

Machine learning

1.3 Anomaly detection 异常检测 (无监督学习)

1.3-1 Problem motivation 问题动机

- 对于给定的无异常 Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, \rightarrow Is x_{test} anomalous?

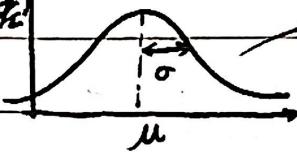
\Rightarrow 建立 x 的概率模型, 获取阈值 ϵ \rightarrow $P(x_{\text{test}}) < \epsilon \rightarrow \text{Anomaly}$

• 常见应用:

$$P(x_{\text{test}}) \geq \epsilon \rightarrow \text{OK}$$

欺诈检测, 工程 (材料检测), Monitoring computers in a data center.

13 $x \sim N(\mu, \sigma^2)$ (均值 / 方差) / (标准差)



$$p(x; \mu, \sigma^2) =$$

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

13-2 Gaussian distribution

$$\mu = \frac{1}{m} \sum_{j=1}^m x^{(j)}$$

$$\sigma^2 = \frac{1}{m} \sum_{j=1}^m (x^{(j)} - \mu)^2$$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) \cdots p(x_n; \mu_n, \sigma_n^2)$$

13-3 Algorithm

- If we have Training set: $\{x^{(1)}, \dots, x^{(m)}\}$, Each example is $x \in \mathbb{R}^n$.

每个特征 $x_j \sim N(\mu_j, \sigma_j^2) \rightarrow$ [独立分布] (实际上在算法中, 不独立看成也很好)

- Anomaly detection algorithm

examples

1. Choose feature x_i that you think might be indicative of anomalies

2. fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} : \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3. Given new example x , Compute $p(x)$

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

Anomaly if $p(x) < \epsilon$.

13

13-4 Developing and evaluating an anomaly detection system

- [Real-number evaluation] (假设有一个正常/异常混合带标签的样本集)

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \leftarrow$ 全为正常样本, $y=0$

Gross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m)}, y_{cv}^{(m)}) \quad \} \text{其中有一些异常样本.}$

Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m)}, y_{test}^{(m)}) \quad \} \quad y=1$

[Algorithm evaluation]

- fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$

- in a cross-validation/test example x , predict $y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (Anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (Normal)} \end{cases}$

→ possible evaluation metrics:

- TP, FP, FN, TN • Precision / Recall

- f1-score # Can also use cross-validation set to choose parameters ϵ .

13-5 Anomaly detection vs. supervised learning

· for Anomaly, 适用于:

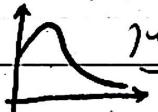
- Very small number of positive examples and large number of negative.
- 分出不同类型的 anomalies. 很难去学习/预测新异常.

· for supervised,

- large number of positive and negative examples
- future positive examples likely to be similar to ones in training set.

13-6 Choosing what features to use.

· For Non-gaussian feature, 可以通过处理让其接近高斯分布

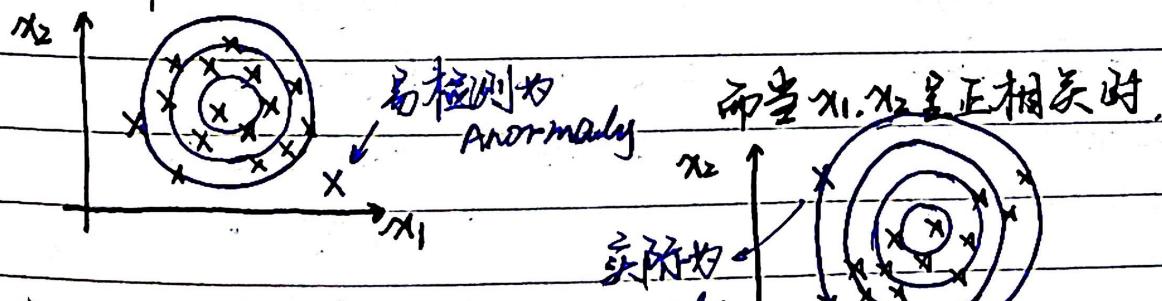
e.g. $x \rightarrow \log(x+c)$ (适合处理“拖尾”式:  )
 $x \rightarrow x^a$ (like $x^{\frac{1}{2}}, x^{\frac{1}{3}} \dots$)

· Error analysis: 常见问题是 prior is comparable (both large) for normal and anomaly examples \Rightarrow 需要建立新特征

13-7 Multivariate Gaussian distribution 多元高斯分布

· Motivation: 当特征不独立 (e.g. 正相关) 时, 预测出错.

E.g. 正常假设下 x_1 和 x_2 是相互独立的特征, 则



为解决这一问题, 开发改良版 anomaly.
并带检测算法. 即 多元高斯分布

Multivariate Gaussian (normal) distribution

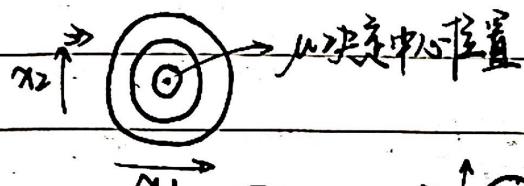
$x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately

Model $p(x)$ all in one go. Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

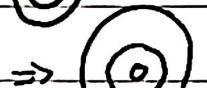
$$\Rightarrow p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu))$$

$\downarrow \det(\Sigma)$

Ex. 1. $\mu = [0]$, $\Sigma = [1]$ 2. $\mu = [0]$, $\Sigma = [0.6]$ 3. $\mu = [0]$, $\Sigma = [0.6]$ \Rightarrow



4. $\mu = [0]$, $\Sigma = [0.1]$ \Rightarrow 5. $\mu = [0]$, $\Sigma = [2]$ \Rightarrow



6. $\mu = [0]$, $\Sigma = [0.1 \ 0.5]$ \Rightarrow 7. $\mu = [0]$, $\Sigma = [1 \ 0.8]$ \Rightarrow 8. $\mu = [0]$, $\Sigma = [-0.5 \ 1]$ \Rightarrow



13-8 Anomaly detection using the multivariate Gaussian distribution

• 流程: 1. Fit model $p(x)$ by setting $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

2. Given a new example x , compute $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu))$$

flag an anomaly if $p(x) < \epsilon$

multivariate Gaussian

• 对比
Original model (单因素)

$$(p(x) = p(x_1; \mu_1, \sigma_1^2) \cdots p(x_n; \mu_n, \sigma_n^2))$$

实际上 $\Sigma = [\sigma_1^2 \ 0 \ \dots \ 0]$ (对称特殊形式)

→ Automatically captures Correlation between features.

缺点 } . computationally more expensive
. must have $m > n$

△ 如果要解决 13-7 问题 → manually

Create features to capture anomalies where x_1, x_2 take unusual combinations of values (e.g. 比如 $x_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$)

- 优点 } . Computationally cheaper
. OK even if m (training set size) is small.

1.14 Recommender Systems 推荐系统

1.4-1 Problem formulation 问题规划

对于一些问题，算法能够自动学习一些特征。

推荐系统 E.g. 用户给电影打分。

1.4-2 Content-based Recommendations 基于内容的推荐算法 (romance)

Movie	Alice ⁽¹⁾	Bob ⁽²⁾	Carol ⁽³⁾	Dave ⁽⁴⁾	x_1	x_2 (action)
love at least	5	5	0	0	0.9	0
Romance forever	5	?	?	0	1.0	0
cute puppies of love	?	4	0	?	0.99	0
Nonstop car chases	0	0	5	4	0.1	1.0
Swords vs karate	0	0	5	?	0	29

- $r_{i,j} = 1$ if user j has rated movie i ($= 0$ otherwise)

- $y_{(i,j)}$ = rating by user j on movie i (if defined)

- $\theta^{(i)}$ = parameter vector for movie i / parameter vector for user j

- $x_{(i)}^{(j)}$ = feature \downarrow

For user j , movie i . predicted rating: $(\theta^{(i)})^T x_{(i)}^{(j)}$

E.g. for Alice, $x_{(3)}^{(1)} = \begin{bmatrix} x_{10} \\ x_{11} \\ x_{12} \end{bmatrix} = \begin{bmatrix} 0 \\ 0.99 \\ 0 \end{bmatrix}$, $\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$.

则 Alice 对于 "cute puppies of love" 的预测为: $(\theta^{(1)})^T x_{(3)}^{(1)} = 4.95$

- $n^{(j)}$ = no. of movies rated by user j . To learn $\theta^{(j)}$.

→ Optimization objective

$$\text{To learn } \theta^{(j)} \rightarrow \min_{\theta^{(j)}} \frac{1}{2} \sum_{i: r_{i,j} \geq 1} ((\theta^{(j)})^T x_{(i)}^{(j)} - y_{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

To learn $\theta^{(1)} \dots \theta^{(n)}$ (n↑ user)

$$\rightarrow \min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j=1}^n \sum_{i: r_{i,j} \geq 1} ((\theta^{(j)})^T x_{(i)}^{(j)} - y_{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2$$

→ Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i: r_{i,j} \geq 1} ((\theta^{(j)})^T x_{(i)}^{(j)} - y_{(i,j)}) x_{ik}^{(j)} \quad (\text{for } k=0)$$

θ₀ 不变

正则化项

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i: r_{i,j} \geq 1} ((\theta^{(j)})^T x_{(i)}^{(j)} - y_{(i,j)}) x_{ik}^{(j)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

14-3 Collaborative filtering

Optimization algorithm

Given $\theta^{(1)}, \dots, \theta^{(n)}$, to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

Given $\theta^{(1)}, \dots, \theta^{(n)}$, to learn $x^{(1)}, \dots, x^{(n)}$:

$$\min_{x^{(1)}, \dots, x^{(n)}} \frac{1}{2} \sum_{i=1}^m \sum_{j:r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2$$

→ Collaborative filtering (and movie ratings).

Given $x^{(1)}, \dots, x^{(n)}$, can estimate $\theta^{(1)}, \dots, \theta^{(n)}$.

Given $\theta^{(1)}, \dots, \theta^{(n)}$, can estimate $x^{(1)}, \dots, x^{(n)}$.

逐行随机选取一些 θ 值 → $x \rightarrow \theta \rightarrow x \rightarrow \dots$

14-4 Collaborative filtering Algorithm

Minimizing $x^{(1)}, \dots, x^{(n)}$ and $\theta^{(1)}, \dots, \theta^{(n)}$ simultaneously (同时梯度下降)

1. Initialize $x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}$ to small random values

2. Minimize $J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$

$$\min_{\substack{x^{(1)}, \dots, x^{(n)}, \\ \theta^{(1)}, \dots, \theta^{(n)}}} J(\dots) = \min_{\substack{i, j: \\ r(i,j)=1}} \frac{1}{2} \sum ((\theta^{(i)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2$$

⇒ Using gradient descent (or other algorithms) E.g. for every $j = 1, \dots, n$

$$x_k^{(j)} := x_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(i)} + \lambda x_k^{(j)}$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)}$$

3. for a user with parameters θ and a movie with (learned) features x . predict a rating of $\theta^T x$.

(注: $x \in R^n$, $\theta \in R^n$, $\theta^T x$, 因为需要算法自动学习特征)

14-5 Vectorization: low rank matrix factorization 矩量化: 低秩矩阵分解

- Predicted rating: $\hat{Y} = \Theta^T X^T$

$$\hat{Y} = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \Rightarrow X = \begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ -x^{(m)} \end{bmatrix} \quad \Theta = \begin{bmatrix} -\theta^{(1)} \\ -\theta^{(2)} \\ \vdots \\ -\theta^{(m)} \end{bmatrix}$$

finding related movies

for product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$

each \rightarrow 5 most similar movies to movie i :

\rightarrow find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$

14-6 Implementational detail: Mean normalization 均值规范化

假设用户 Eve(5) 对前文 5 个电影的评分。 \rightarrow 预测 Eve(5) 对电影评分。

\rightarrow 求 $\theta^{(5)}$ 。因为对 Eve 而言, $r(i,j)$ 均为 0, $\min_j \dots \geq \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta_k^{(j)})^2 - 2\lambda$ 起作用, \rightarrow 使 $\theta^{(5)}$ 不可解, i.e. $\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

\Rightarrow Eve(5) 对评分均为 0, 没有意义

Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \quad \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.5 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.5 & -2.5 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

for user j . on movie i predict:

$$(\theta^{(j)})^T (x^{(i)}) + \mu_i$$

User 5 (Eve):

$$(\theta^{(5)})^T (x^{(i)}) + \mu_i = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.5 \\ 1.25 \end{bmatrix}$$

1.15 Large scale machine learning 1.15 large scale machine learning

1.5-1 Learning with large datasets

* **预览检查**: 假设 train set $m = 100,000,000$, 需要考虑选择其中 1000 个样本训练结果会不会相似 (参见 8-6)

1.5-2 Stochastic gradient descent 随机梯度下降.

对于训练样本过大的情况, 例如 $m = 300,000,000$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(i)$$

Batch gradient descent (批量)

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

→ Repeat {

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

{ (for every $j = 0, \dots, n$) }

}

$$\frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta)$$

每次迭代都用 m 个样本

Stochastic gd

$$\text{Cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

1. Randomly shuffle (reorder) training example

2. Repeat {

for $i = 1, \dots, m$ }

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

{ (for every $j = 0, \dots, n$) }

一次迭代针对单个样本

1.5-3 Mini-batch gradient descent.

- Batch gd: Use all m examples in each iteration
- Stochastic gd: Use 1 example in each iteration
- Mini-batch gd: Use b example in each iteration

E.g. say $b = 10$. $m = 1000$

Repeat { for $i = 1, 11, 21, 31, \dots, 991$

$$\theta_j := \theta_j - \alpha \frac{1}{b} \sum_{k=i}^{i+9} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

for every $j = 0, \dots, n$. } }

如何保证学习率α

15-4 Stochastic gradient descent convergence 收敛

- Cost(θ , $(x^{(i)}, y^{(i)})$) = $\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$ ← 可以用来衡量算法表现
- During learning, compute Cost(θ , $(x^{(i)}, y^{(i)})$) before updating θ using $(x^{(i)}, y^{(i)})$
- (-般) Every 1000 iterations, plot Cost(θ , $(x^{(i)}, y^{(i)})$) averaged over the last 1000 examples processed by algorithm.

• 一般地，学习率α是设置好的不变常数。但，为了更好地收敛到全局最小值，
可以让学习率α随时间逐渐减小。e.g. $\alpha = \frac{\text{Constant-1}}{\text{no.of iteration} + \text{Constant-2}}$
(注：Constant 1 和 2 是需要额外设置的参数) 迭代次数

15-5 Online learning

E.g. Product Search (learning to search)

User searches for "Android phone"

Have 100 phones in store. Will return 10 results.

x = features of phone, how many words in user query match name of phone,
how many words in query match description of phone, etc. (predicted CTR)

y = 1 If users click on link, $y=0$. otherwise. → learn $p(y=1|x; \theta)$

Use to show user the 10 phones they're most likely to click on.

实际上每个用户访问网站，就得到十个样本数据。从第一部手机得到一个
特征矢量 x ，以及对应 y 值 \Rightarrow 如何运行该网站？

→ 不停地给用户提供他们可能会喜欢的十部手机的预测，那么
每次一个用户访问，你将会得到十个样本- (x, y) 数据对。然后利用一个在线学习
算法来更新你的参数，i.e. $\theta_j := \theta_j - \alpha (h_{\theta}(x) - y) \cdot x_j$
对这十个样本利用梯度下降法来更新参数。→ 然后可以丢弃这些数

(因为具有大量用户流
所以要使用重复数据)

→ 数据过于庞大以至于无法使用单个设备

15-6 Map-reduce and data parallelism 减少映射与数据并行

E.g. 假设 $m=400$

Training Set → 4份

Computer 1 : $\text{temp}_j^{(1)} = \sum_{i=1}^{100} (h_0(x_i^{(1)}) - y_i^{(1)}) \cdot x_j^{(1)}$

Computer 2 :
3 : Use $(x_1^{(3)}, y_1^{(3)}), \dots, (x_{100}^{(3)}, y_{100}^{(3)})$

4 : Use $(x_1^{(4)}, y_1^{(4)}), \dots, (x_{100}^{(4)}, y_{100}^{(4)})$

$\text{temp}_j^{(4)} = \sum_{i=1}^{100} (h_0(x_i^{(4)}) - y_i^{(4)}) \cdot x_j^{(4)}$

Combine results : $\theta_j := \theta_j - \alpha \frac{\text{temp}_j}{100}$

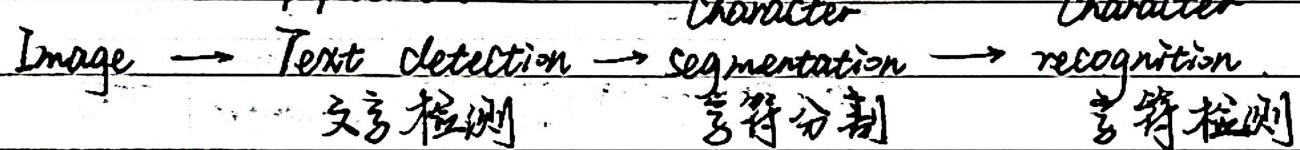
$\text{temp}_j = \text{temp}_j^{(1)} + \dots + \text{temp}_j^{(4)}$

$j = 0, \dots, n$

16 Application Example : Photo OCR 应用实例：图片文字识别

16-1 Problem description and pipeline 问题描述和流程图

• Photo OCR pipeline :



16-2 Sliding window 滑动窗口

• Text detection : 利用一系列包含正样本(对应正确文字)和负样本(无文字)
→ 训练一个固定比例的滑动窗口分类器。
→ 应用到新图片上 ⇒ 得到一张通过像素值大小反映某处有文字
的概率 → 进行一定的形态学处理 → 将含有文字的区域进行分离。

• Character segmentation : 再次使用监督学习算法，
→ 滑动扫描直到字符分割完成

E.g.

16-3 Getting lots of data: Artificial data synthesis 获取大量数据

一个可靠的方法得到高性能机器学习系统的方法：

和人工数据合成

使用一个高偏差机器学习算法，并且使用庞大训练集训练你。

→ 如何获得大量数据集？对于特定问题可以采用人工数据合成

一般两种形式 1. 重新创造新数据

II. 已经有小的标注训练集，以某种方式扩充训练集

I. E.g. for Photo OCR, 可利用免费字体库，将不同字体和不同背景结合来创建一个接近真实数据的人工合成训练集。

II. 引入失真 (Introducing distortions)

E.g. 简符拉伸；Audio: Background noise, bad cellphone connection.

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

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

(上层分析)

16-4 Ceiling analysis: What part of the pipeline to work on next.

• Image → Text detection → Character segmentation → Character recognition.

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

整个系统

Component	Accuracy	整系统	依次将流程2.3.4通过人工
Overall system	72%	17%	控制调试为 100% accuracy
Text ... 100% accuracy	89%		观察整个系统的 accuracy
Ch... seg...	90%	1%	增加
Ch... rec...	100%	1%	
2.3.4 100% accuracy			