# Classifier Explanation
## Introduction to the Algorithms LIME and SP-LIME

Leonhard Applis

TH Nürnberg

21.1.2019

# Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

**Me:** Hey Siri, order me a Pizza

**Siri:** *(After a short break that nearly drains your whole battery)* Ok, I'm calling your mother...

**Me:** Wait! Why would you do this!?

**Siri:** This is the 5th time you ordered Pizza this week.

What do we want from our model?

1. Why did failed predictions fail?
2. Why did correct predictions succeed?
3. Why is my model uncertain about a prediction?

**special importance:**
**setting a model *live*, where it's not *prelabeled***

Interpretations must be ...

- *human-readable*
- reproducable (same input + same model $\rightarrow$ same output)
- **model agnostic**, meaning they can work with any (black-box) model

Difficulties:

- Models can be huge (millions of weights)
- Inputvectors can be huge (e.g. images)
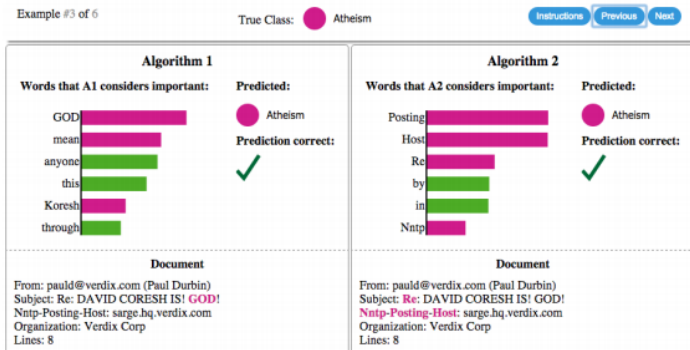- Some models are to complex by it's structure to be readable, (e.g. neural networks)

Figure: LIME-Text: predicting "Atheism" for given text

Both algorithms predict correct - yet Algorithm 2 has strange reasons.

trusting predictions $\neq$ trusting a model

What do we want?

1. get an *overview* of our Model
2. compare models in reasonable time
3. proove correctness & flaws of a model
4. improve our models

Several topics which benefit from machine learning, but need special care:

1. Terrorism-detection
2. Medical diagnosis & prescriptions
3. Fraud-detection

Noone will buy a model, if you can't prove that it's performing reasonable predictions.

There are several issues, at which explanations can help you improve your models:

1. Filtering of Features
2. Find overfitted weighting of features
3. Find Links in Classification (Similiar Classes and Features)

Gaining insights from explanations can help you improve your model!

# Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

## What do we want:

- Human Readable Model Explanation
- For Every Classifier
- For Every Input

### features $\neq$ human readable

To gain *readability*:

- show influence relative to each other, not as numbers
- only show most important features
- use *superpixels* instead of pixels

Let:

1. $G$ be any possible explanation model
2. $g$ be our explanation Model
3. $\Omega(g)$ the complexity of our Model
   - Weights in a regressions model
   - Depth of an decisiontree
   - Number of trees in a random forest
4. $f : Features -> Class$ be the real classification
5. $\Pi_x(z)$ as proximity-measure from $x$ to $z$
6. $\mathcal{L}(f, g, \Pi_x)$ measure of un-faithfullness of $g$ compared to $f$ given the proxmity $\Pi_x$

Wanted:

$$\xi(x) = argmin_{g \in G} \ \mathcal{L}(f, g, \Pi_x) + \Omega(g)$$

Read:

- We want for every input $x$
- an explanation(-model)
- where complexity of $g$ and the failure of $g$ are minimal
- given a set of possible explanations $G$

We do so by picking samples $x$' as subsets from an input $x$ and optimizing our model $g$

Here is the Pseudocode.

Put the funky red-blue image with the red-crosses from the paper here

maybe: Example

# Table of Contents

Explain that we have to little time to inspect everything Looking for a
new way to pick samples

Here is the Pseudocode

# SPLIME Example

I guess this needs more than 2 Pages, we should add an example

# Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Intro
LIME
SP-
LIME
Example:
Traffic
Sign
Recogni-
tion

Setup Problem, Show Code, Plot Examples, nice This could be left out from the presentation, and just be a live demo
Do both: LIME and ANCHOR and sample with SPLIME