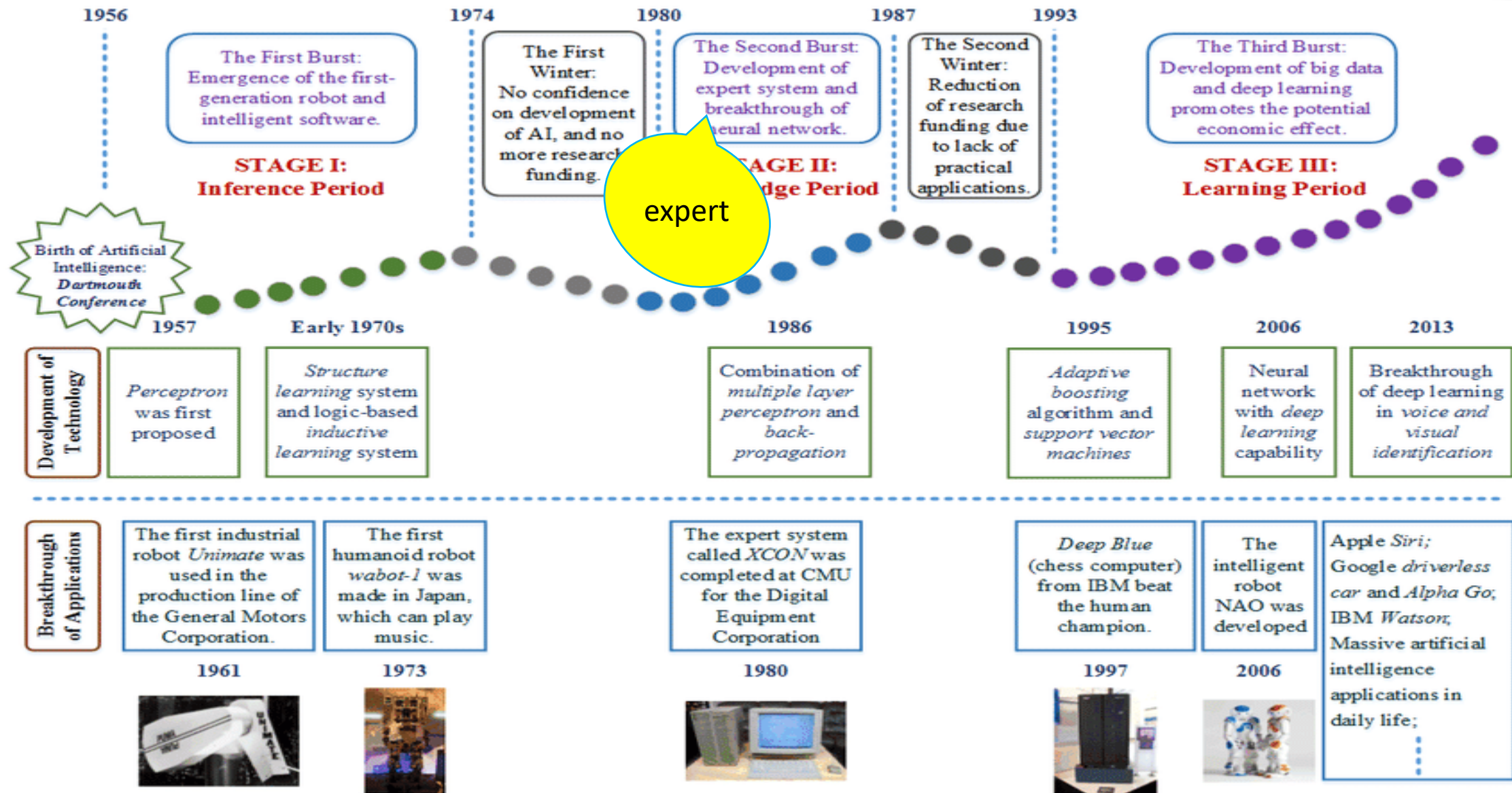


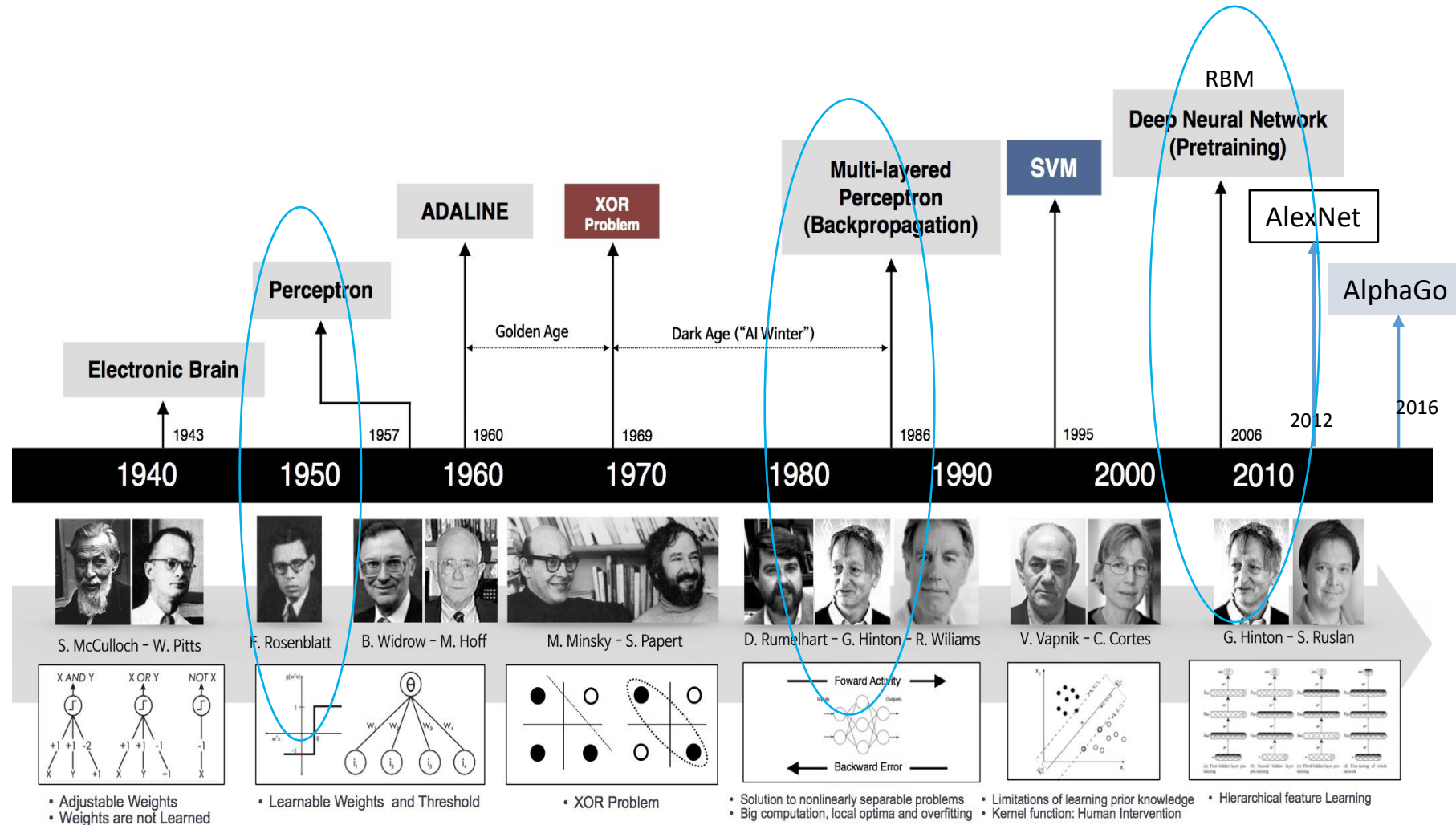


Introduction

Pau-Choo Chung (Julia)

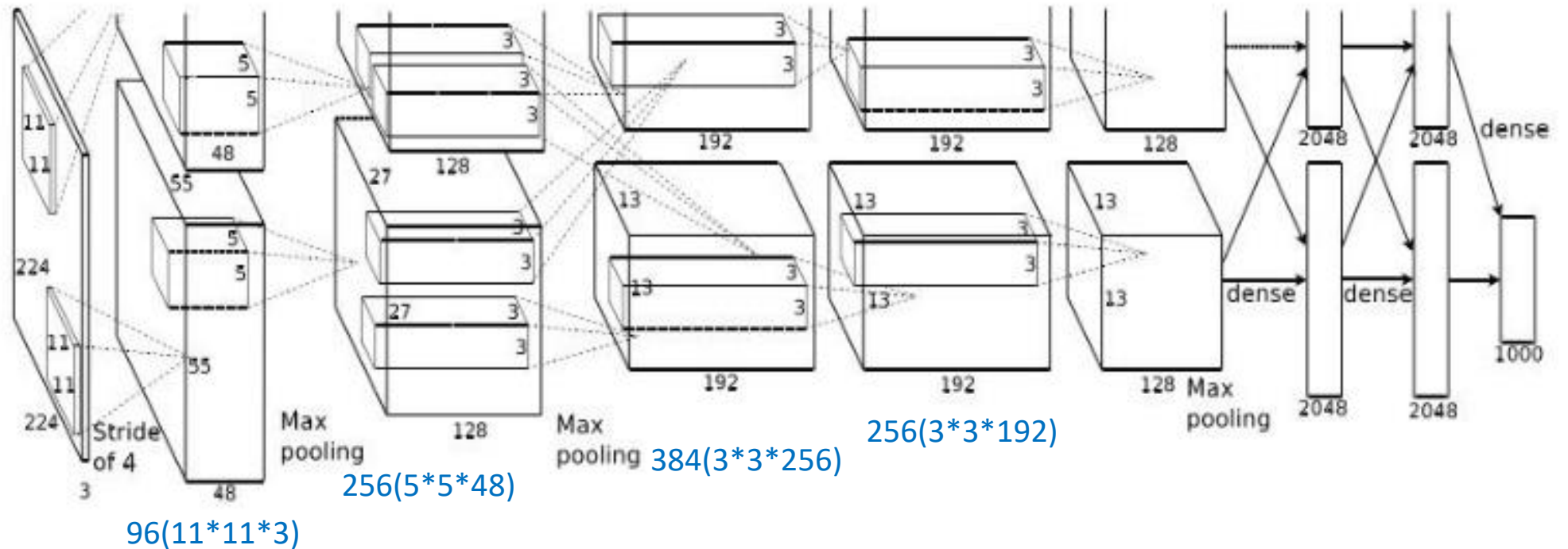






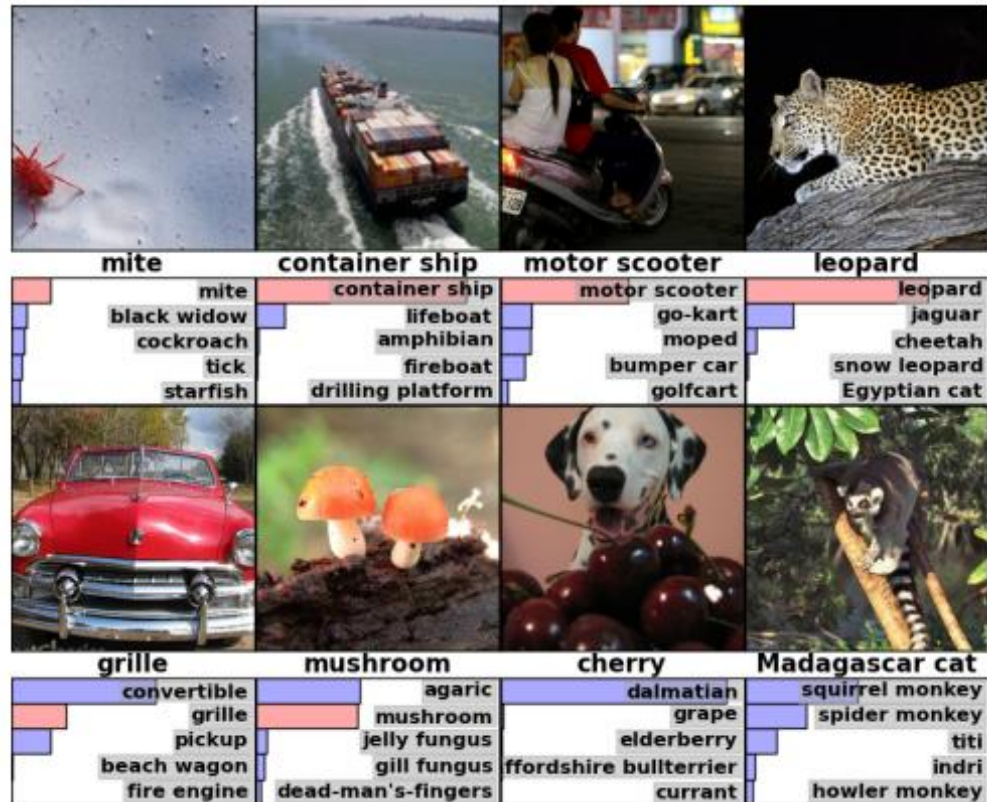
AlexNet

- 5 convolutional layers and 2 fully connected layers
- Max pooling layer after first, second and fifth layers
- Use ReLU as activation function





- 1.3 million high-resolution images in the LSVRC-2010(ImageNet Large Scale Visual Recognition Competition)
- GPU implementation of convolutional nets



Model	Top-1	Top-5
<i>Sparse coding [2]</i>	47.1%	28.2%
<i>SIFT + FVs [24]</i>	45.7%	25.7%
CNN	37.5%	17.0%

Table 1: Comparison of results on ILSVRC-2010 test set. In *italics* are best results achieved by others.

Source: [ImageNet Classification with Deep Convolutional Neural Networks](#) (NIPS 2012)



VGG

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

11, 13, 16 and 19 layers implementation
3*3 and 1*1 kernel sizes for deep model

Using smaller masks
The first using 1*1 kernel size

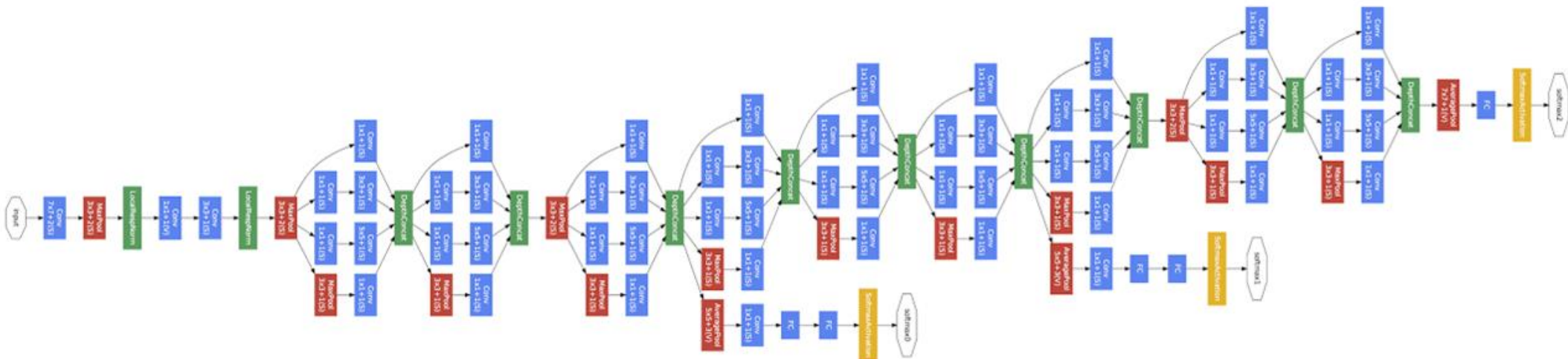
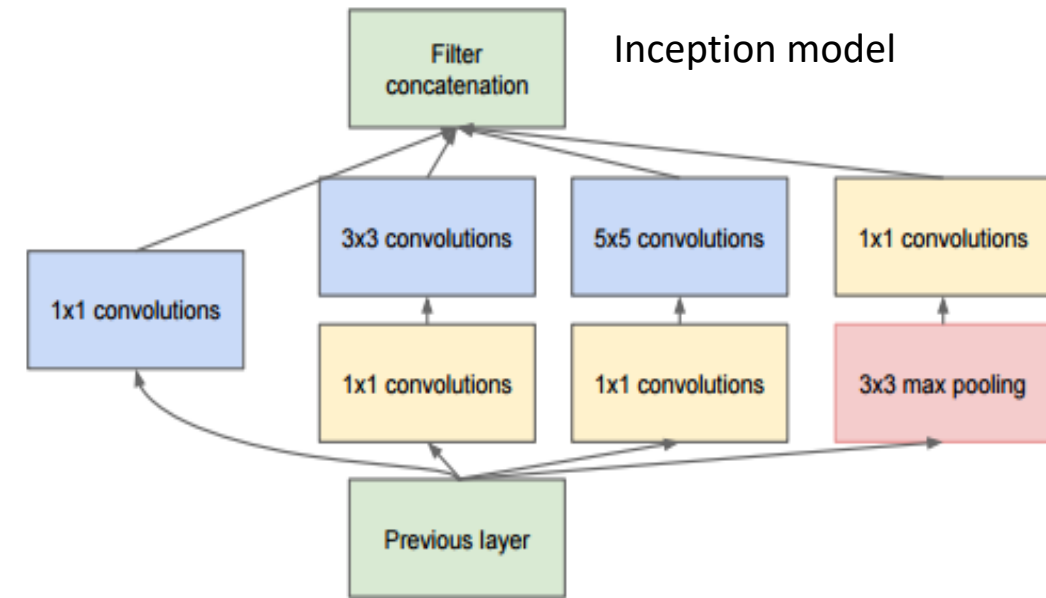


Simplifying networks to afford deeper



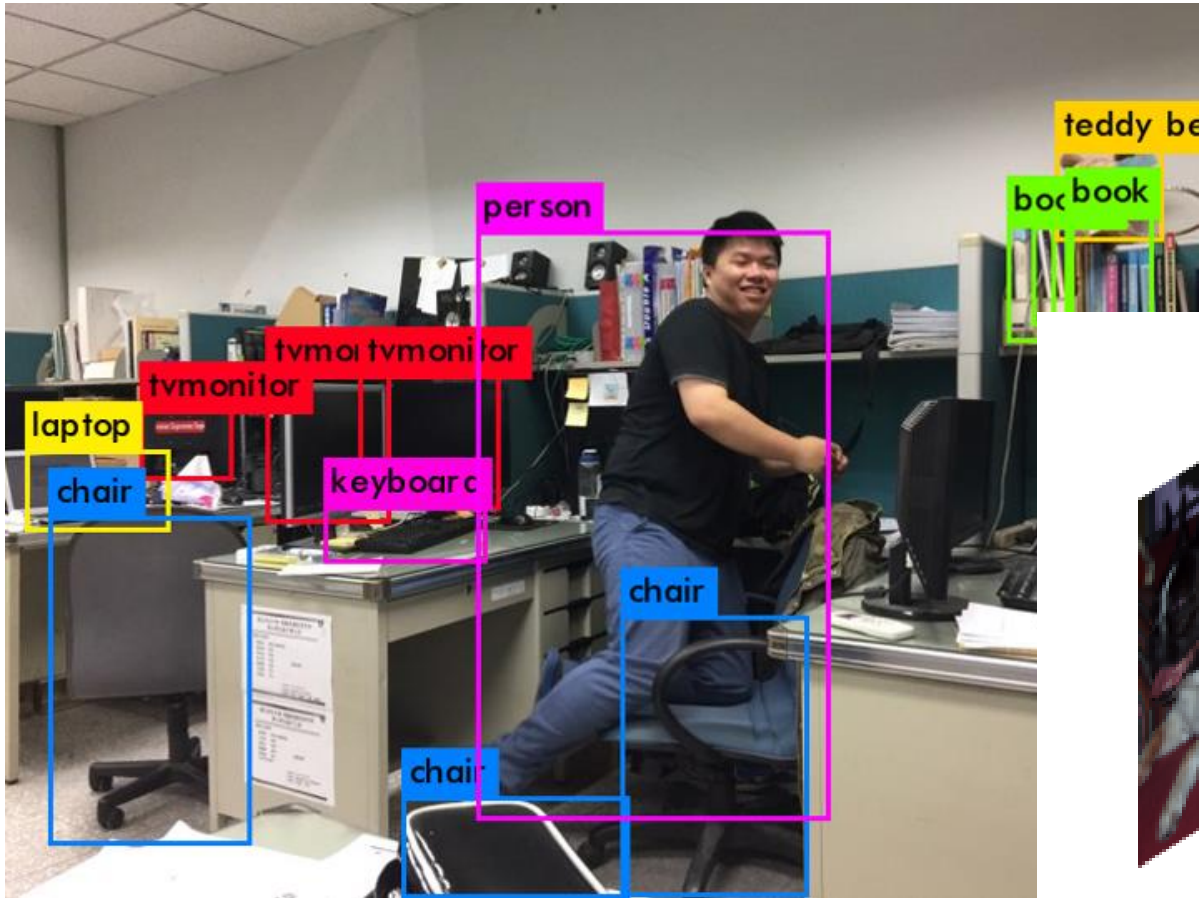
GoogleNet

- Inception
- compute various convolution kernel sizes in same network
- 1×1 convolutions for reducing the number of maps





Object detection (Yolo V2)



Segmentation

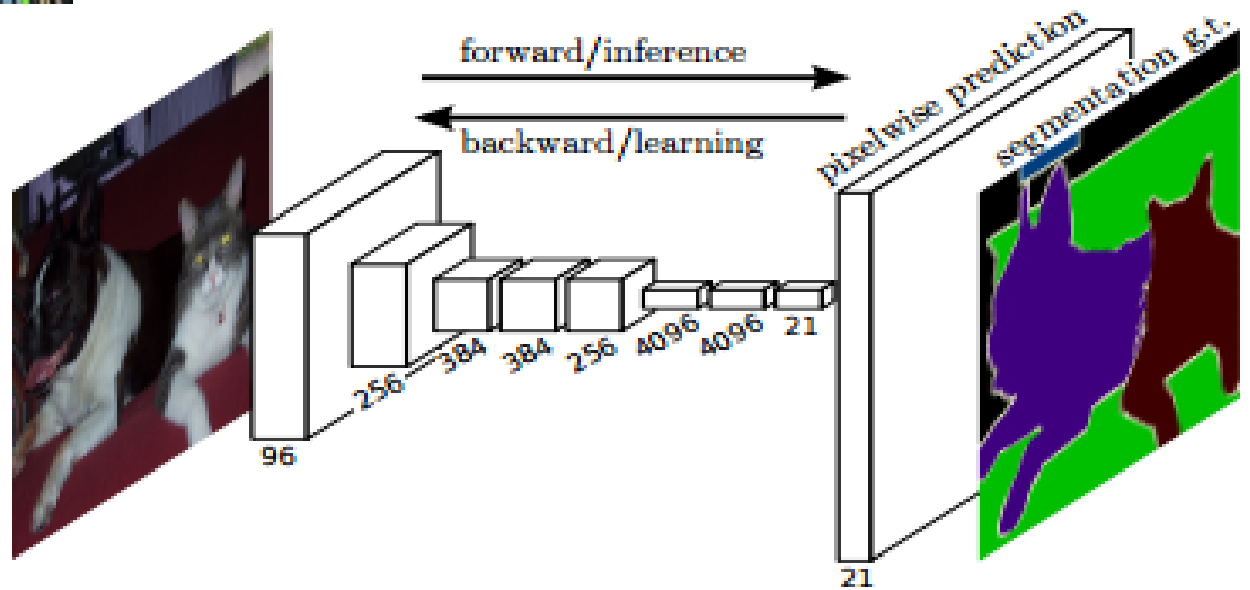
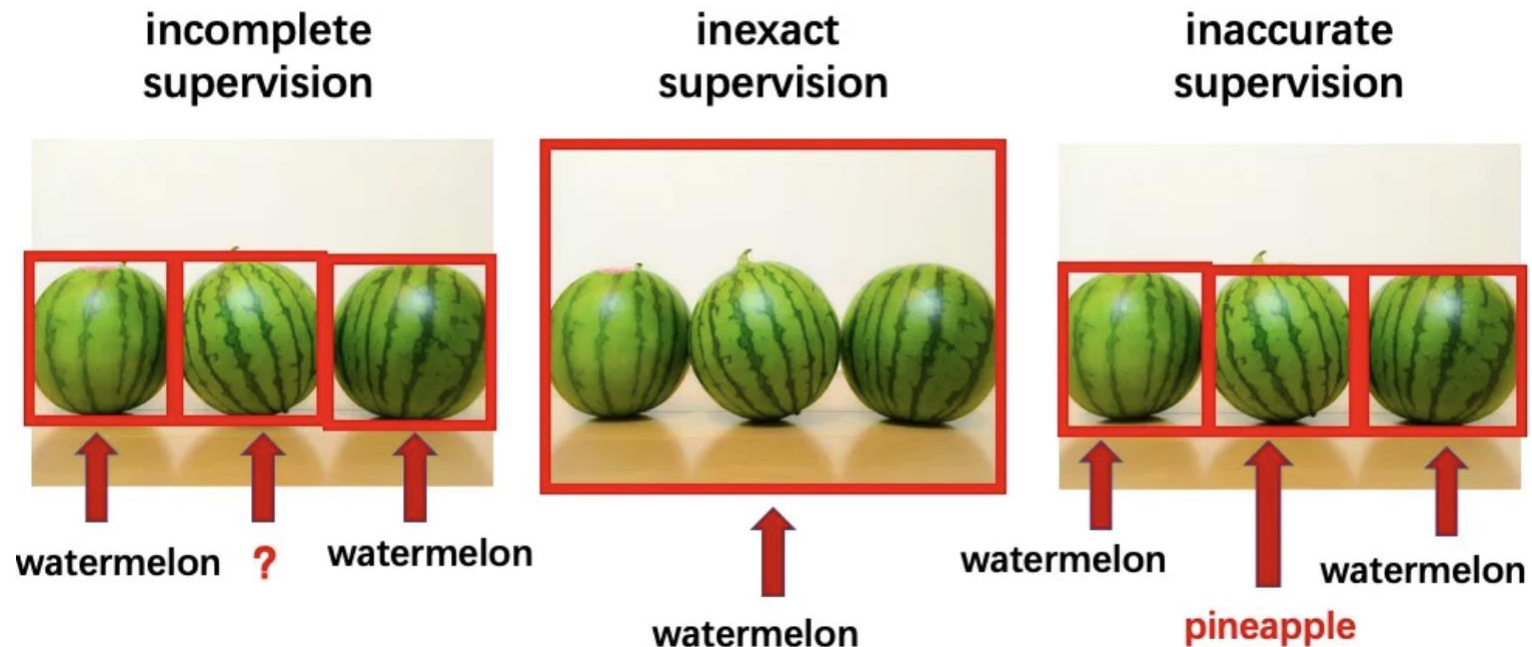


Figure 1. Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.



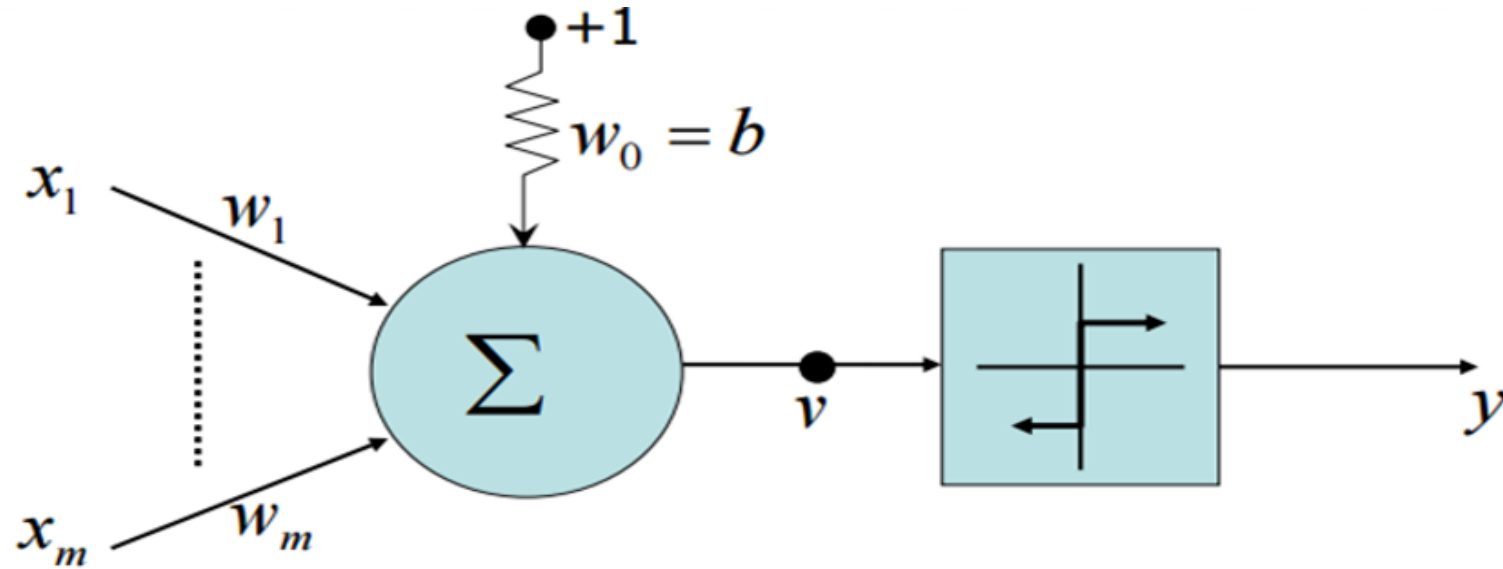
Learning approaches

- Supervised : need labeled ground truth
- Unsupervised: no labels
- Weakly supervised:
 - Incomplete
 - Inexact
 - inaccurate





Perceptron



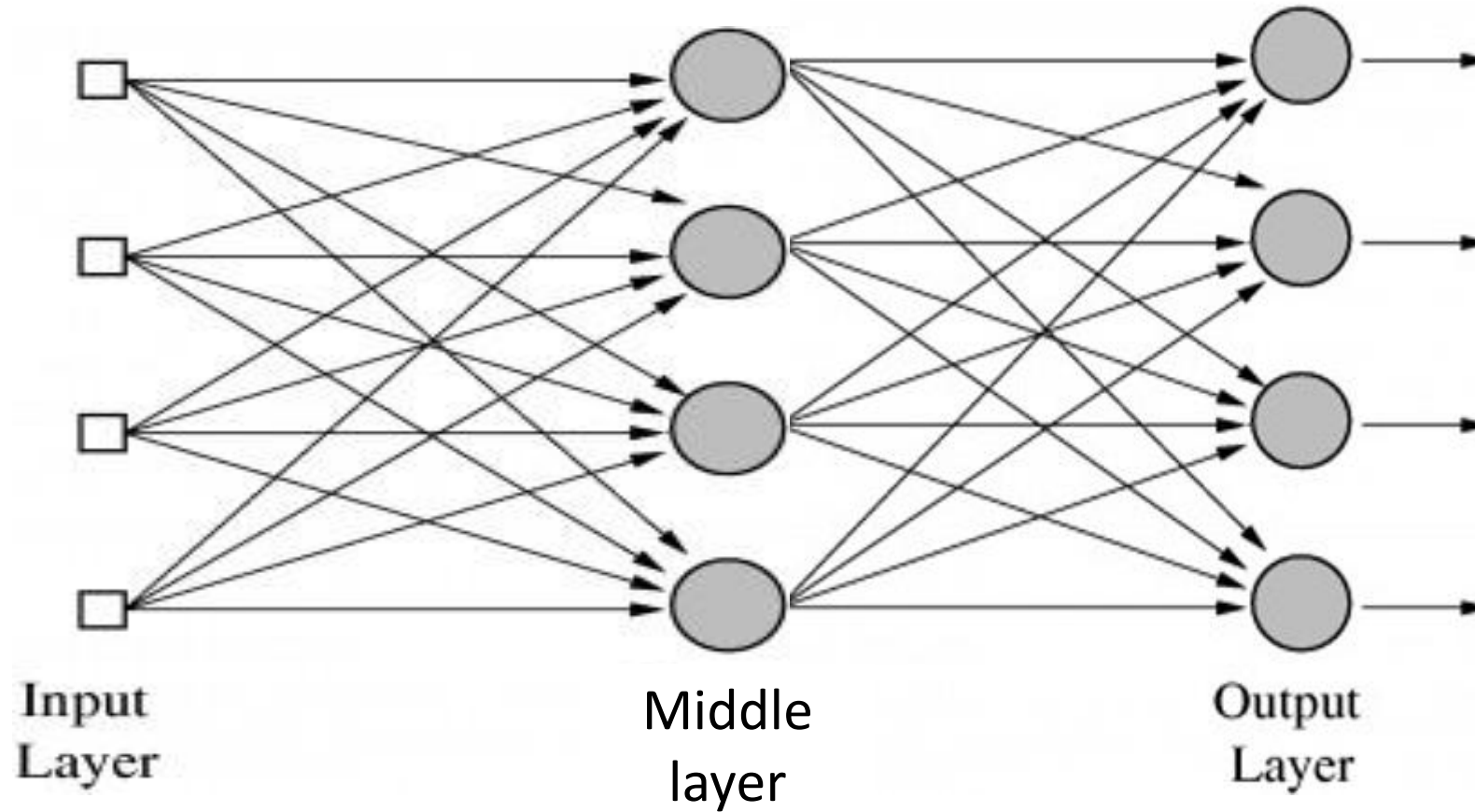
$$v = \sum_{i=0}^m x_i w_i = \vec{w}^T(n) \vec{x}(n)$$

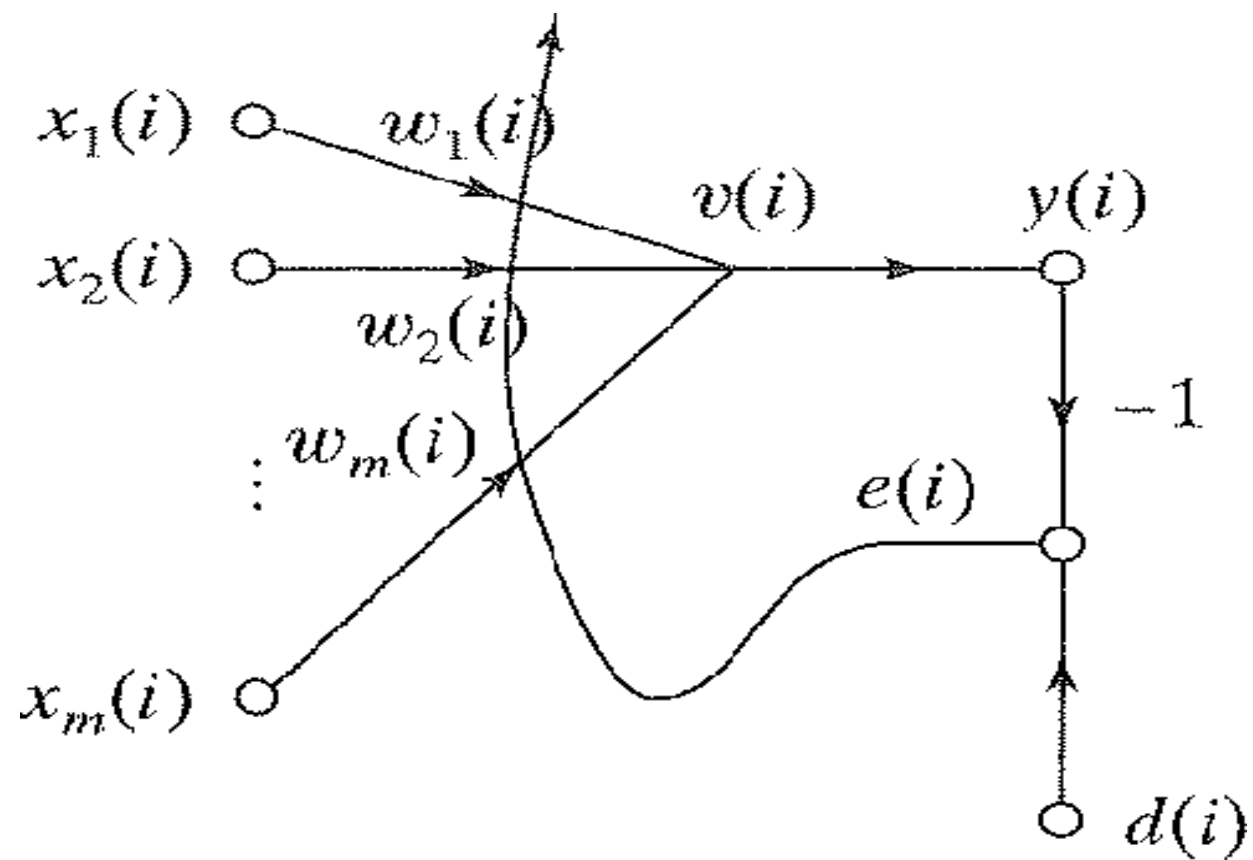
$$y = f(v)$$

$f(v)$: non-linear activation function



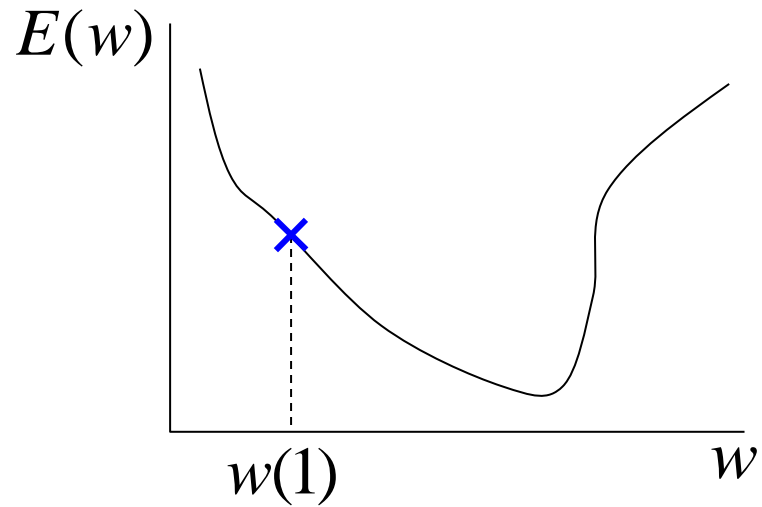
Multi-layer Perceptron







Gradient-Decent



設目前 $w = w(1)$, 請將 w 加一值

使得 $E(w(2)) < E(w(1))$

其中 $w(2) = w(1) + \text{此值}(\Delta w)$

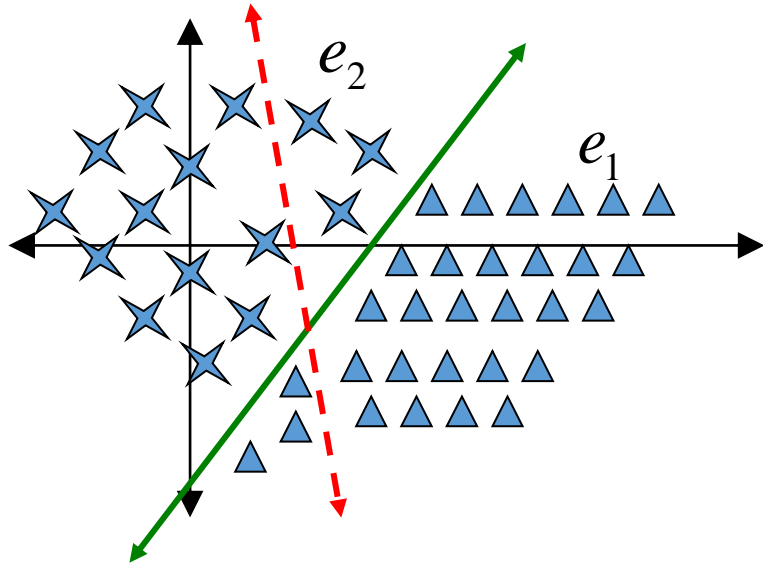
令 $E'(w(1))$ 為 $E(\cdot)$ 在 $w(1)$ 之微分

由圖知 $E'(w(1)) < 0$, 但明顯地 $\Delta w > 0$

$$\text{因此 } \Delta w = -\eta E'(w(1)) = -\eta \left. \frac{\partial f(w)}{\partial w} \right|_{w=w(1)}$$



perceptron



Given
$$\underline{\underline{E = \sum_{k=1}^{\phi} (e_k(n))^2}} \quad \geq 0$$

當 pattern $\vec{x}(n)$ 代入

若 $d(n) \neq y(n)$

→ 希望改 \vec{w} 使 $E(n)$ 小

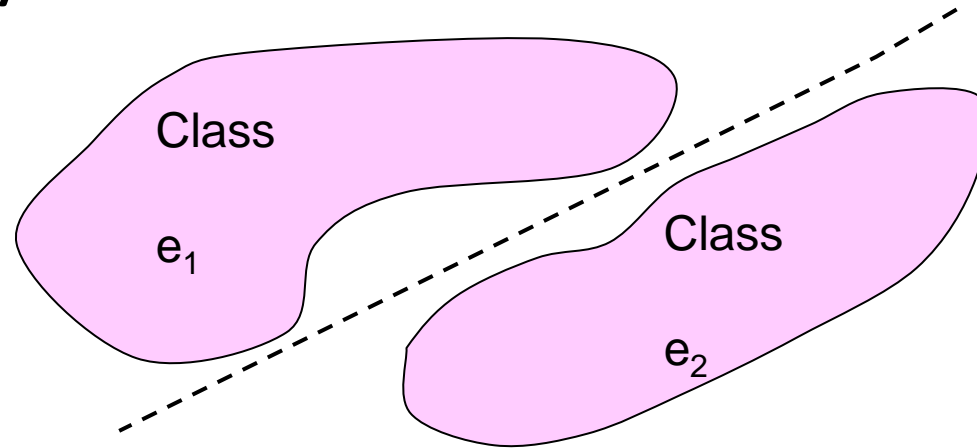
$$\Delta \vec{w} = -\eta \frac{\partial E(n)}{\partial \vec{w}} = +\eta [d(n) - y(n)] \vec{x}(n)$$

若 $d > 0$ 則 $\Delta \vec{w} = \eta_1 \vec{x}(n)$

若 $d < 0$ 則 $\Delta \vec{w} = -\eta_1 \vec{x}(n)$



Summary



For every pattern

$$\vec{w}^T \vec{x} > 0 \quad \forall \vec{x} \in e_1$$

$$\vec{w}^T \vec{x} < 0 \quad \forall \vec{x} \in e_2$$

- Adapting weights:

1. If correctly classified

$$\vec{w}(n+1) = \vec{w}(n) \quad \text{if } \vec{w}^T \vec{x} > 0 \text{ and } \vec{x} \in e_1$$

$$\vec{w}(n+1) = \vec{w}(n) \quad \text{if } \vec{w}^T \vec{x} < 0 \text{ and } \vec{x} \in e_2$$

2. If incorrectly classified (also called Hebbian learning)

$$\vec{w}(n+1) = \vec{w}(n) - \eta(n)\vec{x}(n) \quad \text{if } \vec{w}^T \vec{x} > 0 \text{ and } \vec{x} \in e_2$$

$$\vec{w}(n+1) = \vec{w}(n) + \eta(n)\vec{x}(n) \quad \text{if } \vec{w}^T \vec{x} < 0 \text{ and } \vec{x} \in \underline{e_1}$$



Summary of Learning

- I. Initialize the weight $w(0)=0$
- II. at $t=n$ apply the training patterns
- III. Compute the actual responses

$$y(n) = \text{sign}[\vec{w}^T(n)\vec{x}(n)]$$

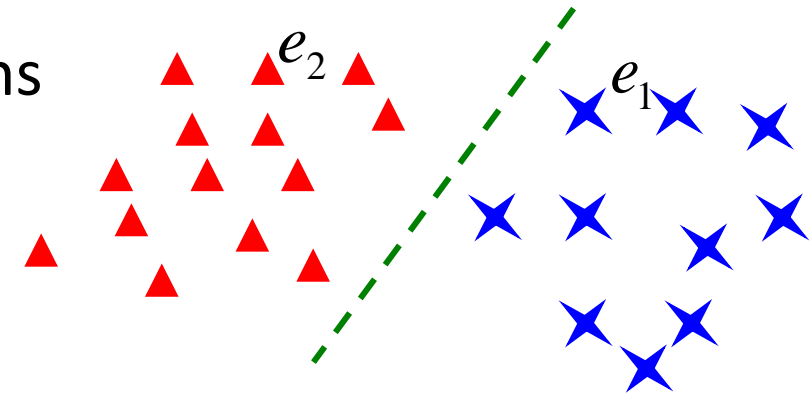
- IV. Update the weight

$$\vec{w}(n+1) = \vec{w}(n) + \eta[d(n) - y(n)]\vec{x}(n)$$

$$\because d(n) = \begin{cases} 1 & \text{if } \vec{x} \in e_1 \\ -1 & \text{if } \vec{x} \in e_2 \end{cases}$$

$$\rightarrow \text{if } d(n) \neq y(n), \text{ then } \vec{w}(n+1) = \begin{cases} \vec{w}(n) + \eta \vec{x}(n) & \text{if } \vec{x} \in e_1 \\ \vec{w}(n) - \eta \vec{x}(n) & \text{if } \vec{x} \in e_2 \end{cases}$$

- V. Increment n , go back to II





Learning Rate Annealing Schedules

1. $\eta = \eta_0$ for all n
2. $\eta(n) = \frac{c}{n}$ stochastic approximation schedule
3. $\eta(n) = \frac{\eta_0}{1 + (n/\tau)}$ Darken & Moody (1992)

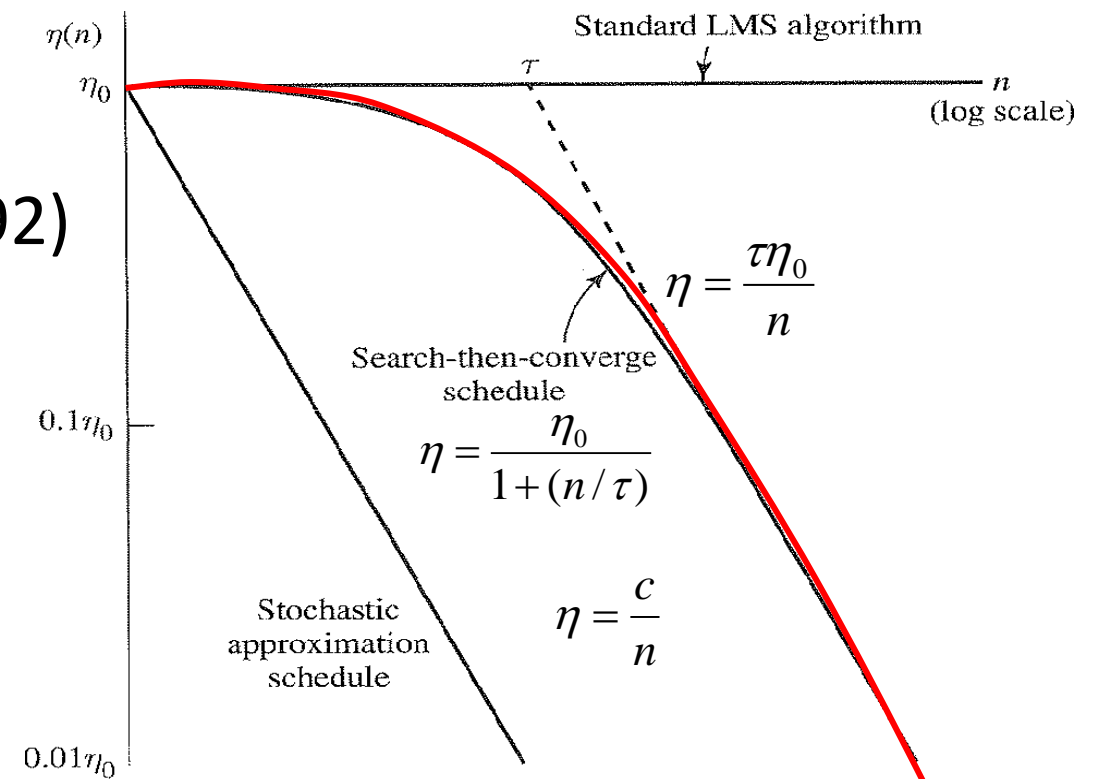


FIGURE 3.5 Learning-rate annealing schedules.