

Deep Learning: Functional View and Features

Dr. Xiaowei Huang

<https://cgi.csc.liv.ac.uk/~xiaowei/>

Up to now,

- Overview of Machine Learning
- Traditional Machine Learning Algorithms
- Deep learning
 - Introduction to Tensorflow
 - Introduction to Deep Learning (history of deep learning, including perceptron, multi-layer perceptrons, why now?, etc)

Today's Topics

- Functional View of DNNs
- Learning Representations & Features

Functional View of DNNs

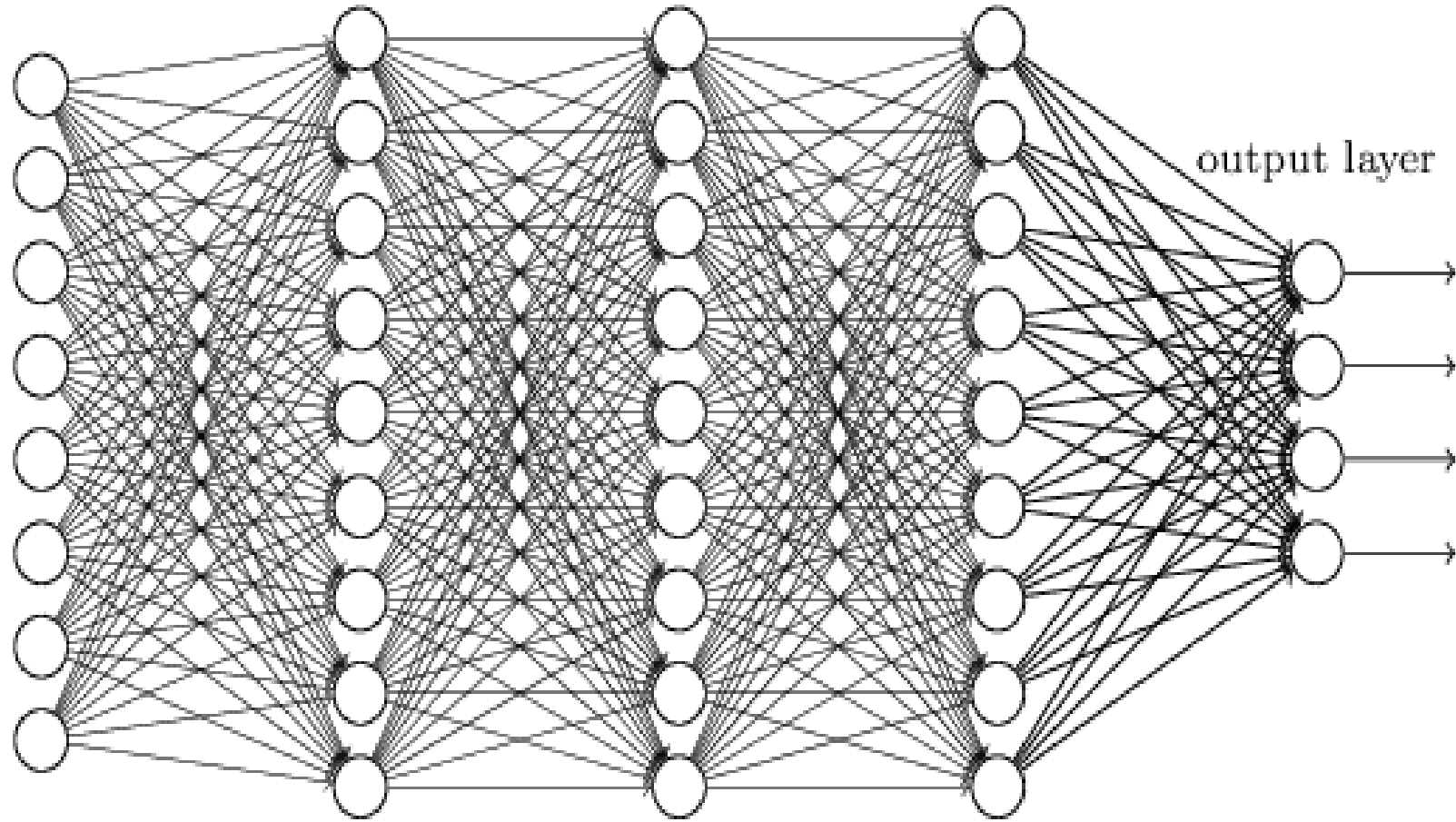
input layer

hidden layer 1

hidden layer 2

hidden layer 3

output layer



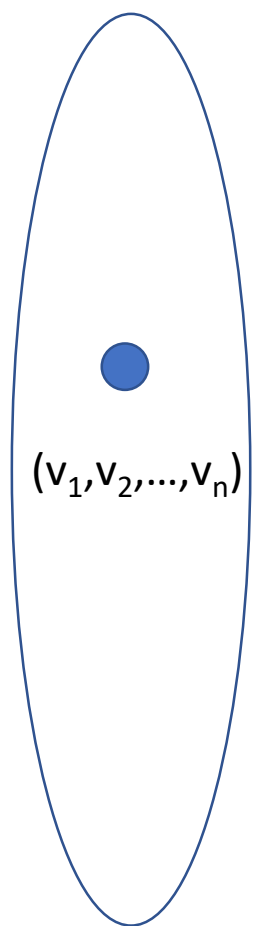
Functional View of DNNs

- A family of parametric, non-linear and hierarchical representation learning functions, which are massively optimized with stochastic gradient descent to encode domain knowledge, i.e. domain invariances, stationarity.

$$a_L(x; \theta_{1,\dots,L}) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

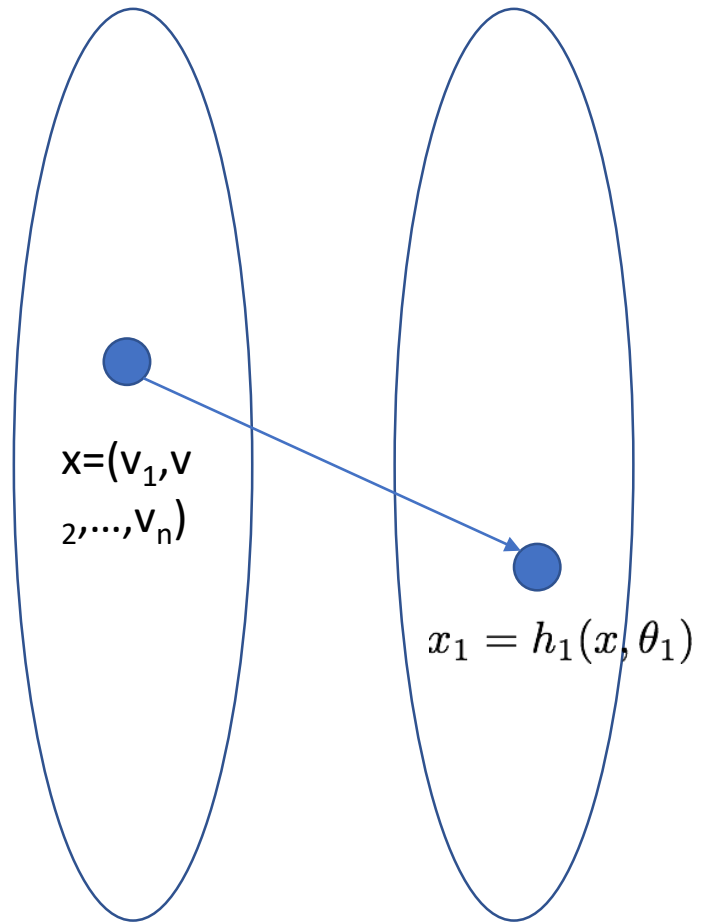
x : input, θ_l : parameters for layer l , $a_l = h_l(x, \theta_l)$: (non-)linear function

Illustration of DNN function



Input space

Illustration of DNN function



Input space

space
of 1st hidden layer

Illustration of DNN function

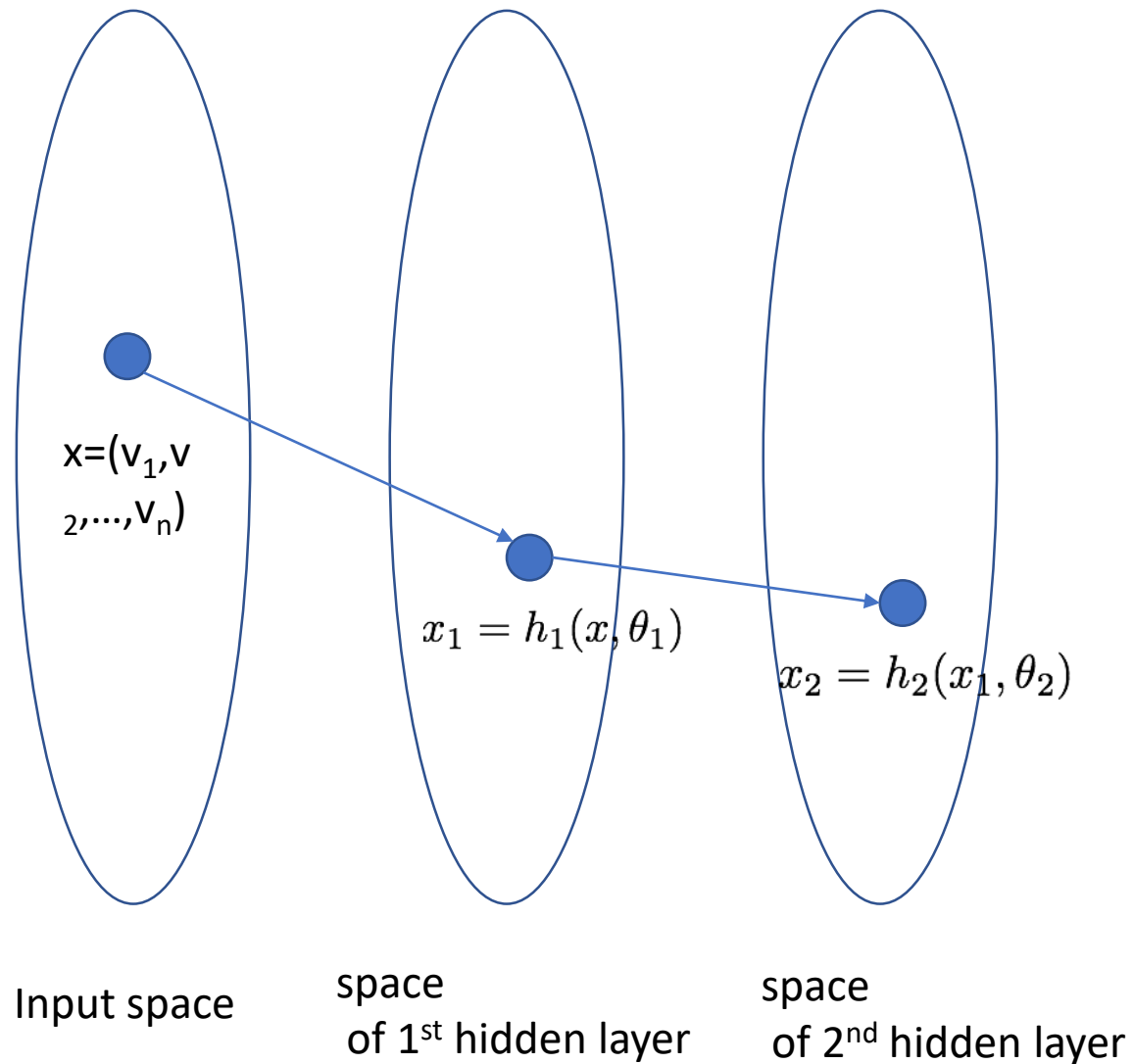
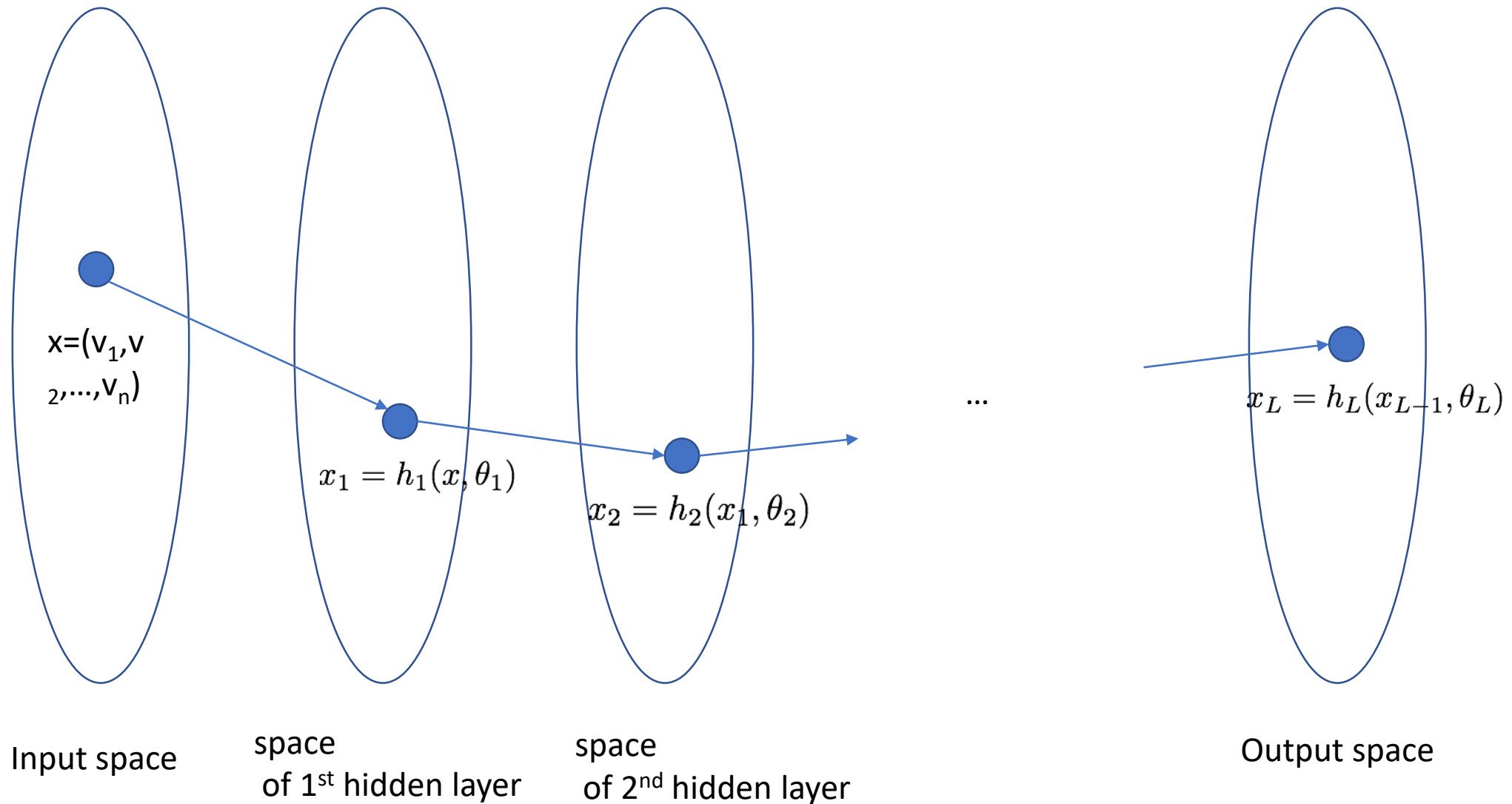


Illustration of DNN function



Note:

- Functions h_1, h_2, \dots, h_L , are usually given.
- Parameters $\theta_1, \dots, \theta_L$ are obtained by learning algorithm

Training Objective

Given training corpus $\{X, Y\}$ find optimal parameters

Find an optimal model
parameterised over θ

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_{1,\dots,L}))$$

Ground truth

Prediction

accumulated loss

Loss function

Learning Representations & Features

Raw digital representation -- Video

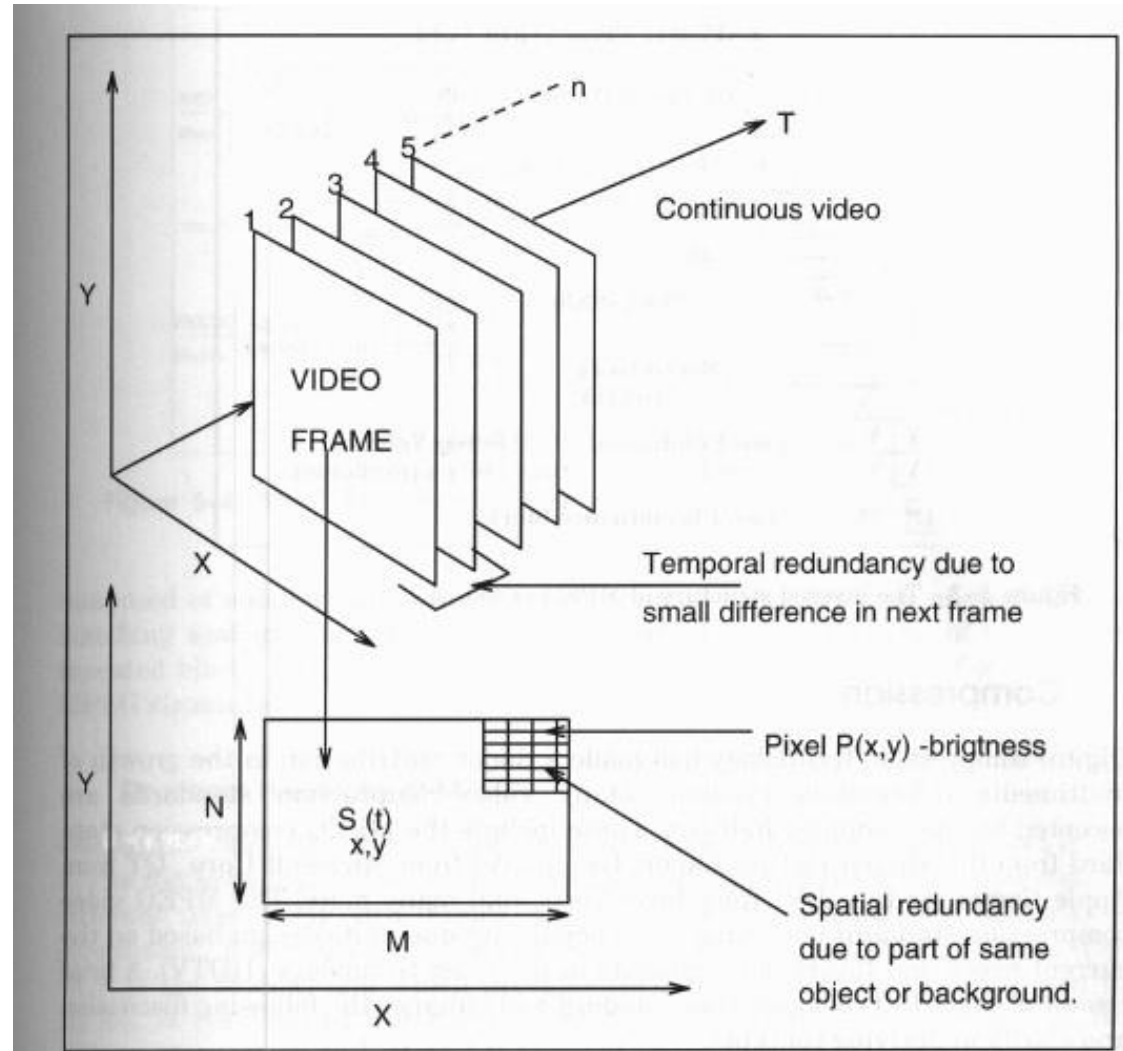
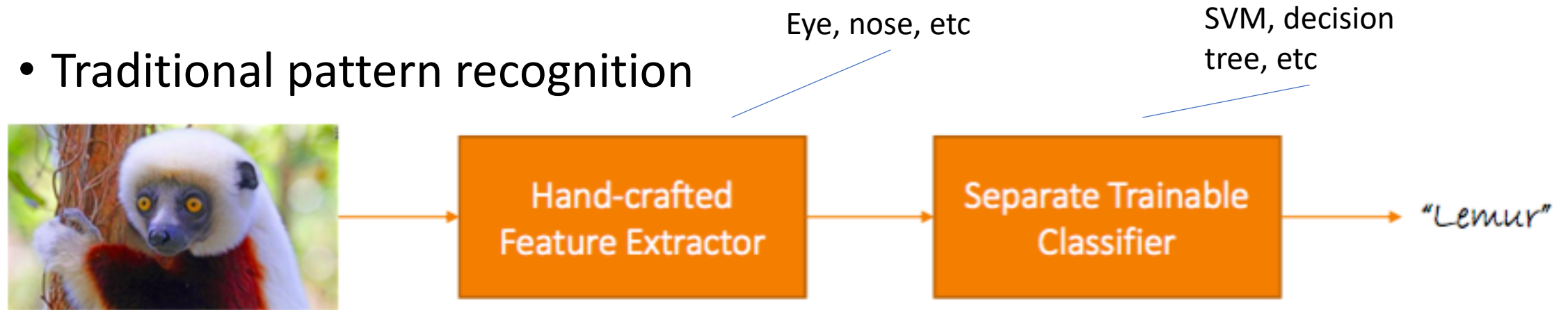


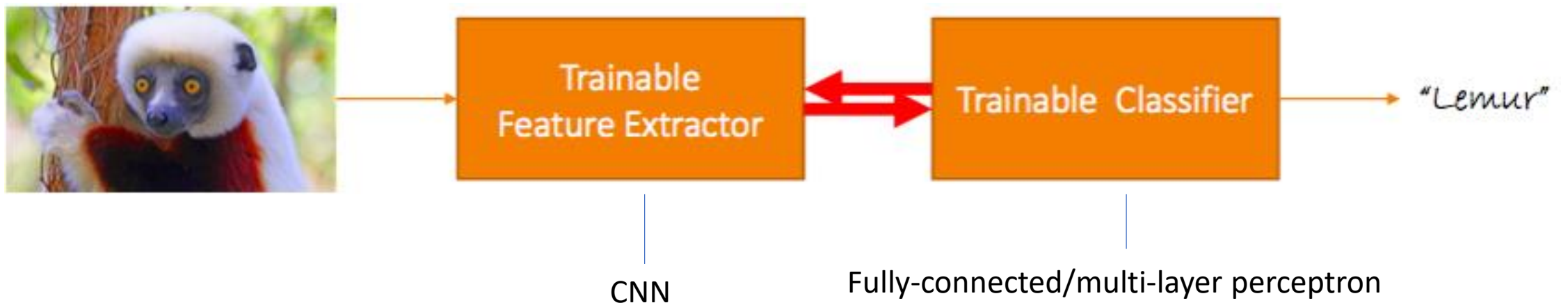
Figure 5-2 Representation of digital video.

Learning Representations & Features

- Traditional pattern recognition

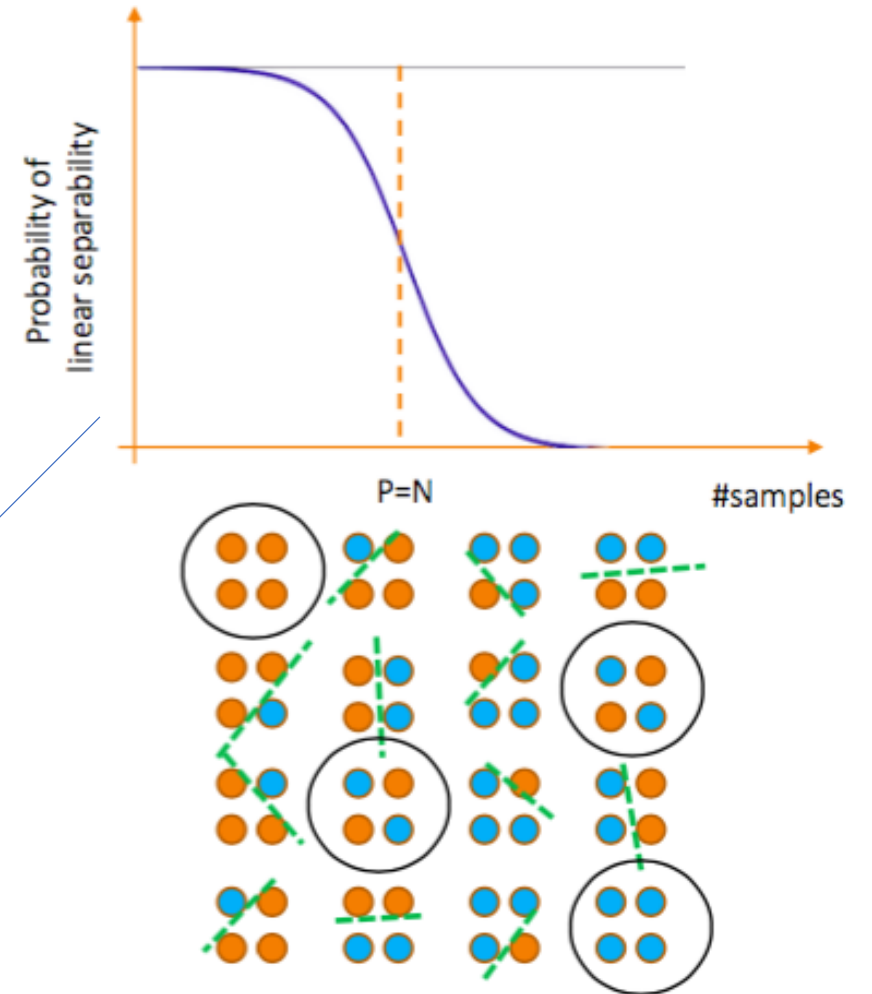


- End-to-end learning Features are also learned from data



Non-separability of linear machines

- $X = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$
- Given the n points there are in total 2^n dichotomies
- Only about d are linearly separable
- With $n > d$ the probability X is linearly separable converges to 0 very fast
- The chances that a dichotomy is linearly separable is very small



Non-linearizing linear machines

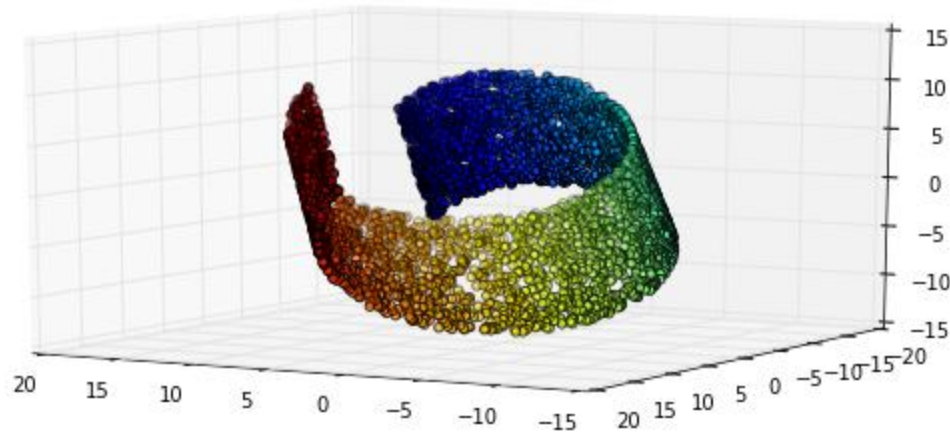
- Most data distributions and tasks are non-linear
- A linear assumption is often convenient, but not necessarily truthful
- Problem: How to get non-linear machines without too much effort?
- Solution: Make features non-linear
- What is a good non-linear feature?
 - Non-linear kernels, e.g., polynomial, RBF, etc
 - Explicit design of features (SIFT, HOG)?

Good features

- Invariant
 - But not too invariant
- Repeatable
 - But not bursty
- Discriminative
 - But not too class-specific
- Robust
 - But sensitive enough

How to get good features?

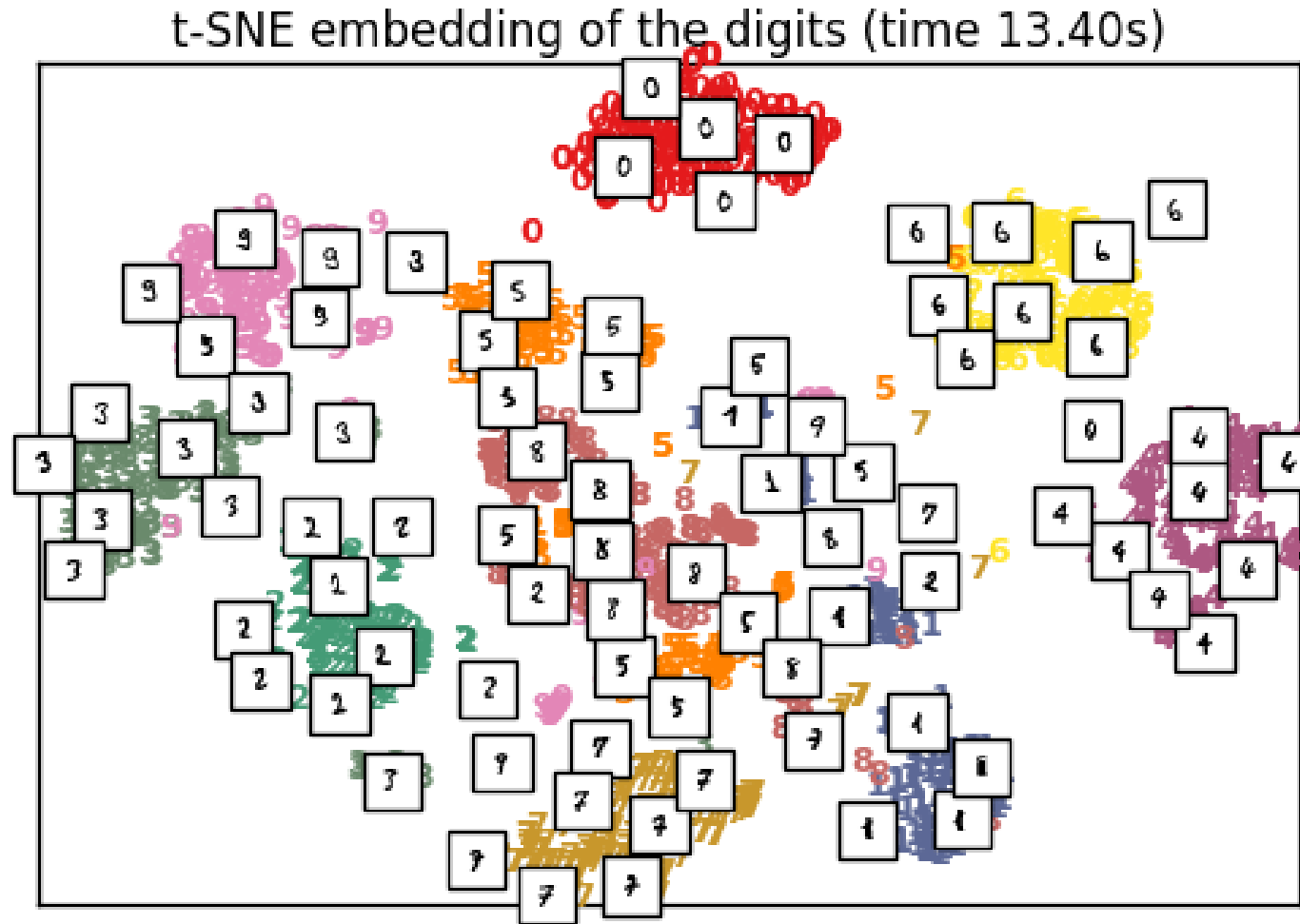
- High-dimensional data (e.g. faces) lie in lower dimensional manifolds



Every point
represents an
input sample.

- This is so-called "swiss roll". The data points are in 3d, but they all lie on 2d manifold, so the dimensionality of the manifold is 2, while the dimensionality of the input space is 3.

Digits lie in a low-dimensional manifold



It is not linear, but can be largely separated with a dimensionality much less than 28×28

How to get good features?

- High-dimensional data (e.g. faces) lie in lower dimensional manifolds
- Although the data points may consist of thousands of features, they may be described as a function of only a few underlying parameters.
 - That is, the data points are actually samples from a low-dimensional manifold that is embedded in a high-dimensional space.
- Goal: discover these lower dimensional manifolds
 - These manifolds are most probably highly non-linear

How to get good features?

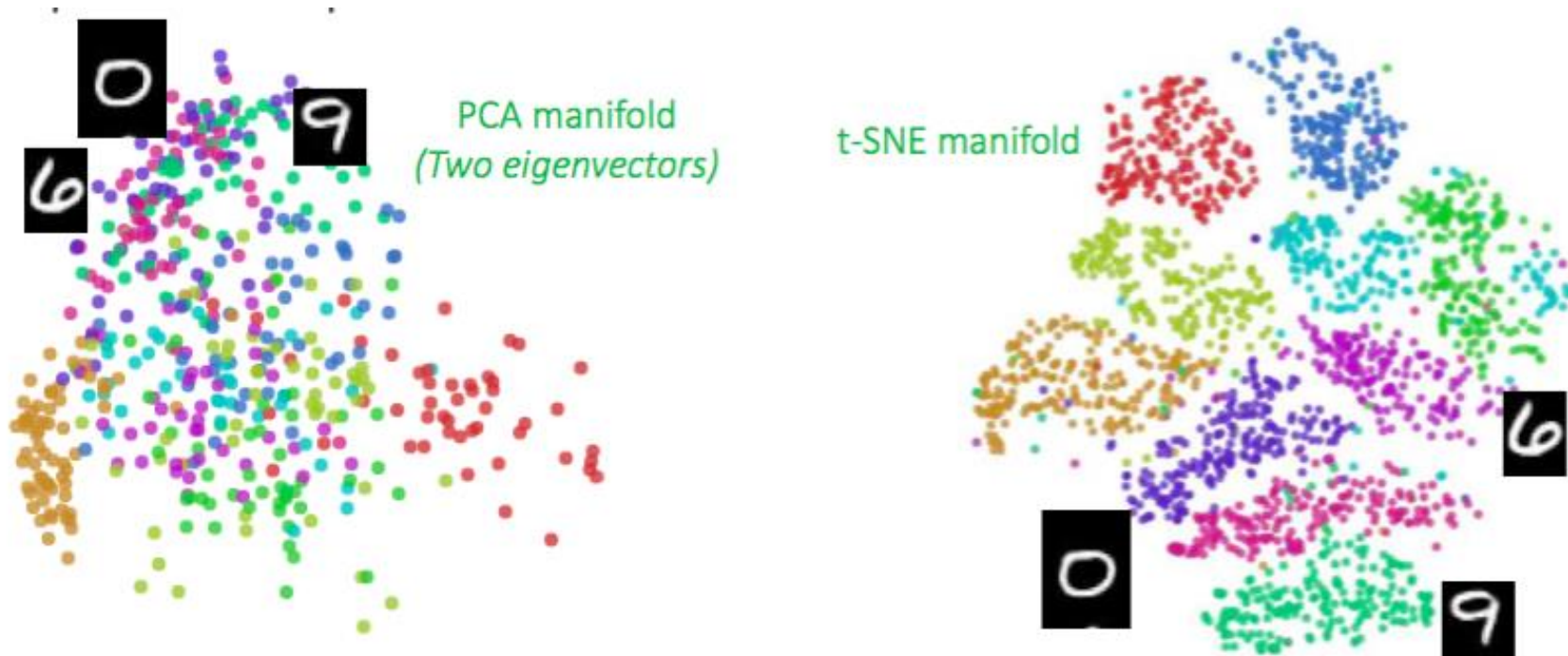
- High-dimensional data (e.g. faces) lie in lower dimensional manifolds
 - Goal: discover these lower dimensional manifolds
 - These manifolds are most probably highly non-linear
- Hypothesis (1): Compute the coordinates of the input (e.g. a face image) to this non-linear manifold -> data become separable
- Hypothesis (2): Semantically similar things lie closer together than semantically dissimilar things

Feature manifold example

- Raw data live in huge dimensionalities
- Semantically meaningful raw data prefer lower dimensional manifolds
 - Which still live in the same huge dimensionalities
- Can we discover this manifold to embed our data on?

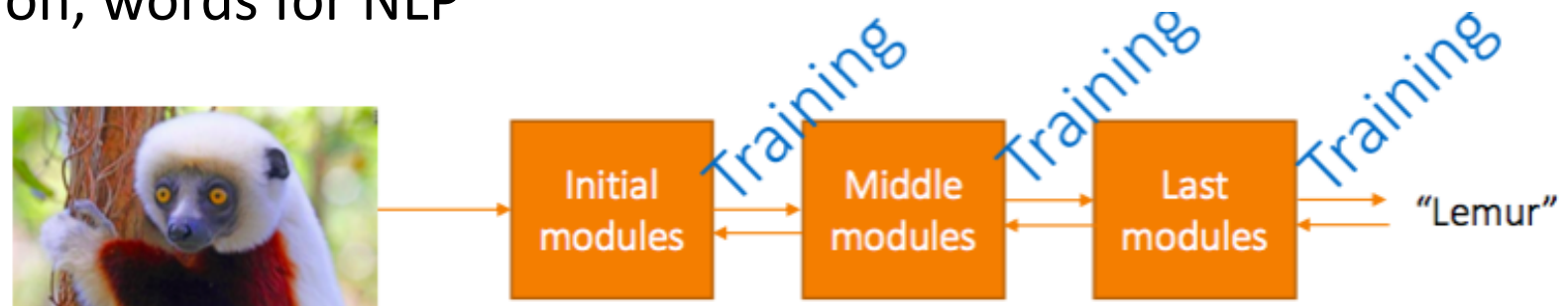
The digits manifold

- There are good features and bad features, good manifold representations and bad manifold representations
- $28 \text{ pixels} \times 28 \text{ pixels} = 784 \text{ dimensions}$



End-to-end learning of feature hierarchies

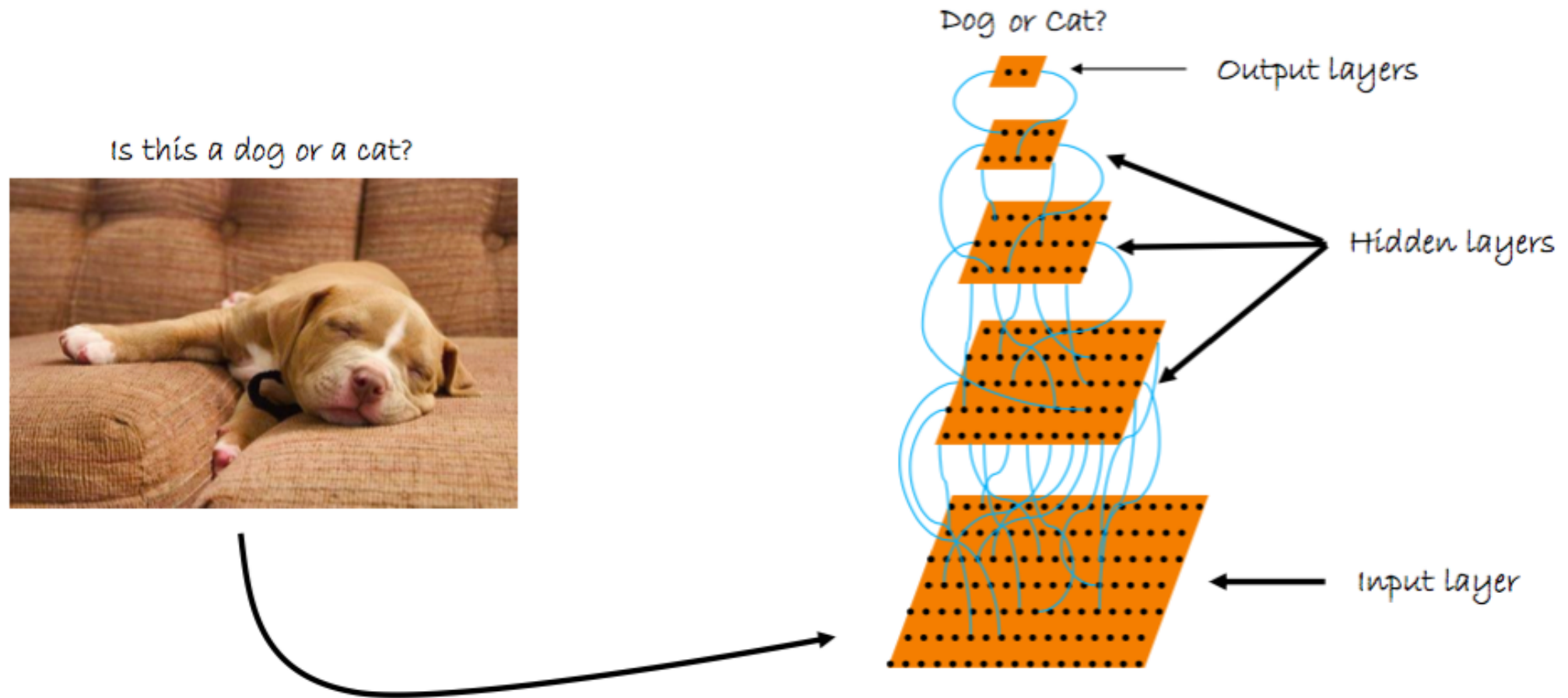
- A pipeline of successive modules
- Each module's output is the input for the next module
- Modules produce features of higher and higher abstractions
 - Initial modules capture low-level features (e.g. edges or corners)
 - Middle modules capture mid-level features (e.g. circles, squares, textures)
 - Last modules capture high level, class specific features (e.g. face detector)
- Preferably, input as raw as possible
 - Pixels for computer vision, words for NLP



Why learn the features?

- Manually designed features
 - Often take a lot of time to come up with and implement
 - Often take a lot of time to validate
 - Often they are incomplete, as one cannot know if they are optimal for the task
- Learned features
 - Are easy to adapt
 - Very compact and specific to the task at hand
 - Given a basic architecture in mind, it is relatively easy and fast to optimize
- Time spent for designing features now spent for designing architectures

Convolutional networks in a nutshell



Feature visualization by CNN

