

Classifier Explanation

Introduction to the Algorithms LIME and ANCHOR

Leonhard Applis , Lisa Branz

TH Nürnberg

21.1.2019

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex I:
Spamfil-
ter

1 Basics

2 LIME

3 SP-LIME

4 Ex I: Spamfilter

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?

- ❶ Why did failed predictions fail?
- ❷ Why did correct predictions succeed?
- ❸ 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

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Maybe just pick the Atheist Christian Example from the paper

Trusting a Model

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Whats the next step after trusting a single prediction
Mention time-required for manually checking a lot of pictures

Example

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Prooving a Model

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Classifier Explanation is required to set anything up IRL, or everyone will hate you very bad

Improving a Model

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

How Classifier Explanation helps you to improve your models
Better Filtering Maybe better weighting of features Finding Links in
Classification (Similiar Classes?)

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex I:
Spamfil-
ter

① Basics

② **LIME**

③ SP-LIME

④ Ex I: Spamfilter

Summary: What do we want

Classifier
Explan-
ation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Human Readable Model Explanation For Every Classifier For Every Input

Definitions

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

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 ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Here is the Pseudocode.

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

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

maybe: Example

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex I:
Spamfil-
ter

1 Basics

2 LIME

3 SP-LIME

4 Ex I: Spamfilter

Problem with Sampling

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

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 ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

Here is the Pseudocode

SPLIME Example

Classifier
Expla-
nation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex 1:
Spamfil-
ter

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

Table of Contents

Classifier
Explan-
ation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex I:
Spamfil-
ter

① Basics

② LIME

③ SP-LIME

④ Ex I: Spamfilter

Spamfilter

Explaining RandomForest for Textclassification

Classifier
Explan-
ation

Leonhard
Applis ,
Lisa
Branz

Basics

LIME

SP-
LIME

Ex I:
Spamfil-
ter

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