

18. Photo OCR

Select text in a Photo

1. Text detection

2. Character segmentation

3. Character classification

⇒ spelling correction

steps to take

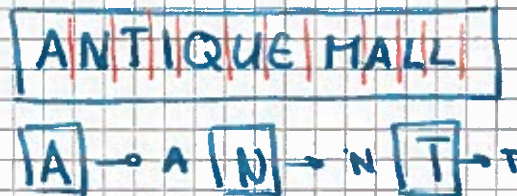


Photo OCR Pipeline



each can be a Machine learning component

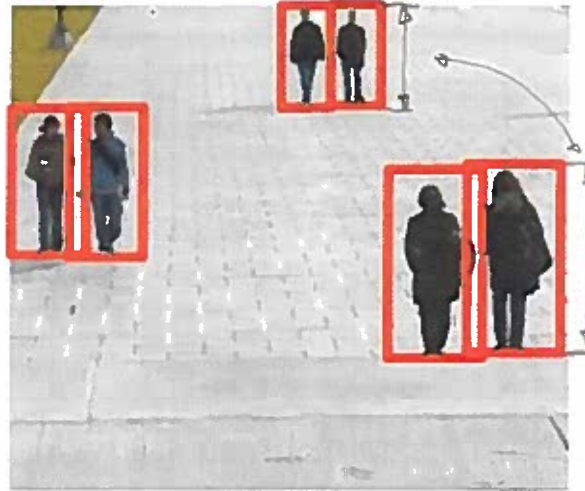
- Designing pipeline
- Artificial data synthesis

- what is the pipeline → how to break down the Pipeline into models

Text detection

diff. aspect ratio

Pedestrian detection



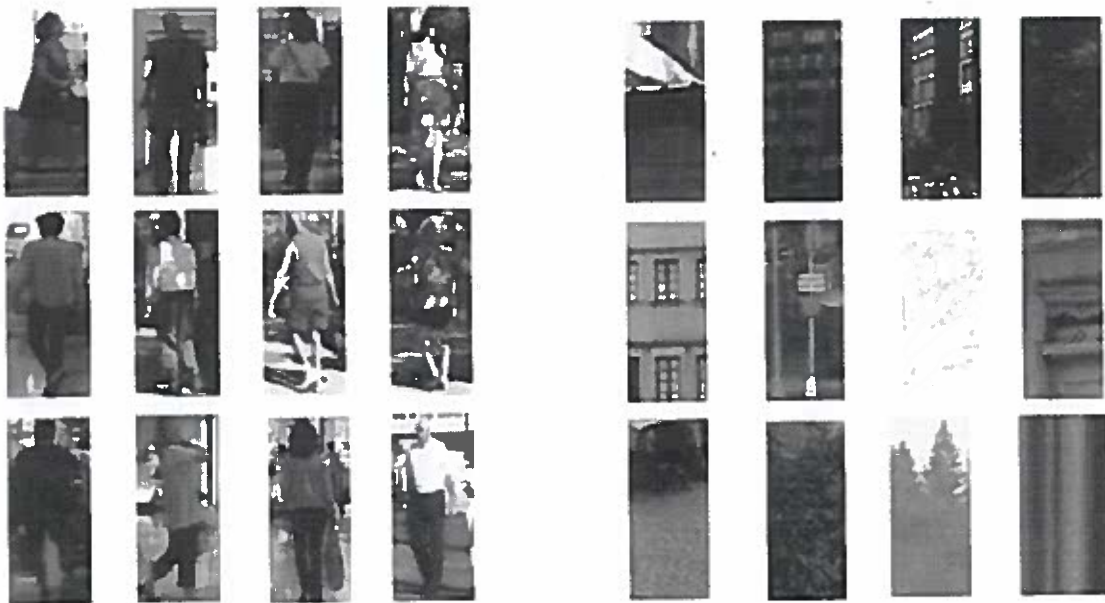
height diff.
but
aspect ratio
the same

try to find the pedestrians \rightarrow aspect ratio (ratio btw. height and width)

Supervised learning for pedestrian detection

x = pixels in 82×36 image patches

to train: 1000 - 10'000 images



positive examples ($y=1$)

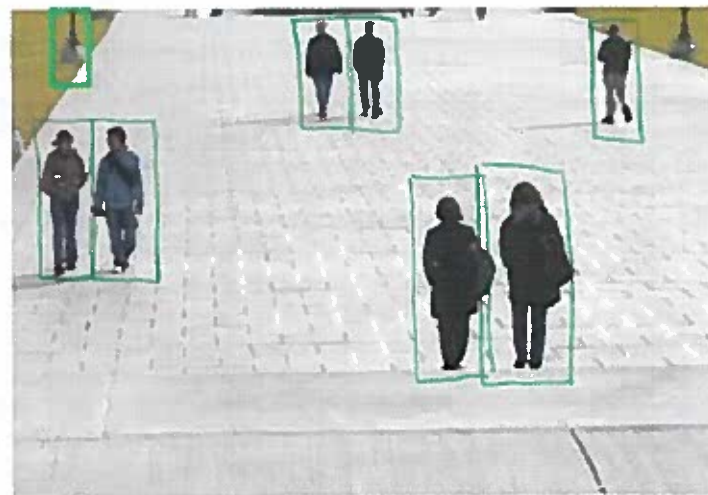
Negative examples ($y=0$)

train a NN to detect a pedestrian.

Sliding window detection:

slide over and send patch to classifier

- step size (stride)
4px. Best is 1px
but is comput.
extensive



- Run image patch
through classifier
to detect ped-
estrians

- start first with smaller image patches and run through image to detect smaller pedestrians.
- start then with a larger image patch to detect larger pedestrians (run through classifier) \Rightarrow take larger patch and resize to 82×36 px

Text detection:



positive examples ($y=1$)

negative examples ($y=0$)

- Instead of detecting pedestrians we detect text, by sliding window approach...

Text detection:

→ run sliding window patches

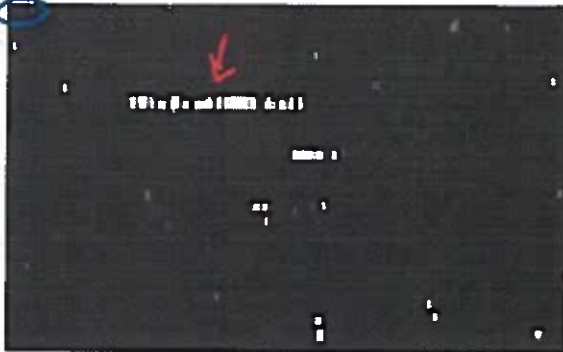
1. step



- white spots show areas where char. have been detected

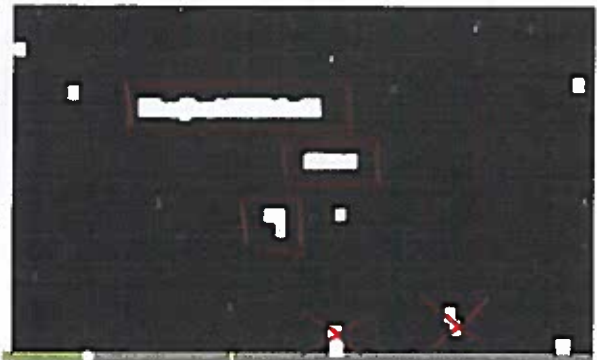
2. step

- black → no text
- white → found text



"expansion" operator

3. step



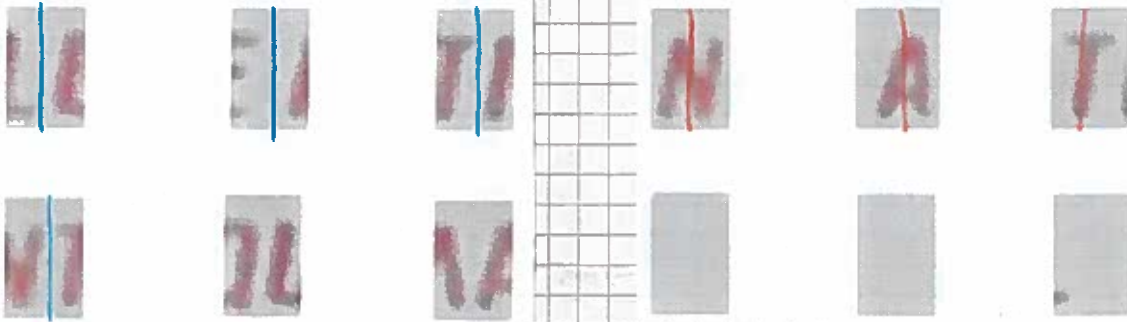
1. step ⇒ sliding window patches through image
2. step ⇒ mark found text as white spots (output of classifier)
3. step ⇒ take each of white blobs and expand that region by verifying if a pixel is white if yes expand by 5 to 10px.
4. step ⇒ connecting boxes with a certain "aspect ratio" mark these as text region
5. step ⇒ cut out the regions and use latter stages of pipeline to read the text.

aspect ratio incorrect

- (light gray → lower probability of text detection
white → high probability of text detection)

is it within some distance of the

1D Sliding window for character segmentation

positive examples ($y=1$)negative examples ($y=0$)

- detect if there is a split between two distinct characters (train a classifier e.g. NN. or other) to detect middle point..
- slide patch over the text as above to detect separation of characters (character segmentation).

Photo OCR Pipeline:

1. Text detection



2. Character segmentation

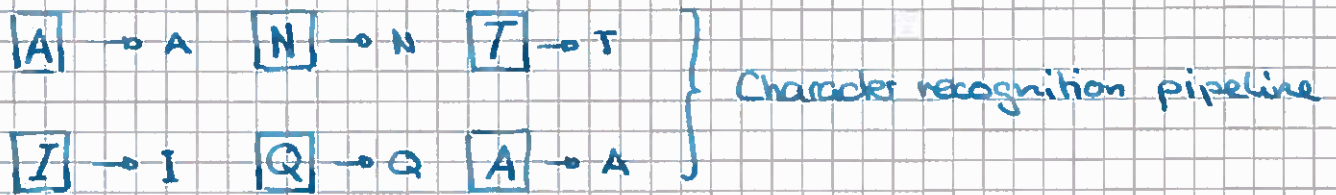


3. Character classification



Character recognition

get a low bias Model
and get lots of data!

Artificial data synthesis for Photo OCR

Real data

Abcdefg
Abcdefg
Abcdefg
 Abcdefg
 Abcdefg

- Use a font library to create data, use characters and place on different backgrounds.

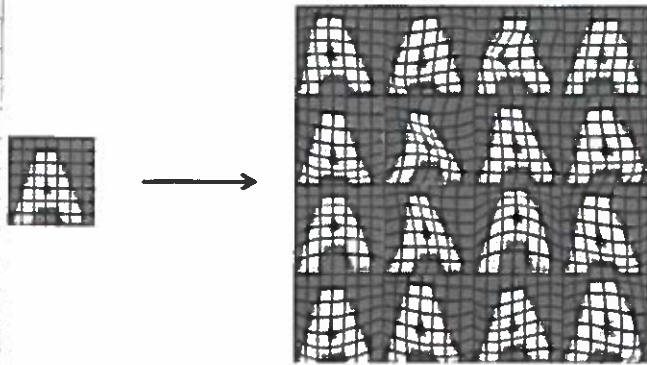


Real data



Synthetic data

- random background
- blurring
- distortion

Synthesizing data by introducing distortions

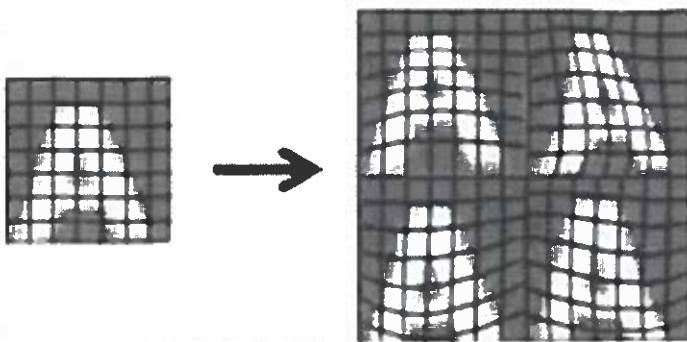
- Introduce artificial warping (distortions)
-

Speech recognition:

- Original audio
- Audio on bad cell phone connection \rightarrow add distortion
- Noisy background - Crowd \rightarrow add distortion
- Noisy background - Machinery \rightarrow add distortion

\Rightarrow multiply the original data by synthesizing data with distortion

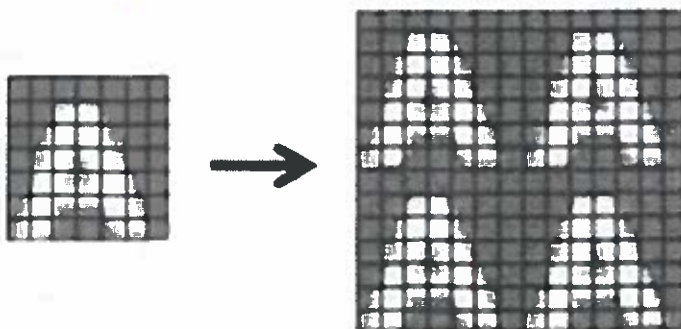
Distortion introduced should be a representation of the type of noise/distortions in the test sets.



Audio:

- Background noise, bad cell phone connection

Usually does not apply to add purely random/meaningless noise to your data



x_i = intensity (brightness) of pixel i

$$x_i \leftarrow x_i + \text{random noise}$$

\Rightarrow look for meaningful distortions
gaussian noise pixel \rightarrow meaning loss

Discussion on getting more data

1. Make sure you have a low bias classifier before expanding the effort. (Plot learning curves). E.g. keep increasing the number of features / number of hidden units in neural network until you have a low bias classifier. \Rightarrow plot learning curves
2. "How much work would it be to get 10x as much data as we currently have?"
 - Artificial data synthesis
 - Collect/label it yourself \Rightarrow how many hours to label and collect.
 - \hookrightarrow 10 sec/example \Rightarrow 10'000 sec \Rightarrow \sim 3h
 - $\frac{1}{10}$ $m=1000 \Rightarrow$ 3h
 - $m=10'000 \Rightarrow$ 30h
 - Crowd source (E.g. Amazon mechanical turk system)
 - sanity check with learning curves
 - how long for creating data

Ceiling analysis: What part of the pipeline to work next

9

Estimating the errors due to each component (ceiling analysis)

feed in perfect text detection data / feed perfect char. segm. data



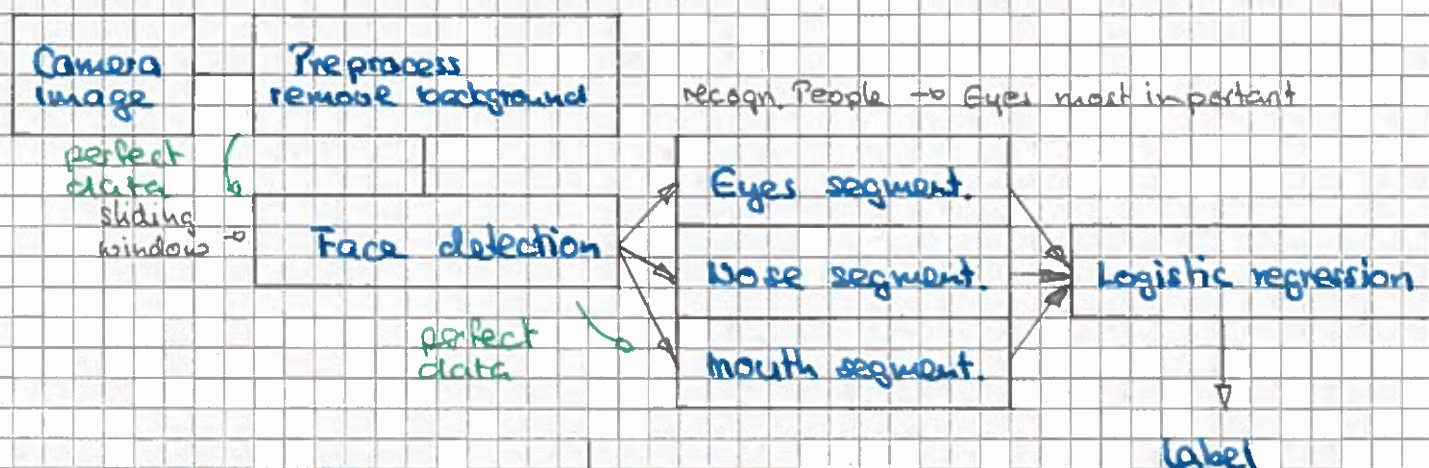
What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy (final output)
Overall system	72%
Text detection	+17% → 89%
Char. segmentation	+1% → 90%
Char. recognition	+10% → 100%

Single row number accuracy
 ↳ with perfect Text detection rate is...
 place resources on problem resolving.

Another ceiling analysis example:

Face recognition from images (Artificial example)



Component	Accuracy
Overall system	85%
Preprocess (remove backg.)	85.1% → 0.1%
Face detection	91% → 5.9%
Eyes	95% → 4.0%
Nose	96% → 1.0%
Mouth	97% → 1.0%
Logistic regression	100% → 3.0%