# Lecture 8. Deep Learning. Convolutional ANNs. Autoencoders

## COMP90051 Statistical Machine Learning

Semester 2, 2018
Lecturer:  Ben Rubinstein

THE UNIVERSITY OF
MELBOURNE

# This lecture

- ## Deep learning
  - * Representation capacity
  - * Deep models and representation learning

- ## Convolutional Neural Networks
  - * Convolution operator
  - * Elements of a convolution-based network

- ## Autoencoders
  - * Learning efficient coding

# Deep Learning and Representation Learning

### Hidden layers viewed as feature space transformation

# Representational capacity

- ANNs with a single hidden layer are universal approximators

- For example, such ANNs can represent any Boolean function

$$OR(x_1, x_2) \qquad u = g(x_1 + x_2 - 0.5)$$

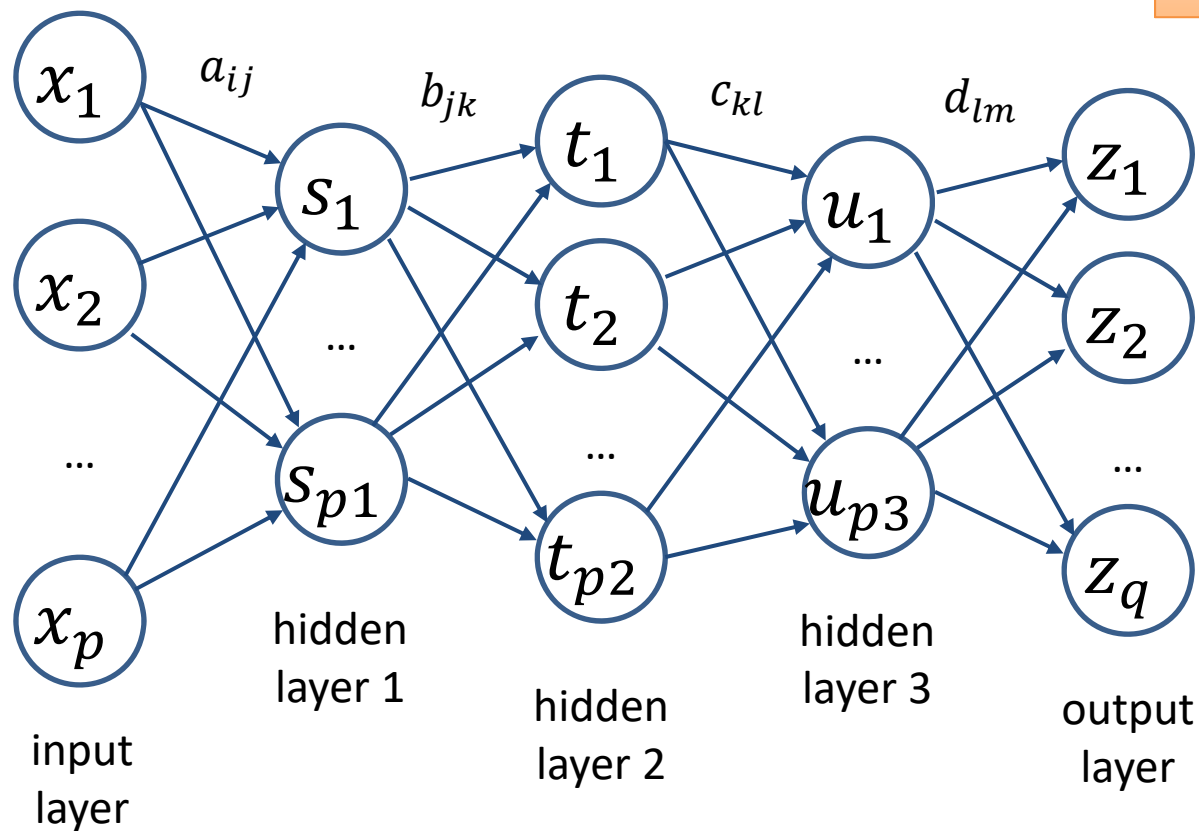$$AND(x_1, x_2) \qquad u = g(x_1 + x_2 - 1.5)$$

$$NOT(x_1) \qquad\qquad u = g(-x_1)$$

$$g(r) = 1 \text{ if } r \geq 0 \text{ and } g(r) = 0 \text{ otherwise}$$

- Any Boolean function over $m$ variables can be implemented using a hidden layer with up to $2^m$ elements

- More *efficient* to stack several hidden layers

# Deep networks

"Depth" refers to number of hidden layers



$$s = \tanh(A'x) \quad t = \tanh(B's) \quad u = \tanh(C't) \quad z = \tanh(D'u)$$
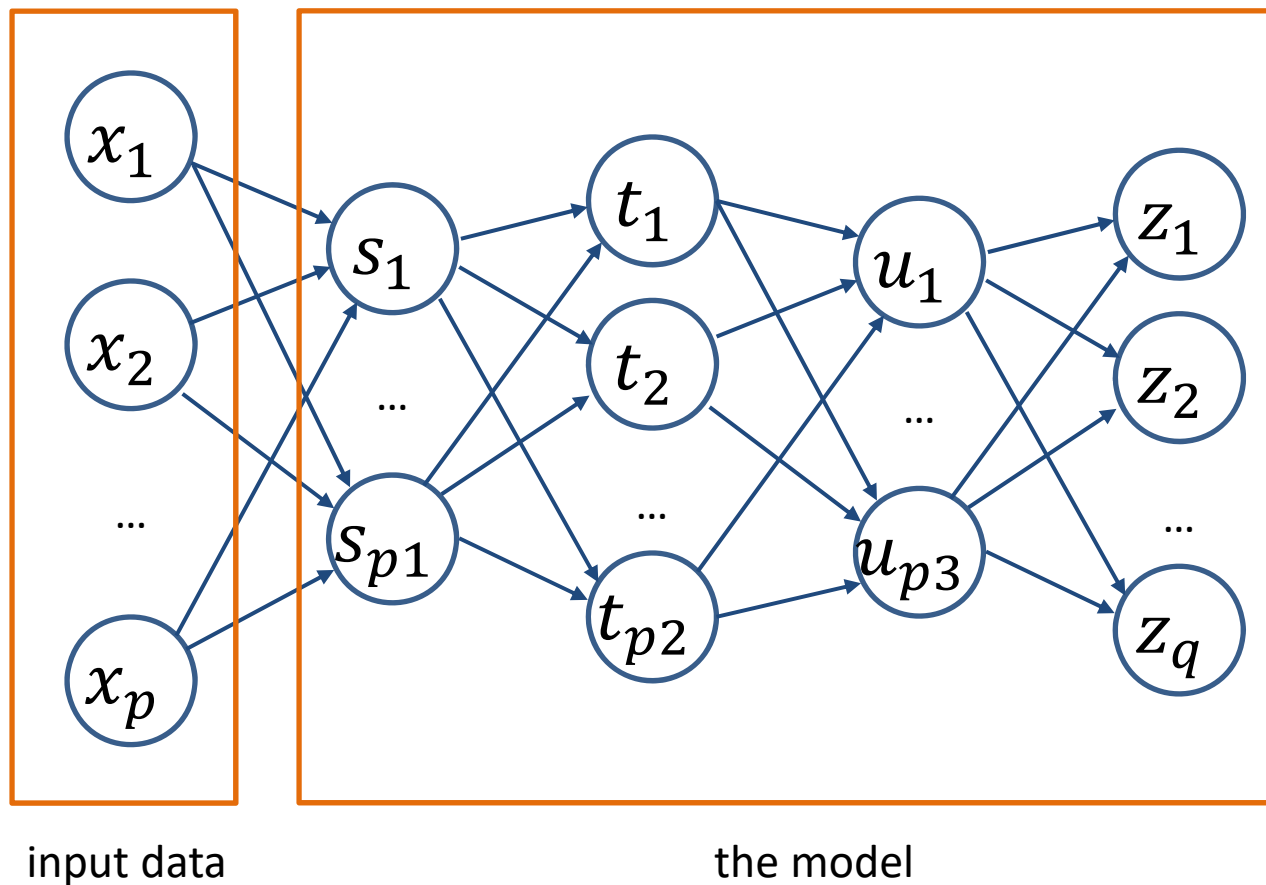
5

# Deep ANNs as representation learning

- Consecutive layers form _representations_ of the input of increasing complexity

- An ANN can have a simple _linear_ output layer, but using complex _non-linear_ representation
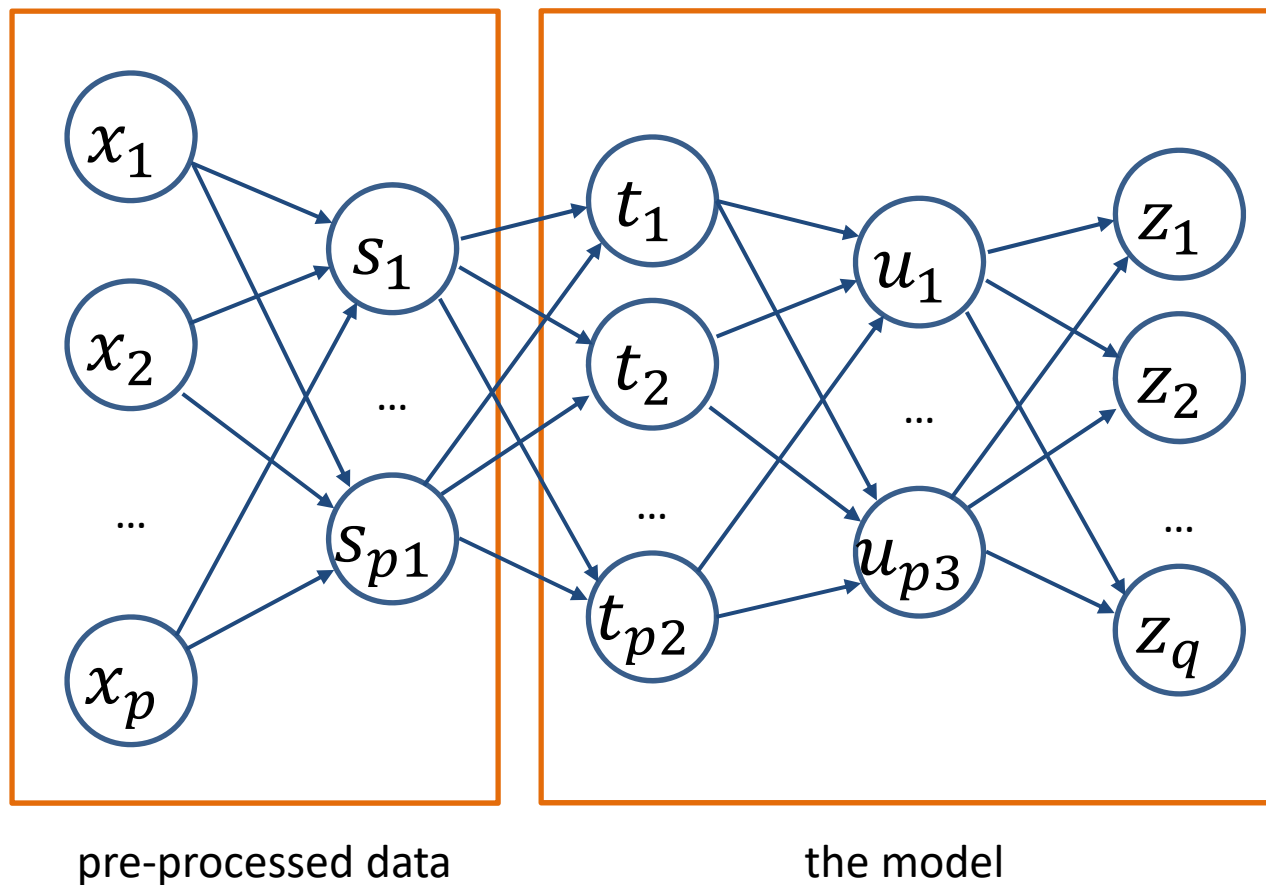
$$z = \tanh\left(D'\left(\tanh\left(C'\left(\tanh\left(B'\left(\tanh(A'x)\right)\right)\right)\right)\right)\right)$$

- Equivalently, a hidden layer can be thought of as the transformed feature space, e.g., $u = \varphi(x)$

- Parameters of such a transformation are learned from data

Bias terms are omitted for simplicity

# ANN layers as data transformation



input data                          the model

# ANN layers as data transformation



pre-processed data                    the model

# ANN layers as data transformation



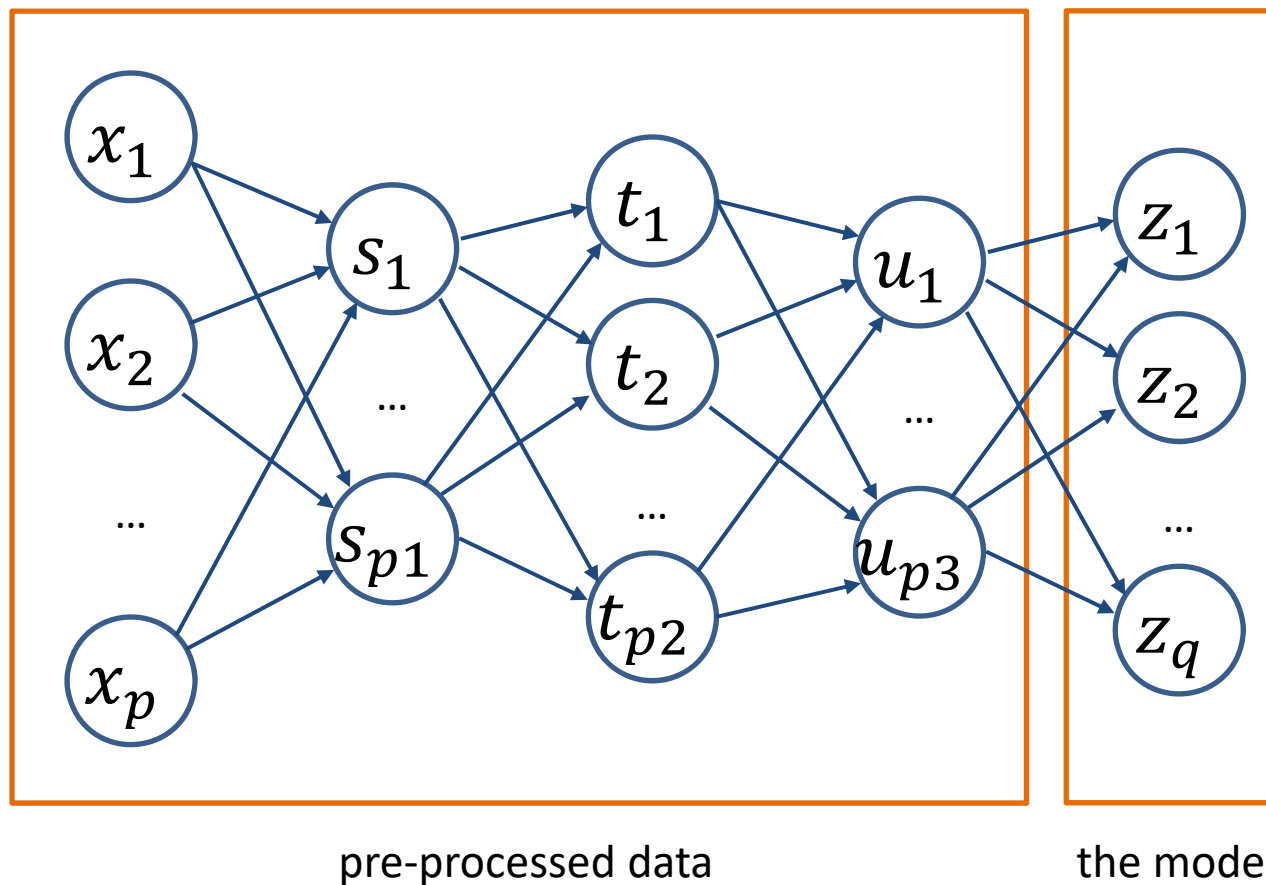pre-processed data                              the model

# ANN layers as data transformation



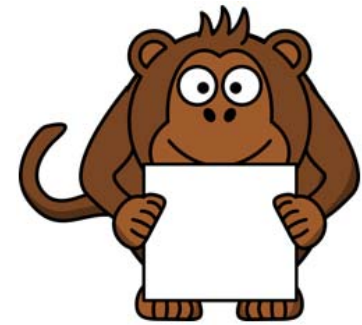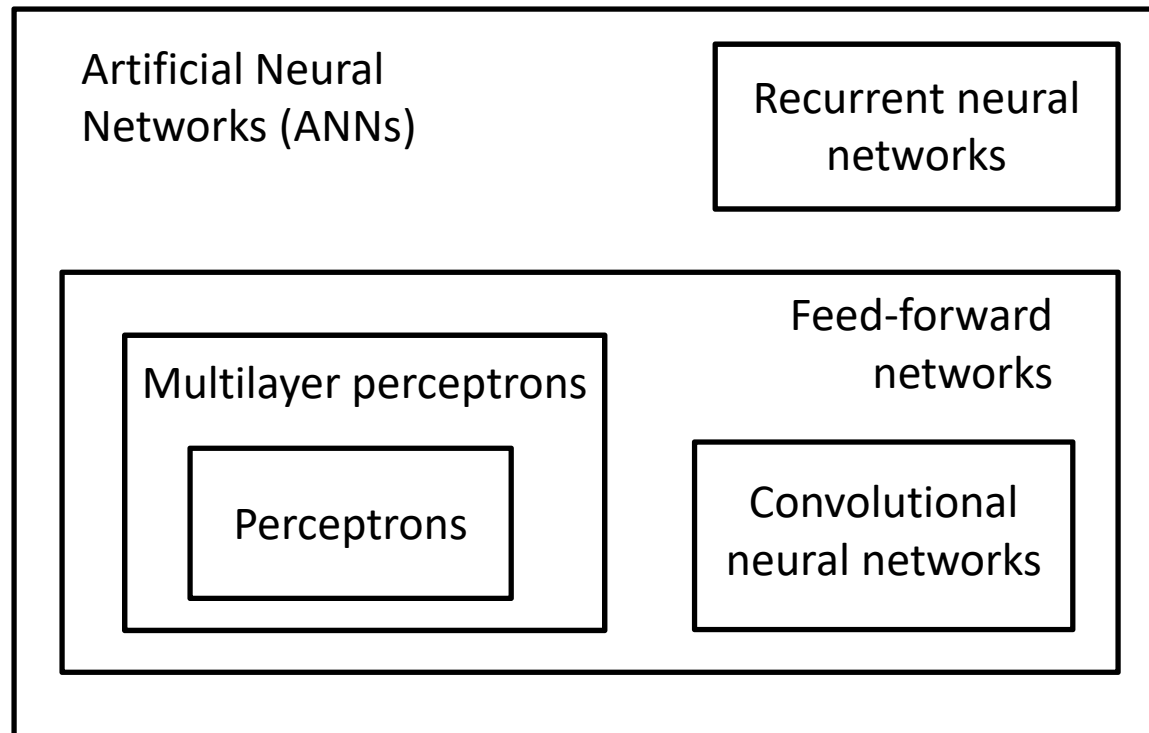pre-processed data                          the model

# Depth vs width

- A single infinitely wide layer in theory gives a universal approximator

- However (empirically) depth yields more accurate models
  Biological inspiration from the eye:
  * first detect small edges and color patches;
  * compose these into smaller shapes;
  * building to more complex detectors, of e.g. textures, faces, etc.

- Seek to mimic layered complexity in a network

- However *vanishing gradient problem* affects learning with very deep models

# Animals in the zoo

Artificial Neural
Networks (ANNs)

Recurrent neural
networks

Feed-forward
networks

Multilayer perceptrons

Perceptrons
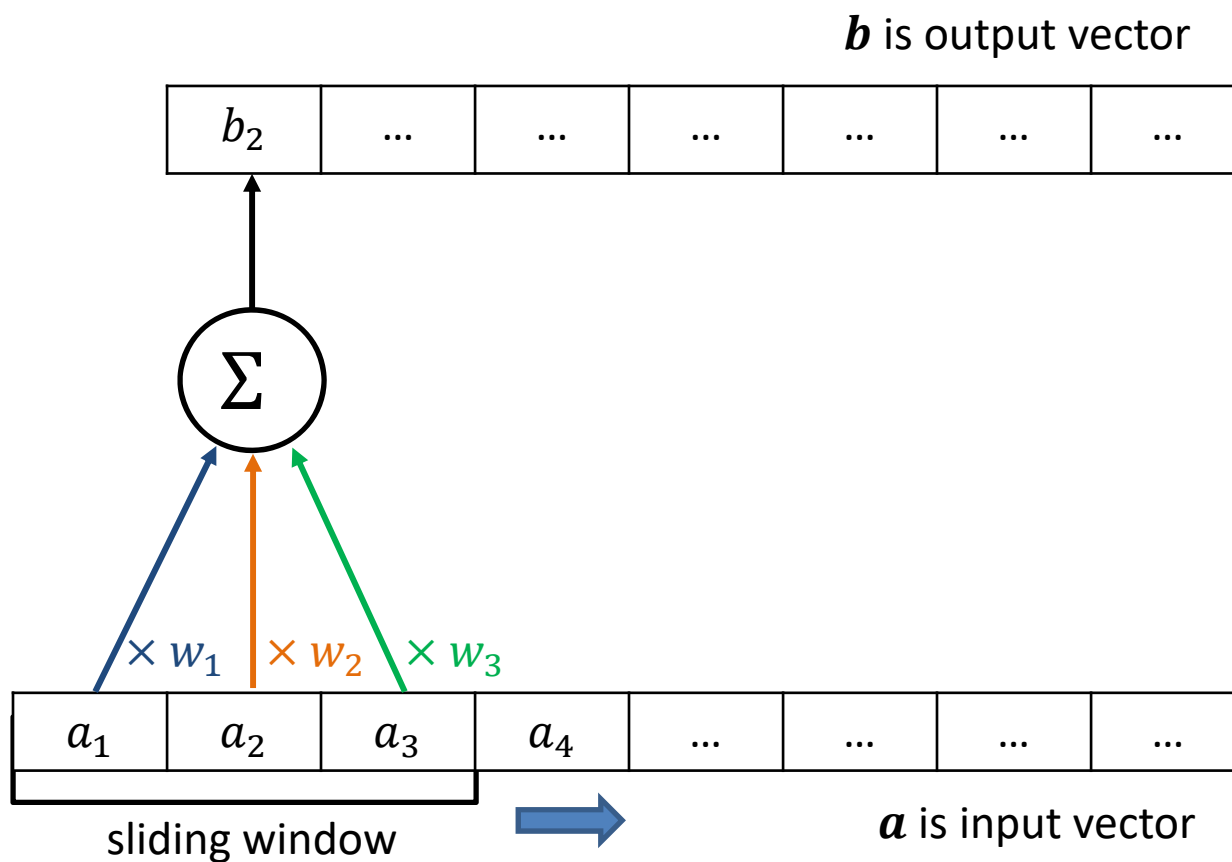
Convolutional
neural networks

art: OpenClipartVectors
at pixabay.com (CC0)

- Recurrent neural networks are not covered in this subject
- An autoencoder is an ANN trained in a specific way.
  * E.g., a multilayer perceptron can be trained as an autoencoder, or a recurrent neural network can be trained as an autoencoder.

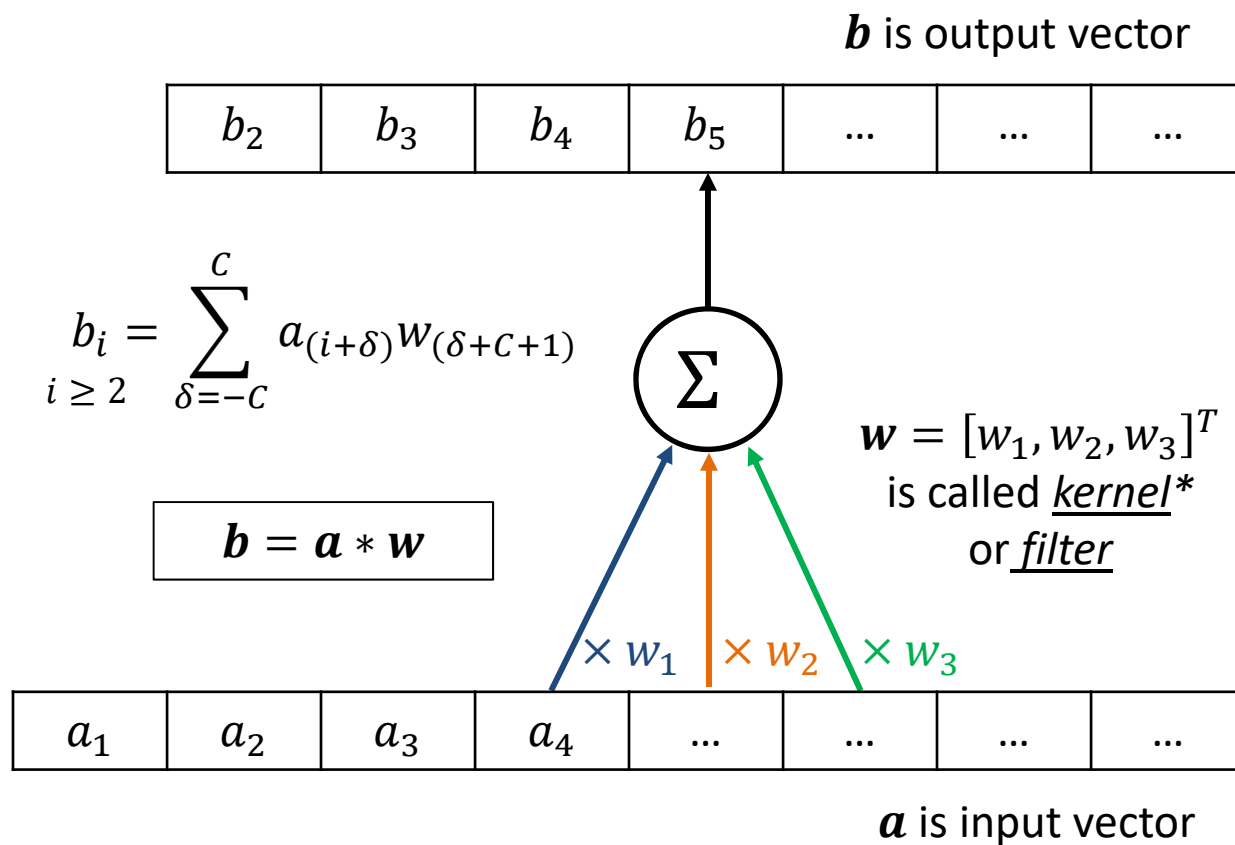12

# Convolutional Neural Networks (CNN)

Based on repeated application of small filters to patches of a 2D image or range of a 1D input
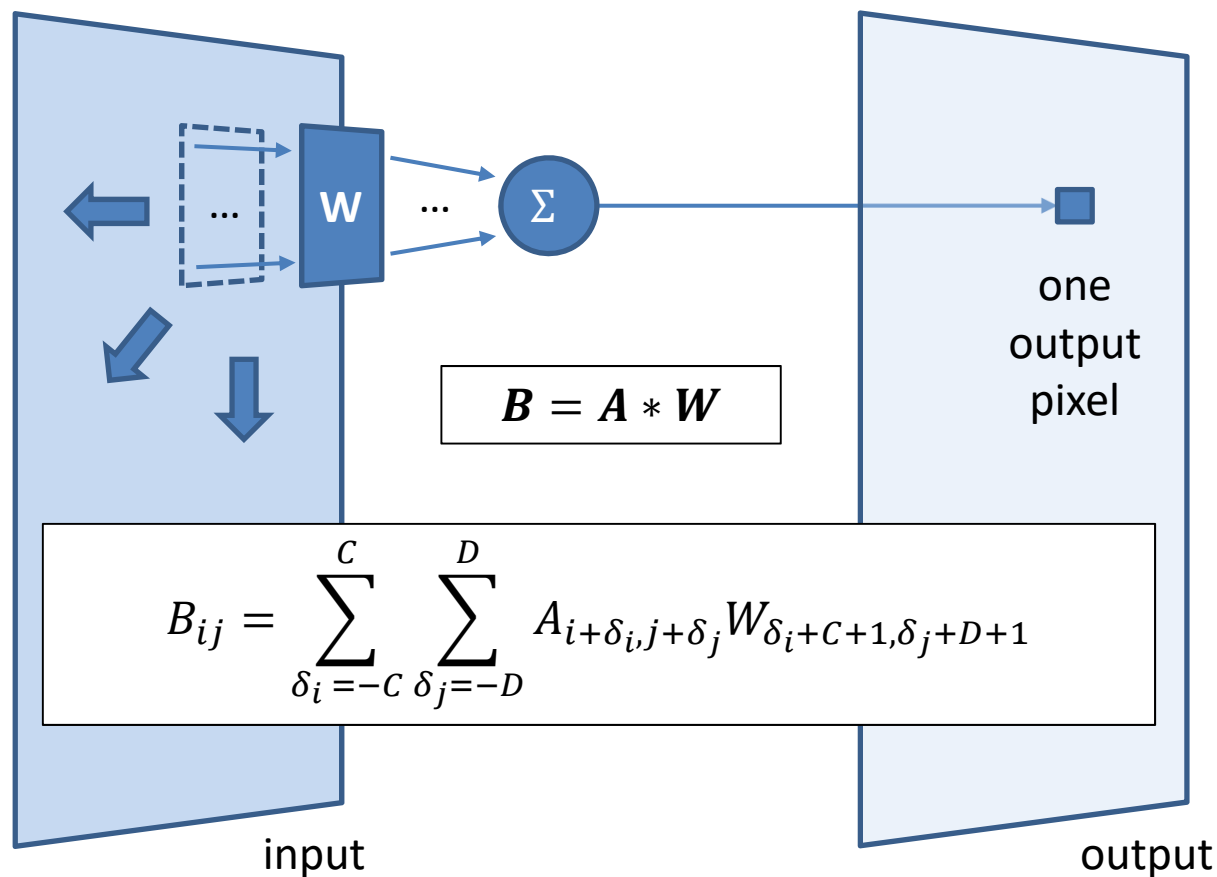
# Convolution

$\boldsymbol{b}$ is output vector



$\boldsymbol{a}$ is input vector

# Convolution

$b$ is output vector

| $b_2$ | $b_3$ | $b_4$ | $b_5$ | ... | ... | ... |
|-------|-------|-------|-------|-----|-----|-----|

$$b_i = \sum_{\delta=-C}^{C} a_{(i+\delta)} w_{(\delta+C+1)}$$
$$i \geq 2$$

$$\boxed{b = a * w}$$

$\Sigma$

$w = [w_1, w_2, w_3]^T$
is called *kernel**
or *filter*

$\times w_1$  $\times w_2$  $\times w_3$

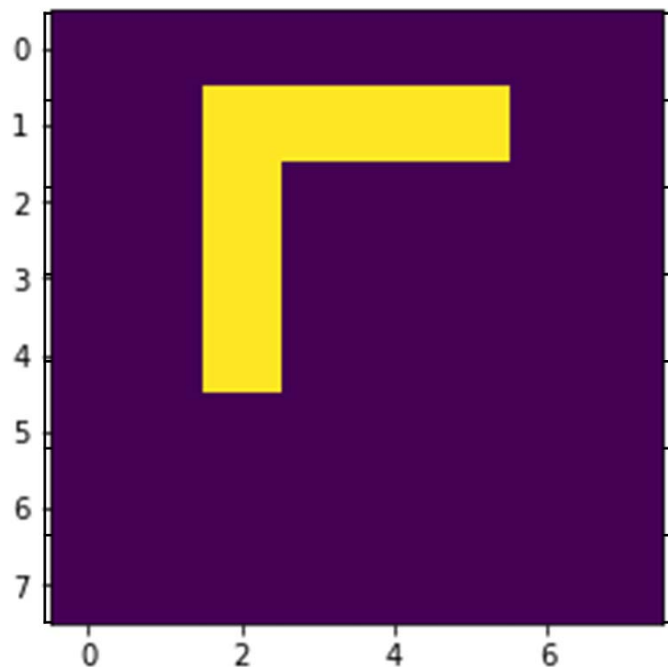| $a_1$ | $a_2$ | $a_3$ | $a_4$ | ... | ... | ... | ... |
|-------|-------|-------|-------|-----|-----|-----|-----|

$a$ is input vector

*Later in the subject, we will also use an unrelated definition of kernel as a function representing a dot product     15

# Convolution on 2D images



$$B = A * W$$

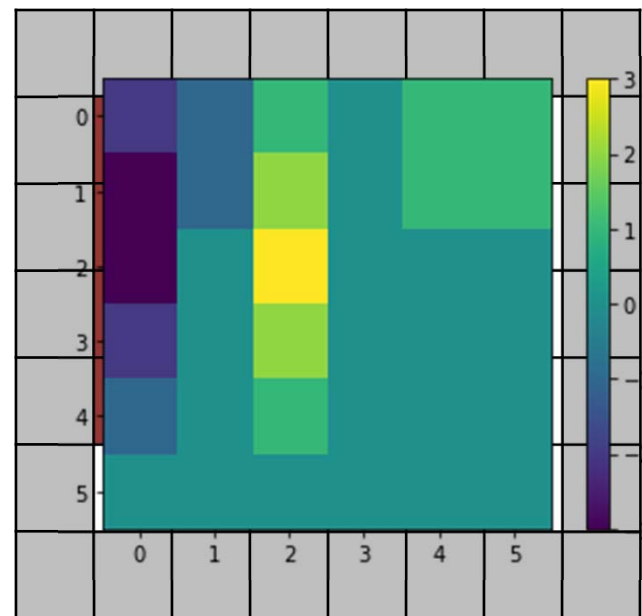$$B_{ij} = \sum_{\delta_i = -C}^{C} \sum_{\delta_j = -D}^{D} A_{i+\delta_i, j+\delta_j} W_{\delta_i + C + 1, \delta_j + D + 1}$$

one output pixel

input                              output

# Filters as feature detectors

convolve with a
vertical edge filter

| -1 | 0 | 1 |
|----|---|---|
| -1 | 0 | 1 |
| -1 | 0 | 1 |

activation
function

$A$ is input image

filtered
image

17

# Filters as feature detectors



$A$ is input image

convolve with a
vertical edge filter

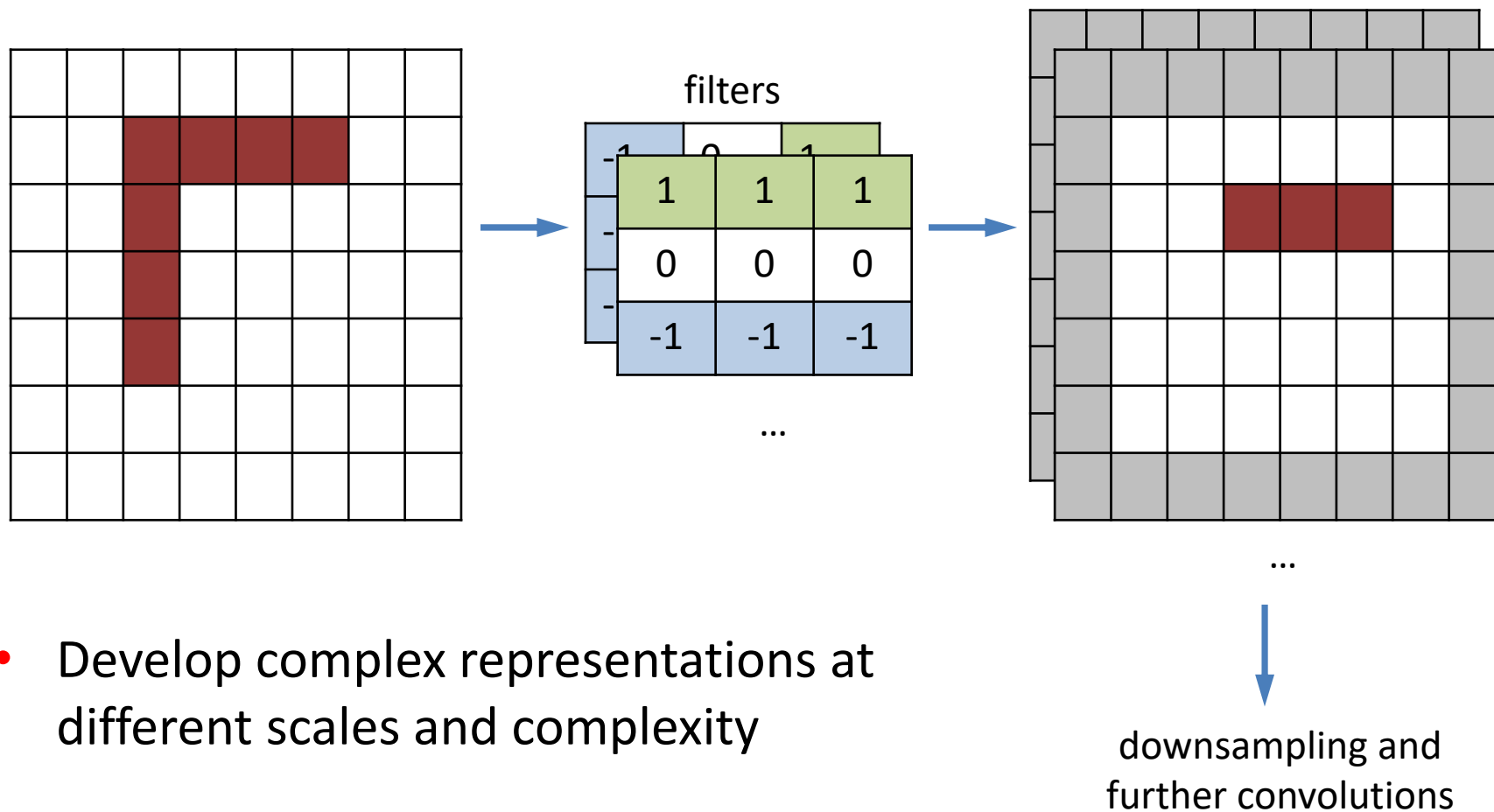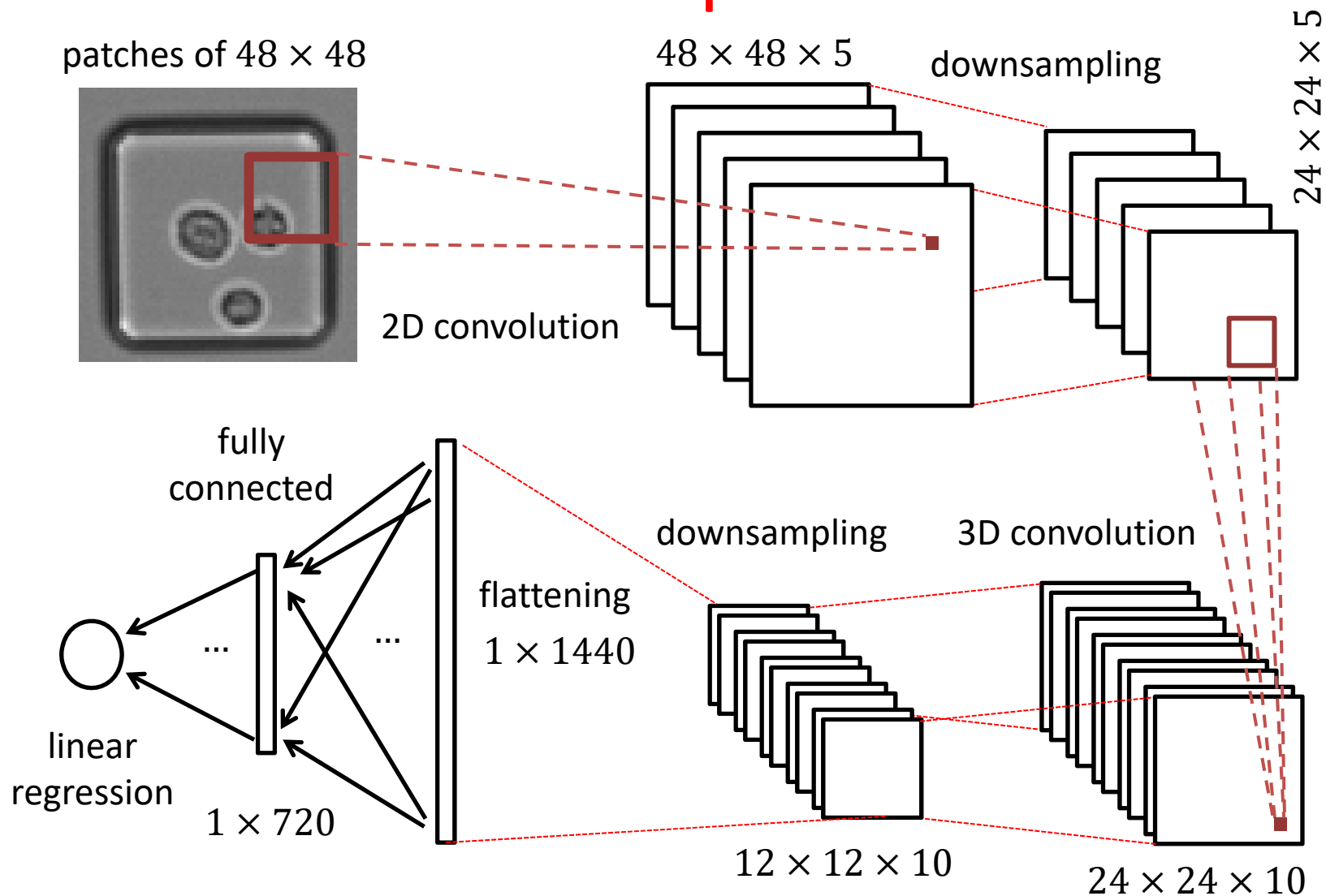| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

activation
function

filtered
image

# Stacking convolutions



filters

...

- Develop complex representations at different scales and complexity

- Filters are learned from training data!

downsampling and further convolutions

# CNN for computer vision

patches of $48 \times 48$

$48 \times 48 \times 5$

downsampling

$24 \times 24 \times 5$

2D convolution

fully connected

downsampling

3D convolution

flattening $1 \times 1440$

linear regression

$1 \times 720$

$12 \times 12 \times 10$

$24 \times 24 \times 10$

Implemented by Jizhizi Li
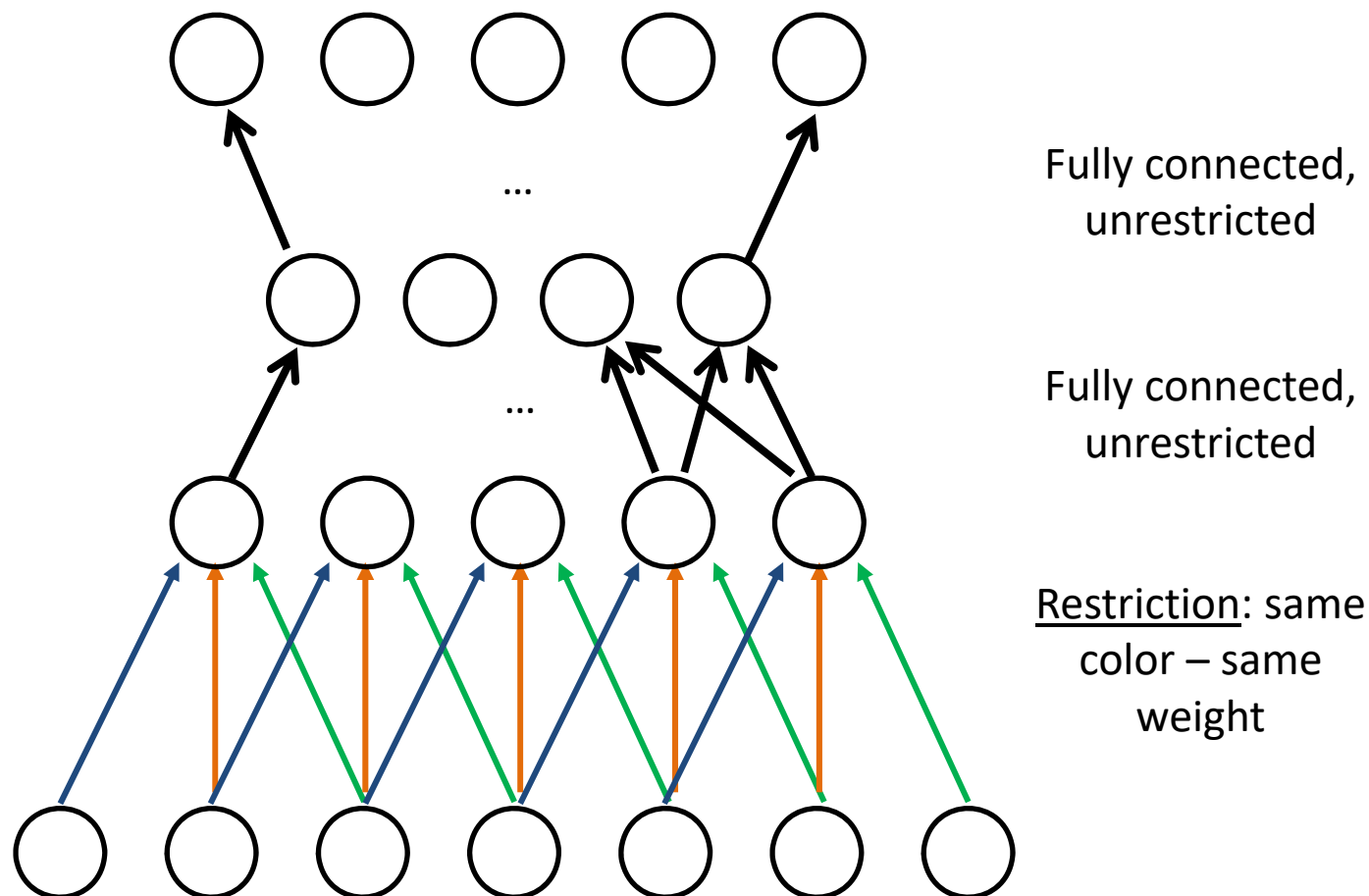based on LeNet5: http://deeplearning.net/tutorial/lenet.html

20

# Components of a CNN

- ## Convolutional layers
  - ∗ Complex input representations based on convolution operation
  - ∗ Filter weights are learned from training data

- ## Downsampling, usually via Max Pooling
  - ∗ Re-scales to smaller resolution, limits parameter explosion

- ## Fully connected parts and output layer
  - ∗ Merges representations together

# Downsampling via max pooling

- Special type of processing layer. For an $m \times m$ patch
$$v = \max(u_{11}, u_{12}, \ldots, u_{mm})$$

- Strictly speaking, not everywhere differentiable. Instead, gradient is defined according to "sub-gradient"

  * Tiny changes in values of $u_{ij}$ that is not max do not change $v$

  * If $u_{ij}$ is max value, tiny changes in that value change $v$ linearly

  * Use $\frac{\partial v}{\partial u_{ij}} = 1$ if $u_{ij} = v$, and $\frac{\partial v}{\partial u_{ij}} = 0$ otherwise

- Forward pass records maximising element, which is then used in the backward pass during back-propagation
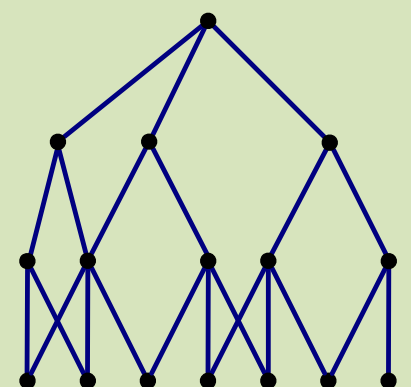
# Convolution as a regulariser



Fully connected,
unrestricted

Fully connected,
unrestricted

Restriction: same
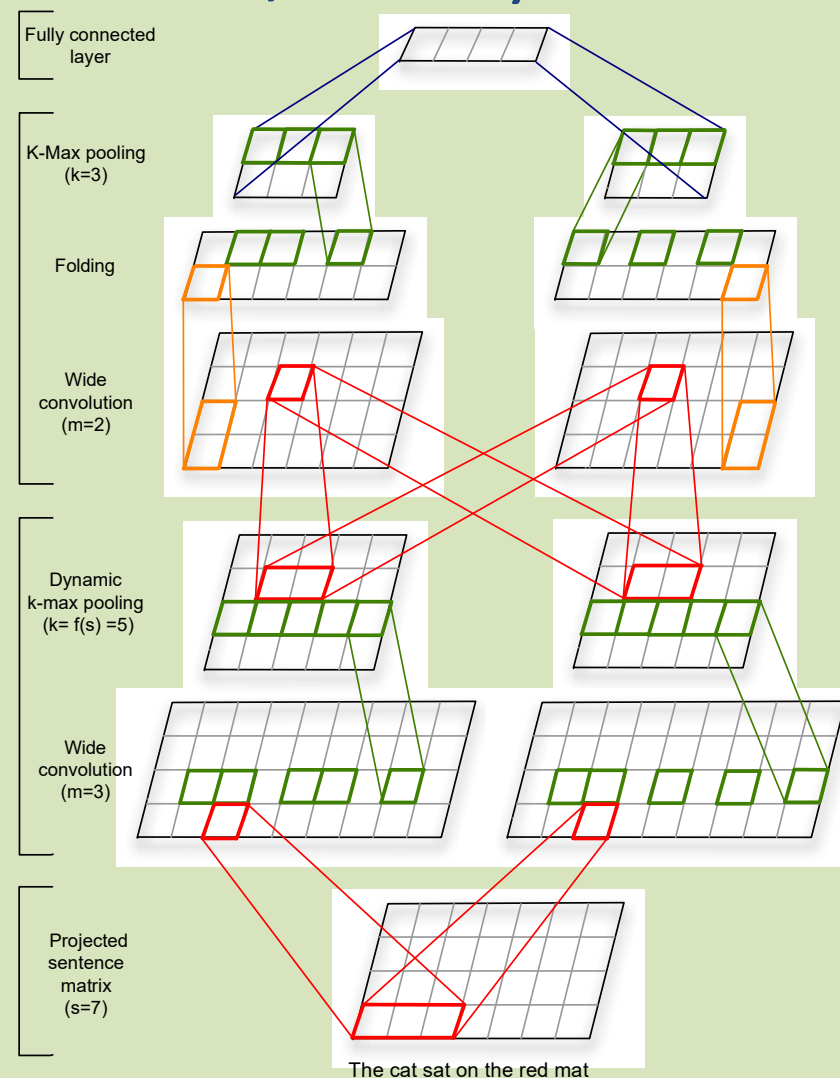color – same
weight

# Document classification (Kalchbrenner et al, 2014)

Structure of text important for classifying documents

Capture patterns of nearby words using 1d convolutions



The  cat  sat  on  the  red  mat

Fully connected layer

K-Max pooling (k=3)

Folding

Wide convolution (m=2)

Dynamic k-max pooling (k= f(s) =5)

Wide convolution (m=3)

Projected sentence matrix (s=7)
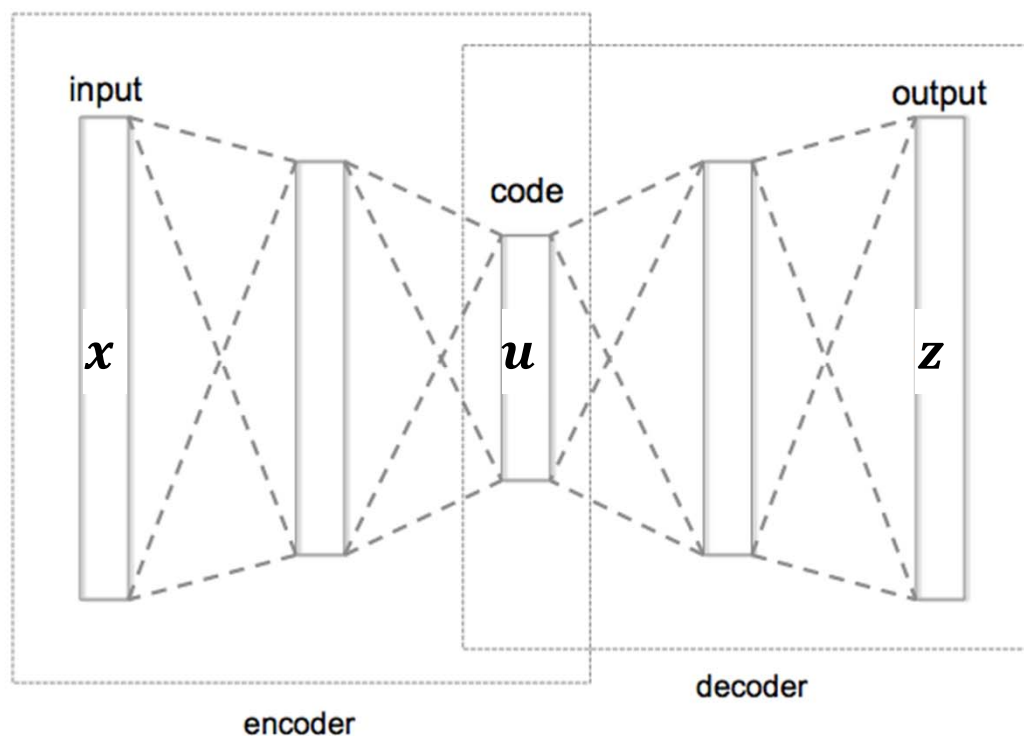
The cat sat on the red mat

24

# Autoencoder

An ANN training setup that can be used
for unsupervised learning, initialisation,
or just efficient coding

# Autoencoding idea

- Supervised learning:
  - Univariate regression: predict $y$ from $\boldsymbol{x}$
  - Multivariate regression: predict $\boldsymbol{y}$ from $\boldsymbol{x}$

- Unsupervised learning: explore data $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$
  - No response variable

- For each $\boldsymbol{x}_i$ set $\boldsymbol{y}_i \equiv \boldsymbol{x}_i$

- Train an ANN to predict $\boldsymbol{y}_i$ from $\boldsymbol{x}_i$

- Pointless?

# Autoencoder topology

- Given data without labels $x_1, \ldots, x_n$, set $y_i \equiv x_i$ and train an ANN to predict $z(x_i) \approx x_i$

- Set hidden layer $u$ in the middle "thinner" than input

# Introducing the bottleneck

- Suppose you managed to train a network that gives a good restoration of the original signal $z(x_i) \approx x_i$

- This means that the data structure can be effectively described (encoded) by a lower dimensional representation $u$



adapted from: Chervinskii at
Wikimedia Commons (CC4)

28

# Dimensionality reduction

- Autoencoders can be used for compression and dimensionality reduction via a non-linear transformation

- If you use linear activation functions and only one hidden layer, then the setup becomes almost that of Principal Component Analysis (stay tuned!)
  - ∗ ANN might find a different solution, doesn't use eigenvalues

# Tools

- Tensorflow, Theano, Torch
  - ∗ python / lua toolkits for deep learning
  - ∗ symbolic or automatic differentiation
  - ∗ GPU support for fast compilation
  - ∗ Theano tutorials at http://deeplearning.net/tutorial/

- Various others
  - ∗ Caffe
  - ∗ CNTK
  - ∗ deeplearning4j …

- Keras: high-level Python API. Can run on top of TensorFlow, CNTK, or Theano

# This lecture

- Deep learning
  * Representation capacity
  * Deep models and representation learning

- Convolutional Neural Networks
  * Convolution operator
  * Elements of a convolution-based network

- Autoencoders
  * Learning efficient coding

- Workshops Week #5: Neural nets

- Next lectures: Kernel methods