

Label-Descriptive Patterns and Their Application to Characterizing Classification Errors

Michael Hedderich*



Jonas Fischer*



Dietrich Klakow



Jilles Vreeken





Blackbox
Classifier

Task, e.g.
Visual Question Answering



Blackbox Classifier



Instance	Correct Prediction?
Are there many ducks playing?	✓
How many ducks are in the picture?	✗
What are the ducks eating?	✗
Do you see ducks in the puddle?	✓
How many roosters are in the puddle?	✗

Task, e.g.
Visual Question Answering



Blackbox
Classifier



Instance	Correct Prediction?
Are there many ducks playing?	✓
How many ducks are in the picture?	✗
What are the ducks eating?	✗
Do you see ducks in the puddle?	✓
How many roosters are in the puddle?	✗

Task, e.g.
Visual Question Answering

Why?

Characterize classification *errors* with *patterns*

- global over all instances
- non-redundant
- easy to interpret
- actionable

Length in bits of sending
data and model

$$\operatorname{argmin}_{M \in \mathcal{M}} L(D, M)$$

Model M over
model class \mathcal{M}

Data is given by instances of correct and wrong classifications with corresponding label

Model is composed of label specific patterns

$$\operatorname{argmin}_{M \in \mathcal{M}} L(D, M)$$

D

How many ducks are in the picture?

X

How many roosters can you see?

X

What colour is the water?

X

What color are the ducks?

X

Do you see a rooster in the puddle?

✓

Are there many ducks playing?

✓

D^-

D^+

M

How and many

What and (color xor colour)

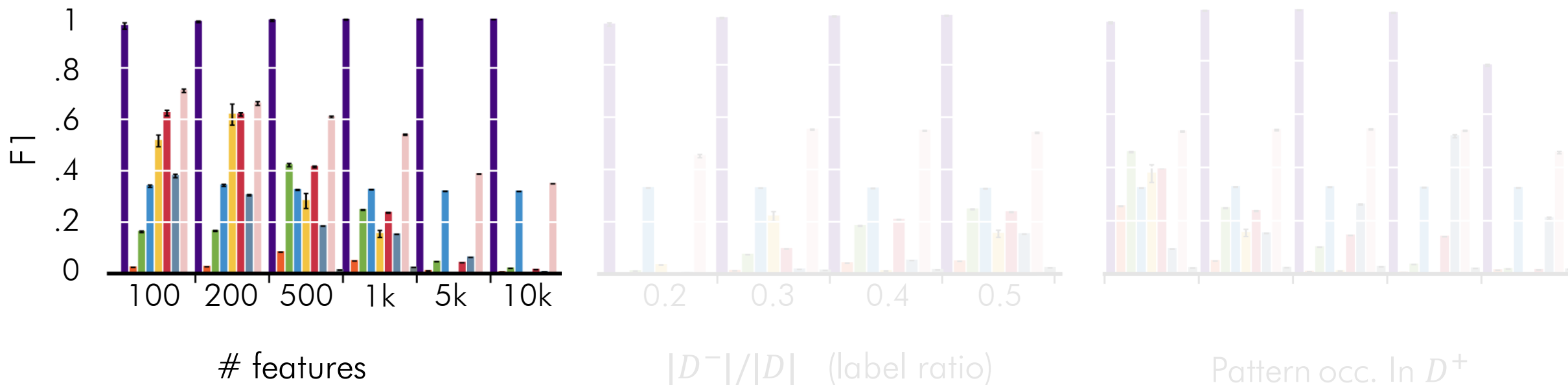
Data is given by instances of correct and wrong classifications with corresponding label

Model is composed of label specific patterns

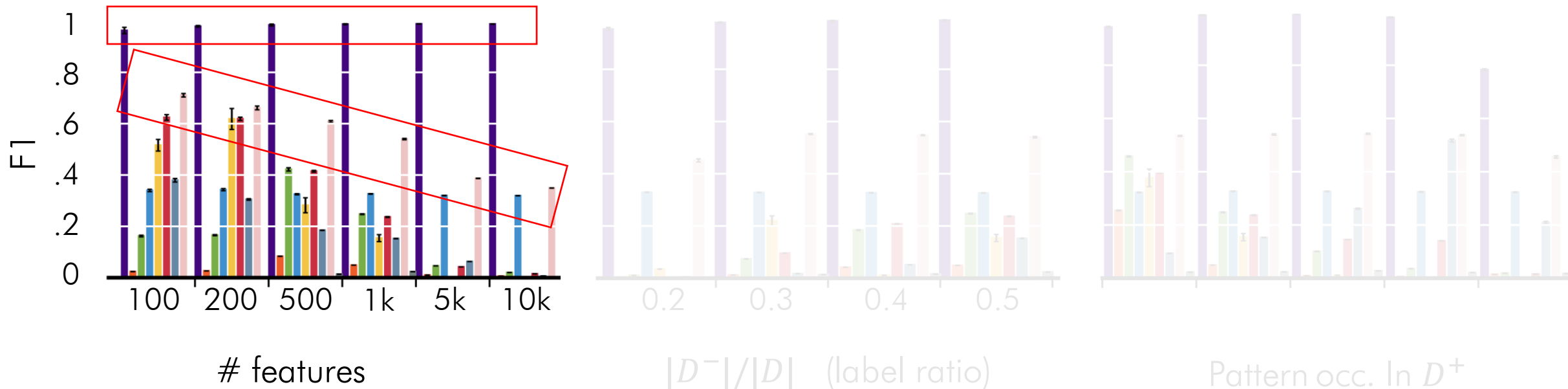

$$\operatorname{argmin}_{M \in \mathcal{M}} L(D, M)$$

For efficient search in practice: **PREMISE**

- iteratively explores the pattern space in bottom-up fashion
- uses word embeddings to explore mutually exclusive patterns
- employs statistical testing for filtering spurious patterns

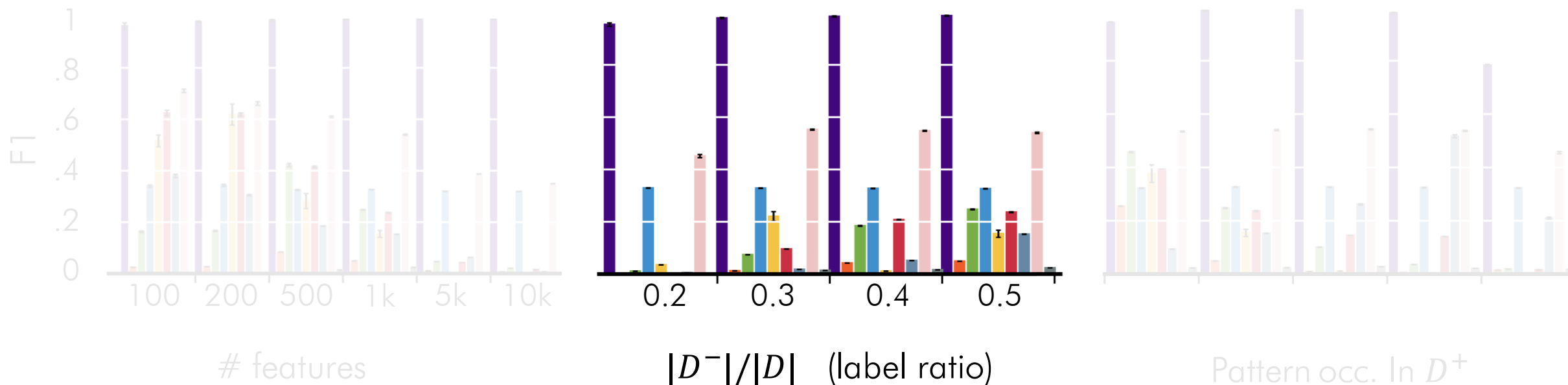


■ PREMISE
 ■ TREE
 ■ RIPPER
 ■ SUBGROUP
 ■ GRAB
 ■ CLASSY
 ■ SPUMANTE
 ■ CORTANA
 ■ C-SALT



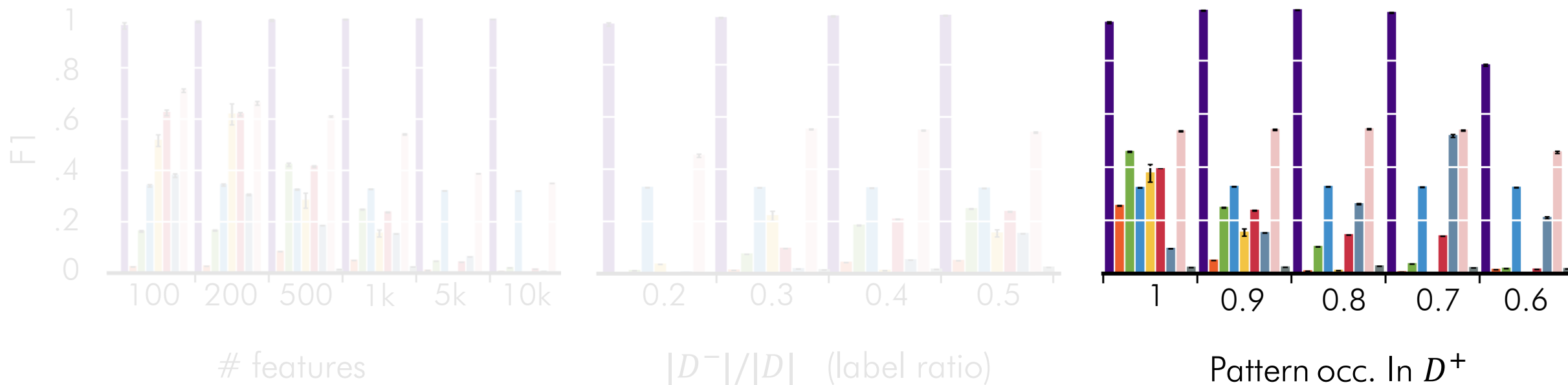
Current methods discover more spurious patterns with increase in features.

- PREMISE
- TREE
- RIPPER
- SUBGROUP
- GRAB
- CLASSY
- SPUMANTE
- CORTANA
- C-SALT



More extreme label ratios are an issue for many methods.

■ PREMISE
 ■ TREE
 ■ RIPPER
 ■ SUBGROUP
 ■ GRAB
 ■ CLASSY
 ■ SPUMANTE
 ■ CORTANA
 ■ C-SALT



Many more experiments in the paper!

■ PREMISE ■ TREE ■ RIPPER ■ SUBGROUP ■ GRAB ■ CLASSY ■ SPUMANTE ■ CORTANA ■ C-SALT

LXMERT (Tan & Bansal, 2019) 70% VQA accuracy



LXMERT (Tan & Bansal, 2019) 70% VQA accuracy

How and many
hanging and from
(kind xor sort) and of
(would xor could xor might xor can) and you
number
letter xor letters

← Basic math
(counting,...)

More experiments on VQA and how to improve an NER classifier in our paper!

LXMERT (Tan & Bansal, 2019) 70% VQA accuracy

How and many

hanging and from

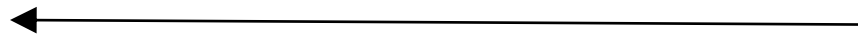
(kind xor sort) and of

(would xor could xor might xor can) and you

number

letter xor letters

Spatial
reasoning



More experiments on VQA and how to improve an NER classifier in our paper!

LXMERT (Tan & Bansal, 2019) 70% VQA accuracy

How and many

hanging and from

(kind xor sort) and of

(would xor could xor might xor can) and you

number

letter xor letters

Ontology
relationships



More experiments on VQA and how to improve an NER classifier in our paper!

LXMERT (Tan & Bansal, 2019) 70% VQA accuracy

How and many

hanging and from

(kind xor sort) and of

(would xor could xor might xor can) and you ← Soft questions

number

letter xor letters

More experiments on VQA and how to improve an NER classifier in our paper!

LXMERT (Tan & Bansal, 2019) 70% VQA accuracy

How and many
hanging and from
(kind xor sort) and of
(would xor could xor might xor can) and you
number
letter xor letters

Recognition
subtask



More experiments on VQA and how to improve an NER classifier in our paper!

Goal:

Discover *easy-to-interpret* patterns that characterize misclassifications.

Method:

Premise, based on Minimum Description Length Principle.

Results:

Consistently discovers non-redundant and descriptive patterns, where state-of-the-art fails. Patterns are easy to understand and to act upon.

Outlook:

Scale to more features to investigate patterns of neuron activations characterizing misclassifications.

<https://github.com/uds-lsv/premise>

 PyPremise



Goal:

Discover *easy-to-interpret* patterns that characterize misclassifications.

Method:

Premise, based on Minimum Description Length Principle.

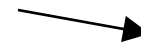
Results:

Consistently discovers non-redundant and descriptive patterns, where state-of-the-art fails. Patterns are easy to understand and to act upon.

Outlook:

Scale to more features to investigate patterns of neuron activations characterizing misclassifications.

Analyze your own classifier!



<https://github.com/uds-lsv/premise>

 PyPremise

