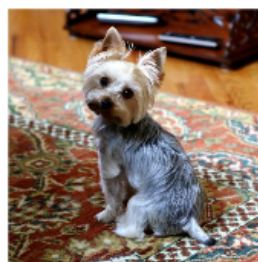


结构化机器学习项目 第二周 机器学习策略

2-1 进行误差分析

深度学习算法对训练集中的随机误差具有相当的鲁棒性。

Look at dev examples to evaluate ideas



90% accuracy
→ 10% error

Should you try to make your cat classifier do better on dogs? ←

Error analysis:

- Get ~100 mislabeled dev set examples. → 5-10 min
- Count up how many are dogs.

→ 5%
5/100

10%
↓
9.5%

"ceiling"

→ 50%
50/100

10%
↓
5%

收集错误样例：

在开发集（测试集）中，获取大约100个错误标记的例子，并统计其中有多少个是狗。

- 假设一种情况是100个数据中，有5个样例是狗，那么如果我们对数据集的错误标记做努力去改进模型的精度，那么可以提升的上限就是5%5%，即仅仅可以达到9.5%9.5%的错误率，这有时称为**性能上限**。那么这种情况下，可能这样耗时的努力方向就不是很值得的一件事情。
- 另外一种假设是100个数据中，有50多个样例是狗，那么这种情况下，我们去改进数据集的错误标记，就是一个比较值得的改进方向，可以将模型的精确度提升至95%95%。

并行分析：

- 修改那些被分类成猫的狗狗图片标签；
- 修改那些被错误分类的大型猫科动物，如：狮子，豹子等；
- 提升模糊图片的质量。

为了并行的分析，建立表格来进行。以单个错误分类样本为对象，分析每个样本错误分类的原因。

Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ←
- Fix great cats (lions, panthers, etc..) being misrecognized ←
- Improve performance on blurry images ←

Image	Dog	Great Cats	Blurry	Instagram	Comments
1	✓			✓	Pitbull
2			✓	✓	
3		✓	✓		Rainy day at zoo
⋮	⋮	⋮	⋮	⋮	
% of total	8%	43%	61%	12%	

Andr

最后，统计错误类型的百分比，这个分析步骤可以给我们一个粗略的估计，让我们大致确定是否值得去处理每个不同的错误类型。

2.2 清楚标记错误的数据

以猫和狗分类问题为例子。

Incorrectly labeled examples



DL algorithms are quite robust to random errors in the training set.

情况一：

深度学习算法对训练集中的随机误差具有相当的鲁棒性。

只要我们标记出错的例子符合随机误差，如：做标记的人不小心错误，或按错分类键。那么像这种随机误差导致的标记错误，一般来说不管这些误差可能也没有问题。

所以对于这类误差，我们可以不去用大量的时间和精力去做修正，只要数据集足够大，实际误差不会因为这些随机误差有很大的变化。

Error analysis

	Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
	...					
	98				✓	Labeler missed cat in background
	99		✓			
	100				✓	Drawing of a cat; Not a real cat.
	% of total	8%	43%	61%	6%	
Overall dev set error				10%	2%
Errors due incorrect labels				0.6% ←	0.6%
Errors due to other causes				9.4% ←	1.4%
					↑	2.1% 1.9%

Goal of dev set is to help you select between two classifiers A & B.

Andrew Ng

方法二：

训练集均是来自网上下载的20万张高清图片，当然也可以加上5000张手机非高清图片；对于开发和测试集都是手机非高清图片。

- 好处：开发集全部来自手机图片，瞄准目标；
- 坏处：训练集和开发、测试集来自不同的分布。

从长期来看，这样的分布能够给我们带来更好的系统性能。

Correcting incorrect dev/test set examples

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution
- Consider examining examples your algorithm got right as well as ones it got wrong.
- Train and dev/test data may now come from slightly different distributions.

收集正确的开发/测试集样本

- 确保使用来自同一分布的数据处理你的开发集和测试集。
- 考虑检查你的算法是否正确

2.3 快速搭建你的第一个系统并进行迭代

Speech recognition example



- • Noisy background
 - • Café noise
 - • Car noise

- • Accent
- • Far from
- • Young
- • Stutter
- • ...

Guideline:
**Build your first
system quickly,
then iterate**

- • Set up dev/test set and metric
- Build initial system quickly
- Use Bias/Variance analysis & Error analysis to prioritize next steps.

Andrew Ng

- 设置开发、测试集和优化指标（确定方向）；
- 快速地建立基本的系统；
- 使用偏差方差分析、误差分析去确定后面步骤的优先步骤。

例子2：

Speech recognition example

Speech activated rearview mirror



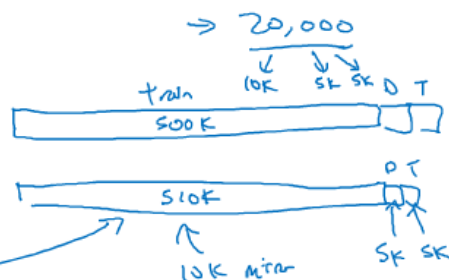
Training

Purchased data X, y
Smart speaker control
Voice keyboard
...

500,000 utterances

Dev/test

Speech activated rearview mirror



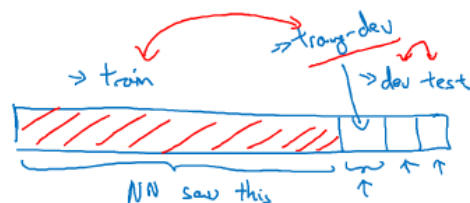
2.5 不匹配数据划分的偏差和方差

Cat classifier example

Assume humans get $\approx 0\%$ error.

Training error 1%
Dev error 10%

Training-dev set: Same distribution as training set, but not used for training



Training error	1%	1%
→ Training-dev error	9%	1.5%
→ Dev error	10%	10%
	Variance	Data mismatch
Human error	0%	0%
Training error	10%	10%
Training-dev error	11%	11%
Dev error	12%	20%
	Bias	Bias + Data mismatch

Andrew Ng

以猫分类为例，假设以人的分类误差0%0%作为贝叶斯误差。若我们模型的误差为：

- Training error : 1%1%
- Dev error : 10%

如果我们的训练集和开发、测试集来自相同的分布，那么我们可以说模型存在很大的方差问题。但如果数据来自不同的分布，那么我们就不能下这样的定论了。

那么我们如何去确定是由于分布不匹配的问题导致开发集的误差，还是由于算法中存在的方差问题所致？

设立“训练开发集”

训练开发集，其中的数据和训练数据来自同一分布，但是却不用于训练过程。

如果最终，我们的模型得到的误差分别为：

- Training error : 1%1%
- Training-dev error : 9%9%
- Dev error : 10%

那么，由于**训练开发集**尽管和训练集来自同一分布，但是却有很大的误差，模型无法泛化到同分布的数据，那么说明我们的模型存在**方差问题**。

但如果我们的模型得到的误差分别为：

- Training error : 1%1%
- Training-dev error : 1.5%1.5%
- Dev error : 10%10%

那么在这样的情况下，我们可以看到，来自同分布的数据，模型的泛化能力强，而开发集的误差主要是来自于**分布不匹配**导致的。

分布不同的偏差方差分析

通过：Human level、Training set error、Training-dev set error、Dev error、Test error 之间误差的大小，可以分别得知我们的模型，需要依次在：可避免的偏差、方差、数据分布不匹配、开发集的或拟合程度，这些方面做改进。

Bias/variance on mismatched training and dev/test sets

Human level	4%	↓ avoidable bias	4%
Training set error	7%	↑ variance	7%
Training-dev set error	10%	↓ data mismatch	10%
→ Dev error	12%	↑ degree of overfitting to dev set.	6%
→ Test error	12%		6%

通常情况来说，通过不同的集合上的误差分析，我们得出的结果会是中间一列误差由小变大，即误差上升的情况。但是也有一定的可能会出现右边一列误差在开发测试集上又表现的好的情况。

下面通过一个后视镜语音检测的例子来说明。我们以该例子建立更加一般的表格。

More general formulation

Rearview mirror

	General speech recognition	Rearview mirror speech data.	
Human level	"Human level" 4%	6%	↑ Avoidable bias
Error on examples trained on	"Training error" 7%	6%	
Error on examples <u>not</u> trained on	"Training-dev error" 10%	"Dev/Test error" 6%	↑ Variance

↔ data mismatch

Andrew Ng

其中，横向分别是：普通语音识别数据、后视镜语音识别数据；纵向分别是：Human level、训练数据误差、未训练数据误差。表格中不同的位置分别代表不同的数据集。

通常情况下，我们分析误差会是一个递增的情况，但是可能对于我们的模型，在后视镜语音识别的数据数据上，已经达到人类水平误差的6%6%了，所以最终的开发测试集也会是6%6%的误差，要比训练误差和训练开发误差都要小。所以如果遇到这种情况，就要利用上表进行分析。

2.6 定位数据不匹配问题

Addressing data mismatch

- • Carry out manual error analysis to try to understand difference between training and dev/test sets

E.g. noisy - car noise

street numbers

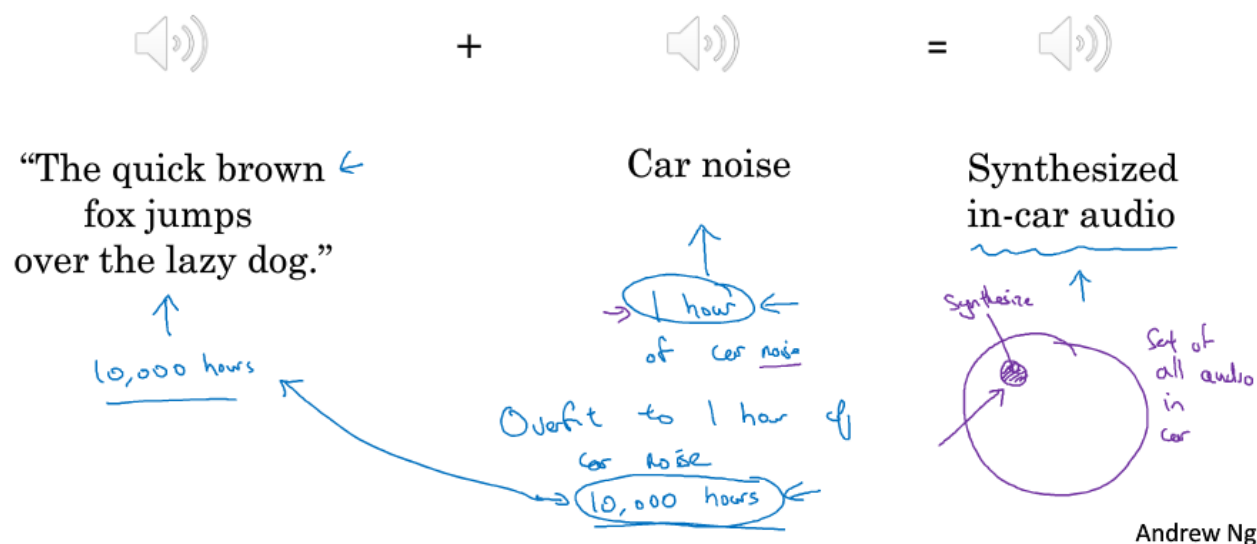
- • Make training data more similar; or collect more data similar to dev/test sets

E.g. Simulate noisy in-car data

如果通过上一节的误差分析，我们可以得知，模型最终在开发和测试集上的误差最终是由于数据分布不匹配而导致。那么这样的情况下如何解决？

- 进行人工误差分析，尝试去了解训练集和开发测试集的具体差异在哪里。如：噪音等；
- 尝试把训练数据变得更像开发集，或者收集更多的类似开发集和测试集的数据，如增加噪音；

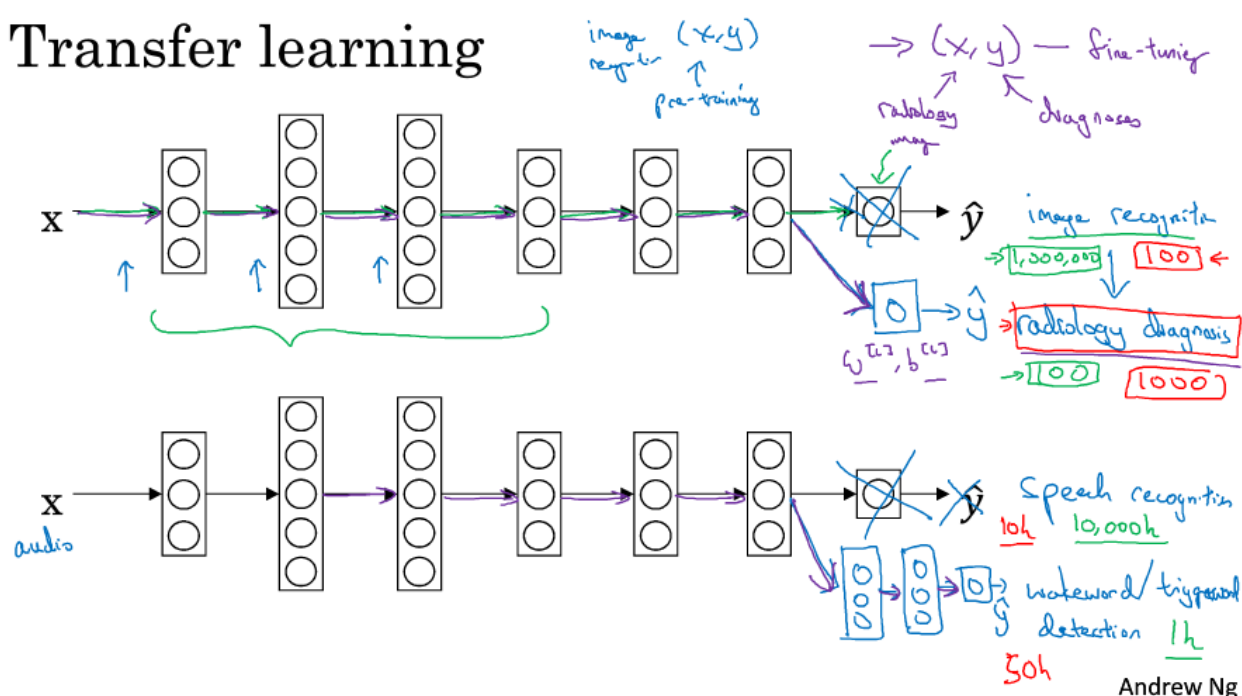
Artificial data synthesis



2.7 迁移学习

将从一个任务中学到的知识，应用到另一个独立的任务中。

Transfer learning



迁移学习的意义：

迁移学习适合以下场合：迁移来源问题有很多数据，但是迁移目标问题却没有那么多的数据。

同样一个例子是语音识别，可能在普通的语音识别中，我们有庞大的数据量来训练模型，所以模型从中学到了很多人类声音的特征。但是对于触发字检测任务，可能我们拥有的数据量很少，所以对于这种情况下，学习人类声音特征等知识就显得相当重要。所以迁移学习可以帮助我们建立一个很好的唤醒字检测系统。

Trans from $A \rightarrow B$

- 任务A和任务B有着相同的输入；
- 任务A所拥有的数据要远远大于任务B（对于更有价值的任务B，任务A所拥有的数据要比B大很多）；
- 任务A的低层特征学习对任务B有一定的帮助；

与迁移学习的串行学习方式不同，在多任务学习中，多个任务是并行进行学习的，同时希望各个任务对其他任务均有一定的帮助。

Simplified autonomous driving example



$x^{(i)}$

$y^{(i)}$
 $(4,1)$
 Pedestrians
 Cars
 Stop signs
 Traffic lights
 \vdots

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(m)} \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}$$

$(4, m)$

自动驾驶的例子

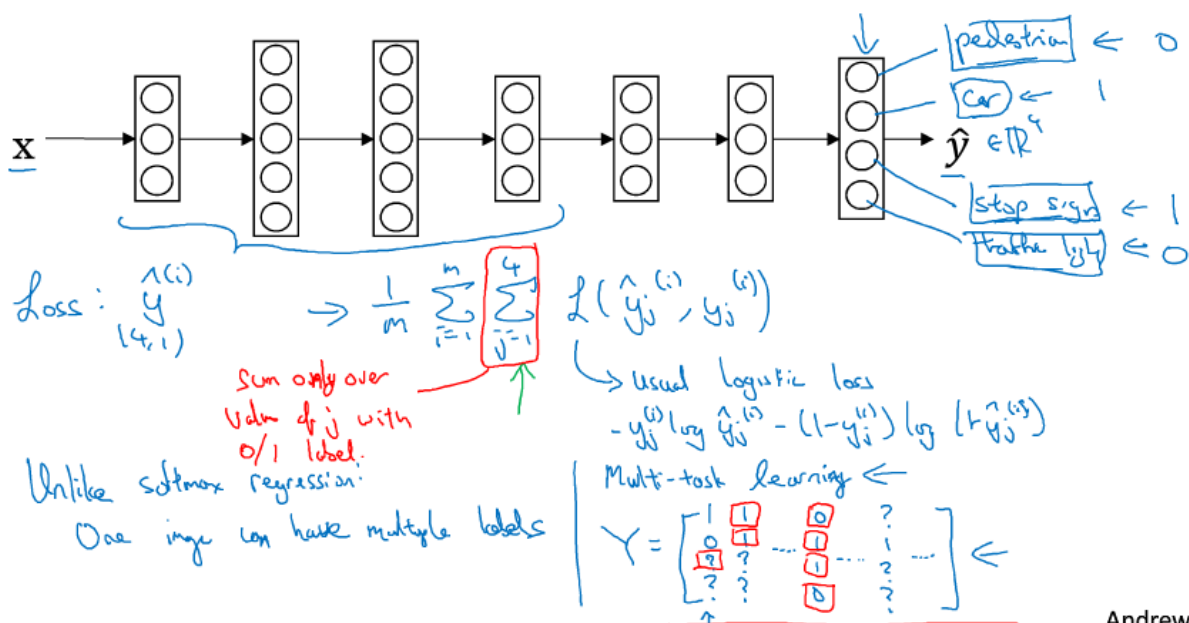
假设在自动驾驶的例子中，我们需要检测的物体很多，如行人、汽车、交通灯等等。

对于现在的任务，我们的目标值变成了一个向量的形式向量中的每一个值代表检测到是否有如行人、汽车、交通灯等，一张图片有多个标签。

$$y^{(i)} = [11110101] \text{ Pedestrians Cars Road signs - Stop Traffic lights}$$

模型的神经网络结构如下图所示：

Neural network architecture

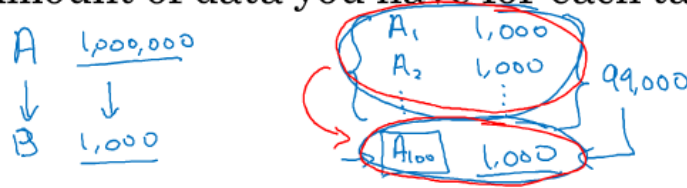


什么时候使用多任务学习有效？

- 如果训练的一组任务可以共用低层特征；
- 通常，对于每个任务大量的数据具有很大的相似性；（如，在迁移学习中由任务A“100万数据”迁移到任务B“1000数据”；多任务学习中，任务 A_1, \dots, A_n ，每个任务均有1000个数据，合起来就有 $1000n$ 个数据，共同帮助任务的训练）
- 可以训练一个足够大的神经网络并同时做好所有的任务。

When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.



- Can train a big enough neural network to do well on all the tasks.

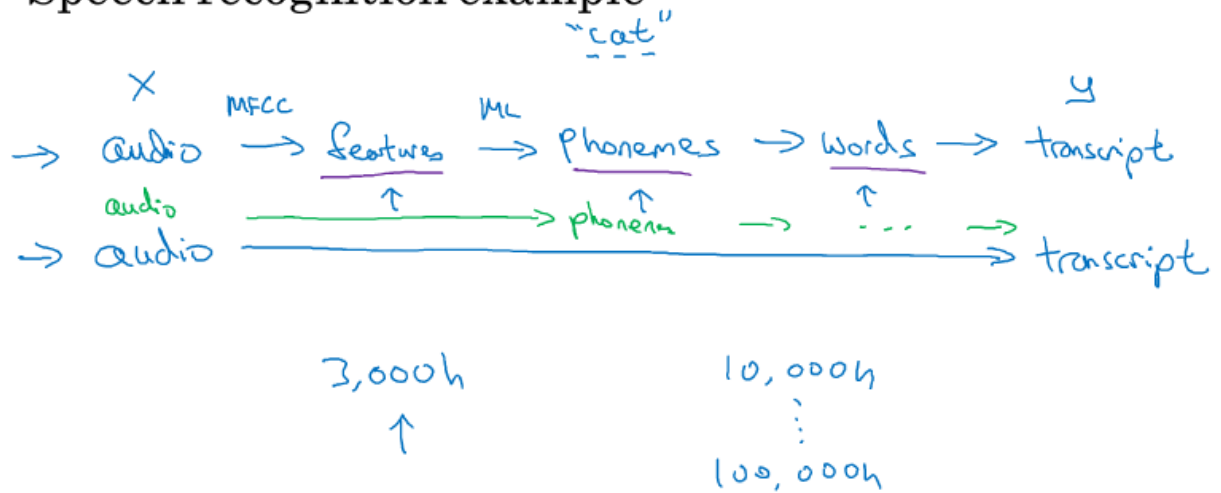
2.9 什么是端到端的深度学习？

定义：

相对于传统的一些数据处理系统或者学习系统，它们包含了多个阶段的处理过程，而端到端的深度学习则忽略了这些阶段，用单个神经网络来替代。

What is end-to-end learning?

Speech recognition example

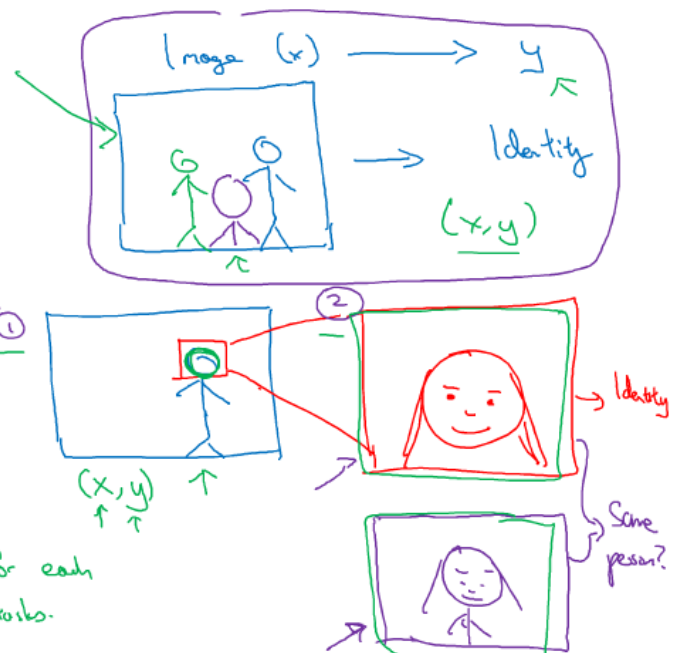


Face recognition



[Image courtesy of Baidu]

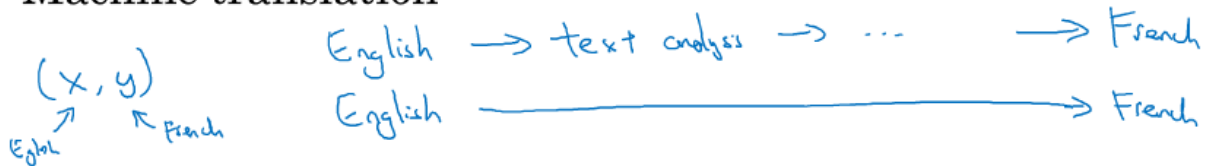
Have data for each of 2 sub-tasks.



Andrew Ng

More examples

Machine translation



Estimating child's age:



Andrew Ng

优缺点：

- 优点：

端到端学习可以直接让数据“说话”；所需手工设计的组件更少。

- 缺点：

需要大量的数据；排除了可能有用的手工设计组件。

Pros and cons of end-to-end deep learning

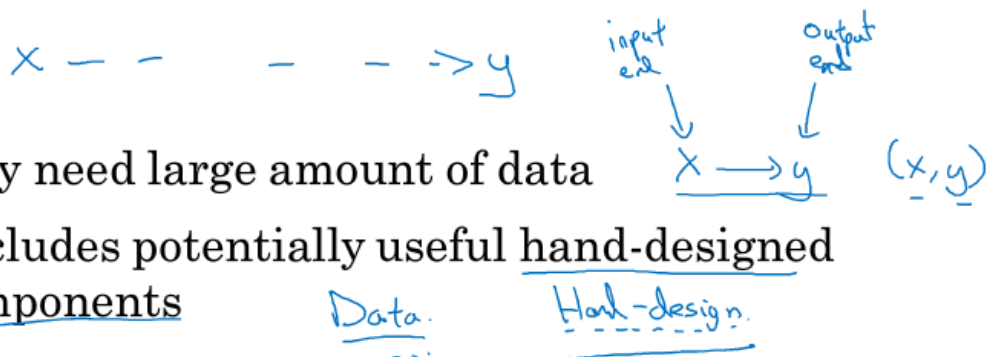
Pros:

- Let the data speak $x \rightarrow y$
- Less hand-designing of components needed



Cons:

- May need large amount of data
- Excludes potentially useful hand-designed components

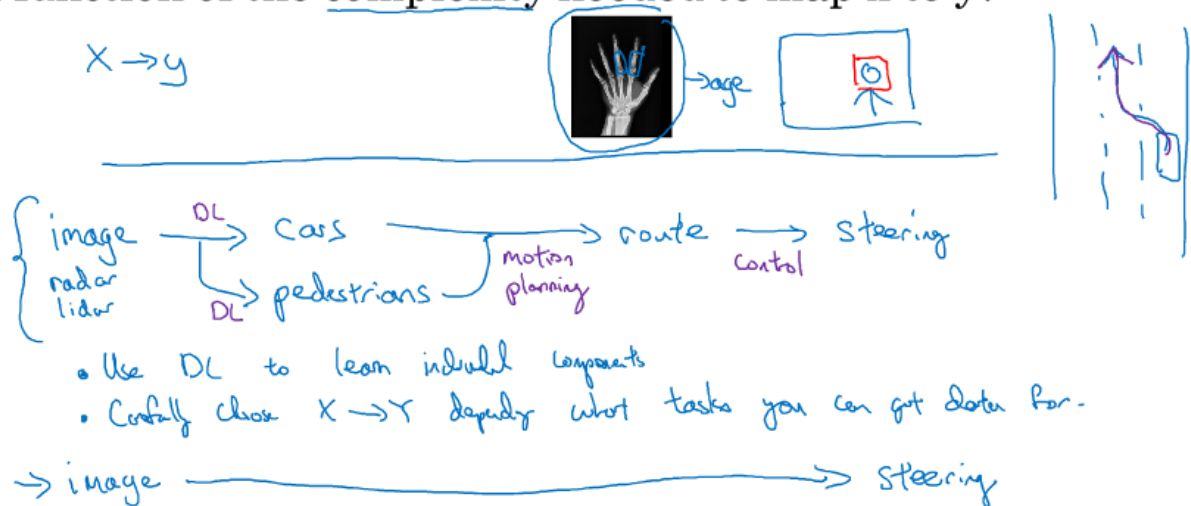


Andrew

关键问题

Applying end-to-end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map x to y ?



Andrew N

例如在无人驾驶中。