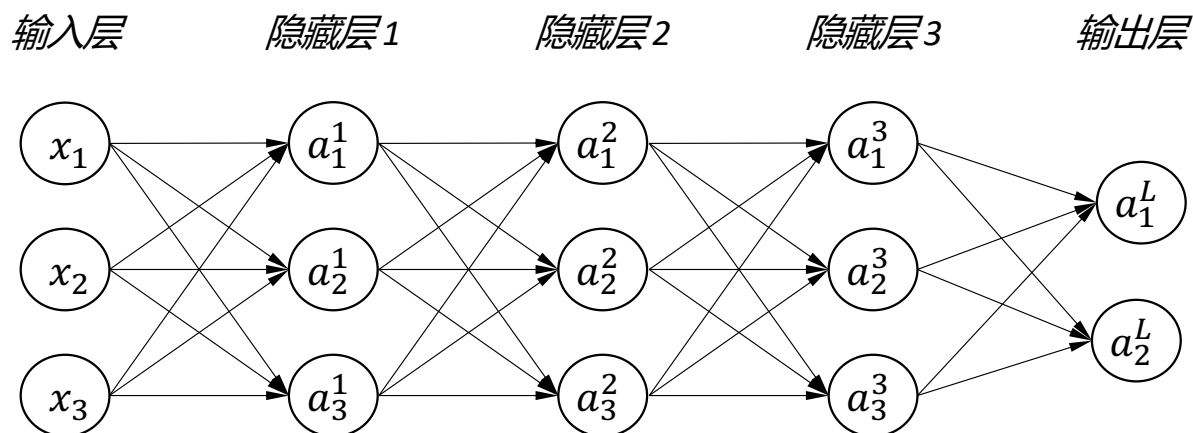# 复习：深度学习部分

主讲人：董豪　讲义：董豪

- 单个神经元 Single Neuron

- 激活函数 Activation Functions

- 多层感知器 Multi-layer Perceptron

- 损失函数 Loss Functions

- 优化 Optimisation
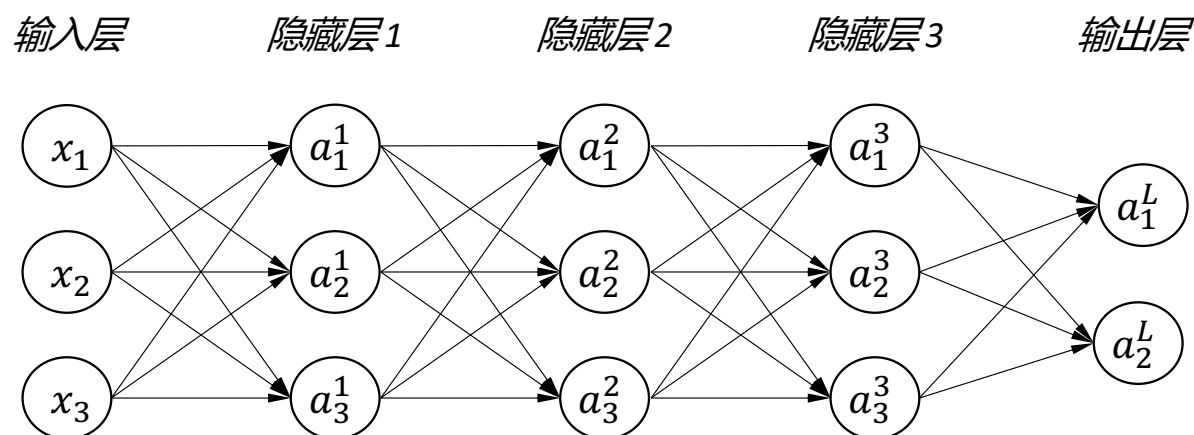
- 正则化 Regularisation

- 实现 Implementation

- 误差反向传播（Error Back-Propagation）

误差反向传播是用来计算网络中所有参数的梯度 $\frac{\partial \mathcal{L}}{\partial \theta}$ 的方法。计算梯度时，引入对 $\mathcal{L}$ 求输出值 $z$ 的偏导 $\delta = \frac{\partial \mathcal{L}}{\partial z}$，作为中间结果（intermediate result)，基于这个中间结果来计算出每一层的梯度 $\frac{\partial \mathcal{L}}{\partial \theta}$。

- 误差反向传播（Error Back-Propagation）

输入层　　　　隐藏层1　　　　隐藏层2　　　　隐藏层3　　　　输出层



层索引（layer index）$l = 1 \dots L$ 代表第一层隐藏层（1）到输出层（L）

输入层用 $\boldsymbol{x} = \boldsymbol{a}^0$ 来表示

$$\boldsymbol{a}^l = f(\boldsymbol{z}^l) = \frac{1}{1 + e^{-\boldsymbol{z}^l}}$$

我们通过这个简单模型和损失函数来讲解

$$\boldsymbol{z}^l = \boldsymbol{W}^{l^T} \boldsymbol{a}^{l-1} + \boldsymbol{b}^l$$

$$\mathcal{L} = \frac{1}{2}(\boldsymbol{y} - \boldsymbol{a}^L)^2$$

- 误差反向传播（Error Back-Propagation）：列格式（column format）教材常用

**1. 已知**

- $a^l = f(z^l) = \frac{1}{1+e^{-z^l}}$
- $z^l = W^{l^T} a^{l-1} + b^l$
- $\mathcal{L} = \frac{1}{2}(y - a^L)^2$

**2. 则有如下求导**

- $\frac{\partial a^l}{\partial z^l} = f'(z^l) = a^l \circ (1 - a^l)$
- $\frac{\partial \mathcal{L}}{\partial a^L} = (a^L - y)$

  逐点相乘
  Hadamard (element-wise) product

- $\frac{\partial z^l}{\partial W^l} = a^{l-1}$ and $\frac{\partial z^l}{\partial b^l} = 1$

**3. 输出层的中间结果 $l = L$**

- $\delta^L = \frac{\partial \mathcal{L}}{\partial z^L} = \frac{\partial \mathcal{L}}{\partial a^L} \frac{\partial a^L}{\partial z^L} = (a^L - y) \circ (a^L \circ (1 - a^L))$

  链式法则（chain rule）

**4. 其他层的中间结果 $l = 1 \dots L-1$**

- $\delta^l = \frac{\partial \mathcal{L}}{\partial z^l} = \frac{\partial \mathcal{L}}{\partial z^{l+1}} \frac{\partial z^{l+1}}{\partial z^l} = \delta^{l+1} \frac{\partial z^{l+1}}{\partial z^l}$

  - $z^{l+1} = W^{l+1^T} a^l + b^{l+1}$
  - $\frac{\partial z^{l+1}}{\partial z^l} = W^{l+1^T} f'(z^l) = W^{l+1^T} \circ (a^l \circ (1 - a^l))$

- $\delta^l = \frac{\partial \mathcal{L}}{\partial z^l} = \frac{\partial \mathcal{L}}{\partial z^{l+1}} \frac{\partial z^{l+1}}{\partial z^l} = W^{l+1^T} \delta^{l+1} \circ (a^l \circ (1 - a^l))$

**5. 则有梯度**

- $\frac{\partial \mathcal{L}}{\partial W^l} = \frac{\partial \mathcal{L}}{\partial z^l} \frac{\partial z^l}{\partial W^l} = \delta^l \frac{\partial z^l}{\partial W^l} = \delta^l a^{l-1^T}$

- $\frac{\partial \mathcal{L}}{\partial b^l} = \frac{\partial \mathcal{L}}{\partial z^l} \frac{\partial z^l}{\partial b^l} = \delta^l \frac{\partial z^l}{\partial b^l} = \delta^l$

**6. 更新参数**

$$W^l := W^l - \alpha \frac{\partial \mathcal{L}}{\partial W^l} \qquad b^l := b^l - \alpha \frac{\partial \mathcal{L}}{\partial b^l}$$

- 误差反向传播（Error Back-Propagation）：<span style="color:red">行格式（row format）编程常用</span>

**1. 已知**

- $a^l = f(z^l) = \frac{1}{1+e^{-z^l}}$
- $z^l = a^{l-1}W^l + b^l$
- $\mathcal{L} = \frac{1}{2}(y - a^L)^2$

**2. 则有如下求导**

- $\frac{\partial a^l}{\partial z^l} = f'(z^l) = a^l \circ (1 - a^l)$
- $\frac{\partial \mathcal{L}}{\partial a^L} = (a^L - y)$
- $\frac{\partial z^l}{\partial W^l} = a^{l-1}$ and $\frac{\partial z^l}{\partial b^l} = 1$

**3. 输出层的中间结果 $l = L$**

- $\delta^L = \frac{\partial \mathcal{L}}{\partial z^L} = \frac{\partial \mathcal{L}}{\partial a^L}\frac{\partial a^L}{\partial z^L} = (a^L - y) \circ (a^L \circ (1 - a^L))$

**4. 其他层的中间结果 $l = 1 \dots L-1$**

- $\delta^l = \frac{\partial \mathcal{L}}{\partial z^l} = \frac{\partial \mathcal{L}}{\partial z^{l+1}}\frac{\partial z^{l+1}}{\partial z^l} = \delta^{l+1}\frac{\partial z^{l+1}}{\partial z^l}$

  - $z^{l+1} = a^l W^{l+1} + b^{l+1}$
  - $\frac{\partial z^{l+1}}{\partial z^l} = W^{l+1} \circ f'(z^l) = W^{l+1} \circ (a^l \circ (1 - a^l))$
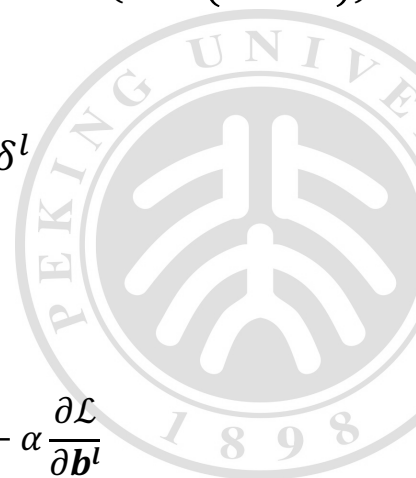
- $\delta^l = \frac{\partial \mathcal{L}}{\partial z^l} = \frac{\partial \mathcal{L}}{\partial z^{l+1}}\frac{\partial z^{l+1}}{\partial z^l} = \delta^{l+1}W^{{l+1}^T} \circ (a^l \circ (1 - a^l))$

**5. 则有梯度**

- $\frac{\partial \mathcal{L}}{\partial W^l} = \frac{\partial \mathcal{L}}{\partial z^l}\frac{\partial z^l}{\partial W^l} = \delta^l\frac{\partial z^l}{\partial W^l} = a^{{l-1}^T}\delta^l$

- $\frac{\partial \mathcal{L}}{\partial b^l} = \frac{\partial \mathcal{L}}{\partial z^l}\frac{\partial z^l}{\partial b^l} = \delta^l\frac{\partial z^l}{\partial b^l} = \delta^l$

**6. 更新参数**

$$W^l := W^l - \alpha\frac{\partial \mathcal{L}}{\partial W^l} \qquad b^l := b^l - \alpha\frac{\partial \mathcal{L}}{\partial b^l}$$

- 梯度消失（Gradient Vanish）问题

$$\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \mathbf{W}^{l+1}{}^T \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$$

刚刚的例子中，中间结果 $\delta$ 中有一项 $\left(\mathbf{a}^l \circ \left(1 - \mathbf{a}^l\right)\right)$，当激活输出 $\mathbf{a}$ 接近0或者1时，中间结果 $\delta$ 会变得很小。由于 $\delta^l$ 又跟 $\delta^{l+1}$ 有关，当反向传播时 $\delta$ 比较小的话，会使得 $\delta$ 越传播越小，使得靠近输入层的参数无法被更新，影响网络训练。

- 解决方法 1: 用 ReLU 代替 Sigmoid 函数（常见的方法）
- 解决方法 2: 逐层训练法（已经很少使用了）
- ...

- 卷积算法 Convolutional Algorithm
- 池化算法 Pooling Algorithm
- 分层表示学习Hierarchical Representation Learning
- 卷积神经网络结构 Convolutional Architectures
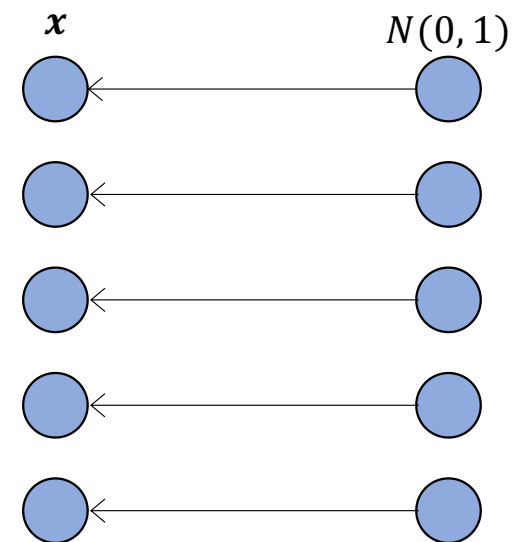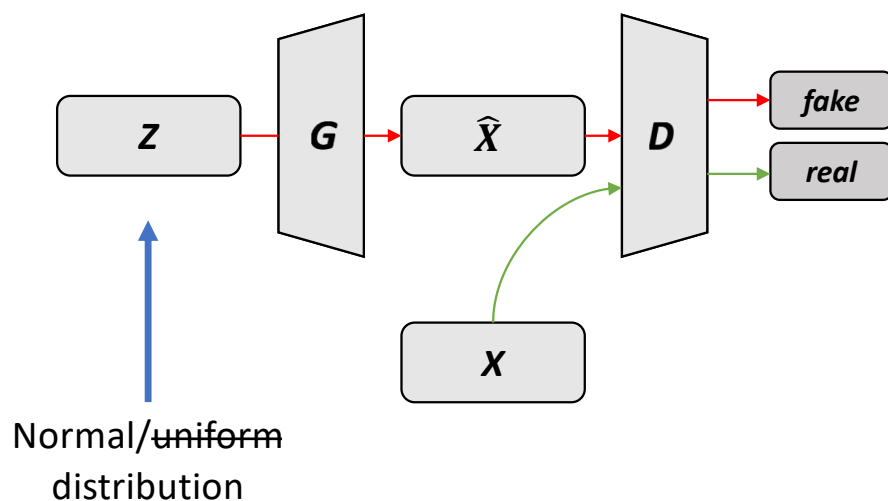- 转置卷积（反卷积）Transposed Convolutional Algorithm

- 引入：生成式模型 Generative Models
- 朴素GAN
- 有条件GAN
- 对抗损失函数 VS 均方差 Adversarial Loss vs. MSE
- GAN面临的挑战 Challenges of GAN

- 朴素GAN



Normal/~~uniform~~ distribution

$x$ $\quad$ $N(0,1)$

Unidirectional Mapping

**GAN: map a distribution to another distribution**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_z}[\log(1 - D(G(\boldsymbol{z})))]$$

$$\mathcal{L}_D = -\mathbb{E}_{\boldsymbol{x} \sim p_{data}}[\log D(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{z} \sim p_z}[\log(1 - D(G(\boldsymbol{z})))]$$

$$\mathcal{L}_G = -\mathbb{E}_{\boldsymbol{z} \sim p_z}[\log D(G(\boldsymbol{z}))]$$

对抗

- 一个简单的例子：辅助分类器生成对抗网络



monarch butterfly     goldfinch     daisy

$$\mathcal{L}_D = -\mathbb{E}_{\boldsymbol{x}\sim p_{data}}[\log D_{\mathrm{x}}(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{z}\sim p_z, c\sim p_c}[\log(1 - D_{\mathrm{x}}(G(\boldsymbol{z},\boldsymbol{c})))]$$

$$- \mathbb{E}_{\boldsymbol{x}\sim p_{data}}[\log D_c(\boldsymbol{x})] - \mathbb{E}_{\boldsymbol{z}\sim p_z, c\sim p_c}[\log(1 - Dc(G(\boldsymbol{z},\boldsymbol{c})))]$$

$$\mathcal{L}_G = -\mathbb{E}_{\boldsymbol{z}\sim p_z, c\sim p_c}[\log D_{\mathrm{x}}(G(\boldsymbol{z},\boldsymbol{c}))] - \mathbb{E}_{\boldsymbol{z}\sim p_z, c\sim p_c}[\log D_c(G(\boldsymbol{z},\boldsymbol{c}))]$$
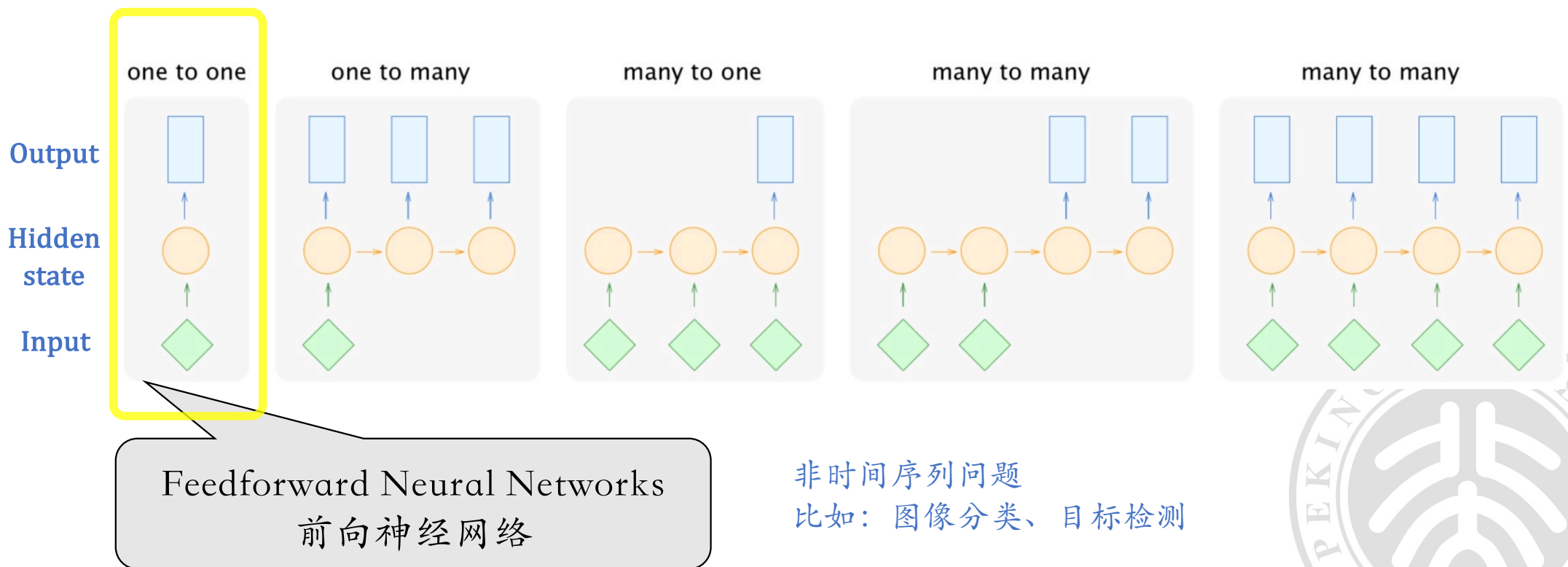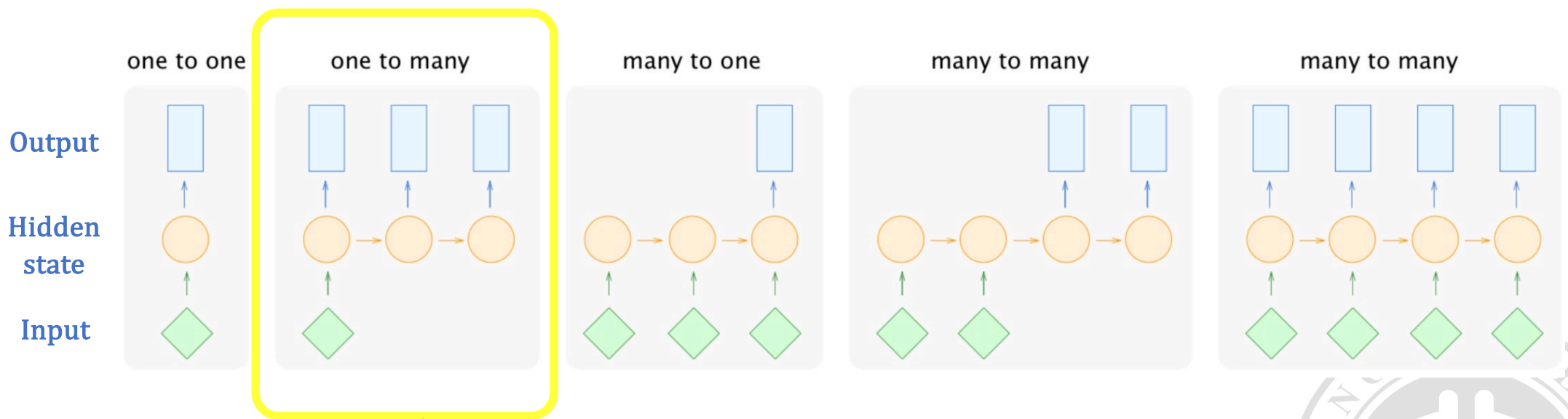
- 词的表示 Word Representation
- 序列数据 Sequential Data
- 朴素循环神经网络 Vanilla Recurrent Neural Network
- 长短期记忆网络 LSTM Long Short-Term Memory
- 序列生成模型 RNNs are Generative Models
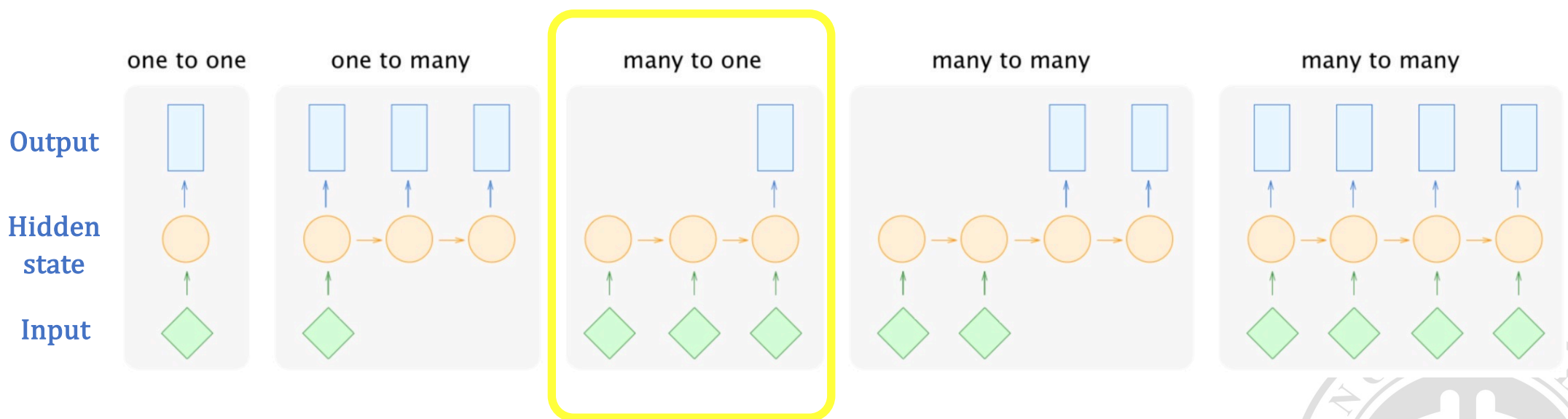- 时间序列应用 Time-series Applications

（以及后续的Transformer）

**Output**

**Hidden state**

**Input**

one to one | one to many | many to one | many to many | many to many

Feedforward Neural Networks
前向神经网络

非时间序列问题
比如：图像分类、目标检测

Output

Hidden state

Input

one to one　one to many　many to one　many to many　many to many

输入一个数据，输出多个数据

图片描述：输入一张图片，生成一句话的描述

Output

Hidden state

Input

one to one    one to many    many to one    many to many    many to many

输入多个数据，输出一个数据

情感分类任务：输入一个有序的句子，输出表示幸福概率的数值。

异步的
(Seq2Seq)

Hidden
state

Input

one to one

one to many

many to one

many to many

many to many

多数据输入和多数据输出

语言翻译：在开始生成翻译句子之前，将整个句子输入到
模型中。

# 序列数据

异步的
(Seq2Seq)

同步的
(simultaneous)



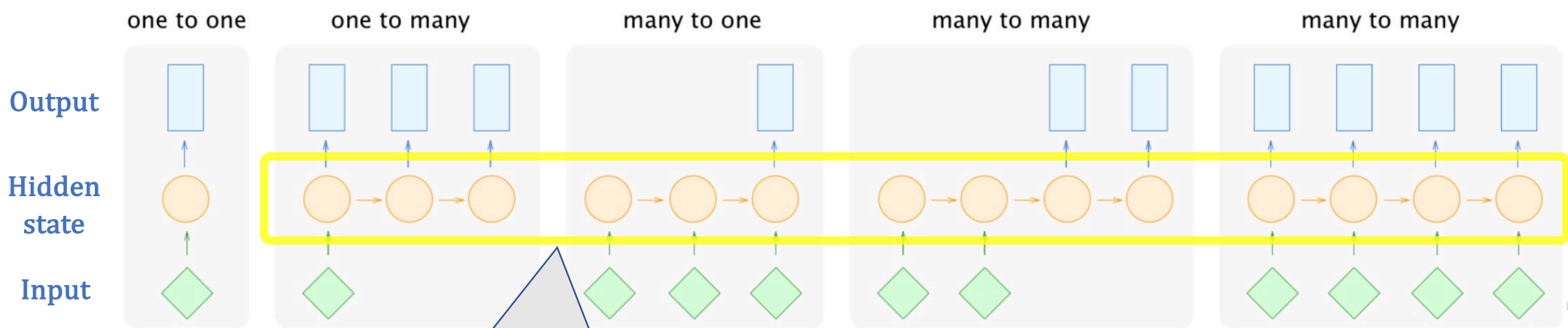| one to one | one to many | many to one | many to many | many to many |
|---|---|---|---|---|

Output

Hidden state

Input

多个数据输入和多个数据输出

天气预测：在每个时间步（time-step）输入信息到模型中，并输出预测的天气状况。

**循环神经网络（Recurrent Neural Nets）** 存储和处理时序信息。

考试加油