

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

Top-Down and Bottom-up Cues for Scene Text Recognition, A. Mishra, K. Alahari, C. V. Jawahar

Vincent BODIN & Thomas MOREAU

December 17, 2013

Introduction

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first



Figure: Task of the project (image from SVT)

Sommaire

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

1 Learning preliminaries

- Learning characters
- Learning words

2 Character detection

3 Recognizing words

4 Implementation and first results

Learning characters

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

We need to build classifiers to recognize characters in a natural picture:

- ① Character databases: ICDAR 2003 [1], Chars74K [2];
- ② features: Histogram Of Gradient (HOG) [3]¹;
- ③ build K SVMs ($K = 62$) with RBF Kernel, Fig.(2):

$$\exp(-\gamma|x - x'|^2), \gamma > 0 \quad (1)$$

method one-versus-all, we used a Python library scikit-learn [4]: K coefficients γ and regularization C optimized by cross-validation => test error without optimization 25%, with optimization 16%.

¹N. Dalal, B. Triggs, Histogram of oriented gradients for human detection

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

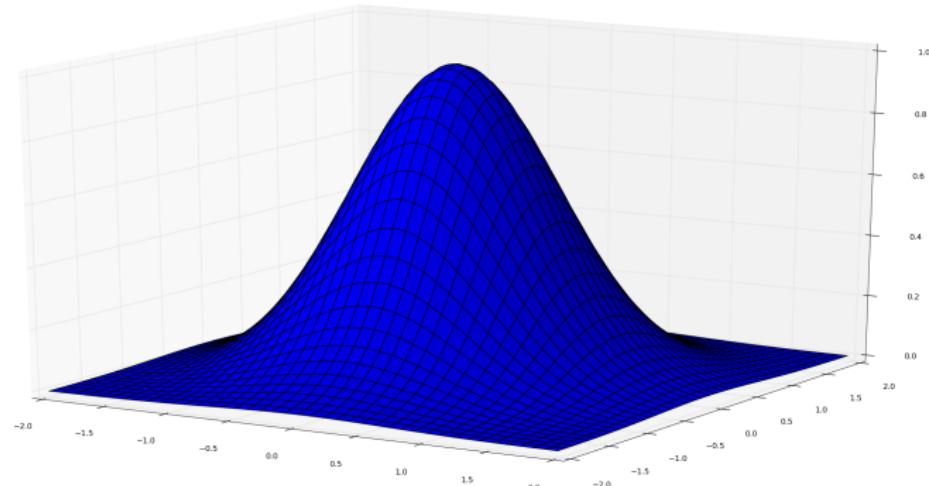


Figure: RBF Kernel

Learning words

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

We have to build a prior lexicon that contains how characters interact to create words.

- ① Word database ICDAR 2003 word [5]²;
- ② compute for each pair of characters (c_i, c_j) their frequency of occurrence $p(c_i, c_j)$ in the database (bi-gram model);
- ③ will be used as an energy term to be minimized afterward.

²<http://algoval.essex.ac.uk/icdar/Datasets.html>

Sommaire

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

1 Learning preliminaries

2 Character detection

3 Recognizing words

4 Implementation and first results

Sliding window and pruning

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

- Scan the image with sliding windows;
- for each window I_i , compute ϕ_i , feature of HOG with 12 orientations;
- run the 62 SVMs on ϕ_i and compute the goodness score:

$$GS(I_i) = \max_c p(c|\phi_i) \exp \left(-\frac{(a_i - \mu_{aj})^2}{2\sigma_{aj}^2} \right) \quad (2)$$

a is the aspect ratio and j is the class that reaches the maximum.

Sliding window and pruning

Top-Down
and

Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

- if $GS(I_i) > \text{threshold}$, keep the window **and** all the probabilities of each character.
- apply Non-Maximum Suppression algorithm (NMS) to merge overlapping windows:
 - ① find I_i with maximum confidence.
 - ② for all others I_j compute:

$$\text{criterion} = \frac{I_i \cap I_j}{I_i \cup I_j} \quad (3)$$

- ③ if criterion $> \text{threshold}$ **and** highest character is the same then merge the windows.
- ④ the ending merged window has the highest confidence and is located at the barycenter of all the windows merged.

Sommaire

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

1 Learning preliminaries

2 Character detection

3 Recognizing words

4 Implementation and first results

Graph construction and retrieval of the word

Top-Down
and

Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

PGM part

Technical part of graphical model implemented for PGM course,
here is a (very) short summary.

Input. set of n windows with possible characters.

- ➊ for each window, assign a node taking values in \mathcal{K}_ϵ (ϵ void label);
- ➋ build edges if the windows are 'close enough';
- ➌ assign unary energy to nodes, **and** pairwise energy for edges (overlapping terms + lexicon prior):

$$E(\mathbf{x}) = \sum_{i \in \mathcal{V}} E_i(x_i) + \sum_{(x_i, x_j) \in \mathcal{E}} E_{i,j}(x_i, x_j) \quad (4)$$

- ➍ minimize the discrete energy on the graph (NP-hard) with TRW-S algorithm [6].

Sommaire

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries
Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

1 Learning preliminaries

2 Character detection

3 Recognizing words

4 Implementation and first results

Implementation and first results



Figure: Image for testing the algorithm

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

Implementation and first results

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

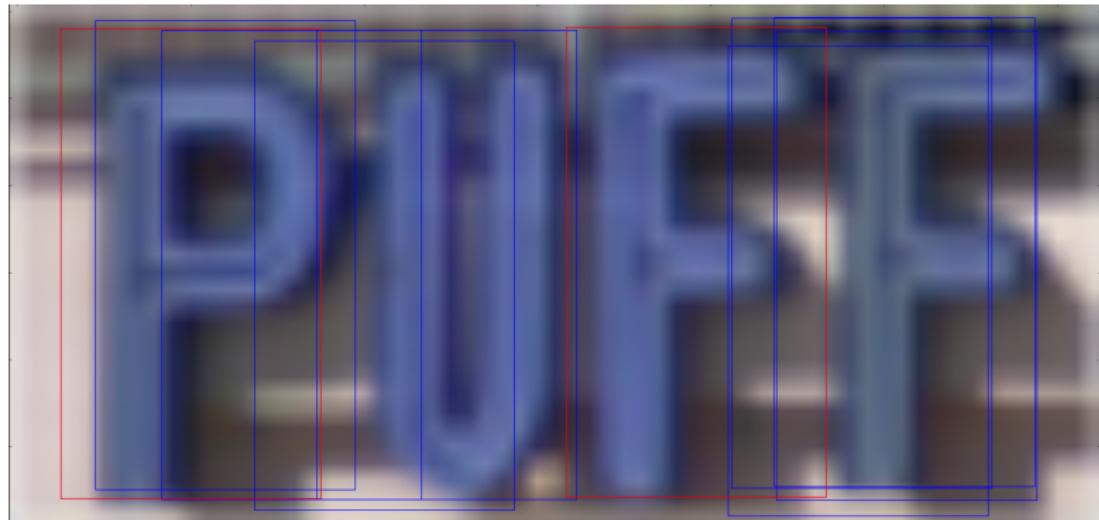


Figure: Word retrieved: 'PE'

Implementation and first results

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

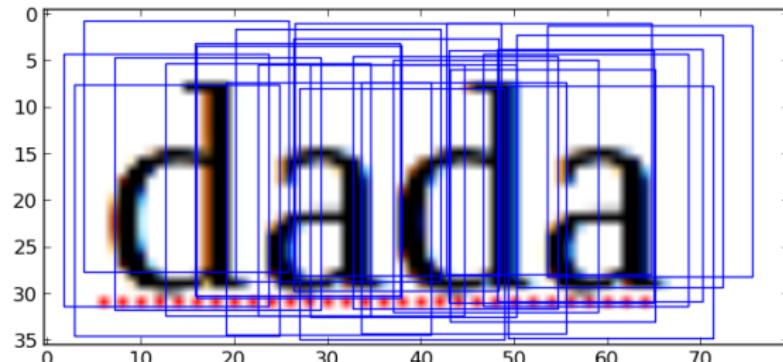


Figure: All the sliding windows kept to build graphical model

Implementation and first results

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

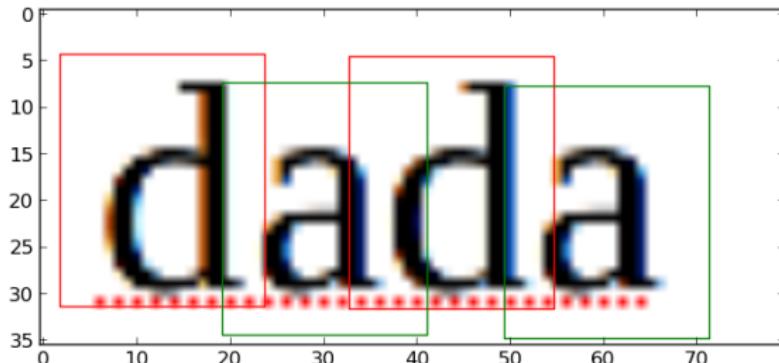


Figure: Presence of the letters 'a' (green) and 'd' (red)

Implementation and first results

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

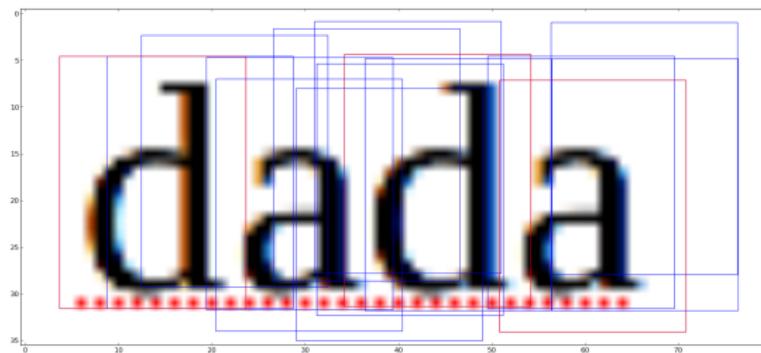


Figure: Word retrieved 'ddaa' -> the second a has a too big overlap
(penalized by pairwise energy, PGM part)

Improvements

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

- ① **Difficult problem:** for each window, classify K times.
- ② SVM is not good enough (albeit a 14-hours cross-validation): error on the training set is still 16%.
- ③ parameters of the algorithm: $2 \times K$ for SVMs + 2 for HOG (bins, number of orientations) + threshold for GS + threshold in NMS + parameters of energy and of TRW-S
=> complex algorithm.
- ④ (PGM part) The energy tends to explode letters instead of joining them into words (example 'dda').

Improvements

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first

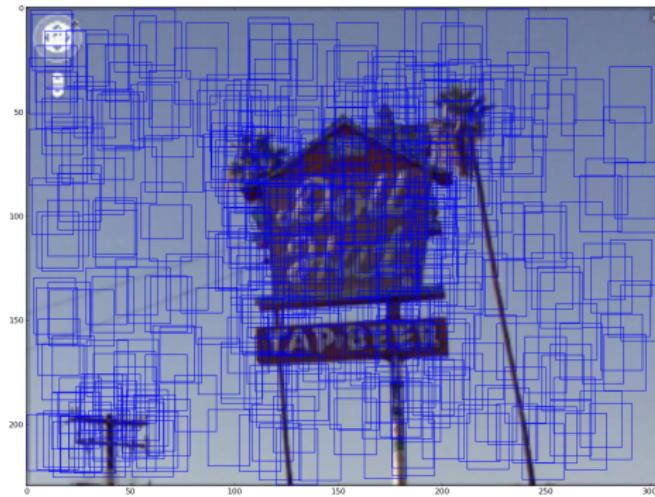


Figure: Flexible GS threshold... too many windows

References

Top-Down
and
Bottom-up
Cues for
Scene Text
Recognition,
A. Mishra, K.
Alahari, C.
V. Jawahar

Vincent
BODIN &
Thomas
MOREAU

Introduction

Learning
preliminaries

Learning
characters
Learning
words

Character
detection

Recognizing
words

Implementation
and first



Robust ocr.

<http://algovat.essex.ac.uk/icdar/Datasets.html>.

Accessed: 2013-12-16.



The chars74k dataset: Character recognition in natural images.

<http://www.ee.surrey.ac.uk/CVSSP/demos/chars74k/>.

Accessed: 2013-12-16.



B. Triggs N. Dalal.

Histogram of oriented gradients for human detection.

CVPR, 2005.



Scikit-learn.

<http://scikit-learn.org/stable/modules/svm.html#parameters-of-the-rbf-kernel>.

Accessed: 2013-12-16.



Robust word recognition.

<http://algovat.essex.ac.uk/icdar/Datasets.html>.

Accessed: 2013-12-16.



V. Kolmogorov.

Convergent tree-reweighted message passing for energy minimization.

IEEE Transactions on Pattern Analysis and Machine Learning (PAMI), 2006.