
OrationesPython Documentation

Release 1.0.1

Timo Mätäsaho

December 13, 2013

CONTENTS

1	Documentation for the Code	3
1.1	Orationes Search	3
1.2	Orationes getBoxes	3
1.3	Utility functions	4
1.4	Helper Functions	10
2	Indices and tables	11
	Python Module Index	13
	Index	15

Contents:

DOCUMENTATION FOR THE CODE

1.1 Orationes Search

`osearch.init(sysargv)`

The first function to be called when this script is run. Checks if the system arguments are correct and returns error when needed.

When calling from PHP use the `$output` and `$return_var` arguments to catch the correct output and possible return error code. If the `$return_var` is 1 then the program has finished without errors.

Refer to PHP [exec](#) command manual for further information of calling Python scripts and programs in general from PHP.

`osearch.osearch(img, switch, txf, sw)`

The main program that only calls the processing methods from OratUtils and HFun classes.

Parameters

- **img** (*string*) – The name of the image.
- **txf** (*string*) – The name of the cleaned XML file.
- **sw** (*string*) – The word or letter that is searched.

Returns error code or 1 and JSON string

The error codes are: 2 - Error in starting parameters 3 - Given image doesn't exist, it's path is faulty or its format is wrong or either the image file is corrupted. 4 - Given string or character to be searched cannot be found from the XML 5 - 6 - 7 - 8 - 9 - Unknown error while opening the image

There are still lots of places that are missing error handling!

The return option have to be chosen between returning the string as a return code or is it just printed out for the calling PHP program. Currently it is being printed.

1.2 Orationes getBoxes

`getBoxes.getBoxesAndLines(img)`

Optional directly callable program that can be used to extract the bounding box and line location information from an image.

Parameters **img** (*string*) – The name of the image.

Returns JSON string

Returns a JSON array containing the possible locations of the text lines and bounding boxes.

1.3 Utility functions

class OratUtils.**OratUtils**

This class contains only static utility methods that are called directly from the main program ‘osearch.py’.

static boundingBox (*image*, *debug*)

This functions tries to determine the bounding boxes for each text line.

Parameters

- **image** (*ndarray*) – the processed image
- **debug** (*bool*) – debug switch

Returns ndarray – bboxes

$$bboxes_{n \times m} = \begin{bmatrix} \text{patch label numbers} \\ \text{starting x-coordinates} \\ \text{starting y-coordinates} \\ \text{ending x-coordinates} \\ \text{ending y-coordinates} \end{bmatrix}$$

debug switch can be used to plot the results of the bounding box founding method and to see whether it is working correctly.

Process pipeline:

1. Calculate the histogram from the image
2. Binarize image with threshold 0.95
3. Label all the patched in on the binarized image
4. Calculate the sizes of the patches
5. Remove unnecessary patches
 - (a) Remove the largest patch. The largest patch is always the patch consisting of the borders and marginals.
 - (b) Remove patches which size is smaller or equal to 50 pixels
 - (c) Remove all the patches which are higher than 70 pixels. This removes the possible remaining marginal patches which weren’t connected to the major marginal and border patch.
6. Perform morpholog operations to clean the image and bind the text lines together
 - (a) Perform erosion with a cross like structure element

$$SE_{e_{5,5}} = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

(b) Perform dilation with a long vertical line. (needs a 70x70 size structure element)

$$SEd_{70,70} = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ & \vdots & & \vdots & \\ 1 & 1 & \dots & 1 & 1 \\ & \vdots & & \vdots & \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

7. Label the morphologically operated image

8. Remove patches which size is less or equal to 4000 pixels

9. Label the image again with new labels

10. Calculate the extreme dimensions of each patch. These values are used as the limiting bounding boxes.

11. Combine the boxes which are horizontally too close as they are thought to be separate boxes on the same textline.

12. Return the bounding boxes

static `contStretch` (*im*, *a*, *h*)

Performs contrast stretching for grayscale images. Pixel intensities are set to differ 'a' times the average intensity from the original intensity values. The new intensity values are sliced to stay between [0, 255].

$$I_{stretched} = I_{old} + a * (I_{old} - I_{average})$$

$$I_{new} = \begin{cases} 0, & I_{stretched} < 0 \\ I_{stretched}, & 0 \leq I_{stretched} \leq 255 \\ 255, & I_{stretched} > 255 \end{cases}$$

Parameters

- **im** (*ndarray*) – The image which contrast is to be stretched
- **a** (*int*) – multiplication coefficient
- **h** (*int*) – The height of image. Used as partial image average switch

Returns *ndarray* – contrast stretched image

Parameter *h* is a switch which could be used to determine if the average intensity is calculated over the whole image or from a small portion of it. Currently it is defaulted in the code to newer happen. Originally the idea was that if the image is very big, the intensity average would be taken from a small sample. To make the function more generic and also because of the nature of the images in Orationes project, it was decided that the average is always calculated over the whole image.

static `findCorr` (*bboxes*, *slines*, *charcount*, *imlines*, *debug*)

Used to find out which bounding box and which line are the same.

Parameters

- **bboxes** (*ndarray*) – Narray containing the coordinates of the bounding boxes
- **slines** (*ndarray*) – A vector containing the y-coordinates of the interesting lines
- **charcount** (*list*) – Contains the lengths of each line

- **imlines** (*ndarray*) – $n*1$ size ndarray containing the lines (or rather their y-position) got from the image by radontransform
- **debug** (*bool*) – Debug switch

Returns ndarray – $m*n$ ndarray containing the starting and ending coordinates of hits

static hfilter (*image, diameter, height, length, n*)

This function performs homomorphic filtering on grayscale images.

Parameters

- **image** (*ndarray*) – 2-dimensional ndarray
- **diameter** (*int*) – filter diameter
- **height** (*int*) – Height of the image
- **length** (*int*) – Length of the image
- **n** (*int*) – Filter order

Returns ndarray – homomorphically filtered image

The image must in ndarray format. In osearch PIL images are converted to scipy images which are in ndarray format. Ndarray format allows easy and fast direct access to the pixel values and this function is written entirely only for the ndarrays.

static packBoxesAndLines (*bboxes, imlines*)

Parameters

- **bboxes** (*ndarray*) – Ndarray containing the coordinates of the bounding boxes
- **imlines** (*ndarray*) – $n*1$ size ndarray containing the lines (or rather their y-position) got from the image by radontransform

Returns jsondata – JSON packed string containing the found bounding boxes and lines

static packCoordsToJson (*slines, origimage, coords, charpos, wordlens, bboxes, debug*)

This function is used to pack the hit coordinates and bounding box coordinates into a JSON string which is returned to the calling PHP site.

Parameters

- **slines** (*ndarray*) – A vector containing the y-coordinates of the interesting lines
- **origimage** (*ndarray*) – Original image. Used when debugging
- **coords** (*ndarray*) – A ndarray containing the coordinates of the hits
- **charpos** (*list of lists*) – List of the character positions got from the XML
- **wordlens** (*list of lists*) – List of wordlengths
- **bboxes** (*ndarray*) – Ndarray containing the coordinates of the bounding boxes
- **debug** (*bool*) – Debug switch

Returns json-string – JSON packed string containing the hits and the bounding boxes

static padlines (*imlines, llines, charlines*)

Parameters

- **imlines** (*ndarray*) – $n*1$ size ndarray containing the lines (or rather their y-position) got from the image by radontransform
- **llines** (*ndarray*) – $n*2$ size ndarray containing the length information of the lines

- **charlines** (*list*) – list of lists telling the position(s) of searched character(s)/word(s) on each line

Returns ndarray – wantedlines

Long: Llines contains the information about the lines got from the XML and also it contains the information if some of the lines are remarkably shorter than other lines. That means that, if there are some lines that are not found from the image, it is assumed that those non-found lines are the shortest lines according to the XML and character count. Those lines are marked as 0 in the second column in llines.

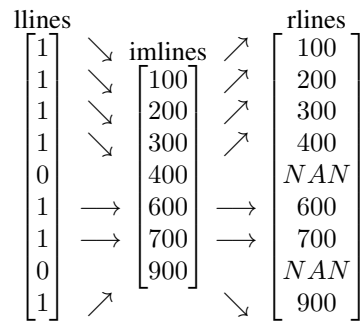
Short: Llines[:,1] contains only 1s and 0s. 1 meaning a line with enough letters to be recognized by poormanradon (pmr) and 0 meaning a line which is probably undetected by pmr

Behavior: Number of lines found from the image using pmr is larger than the number of lines calculated from XML:

TODO! Currently this case is not handled!

Number of lines found from the image using pmr is smaller than the number of lines calculated from XML:

Pad the lines according to the information in *llines[:,1]*:



Number of lines found from the image using pmr equals to the number of lines calculated from XML:

Pick unique lines from imlines and return them as lines the interesting lines.

static poormanradon (*image, iname, height, debug*)

Performs a naive radon-transform and peak detection on the binarized and contrast stretched image and tries to determine where the text lines are in the image.

Parameters

- **image** (*ndarray*) – Image
- **iname** (*string*) – Image name
- **height** (*int*) – Image height
- **debug** (*bool*) – Debug switch

Returns ndarray – Array containing the lines which are found using radon transform

Calculates the intensity sums over each vertical line. The sums are then inverted and peaks are detected from the inverted data. Data inversion wouldn't be necessary in the python code, but this convention comes from the Matlab code that was ported to python.

Before the transform, the image is cleaned so that by using static values (very bad, should be dynamic, but so far there hasn't been time to do that) the marginals and everything outside them is erased and turned to white. Because the distance between the camera and the page differs in each image, the marginals aren't always on the same position. This combined with static values causes inaccuracy in the erasing process and might cause inaccuracy when detecting the peaks and the lines.

In the peak detection, it is assumed that a spike is considered a peak if it's 25 units away from a previous detected peak and also if its value difference is at least 1500 to its previous value.

upLim in the source means upper limit in the image coordinates, which increase when going down in the image. That's why *upLim* is small. Respectively the *downLim* means the bottom limit in the image coordinates and that's why it's bigger than the upper limit.

static processlines (*charcount*, *imlines*)

This functions compares the number of lines found from the image to the number of lines found from the XML file and creates a logical vector telling which lines are probably found and which are not.

Parameters

- **charcount** (*list*) – Contains the lengths of each line
- **imlines** (*ndarray*) – Contains the textlines which are found in 'poormanradon'

Returns ndarray – *l*lines

Returns ndarray – *im*lines

*l*lines is a $n \times 2$ vector, where the *l*lines[:,1] is a logical vector containing the information of the probably found and non-found lines.

$$[1, 1, 1, 1, 0, 1, 0, 1, \dots]^t$$

*im*lines is a ndarray containing the y-coordinates of the textlines found from the image with poormanradon. When padding some of the coordinates are removed (*nofound* < 0, not used), some NAN values are added in between some coordinates (*nofound* > 0) or it is returned unchanged (*nofound* == 0).

We calculate a number 'nofound'

$$\text{nofound} = \text{lines}_{XML} - \text{lines}_{image}$$

Naturally there are three cases.

nofound < 0: This means that there are more lines found from the image than there actually are. Currently nothing's done here to compensate this behavior.

nofound > 0: This means there aren't enough lines found from the image. Usually the non-found lines are assumed to be the very short lines. When padding the indices of the lines, the shortest lines are always set to be the non-found lines.

nofound == 0: It is assumed that all the textlines were found correctly and the *im*lines will be returned unchanged.

static stringparser (*tfile*, *c*)

Performs case sensitive search for text file *tfile* with string or character *c* (char on default). Argument *c* can be any regular expression

Parameters

- **tfile** (*string*) – The string containing the cleaned XML file as a string
- **c** (*string/char/regular expression*) – The letter or string that is searched from the *tfile*

Returns list – charcount

Returns list of lists – charpos

Returns list of lists – charlines

Returns list of lists – wordlens

- *Charcount* is a list containing the lengths of each line.

–[63, 60, 4, 65, 66, 37, 66, ...]

- *Charpos* is a list containing lists including the positions of the found characters or the first letters of the found words.

–[[52], [10, 47, 62], [19, 62], [51], ...]

- *Charlines* is a list of lists where the length of each sublist tells the number of hits on that line and the element values representing the line number from the XML file.

–[[3], [4, 4, 4], [6, 6], [7], ...]

- *Wordlens* is a list containing lists containing the lengths of the words on each line.

–[[3], [3, 3, 3], [3, 3], [3], ...]

static **txtfparser** (*tfile, c*)

Performs case sensitive search for text file *tfile* with string or character *c* (char on default). Argument *c* can be any regular expression

Parameters

- **tfile** (*string*) – The name of the cleaned XML file
- **c** (*string/char/regular expression*) – The letter or string that is searched from the *tfile*

Returns list – charcount

Returns list of lists – charpos

Returns list of lists – charlines

Returns list of lists – wordlens

- *Charcount* is a list containing the lengths of each line.

–[63, 60, 4, 65, 66, 37, 66, ...]

- *Charpos* is a list containing lists including the positions of the found characters or the first letters of the found words.

–[[52], [10, 47, 62], [19, 62], [51], ...]

- *Charlines* is a list of lists where the length of each sublist tells the number of hits on that line and the element values representing the line number from the XML file.

–[[3], [4, 4, 4], [6, 6], [7], ...]

- *Wordlens* is a list containing lists containing the lengths of the words on each line.

–[[3], [3, 3, 3], [3, 3], [3], ...]

1.4 Helper Functions

class `OratUtils.HFun`

Contains helper functions that are used in various places

static `gray2uint8 (image)`

Converts grayscale images to uint8 type

Parameters `image` (`ndarray`) – The grayscale image to converted

Returns `ndarray(uint8)` - I

static `im2float (image)`

Changes uint8 type images to float64 images.

Parameters `image` (`ndarray(uint8)`) – The image to be converted

Returns `ndarray(float64)` – I

INDICES AND TABLES

- *genindex*
- *modindex*
- *search*

PYTHON MODULE INDEX

g

getBoxes, 3

o

OratUtils, 4

osearch, 3

INDEX

`boundingBox()` (OratUtils.OratUtils static method), 4

`contStretch()` (OratUtils.OratUtils static method), 5

`findCorr()` (OratUtils.OratUtils static method), 5

`getBoxes` (module), 3

`getBoxesAndLines()` (in module `getBoxes`), 3

`gray2uint8()` (OratUtils.HFun static method), 10

`hfilter()` (OratUtils.OratUtils static method), 6

`HFun` (class in OratUtils), 10

`im2float()` (OratUtils.HFun static method), 10

`init()` (in module `osearch`), 3

`OratUtils` (class in OratUtils), 4

`OratUtils` (module), 4

`osearch` (module), 3

`osearch()` (in module `osearch`), 3

`packBoxesAndLines()` (OratUtils.OratUtils static method), 6

`packCoordsToJson()` (OratUtils.OratUtils static method), 6

`padlines()` (OratUtils.OratUtils static method), 6

`poormanradon()` (OratUtils.OratUtils static method), 7

`processlines()` (OratUtils.OratUtils static method), 8

`stringparser()` (OratUtils.OratUtils static method), 8

`txtfparser()` (OratUtils.OratUtils static method), 9