

Classifier Explanation

Introduction to the Algorithms LIME and SP-LIME

Leonhard Applis

TH Nürnberg

21.1.2019

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

Table of Contents

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

Trusting a Prediction

Intro

Classifier
Expla-
nation

Leonhard
Applis

Intro

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?

- 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

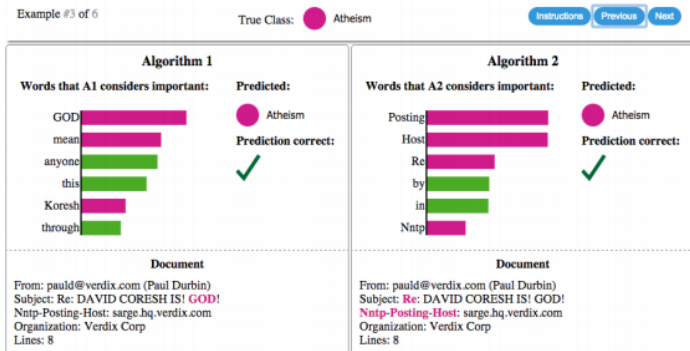


Figure: LIME-Text: predicting "Atheism" for given text

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

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

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
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

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)

²Further Reading:

https://beta.vu.nl/nl/Images/werkstuk-fonti_tcm235-836234.pdf

Table of Contents

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

1 Introduction

2 LIME

3 Example: Traffic Sign Recognition

given neural network :

- 6 layers, first 3 convolutional (complex structure)
- 43 Classes (complex output)
- 64x64 Images (complex input)
- trained with 8k images (rich data-input)
- tested with 2k images reaching 95% accuracy (good?)

The NN was trained with Tensorflow and is shipped with your notebook.

```
from lime import lime_image
from skimage.segmentation import mark_boundaries

#Setup the Explainer
explainer = lime_image.LimeImageExplainer()

#Explain the predictions
explanation = explainer.explain_instance(
    image, model.predict, top_labels=43, hide_color=0,
    num_samples=1000)

#Show the mask for a class
temp, mask = explanation.get_image_and_mask(
    10, positive_only=True, num_features=5, hide_rest=False)

plt.imshow(mark_boundaries(temp / 2 + 0.5, mask))
```

Trafficsign-Recognition

Simple Classification

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

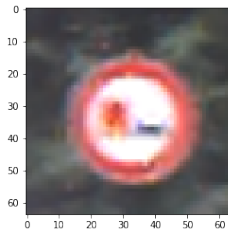


Figure: No Overtaking -
Sample Image from
Test-data

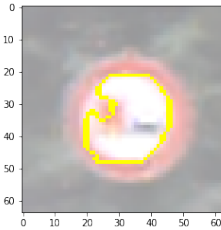


Figure: Prediction -
showing the 5 Superpixels
for *no overtaking*

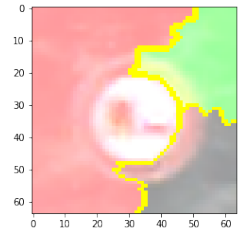


Figure: Prediction -
showing the 4 most
important Superpixels for
right of way crossing

Trafficsign-Recognition

Overfitting

Classifier
Expla-
nation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

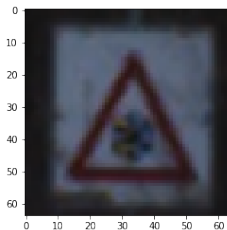


Figure: Frost - Sample
Image from *Training-data*

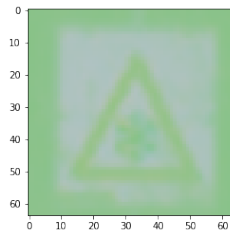


Figure: Prediction for
Frost, a sign for overfitting

Trafficsign-Recognition

Similar Classes I

Classifier
Explain-
ation

Leonhard
Applis

Intro

LIME

Example:
Traffic
Sign
Recogni-
tion

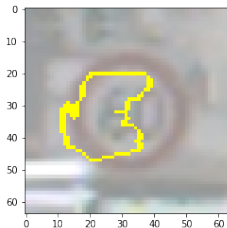


Figure: Prediction: 60 -
only 6 is circled

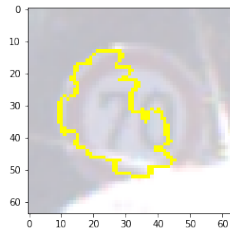


Figure: Prediction: 70 -
only 7 is circled

The model seems to recognize numbers!

Let's have some fun!

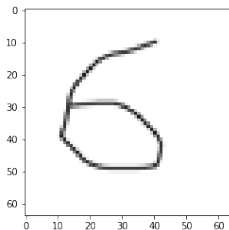


Figure: Only number 6, no
Street Sign

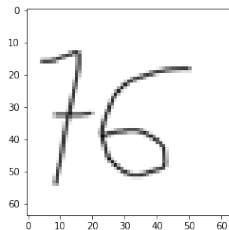


Figure: Number 76 - what
will be predicted?

Let's have some fun!

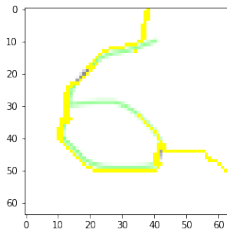


Figure: Number 6 -
78% *Speed Limit 60*

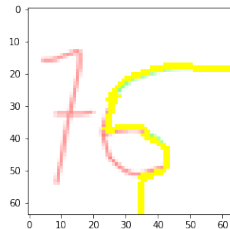


Figure: Number 76 -
99.9% *No Overtaking*

- Accuracy \neq Quality
- most of the predictions *look good*
- trainings-data is heavily overfitted
- everything that is not a streetsign causes trouble
- there can (still) be much more *hidden* problems