Machine Learning

# Application example: <mark>Photo OCR</mark>

# Problem description and pipeline

# The Photo OCR problem

how to get the computer/camera to
read the text/pictures better
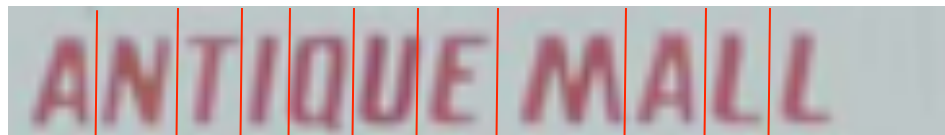


LULA B's ANTIQUE MALL
LULA B's
LULA B's
OPEN
LULA B's

Andrew Ng

# Photo OCR pipeline

1. Text detection

2. Character segmentation

3. Character classification

Cleaning → Cleaning

# Photo OCR pipeline

each of this module can be a machine learning problem

| Image | → | **Text detection** | → | **Character segmentation** | → | **Character recognition** |
|---|---|---|---|---|---|---|

1-5                    1-5                    1-5

Application example: Photo OCR

Sliding windows

Machine Learning

**Text detection**

**Pedestrian detection**

A simpler case

aspect ratio
is similar for
all pedestrians

Andrew Ng

$x =$ pixels in 82x36 image patches

1000
10,000
⋮

Positive examples $(y = 1)$

Negative examples $(y = 0)$

Andrew Ng

# Sliding window detection

step-size /stride

large image patch

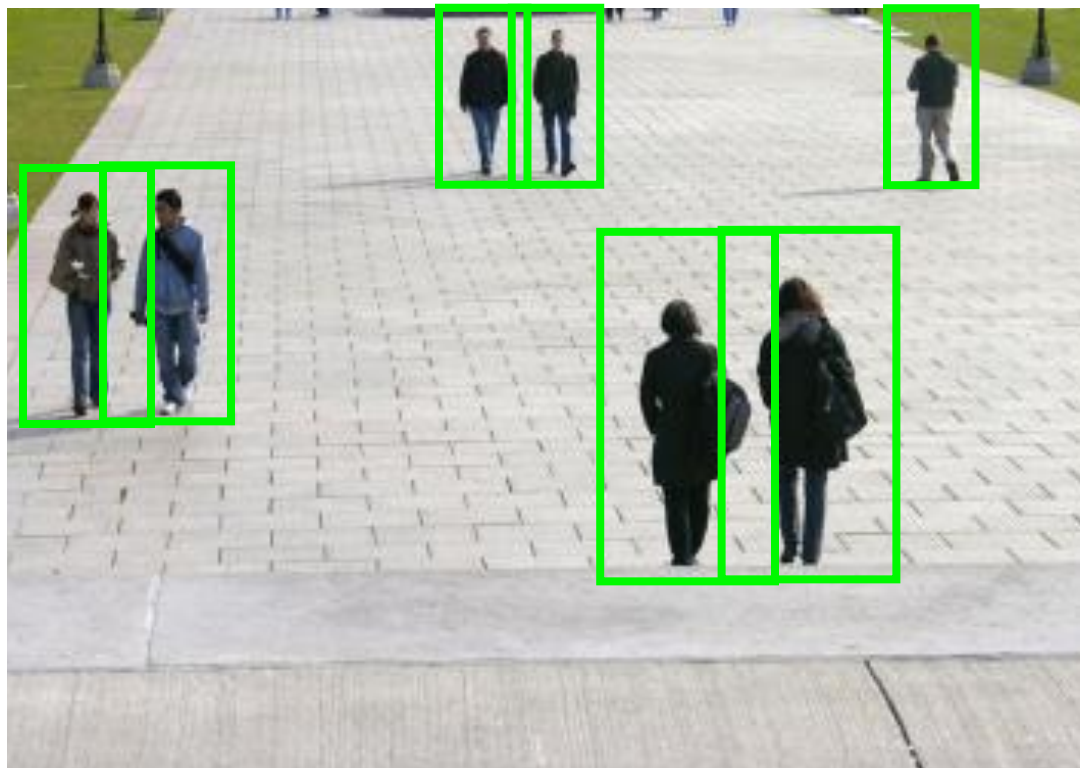we slide the window little bit and check if there is pedestrian in the window



Andrew Ng

82×36

# Sliding window detection

# Sliding window detection

# Text detection
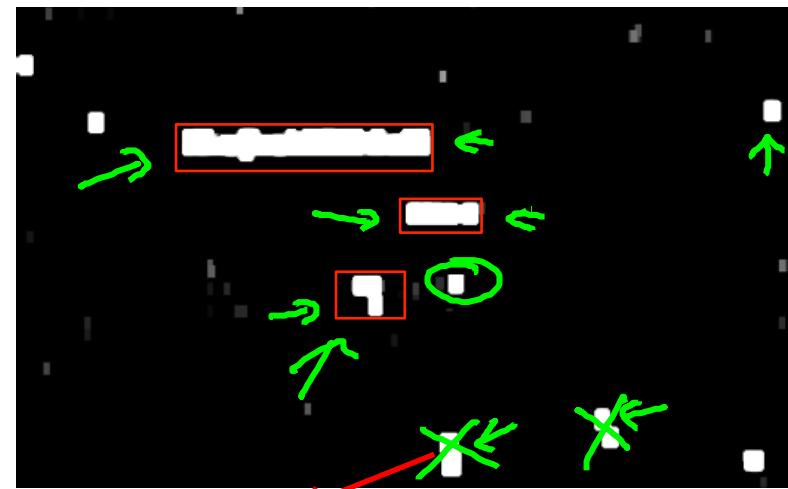
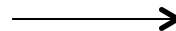# Text detection

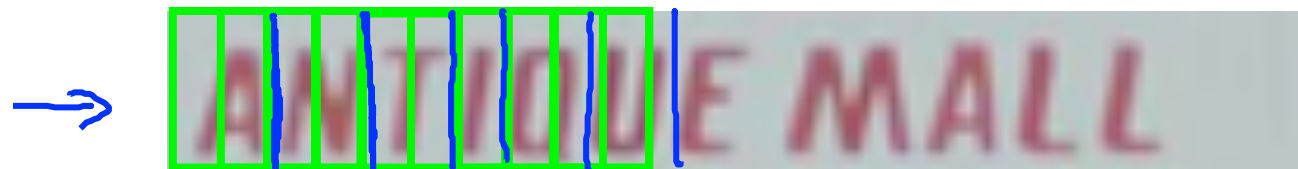

Positive examples $(y = 1)$



Negative examples $(y = 0)$

# Text detection



"expansion"

text region

aspect ratio looks funny for text

[David Wu]

Andrew Ng

# 1D Sliding window for character segmentation



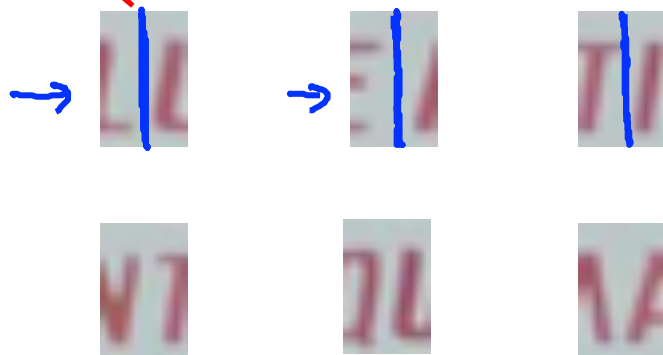Then we use sliding window to check if there is a split between two characters

slide the window over
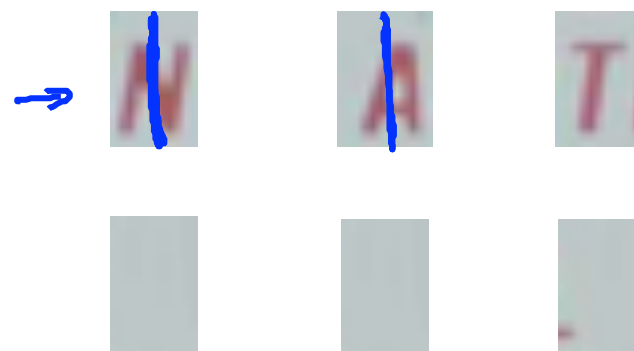
find this split
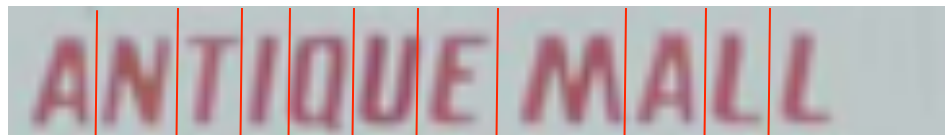
Positive examples $(y = 1)$

Negative examples $(y = 0)$

Andrew Ng

→ 1. Text detection



→ 2. Character segmentation



→ 3. Character classification

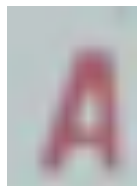just apply some supervised learning algorithm

Machine Learning

Application example: Photo OCR

Getting lots of data: Artificial data synthesis

two main variations

# Character recognition

# Artificial data synthesis for photo OCR



Real data

Abcdefg
Abcdefg
Abcdefg
Abcdefg
Abcdefg

[Adam Coates and Tao Wang]

Andrew Ng

# Artificial data synthesis for photo OCR

Real data → Synthetic data

Creating new data from scratch

[Adam Coates and Tao Wang]

Andrew Ng

# Synthesizing data by introducing distortions

Andrew Ng

# Synthesizing data by introducing distortions: Speech recognition

🔊 Original audio: ⟵

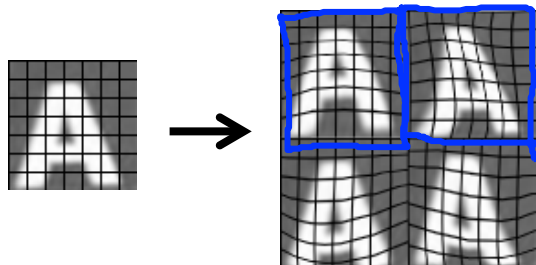🔊 Audio on bad cellphone connection

🔊 Noisy background: Crowd

🔊 Noisy background: Machinery

Andrew Ng

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

Audio:
Background noise,
bad cellphone connection

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

$x_i$ = intensity (brightness) of pixel $i$

$x_i \leftarrow x_i +$ random noise

meaningless noise is less meaningful

[Adam Coates and Tao Wang]

Andrew Ng

## Discussion on getting more data

1.  Make sure you have a low bias classifier before expending 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.

2.  "How much work would it be to get 10x as much data as we currently have?"
    -   Artificial data synthesis
    -   Collect/label it yourself
    -   "Crowd source" (E.g. Amazon Mechanical Turk)

#hours?

$m = 1,000$

$\rightarrow$ 10 secs/example

$m = 10,000$

May be the most popular source

1. Make sure you have a low bias classifier before expending 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.
2. "How much work would it be to get 10x as much data as we currently have?"
   - Artificial data synthesis
   - Collect/label it yourself
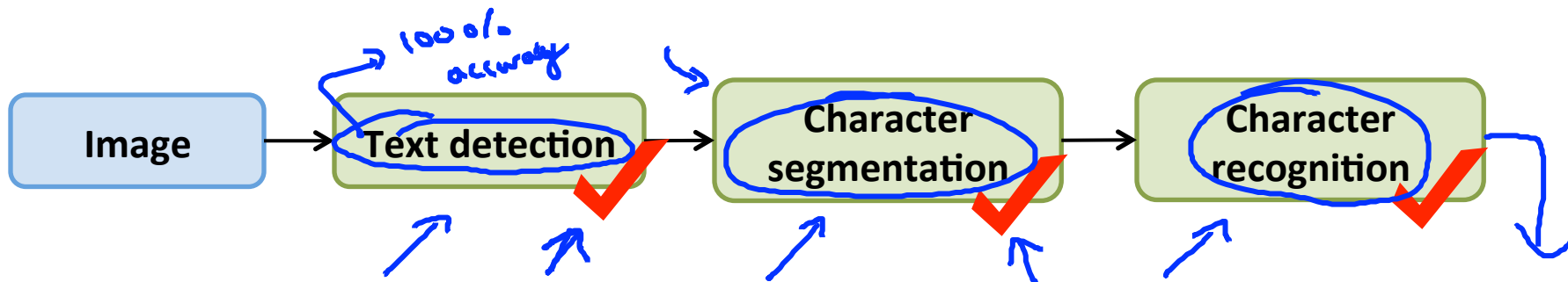   - "Crowd source" (E.g. Amazon Mechanical Turk)

Machine Learning

# Application example: Photo OCR

Ceiling analysis: What part of the pipeline to work on next

**Estimating the errors due to each component (ceiling analysis)**

Image → Text detection → Character segmentation → Character recognition

100% accuracy

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

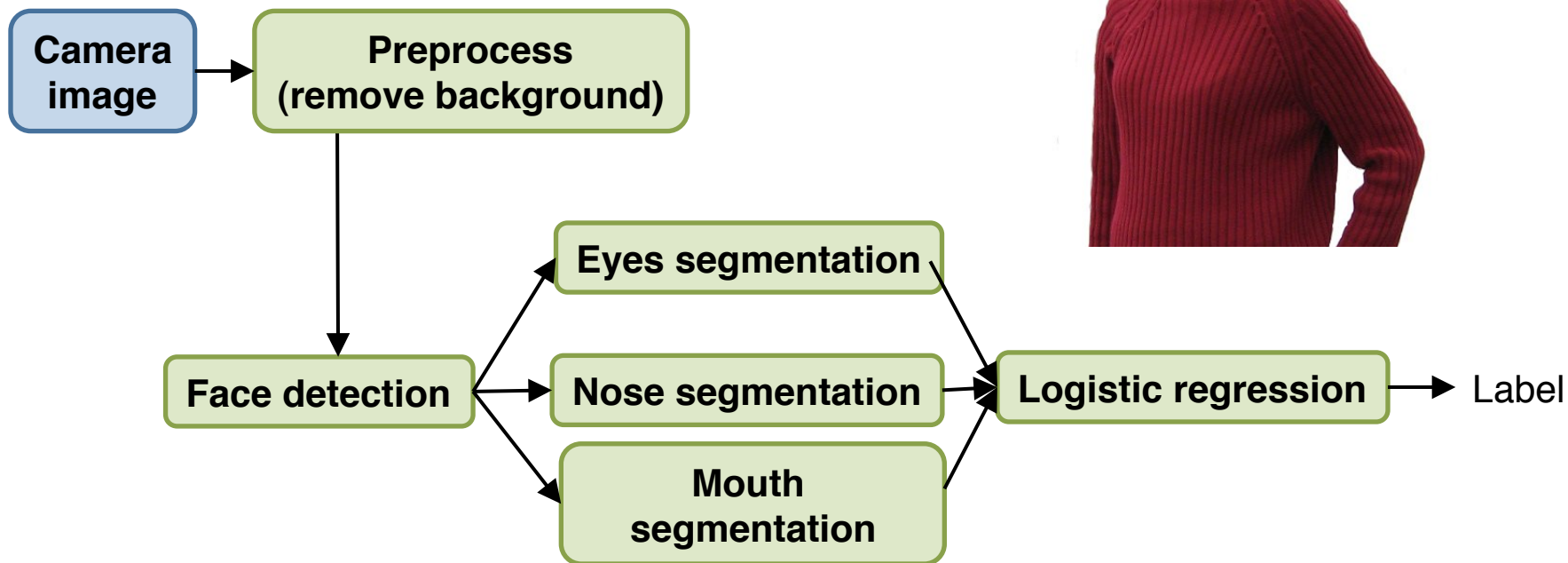| Component | Accuracy |
|---|---|
| Overall system | 72% |
| Text detection (manually set that all text are correctly detected) | 89% |
| Character segmentation (same idea) | 90% |
| Character recognition (same) | 100% |

17%
only 1% improvement (do not spend too much time here!)
1%
10%

Andrew Ng

# Another ceiling analysis example



| Component | Accuracy |
|---|---|
| Overall system | 85% |
| Preprocess (remove background) | 85.1% |
| Face detection | 91% |
| Eyes segmentation | 95% |
| Nose segmentation | 96% |
| Mouth segmentation | 97% |
| Logistic regression | 100% |

Handwritten annotations:
- 0.1% (Preprocess)
- 5.9% (Face detection)
- 4% (Eyes segmentation)
- 1% (Nose segmentation)
- 1% (Mouth segmentation)
- 3% (Logistic regression)

Andrew Ng