



deeplearning.ai

- Transfer learning was a sequential process, learn from A transfer to B, $A \rightarrow B$
 - Multitask learning parallelizes this
- Also \rightarrow these can be thought of as methods to use when you don't have much data for predictions

Learning from multiple tasks

Multi-task learning

Simplified autonomous driving example



An autonomous car needs to identify

- Pedestrians
 - Cars
 - Stop signs
 - Traffic lights
- for proper operation

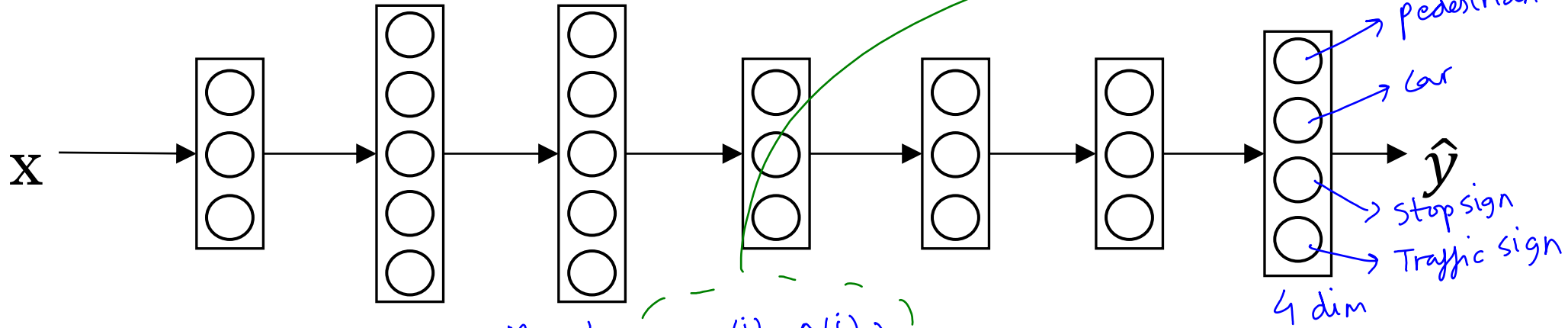
- If Image = $X^{(i)}$

then o/p has 4 labels $\rightarrow y^{(i)} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} \text{pedestrian} \\ \text{car} \\ \text{stop sign} \\ \text{Traffic light} \end{bmatrix}$
(4x1)

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

$$Y = 4 \times M$$

Neural network architecture



If y is only partially labelled
 $y = \begin{bmatrix} 1 & ? & 1 & ? & 1 & ? \\ 1 & 0 & 0 & \dots & 1 & \dots \\ 0 & ? & 0 & \dots & 0 & \dots \\ ? & ? & 1 & ? & ? & ? \end{bmatrix}$ → then just omit the values of y_j that don't have 0/1
 Multi task Learning will still work!

$$\text{Loss} : \frac{1}{M} \sum_{i=1}^M \sum_{j=1}^4 L(y_j^{(i)}, \hat{y}_j^{(i)})$$

$$-y_j^{(i)} \log \hat{y}_j^{(i)} - (1 - y_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$$

(logistic loss)

- Diff b/w this & Softmax in the last layer

Softmax assigned a single label to a single example, this can have multiple labels per example
 either a cat or dog or rabbit or horse

can have car & traffic lights.

This is doing multitask learning
 - solves 4 problems
 - Is there a car? Is there a STOP sign? Is there pedestrian etc.
 - could have used 4 diff NN - 1 to solve each problem, but since here the features may be shared eg stop sign occurs together w/ traffic lights OR pedestrians are often in the Img having traffic light

Hence it may be better to have it in 1 NN.
 - Also more expensive to do 4 v/s 1

When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features. *Stop sign, car, traffic lights are all features of Roads*

- Usually: Amount of data you have for each task is quite similar. *Is not a Hard & fast Rule*

A 1,000,000 examples
↓
B 1000
(Transfer learning)

A₁ 1000
A₂ 1000
A₃ 1000
⋮ ⋮
A₁₀₀ 1000
(Multitask learning 100 tasks)

each have 1000 examples
- If we try to do just task A₁₀₀ in Isolation, we have 1000 examples to train on
- But if we do all of them, it's like we have 99K more samples (as we learn something from others as well)

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