



Try out our code!

Do Machine Learning Models Learn Statistical Rules Inferred from Data?

Aaditya Naik, Yinjun Wu, Mayur Naik, Eric Wong asnaik@seas.upenn.edu

Motivation

Models make mistakes!

Object Detection

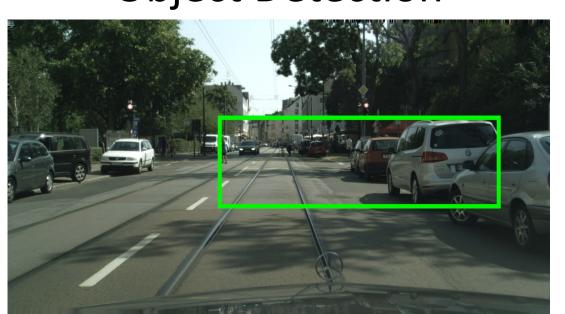
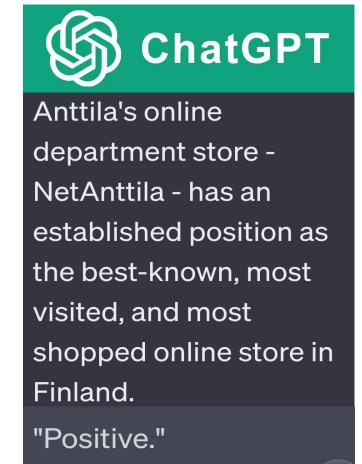


Image Classification





Sentiment Analysis



Rules for estimating errors must be

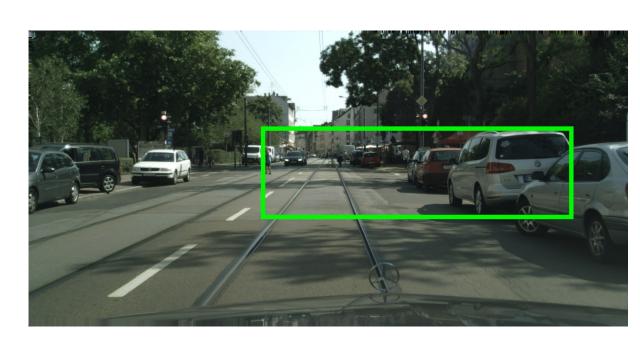
- Valid and hold over a large portion of the data,
- > Expressive to capture complex and interesting patterns,

Statistical Quantile Rules (SQRs)

> Scalable to generate several rules without supervision.

The answer is **Statistical Quantile Rules**:

$$\mathbb{P}(a \le \phi(x) \le b) = 1 - \delta$$

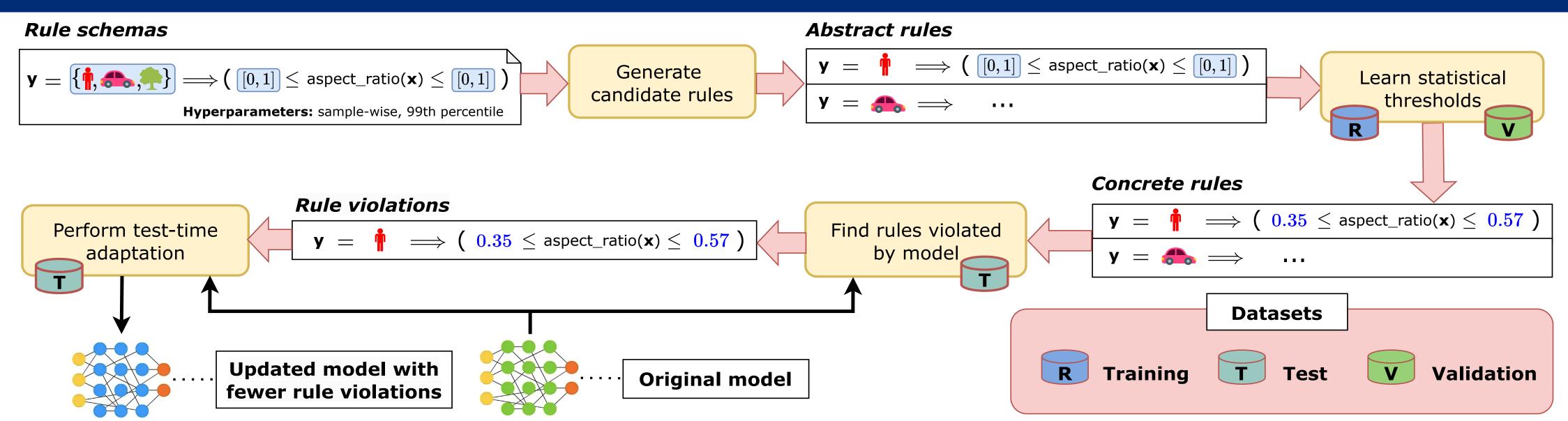


An example of a violation of the rule $0.07 \le aspect_ratio(car) \le 2.77$. 98% of all the ground-truth cars satisfy this rule, making it valid, expressive and scalable.

Observation:

Fundamental errors defy *rules* based on human intuition

The Statistical Quantile Rule Learning (SQRL) Framework



The workflow of SQRL. SQRs are generated from the training and validation data and used to evaluate and improve models.

Examples of Rule Violations

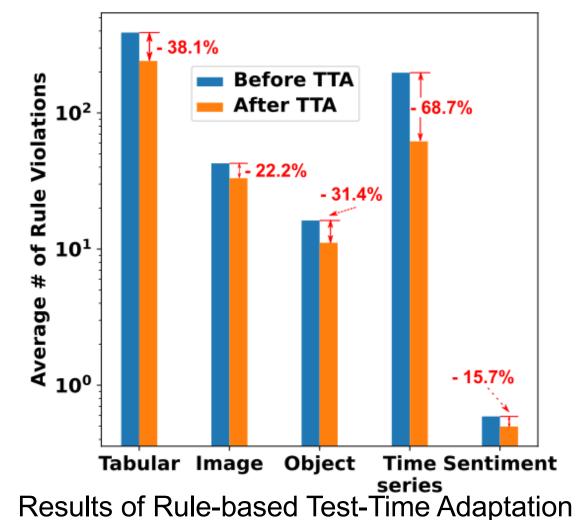
Sentiment Analysis Object Detection Rule Rule $car(x) \Rightarrow 20.22 \le width(x) < 1655.17$ $neutral(x) \Leftarrow 0.02 \leq fitness(x) \leq 0.02$ $\wedge~0.04 \leq \mathrm{news}(x) \leq 0.14$ $\land \ \operatorname{aspect_ratio}(x) < 0.81 \lor \dots$ If the probability that the sentence is about fitness and If an object is a car, then its aspect ratio and width must lie in one of the bounds defined by news is within the above bounds then the sentiment of the rules above. this sentence is **neutral**. **Original Prediction Original Prediction** Anttila's online department store - NetAnttila - has an established position as the best-known, most visited, and most shopped online store in Finland. fitness and health: 0.0210 news and social concern: 0.077 predicted label: positive (wrong) **Prediction after Test Time Adaptation Prediction after Test Time Adaptation** Anttila's online department store - NetAnttila - has an established position as the best-known, most visited, and most shopped online store in Finland.

fitness and health: 0.0210

news and social concern: 0.077

predicted label: neutral (correct)

Examples of Rule Violations



Task	Total	Selected	Violations per sample
Tabular Classification	292,129	400	389.29
Image Classification	73,032	340	42.62
Object Detection	252	252	16.21
Time-Series Imputation	35	35	197.46
Sentiment Analysis	158	158	0.59

Rules

Number of SQRs generated by SQRL.

Examples of Rule Violations

Formalized SQRs to characterize and identify basic errors at scale and proposed the SQRL framework that can find up to 300K rules and up to 158K violations.

We find that models do not always learn statistical rules but can be adapted to correct up to 68.7% rule violations.