

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

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

1 Basics

2 LIME

3 SP-LIME

4 Example: Traffic Sign Recognition

Trusting a Prediction

Intro

Classifier
Expla-
nation

Leonhard
Applis

Basics

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*
- reproducible (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 too complex by its structure to be readable, (e.g. neural networks)

Example

Desired output of a "Atheism"-Classifier

Classifier
Explan-
ation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

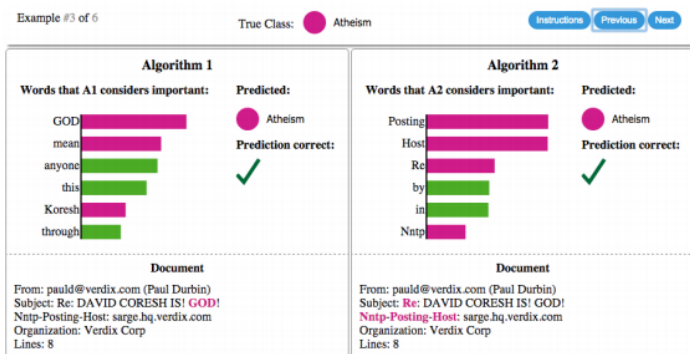


Figure: LIME-Text: predicting "Atheism" for given text

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

trusting predictions \neq trusting a model

What do we want?

- ➊ get an *overview* of our Model
- ➋ compare models in reasonable time
- ➌ proove correctness & flaws of a model
- ➍ 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:

- ❶ Filtering of Features
- ❷ Find overfitted weighting of features
- ❸ 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

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

① Basics

② **LIME**

③ SP-LIME

④ Example: Traffic Sign Recognition

Summary: What do we want

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

Human Readable Model Explanation For Every Classifier For Every Input

All the stuff from the Lime Paper with the fancy ∇ and λ and greeky
cheeky mumbo jumbo

The part talking about making a minimal complex model which is quite accurate

Local Interpretable Model-Agnostic Explanations

The LIME-Algorithm

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

Here is the Pseudocode.

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

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

maybe: Example

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

① Basics

② LIME

③ SP-LIME

④ Example: Traffic Sign Recognition

Problem with Sampling

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

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

Submodular Pick

The SPLIME Algorithm

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

Here is the Pseudocode

SPLIME Example

Classifier
Expla-
nation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

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

Table of Contents

Classifier
Explain-
ation

Leonhard
Applis

Basics

LIME

SP-
LIME

Example:
Traffic
Sign
Recogni-
tion

① Basics

② LIME

③ SP-LIME

④ Example: Traffic Sign Recognition

Trafficsign-Recognition

Explaining RandomForest for Textclassification

Classifier
Expla-
nation

Leonhard
Applis

Basics

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