

Machine Learning

Application example:
Photo OCR 照片 OCR

Problem description
and pipeline

OCR pipeline

The Photo OCR problem

照片无文字符识别以获取照片中的信息

取
汽车



Photo OCR pipeline

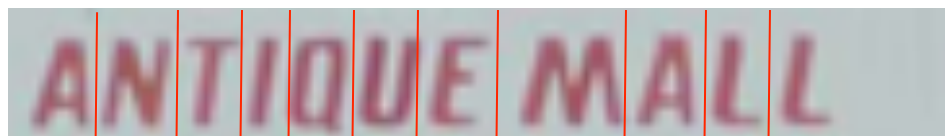
→ 1. Text detection

扫描图像, 找文字信息区域



→ 2. Character segmentation

识别并分离文字区域



→ 3. Character classification

分类器识别字符

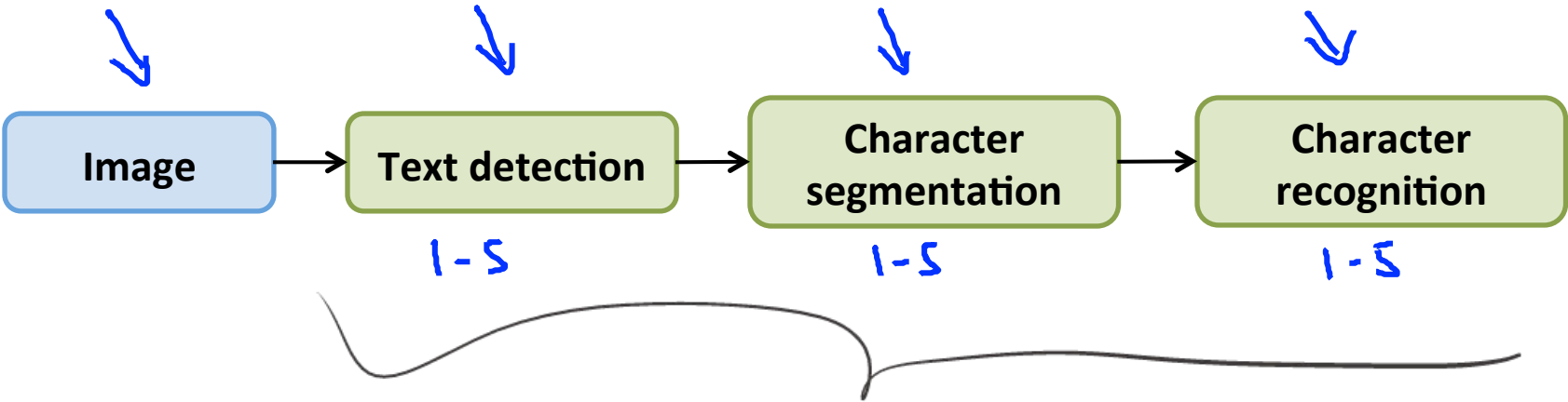


~~Cleaning~~ → Cleaning

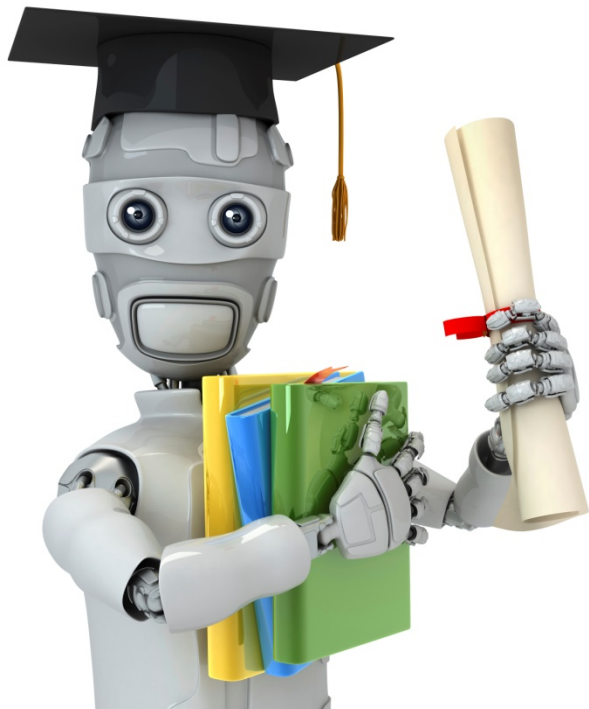
拆字停止系统: 对字符分类算法识别出的字符进行停止

Andrew Ng

Photo OCR pipeline 照片OCR流水线



不同模块组成流水线，便于分工



Machine Learning

Application example: Photo OCR

Sliding windows

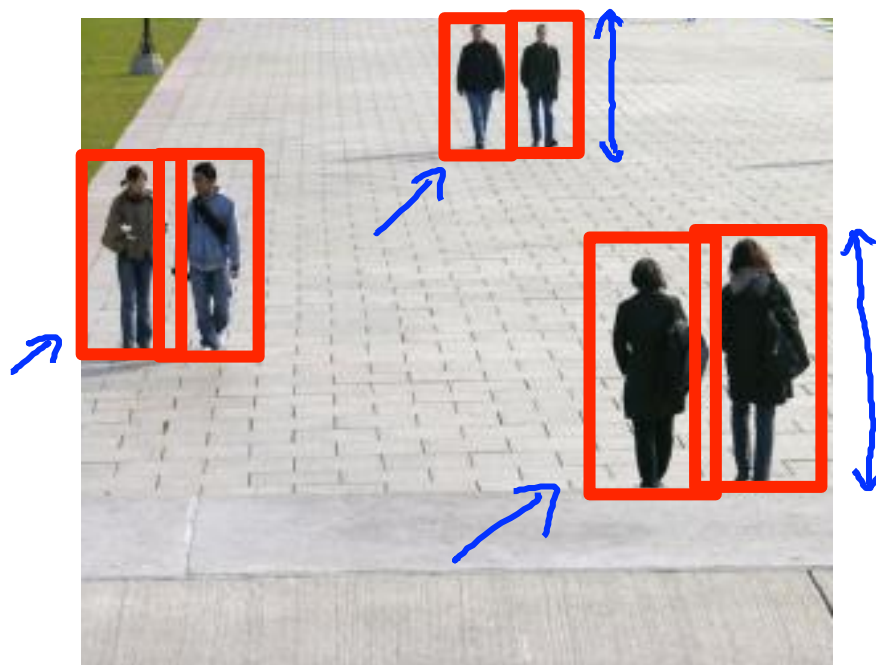
滑动窗口

文字识别 Text detection



文字区域宽高比例不同

行人检测 Pedestrian detection



行人宽高比例类似

即使长、宽不同，但比例类似

Supervised learning for pedestrian detection

x = pixels in 82x36 image patches

比例保持不变

1,000
10,000
...



Positive examples ($y = 1$)

正样本(含行人)



Negative examples ($y = 0$)

负样本(不含行人)

Andrew Ng

监督学习 输入新图像, 进行分类 \Rightarrow 输出

Sliding window detection

step-size / stride



移动窗口/步长/滑动的步数, 一般设为4 or 8 or ...



如果窗口大小, 移动窗口遍历整个图像

修改窗口大小, 移动窗口遍历整个图像

比例不变

(每个窗口输入分类器分类)

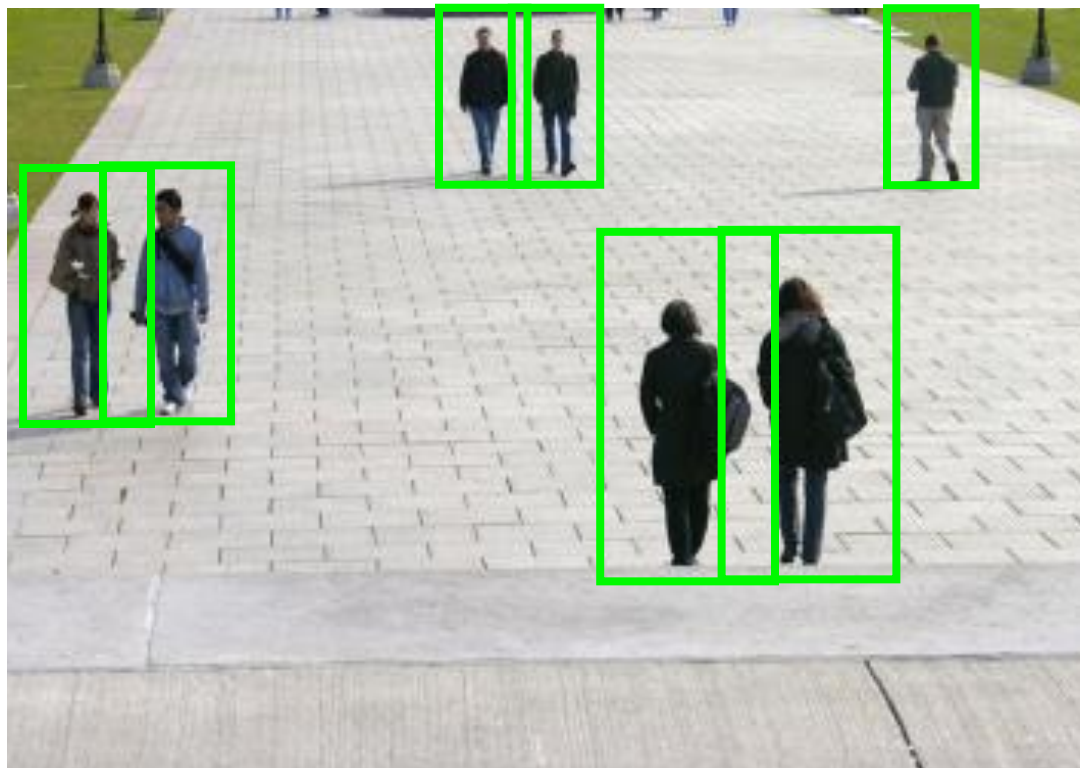
Sliding window detection



Sliding window detection



Sliding window detection



Text detection

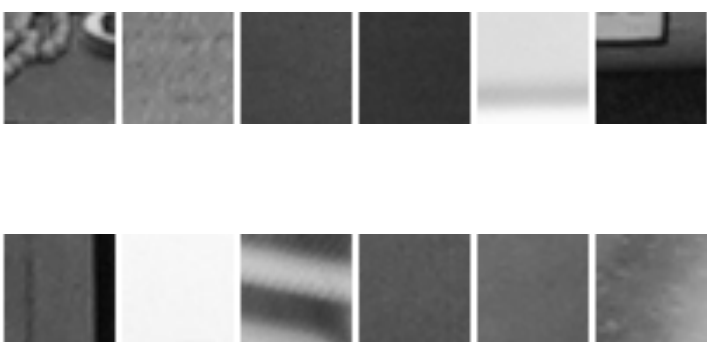


Text detection



Positive examples ($y = 1$)

已样丰(有文字)

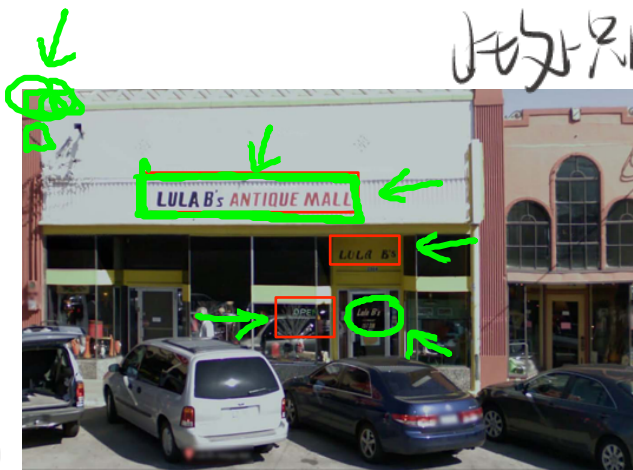


Negative examples ($y = 0$)

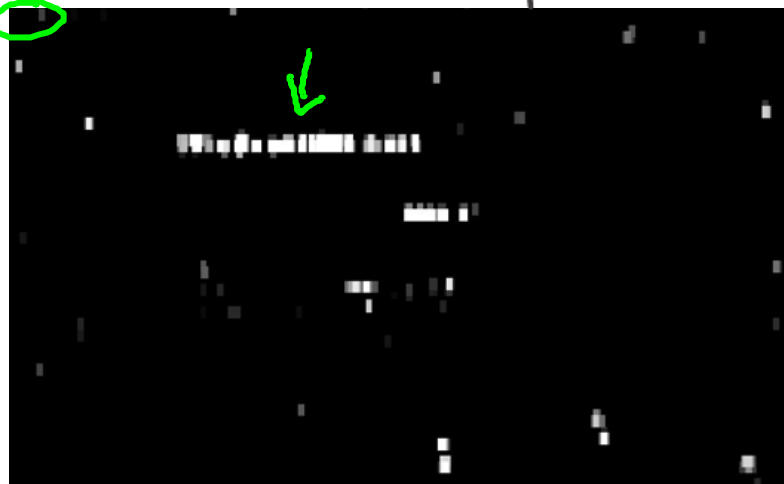
负样丰(无文字)

Text detection

(1) 宽度/高度 $\gamma=0$ 黑无文字
 $\gamma=1$ 白有文字
 (颜色程度: 文字概率)



此处只用了一个阈值比例来检测
 字符文字的比例和统一
 "expansion"



(2) 初大并子
 →



[David Wu]

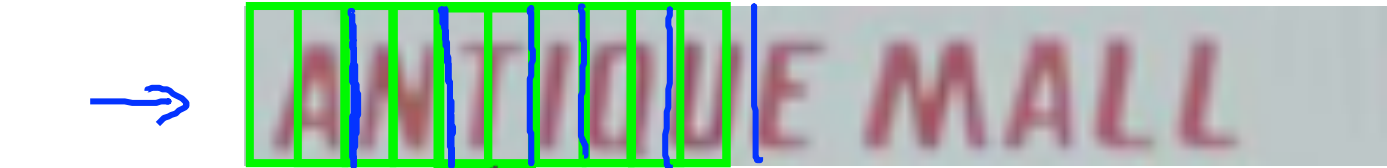
Andrew Ng

判断每个像素是否在周围15-15个像素
 若存在, 则将该范围内的像素均变为白色

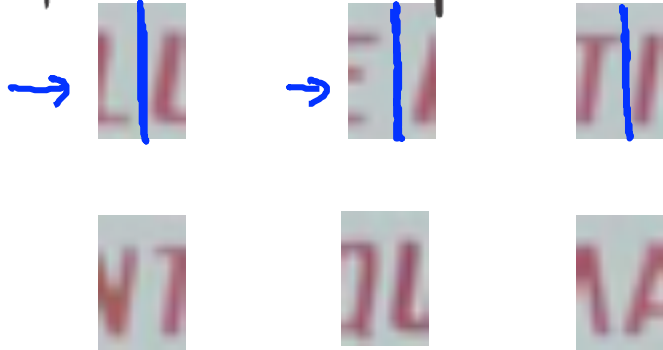
(3) 再在周围绘制边框

若比例奇怪: 宽 \leq 高 (文字是宽 $>$ 高), 则忽略
 < 但可能忽略掉单行短文字也

1D Sliding window for character segmentation 字符分割

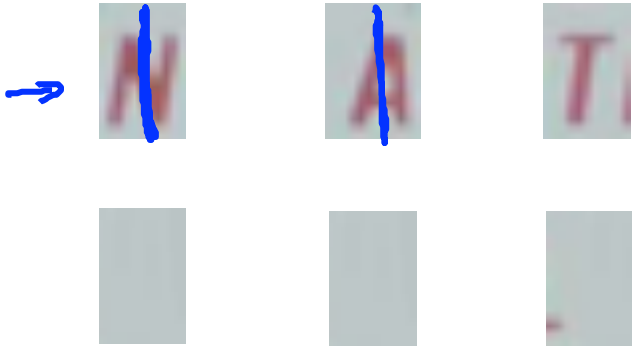


判断是否有分割点的地方



Positive examples ($y = 1$)

正样本(有分割点)



Negative examples ($y = 0$)

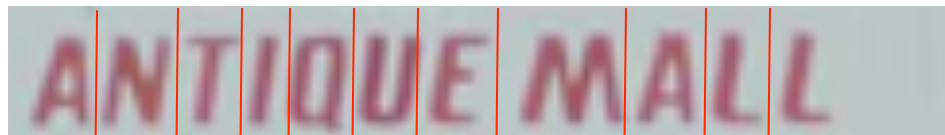
负样本(无分割点)

Photo OCR pipeline

→ 1. Text detection

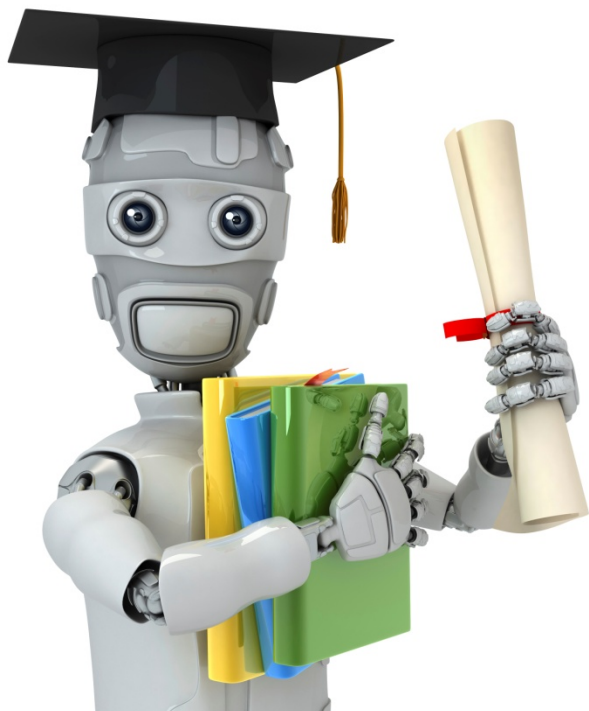


→ 2. Character segmentation



→ 3. Character classification





Machine Learning

Application example: Photo OCR

Getting lots of
data: Artificial
data synthesis

人工数据合成 (生成更多数据)

< 从0开始创造新数据 (1)
已有小数据训练集, 再扩展数据集 (2)

Character recognition



→ **A**



→ **N**



→ **T**



→ **I**

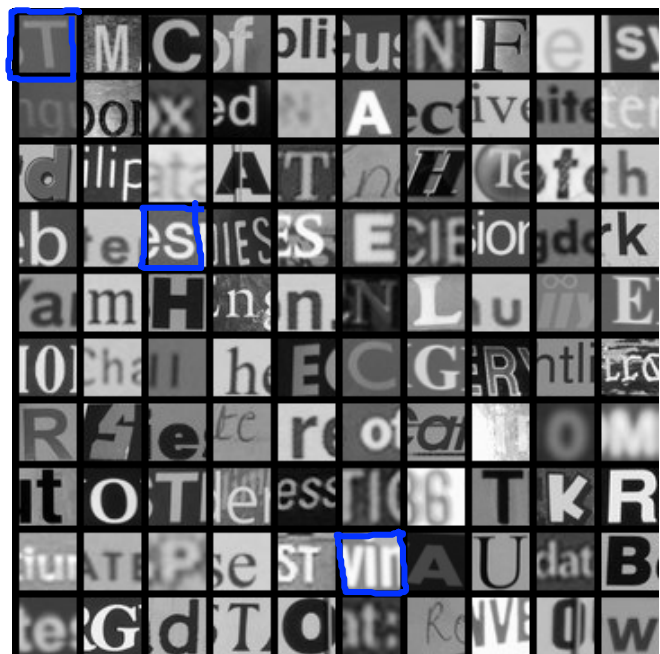


→ **Q**



→ **A**

Artificial data synthesis for photo OCR



Real data

Abcdefg
Abcdefg
Abcdefg
Abcdefg
Abcdefg
Abcdefg

[Adam Coates and Tao Wang]

017 { 用大量字符生成不同字体的字符
↓
再贴到任意不同背景 →
模糊算子
仿射变换 等分 仿射 旋转

Andrew Ng

Artificial data synthesis for photo OCR

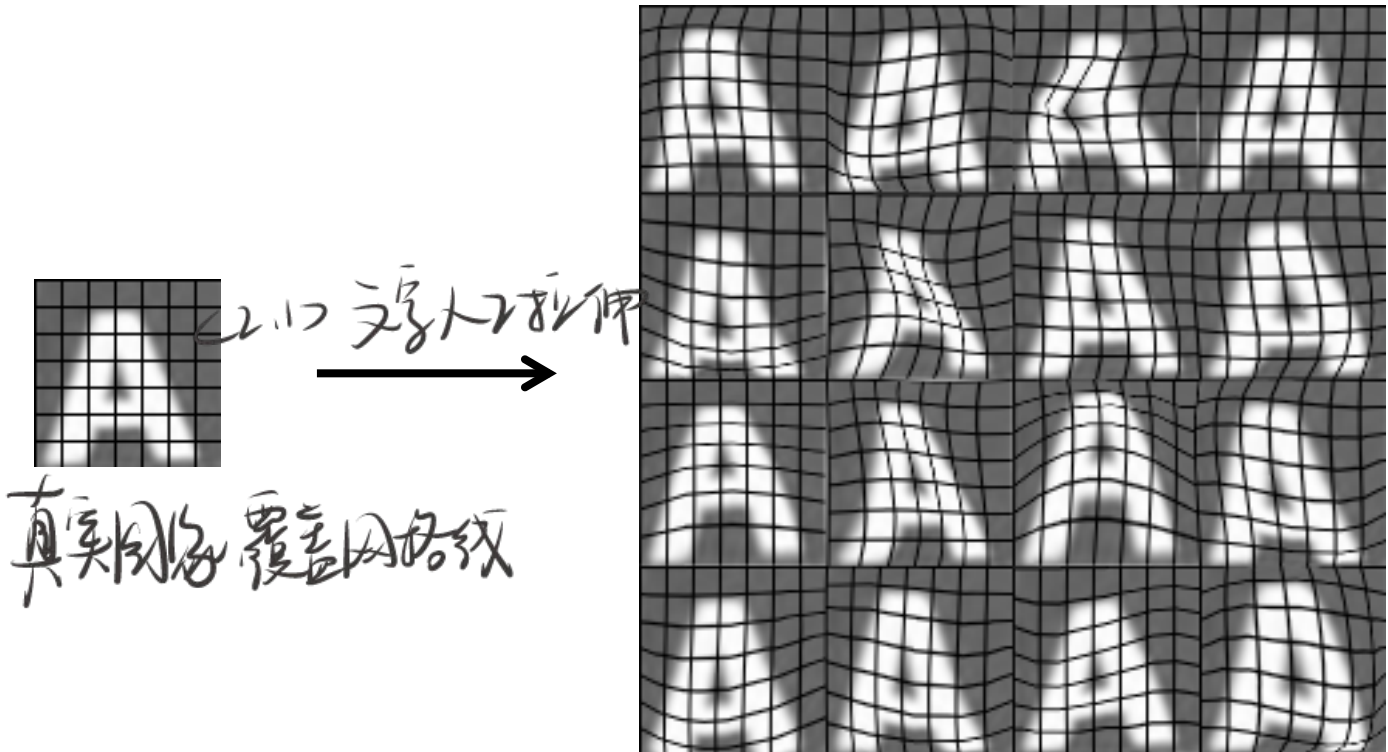


Real data



Synthetic data

Synthesizing data by introducing distortions



[Adam Coates and Tao Wang]

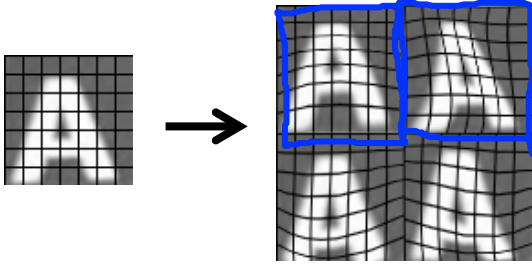
Andrew Ng

(2) 音韻規則出現順序



Synthesizing data by introducing distortions

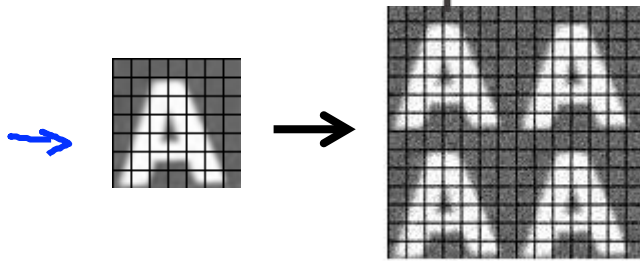
- Distortion introduced should be representation of the type of noise/distortions in the test set. 引入失真要有代表性



→ Audio:

Background noise,
bad cellphone connection

- Usually does not help to add purely random/meaningless noise to your data. 引入无意义的失真(非像素级噪声)



→ x_i = intensity (brightness) of pixel i

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

[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?"

人工合成 Artificial data synthesis

收集、打标签 Collect/label it yourself

- "Crowd source" (E.g. Amazon Mechanical Turk)

外包: 服务商提供网络雇人服务, 以较低价格标注大量数据

→ #Hours?

$n = 1,000$

→ 10 secs/example

$n = 10,000$

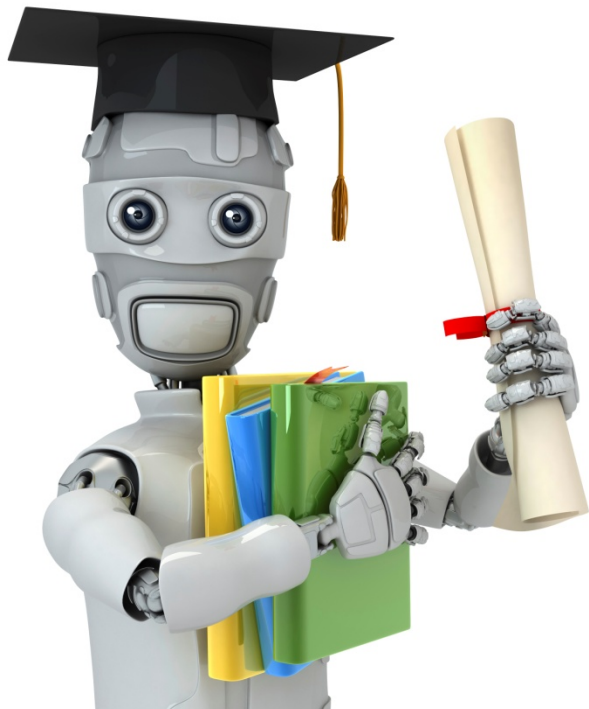
Andrew Ng

绘制一个学习曲线, 得到低偏差高泛化, 这样人工合成或采样才有意义

若泛化误差高, 增加隐藏单元

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)



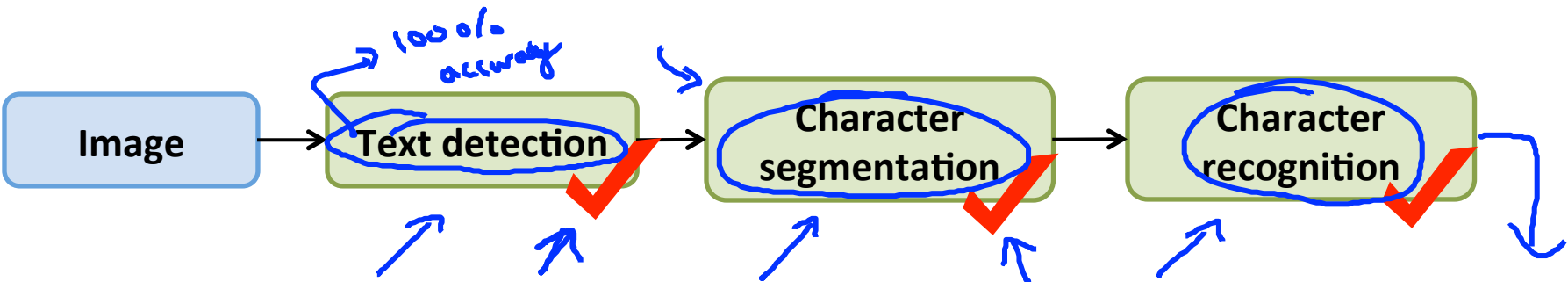
Machine Learning

Application example: Photo OCR

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

上限分析

Estimating the errors due to each component (ceiling analysis)



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

Component	Accuracy
Overall system	72%
Text detection	89%
Character segmentation	90%
Character recognition	100%

再个阶段字符识别

这个测试环节认为是一个环节，这个环节准确率100%

这个环节准确率100%
这个环节准确率100%

这个环节准确率100%
这个环节准确率100%
这个环节准确率100%

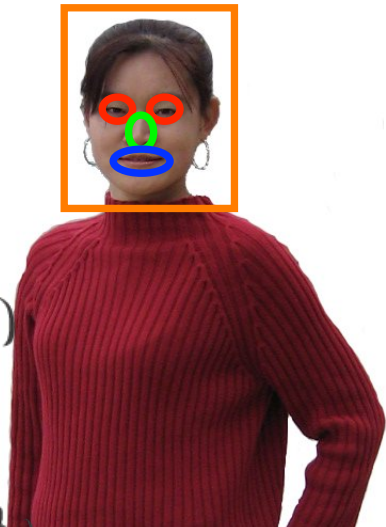
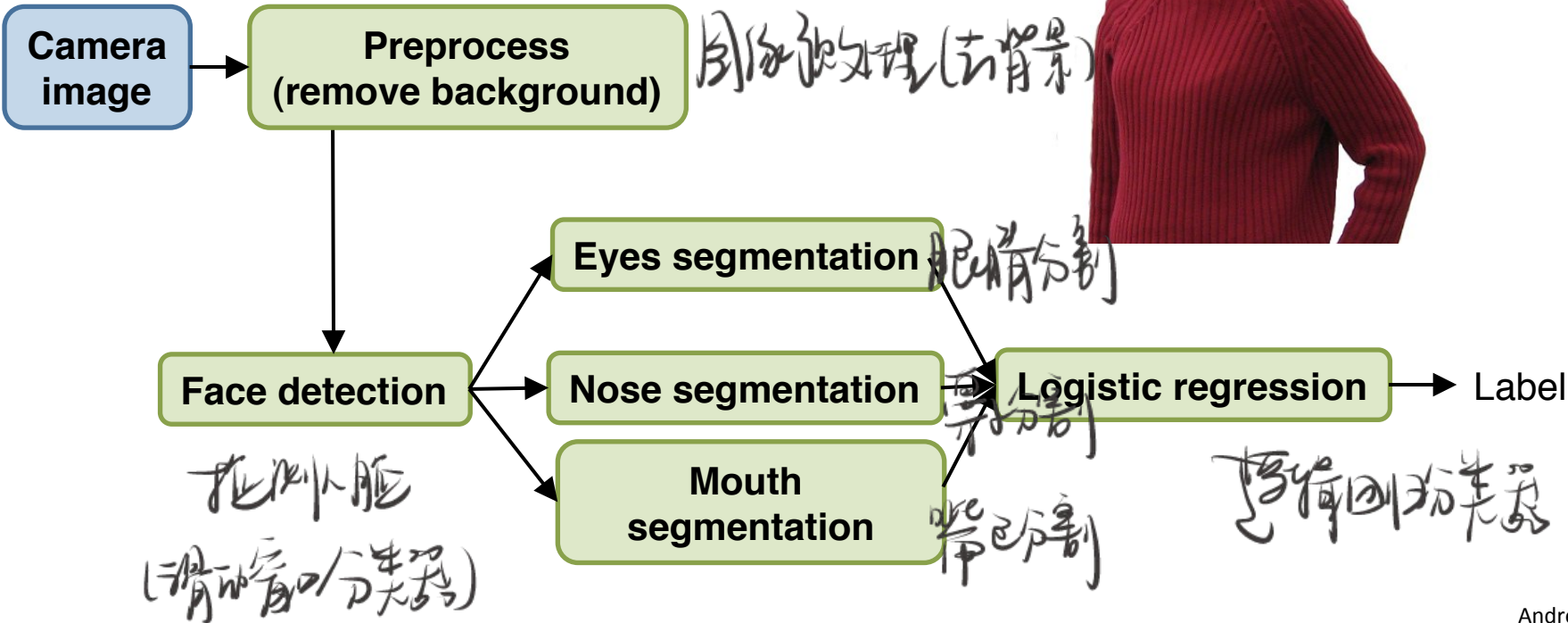
← 17%
← 1%
← 10%

各模块性能提升上限

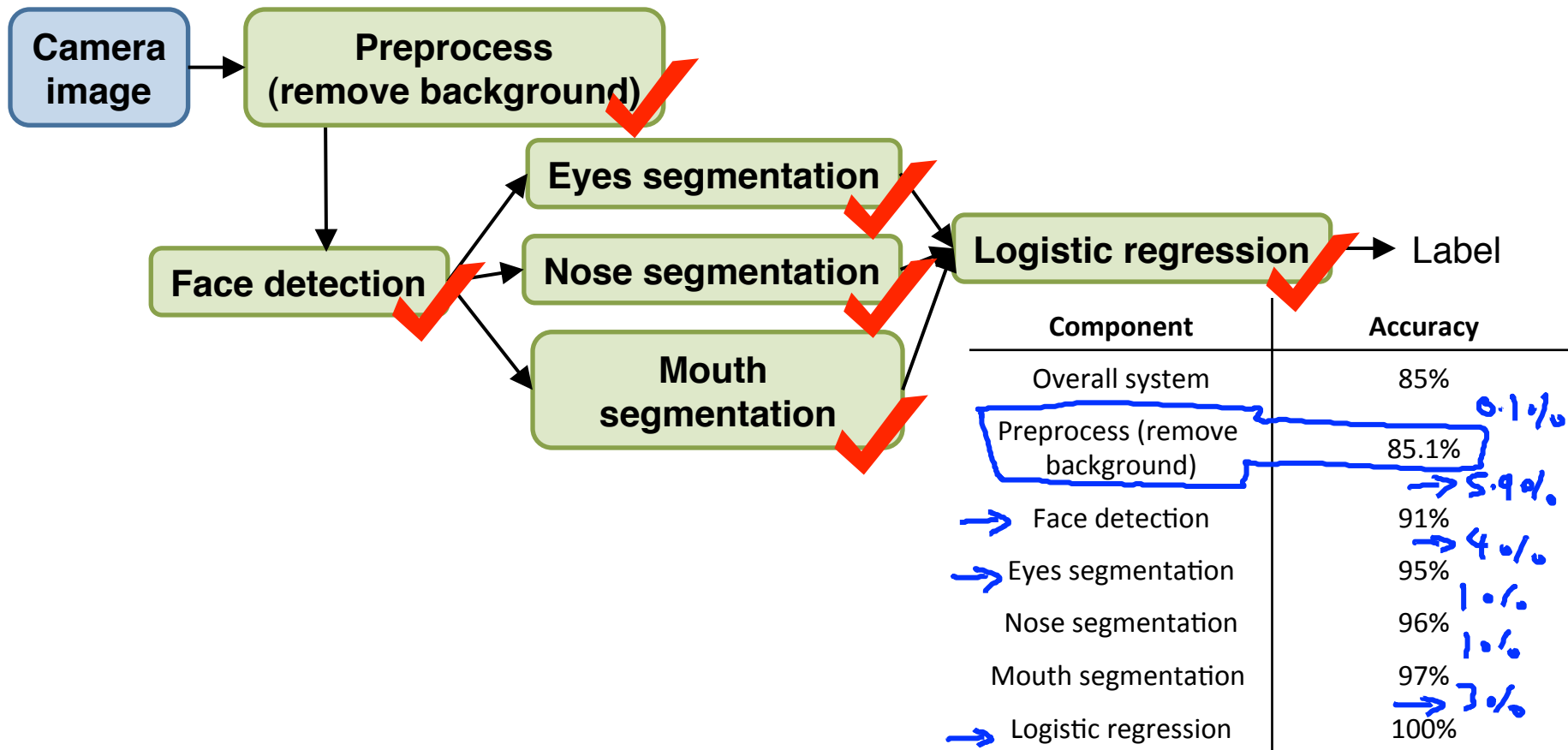
Andrew Ng

Another ceiling analysis example

Face recognition from images
(Artificial example)



Another ceiling analysis example



Summary: Main topics

监督

Supervised Learning

- Linear regression, logistic regression, neural networks, SVMs

$$(x^{(i)}, y^{(i)})$$

无监督

Unsupervised Learning

- K-means, PCA, Anomaly detection

$$x^{(i)}$$

特殊应用

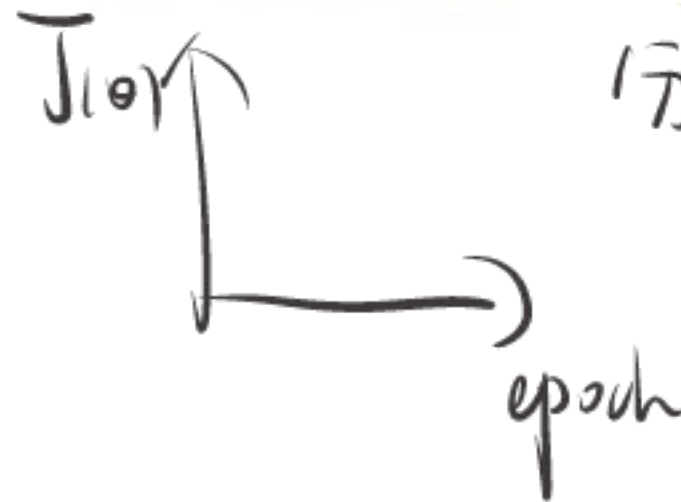
Special applications/special topics

- Recommender systems, large scale machine learning.

→ Advice on building a machine learning system

- Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.

Figure...



分析错误率

