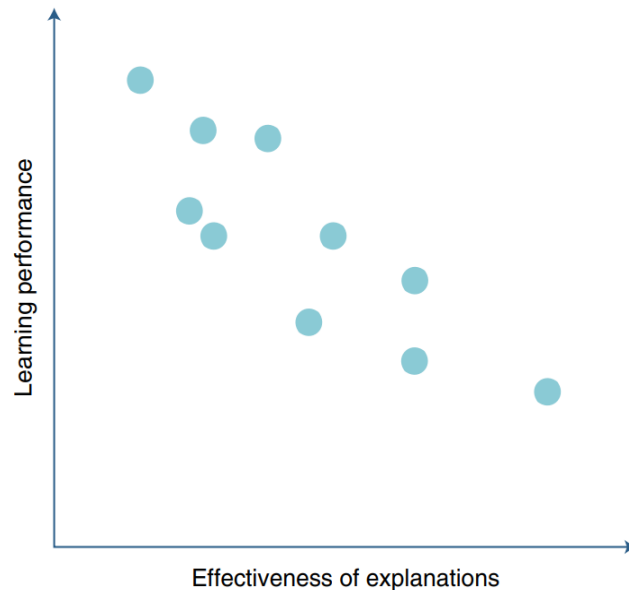


Problem: Irrelevant Features

“2 parents have missing date-of-birth is shown as a significant blue bar which I can’t imagine is protective.” – Child Welfare Screener

Performance and Interpretability

- Interpretability leads to...
 - × More efficient training
 - × Better generalization
 - × Fewer adversarial examples
- The interpretability-performance tradeoff is (mostly) a myth

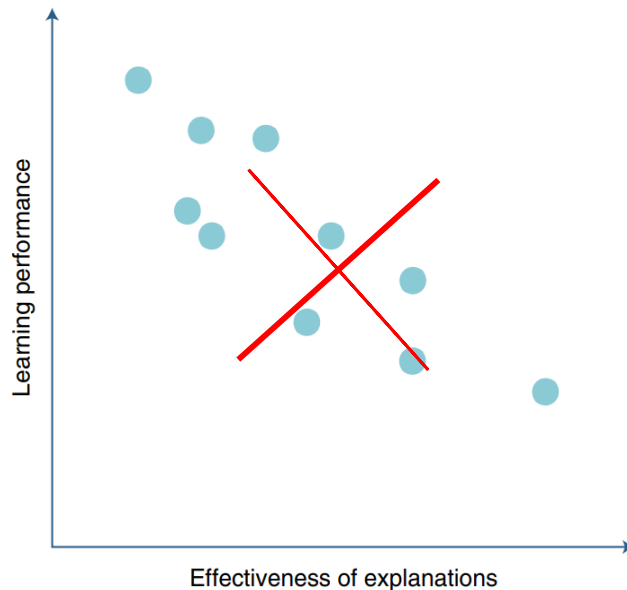


Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell* 1, 206–215 (2019).

Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., & Madry, A. (2019). *Adversarial Examples Are Not Bugs, They Are Features* (arXiv:1905.02175).

Performance and Interpretability

- Interpretability leads to...
 - × More efficient training
 - × Better generalization
 - × Fewer adversarial examples
- The interpretability-performance tradeoff is (mostly) a myth



Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell* 1, 206–215 (2019).

Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., & Madry, A. (2019). *Adversarial Examples Are Not Bugs, They Are Features* (arXiv:1905.02175).

- Introduction to interpretable ML
- **Why** do we care about interpretable features?
- **What are interpretable features *really*?**
- **How** do we get interpretable features?

What are interpretable features *really*?

The features that are most useful and meaningful to the user

Example: Housing Price Prediction



Area Quality (numeric)	Average House Size (numeric)	Most Common House Color (categorical)	Normalized Median Income (numeric)	x12 (numeric)
----------------------------------	--	---	--	-------------------------

	Area Quality (numeric)	Average House Size (numeric)	Common House Color (categorical)	Normalized Median Income (numeric)	x12 (numeric)
Readable	✓	✓	✓	✓	

	Area Quality (numeric)	Average House Size (numeric)	Common House Color (categorical)	Normalized Median Income (numeric)	x12 (numeric)
Readable	✓	✓	✓	✓	
Understandable	✓	✓	✓		

	Area Quality (numeric)	Average House Size (numeric)	Common House Color (categorical)	Normalized Median Income (numeric)	x12 (numeric)
Readable	✓	✓	✓	✓	
Understandable	✓	✓	✓		
Relevant	✓	✓			

	Area Quality (numeric)	Average House Size (numeric)	Common House Color (categorical)	Normalized Median Income (numeric)	x12 (numeric)
Readable	✓	✓	✓	✓	
Understandable	✓	✓	✓		
Relevant	✓	✓			
Abstract Concept	✓				

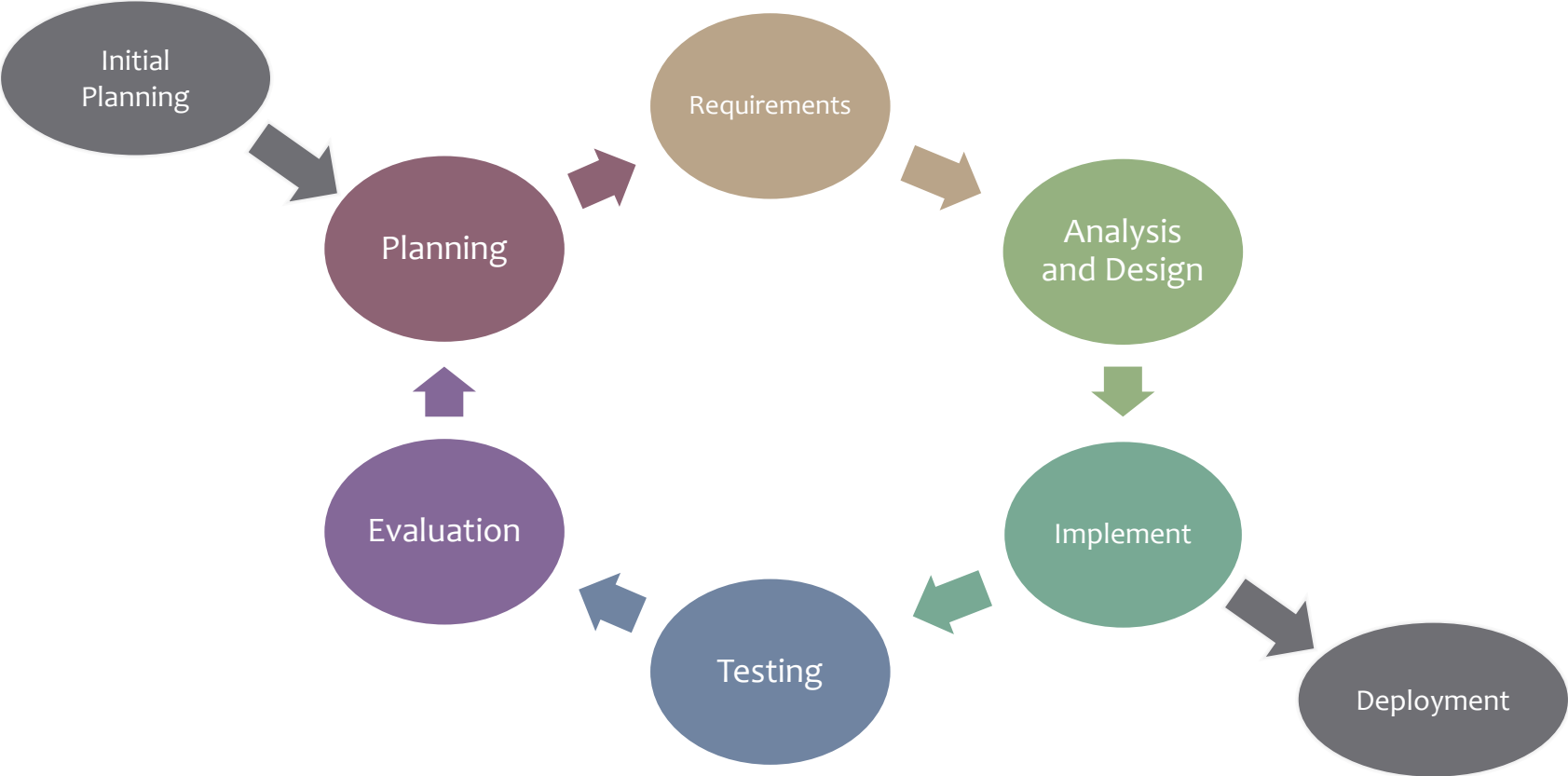
- Introduction to interpretable ML
- **Why** do we care about interpretable features?
- **What** are interpretable features *really*?
- **How** do we get interpretable features?

“[Feature engineering] is the first step to making an interpretable model, even if we don’t have a model yet” – Data Scientist

Methods for Interpretable Features

1. Including the user
2. Using interpretable feature transformations
3. Using interpretable feature generation

Iterative Design Process (for Features)

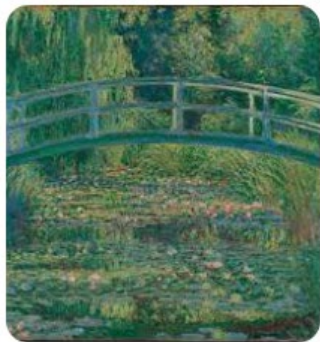


Collaborative Feature Engineering

- Improve the process of crafting *human-generated* features
- Crowd source feature generation
- Allow domain experts to directly participate

Flock: Choosing features through comparisons

1. Machine-generate features for a prediction task
2. Crowd-generate features
3. Cluster crowd-generated features
4. Iterate on inaccurate model nodes





The first painting is probably a Monet because it has lilies in it, and looks like Monet's style. The second probably isn't Monet because Monet doesn't normally put people in his paintings.

Split by conjunctions...

The first painting is probably a Monet because it has lilies in it, and looks like Monet's style. The second probably isn't Monet because Monet doesn't normally put people in his paintings.

Cluster using any clustering algorithm...

**The first painting is probably a Monet because it has lilies in it
It has flowers**

The painting including lilies

There are flowers and lilies in the painting

Crowd-source an aggregated feature label

Does the painting have flowers/lilies?

Results

- Flock outperforms:
 - + Original features/data used directly
 - + Machine-engineered features only
 - + Crowd classifications
- And generates interpretable features, ie.
 - + Contains flowers
 - + Is abstract
 - + Does not contain people

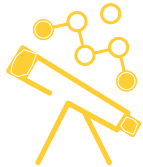
Ballet: Feature Engineering with Feedback

- Abstract away model building/training/evaluating
- Write features with only simple Python

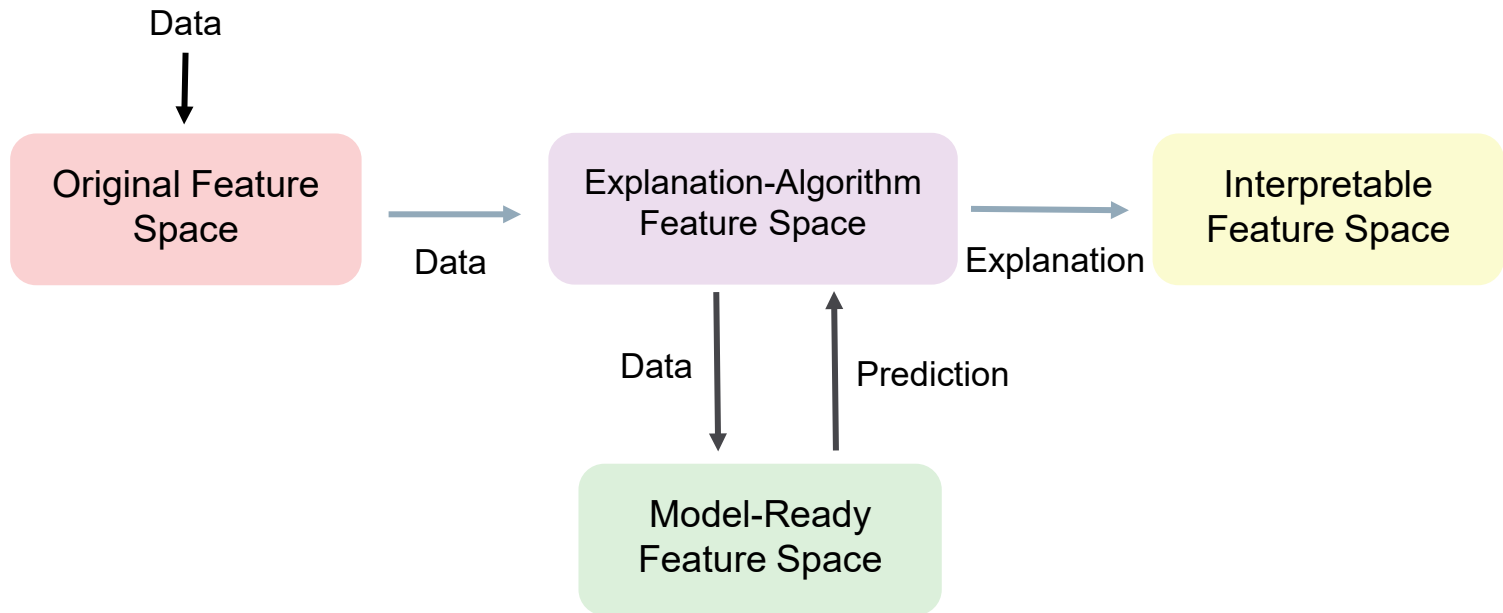
```
def hi_lo_age(dataset):  
    """Whether users are older than 30 years"""  
    from sklearn.preprocessing import binarize  
    threshold = 30  
    return binarize(dataset["users"]["age"]  
                    .values.reshape(-1,1), threshold)
```

Methods for Interpretable Features

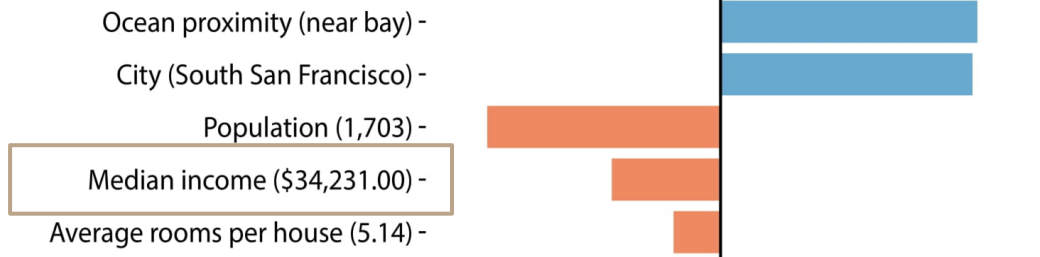
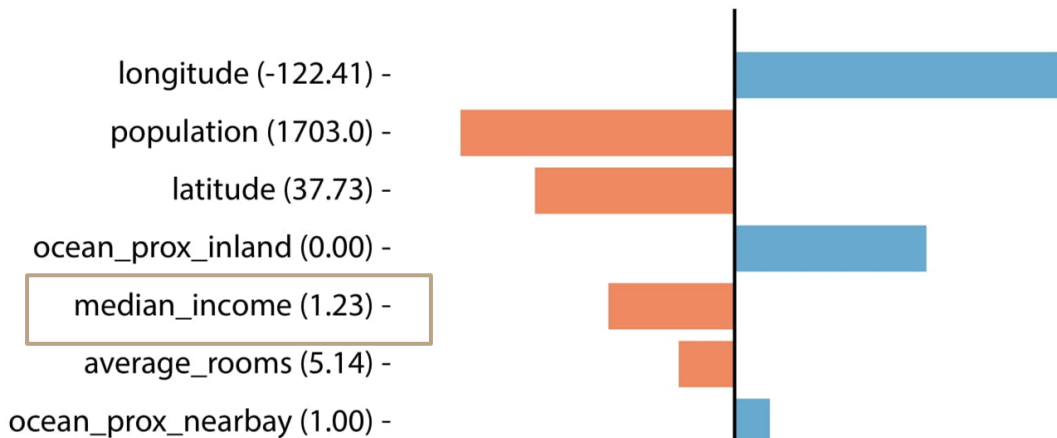
1. Including the user
2. Using interpretable feature transformations
3. Using interpretable feature generation



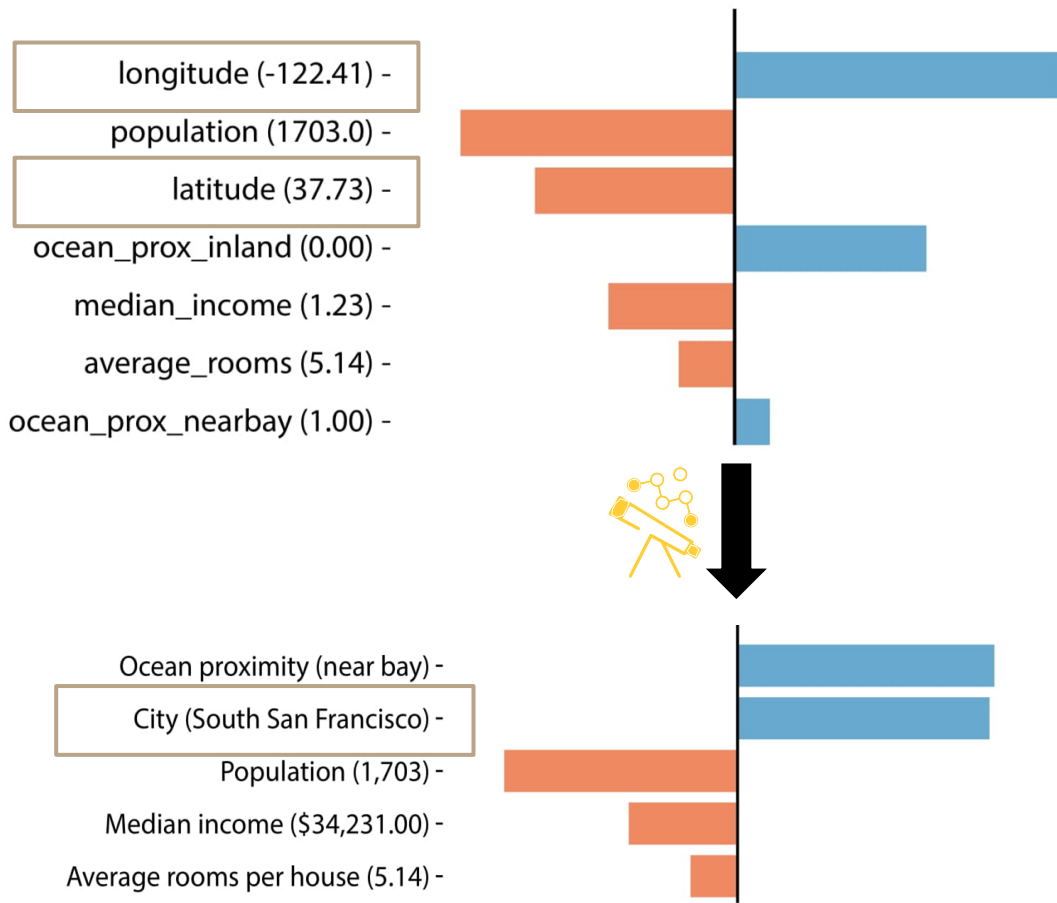
Pyreal: System for Interpretable Transforms



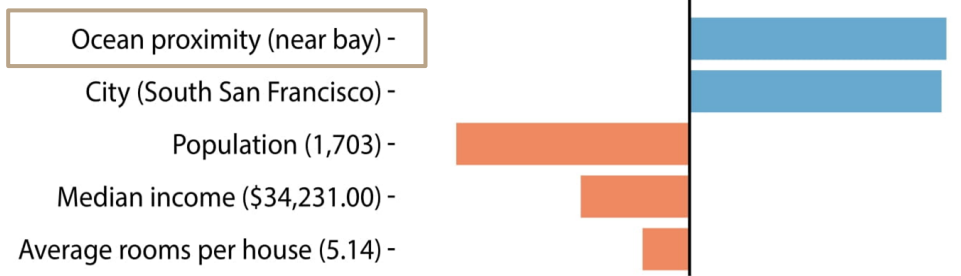
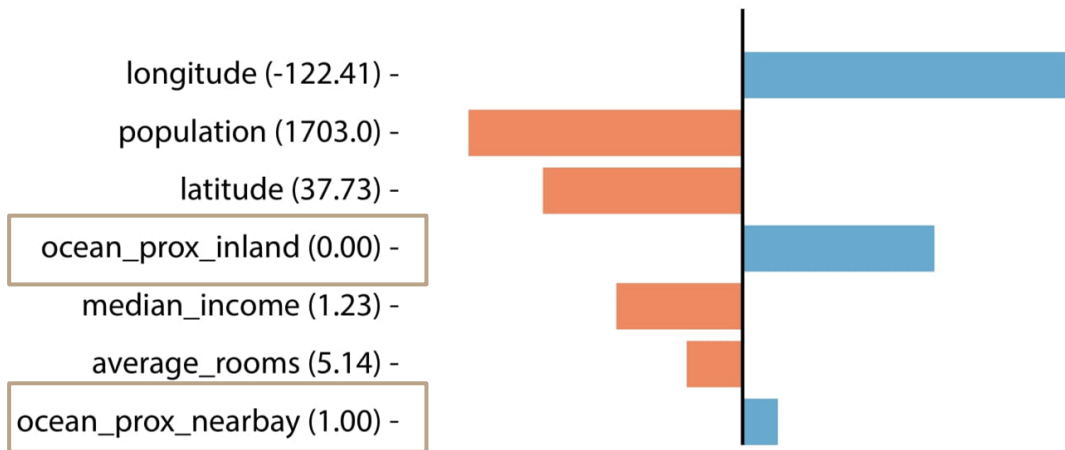
Feature Contributions



Feature Contributions



Feature Contributions



Methods for Interpretable Features

1. Including the user
2. Using interpretable feature transformations
3. Using interpretable feature generation

Mind the Gap Model (MGM)

1. Assign features to groups with AND or OR

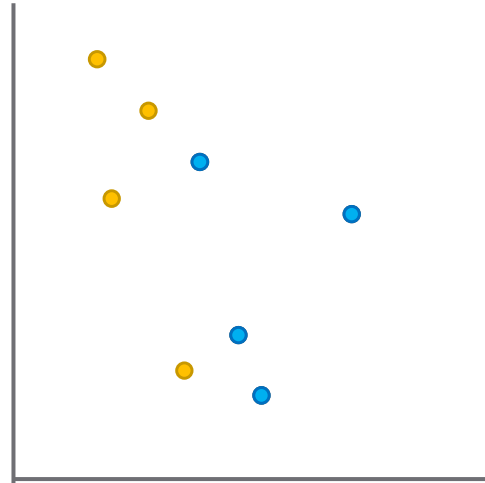
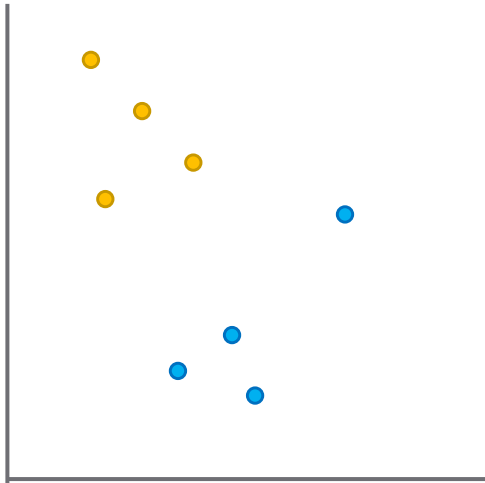
lays eggs

backbone
AND
tail

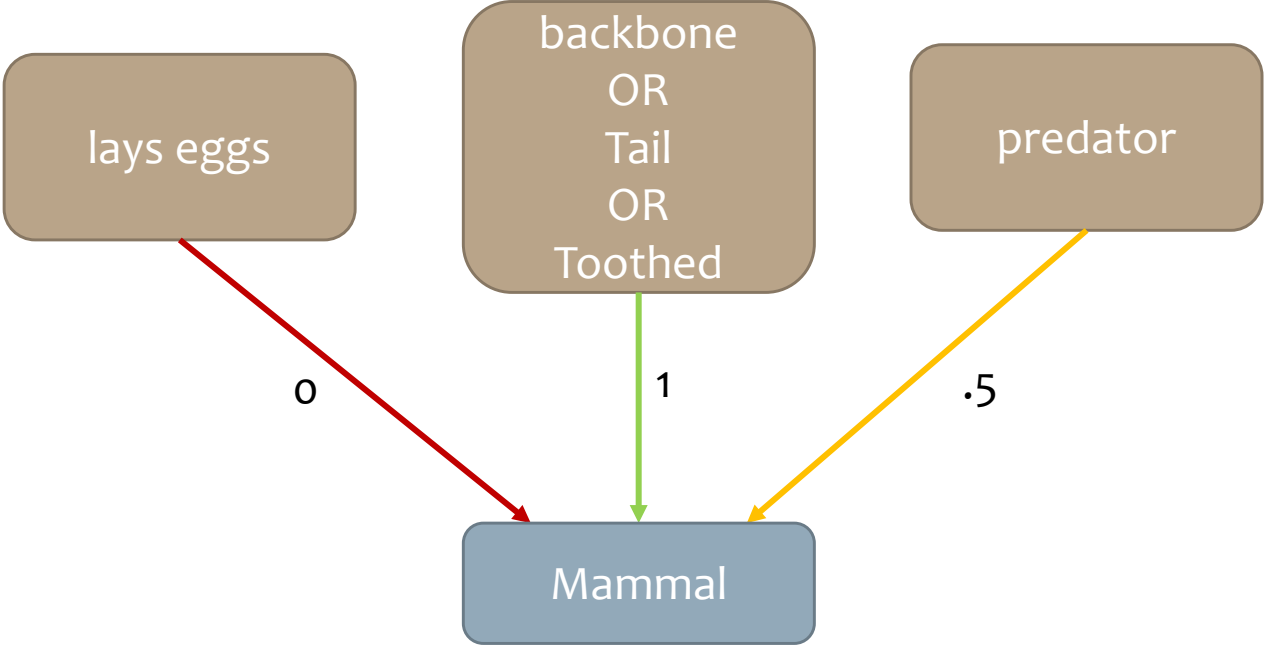
toothed
OR
predator

Mind the Gap Model (MGM)

2. Identify groups maximize separation and iterate



Mind the Gap Model (MGM)



Conclusion

1. ML models are only as interpretable as their features
2. Interpretable features are those that are meaningful to the user
3. Interpretable features are generated by including users, focusing on interpretable transforms, and using feature generation algorithms that consider interpretability

Lab

- Use explanation algorithms to identify flawed data

