

Classifier Explanation

Introduction to the Algorithms LIME and SP-LIME

Leonhard Applis

TH Nürnberg

21.1.2019

Table of Contents

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

4 SP-LIME

Trusting a Prediction

Intro

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

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

Desired output of a "Atheism"-Classifier

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

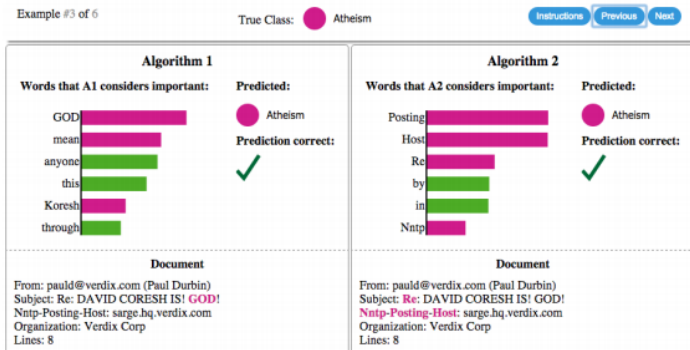


Figure: LIME-Text: predicting "Atheism" for given text

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.

Improving a Model

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

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

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

① Introduction

② LIME

③ Example: Traffic Sign Recognition

④ SP-LIME

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:

- ❶ G be any possible explanation model
- ❷ g be our explanation Model
- ❸ $\Omega(g)$ the complexity of our Model
 - Weights in a regressions model
 - Depth of an decisiontree
 - Number of trees in a random forest
- ❹ $f : Features \rightarrow Class$ be the real classification
- ❺ $\Pi_x(z)$ as proximity-measure from x to z
- ❻ $\mathcal{L}(f, g, \Pi_x)$ measure of un-faithfulness of g compared to f given the proximity Π_x

Wanted:

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

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 ¹

¹We do not really check different models, we train one

Local Interpretable Model-Agnostic Explanations

The LIME-Algorithm

Classifier
Explan-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

Additional Requirements:

LASSO² - *Least Absolute Shrinkage and Selection Operator*

Machine Learning algorithm to select most important features relative to each other.

G are only *sparse linear regression models* (e.g. Decision Trees or simple logistic regression)

Require: Classifier f , Number of samples N

Require: Instance x , and its interpretable version x'

Require: Similarity kernel π_x , Length of explanation K

$Z \leftarrow \{\}$;

foreach $i \in \{1, 2, \dots, N\}$ **do**

$z'_i \leftarrow \text{sample_around}(x')$;

$Z \leftarrow Z \cup z'_i, f(z_i, \pi_x(z_i))$;

end

$w \leftarrow K - \text{Lasso}(Z, K) \triangleright$ with z'_i as features, $f(z)$ as target;

return w ;

²Further Reading:

Table of Contents

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

4 SP-LIME

Trafficsign-Recognition

Explaining RandomForest for Textclassification

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

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

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

① Introduction

② LIME

③ Example: Traffic Sign Recognition

④ SP-LIME

Problem with Sampling

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

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

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

Here is the Pseudocode

SPLIME Example

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

SP-
LIME

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