

Bi-directional learning system

Shikui Tu

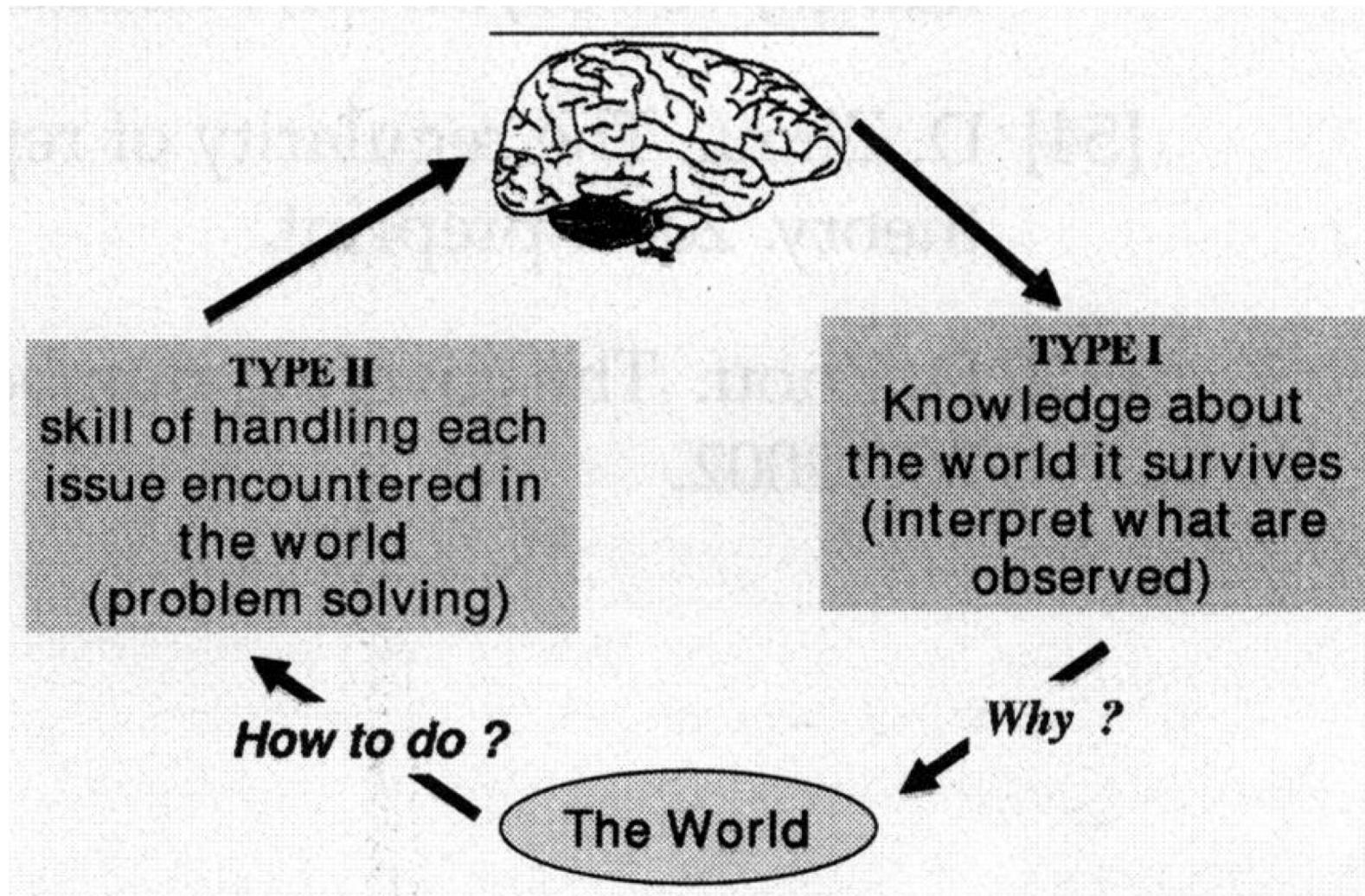
Department of Computer Science and
Engineering, Shanghai Jiao Tong University

2018-05-17

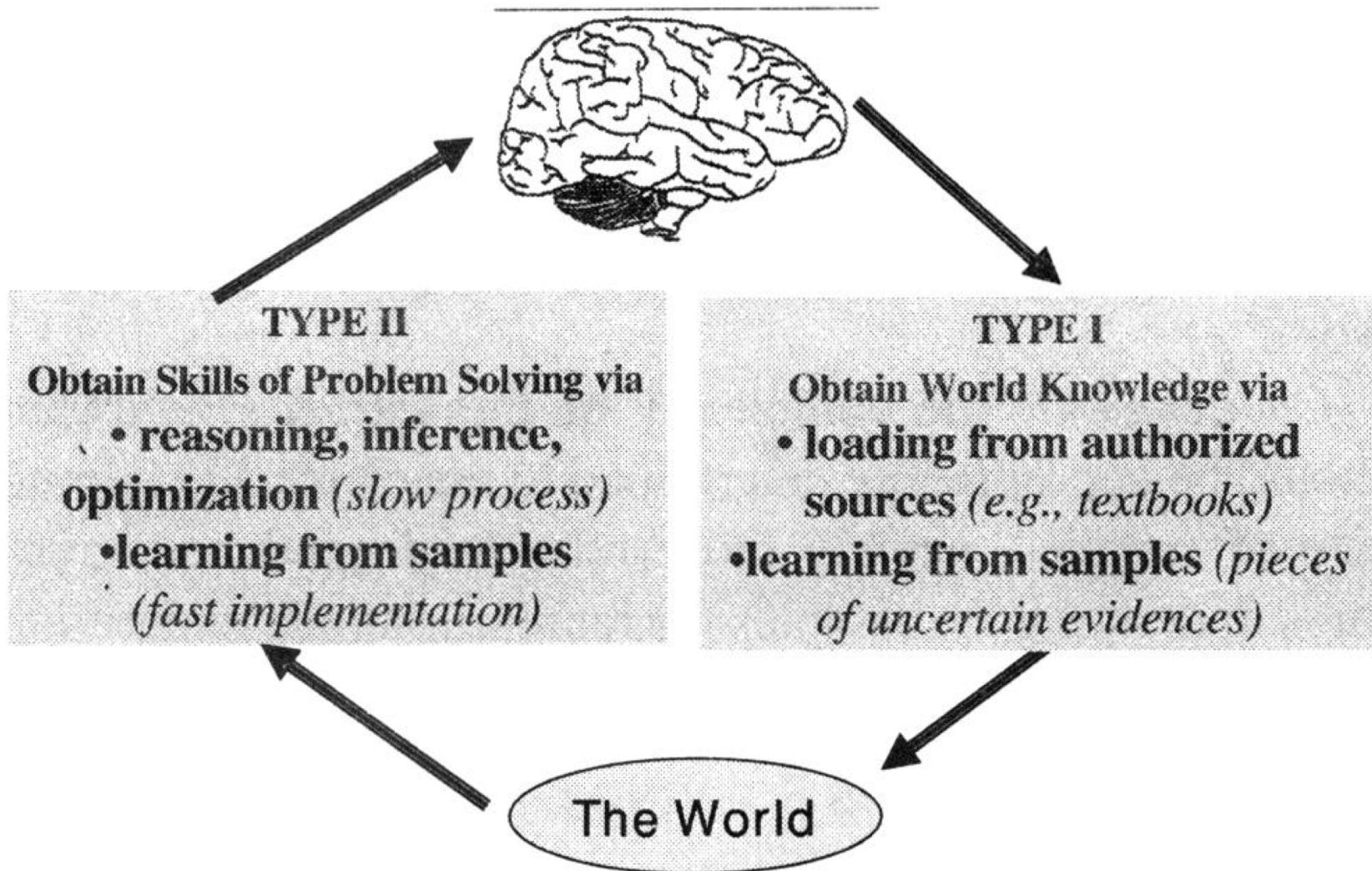
Outline

- **Fundamentals and challenges of learning**
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- Bi-directional learning
 - Inbound learning theory
 - Outbound learning theory
 - Bi-directional architectures
- Early bi-directional deep learning

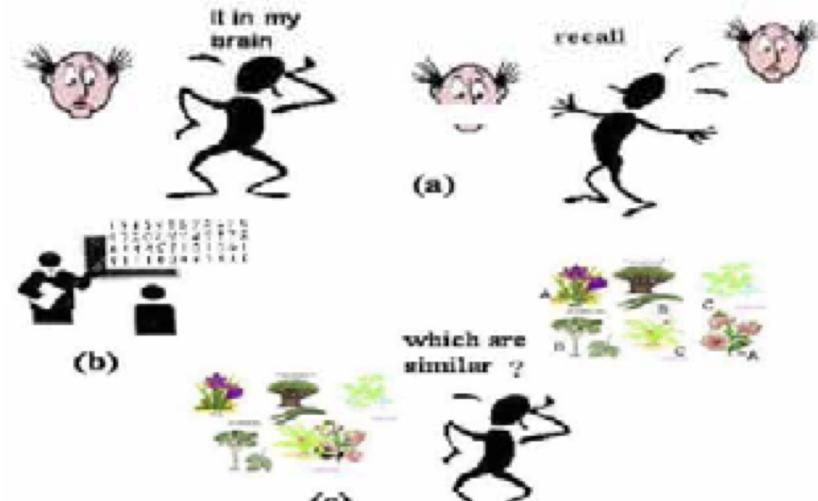
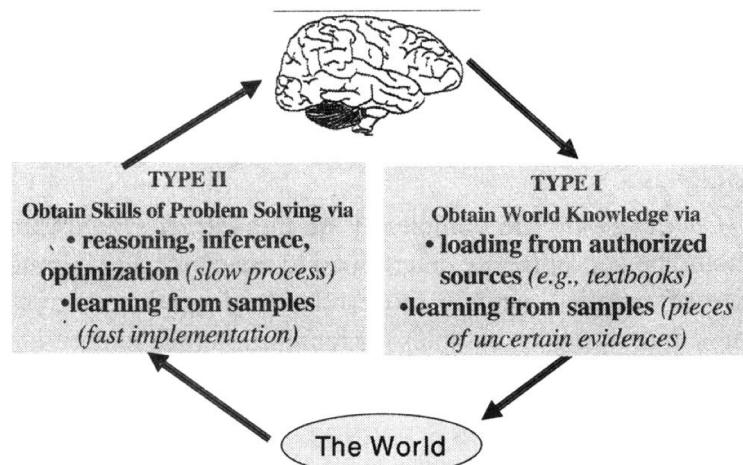
Two types of intelligent ability



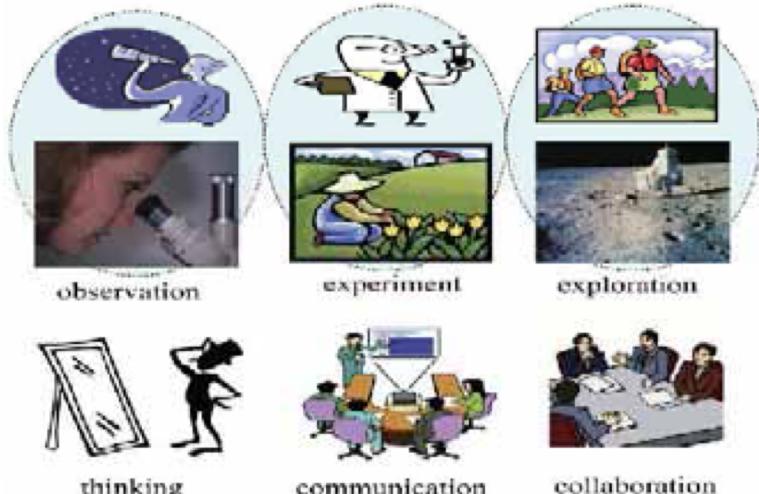
How to get the abilities



Abilities and problem-solving skills



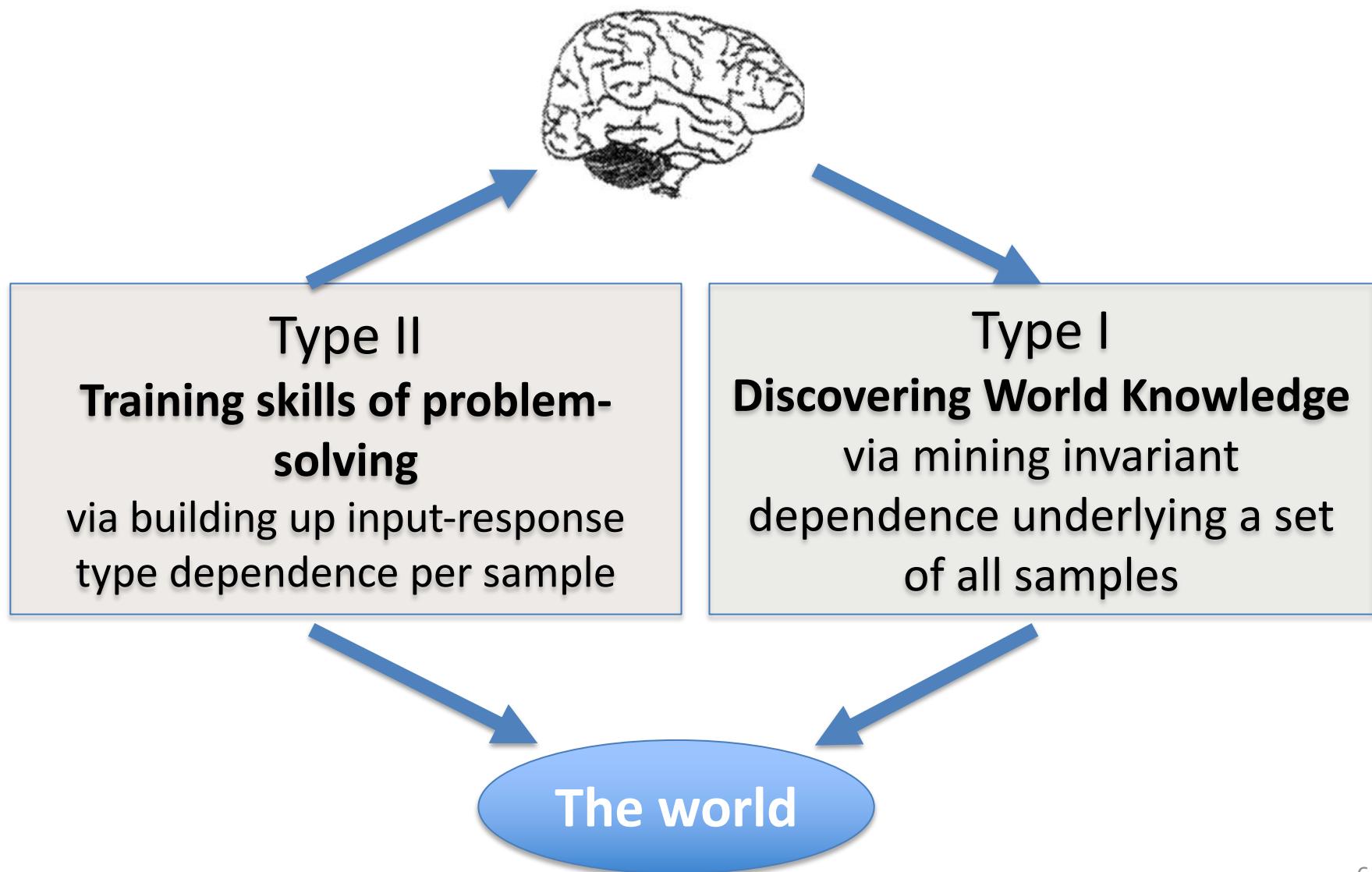
(I) **perception**



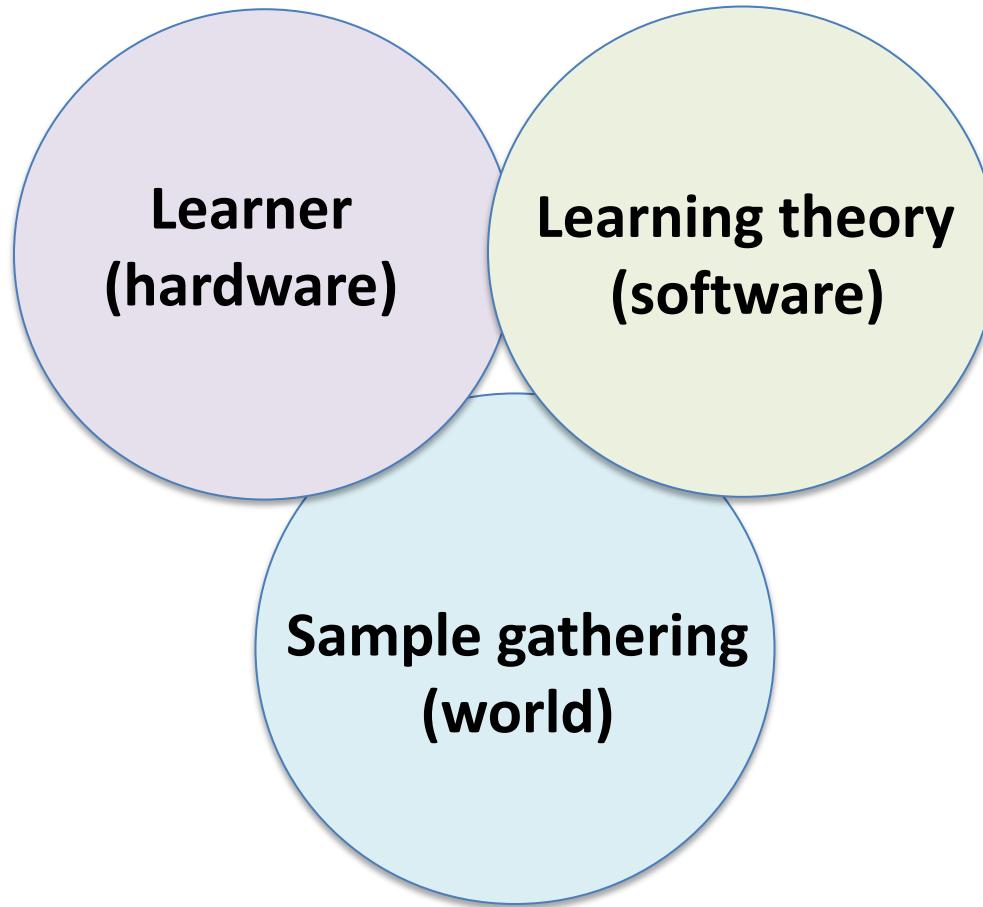
Abilities of knowledge discovery



Two types of learning



Key ingredients of learning

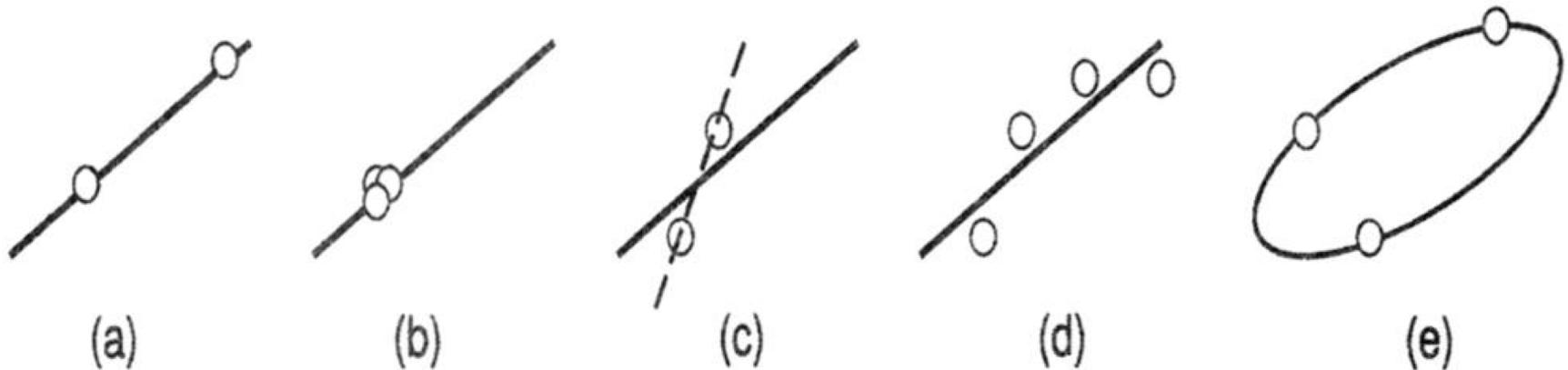


Two key learning challenges

- **Challenge I**
 - The learner's hardware should be designed not only be able to accommodate but also appropriately match the interested dependence structures underlying the world.
- **Challenge II**
 - The complexity of the learner's hardware should be appropriately determined to match usually a finite size samples, namely those reliable dependence structures underlying the samples for representing the underlying world.

Uncertainty in samples and randomly sampling

Consider a simple problem of learning a line from samples:

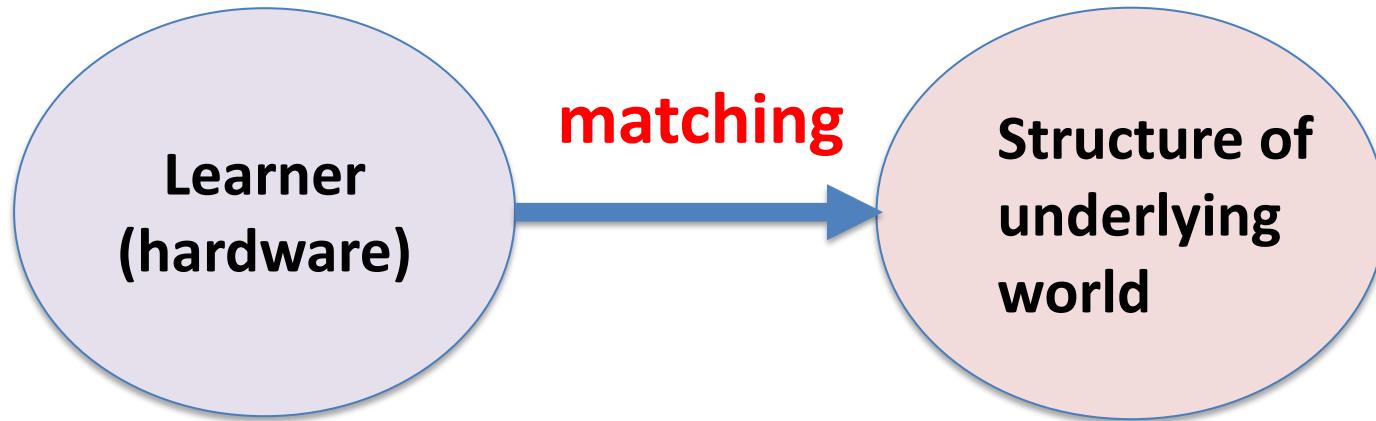


- (a) Conceptual case
- (b) Fail when two samples are the same
- (c) Noise in gathering samples, quantization effects
- (d) More samples to reduce uncertainty
- (e) To find a more complicated structure, but not enough samples



Statistical learning

Learner represents data



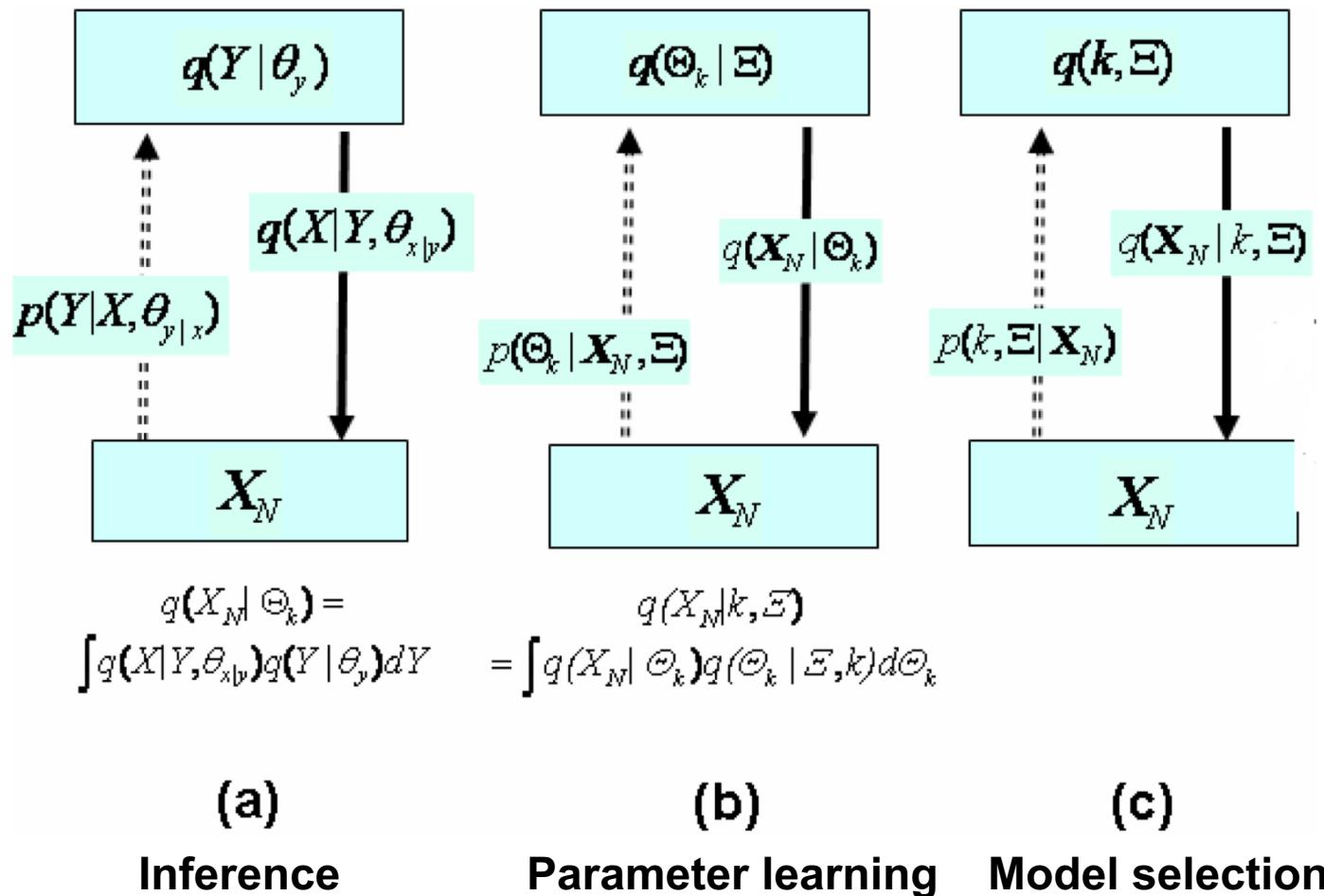
Learner's hardware appropriately represents dependence among data



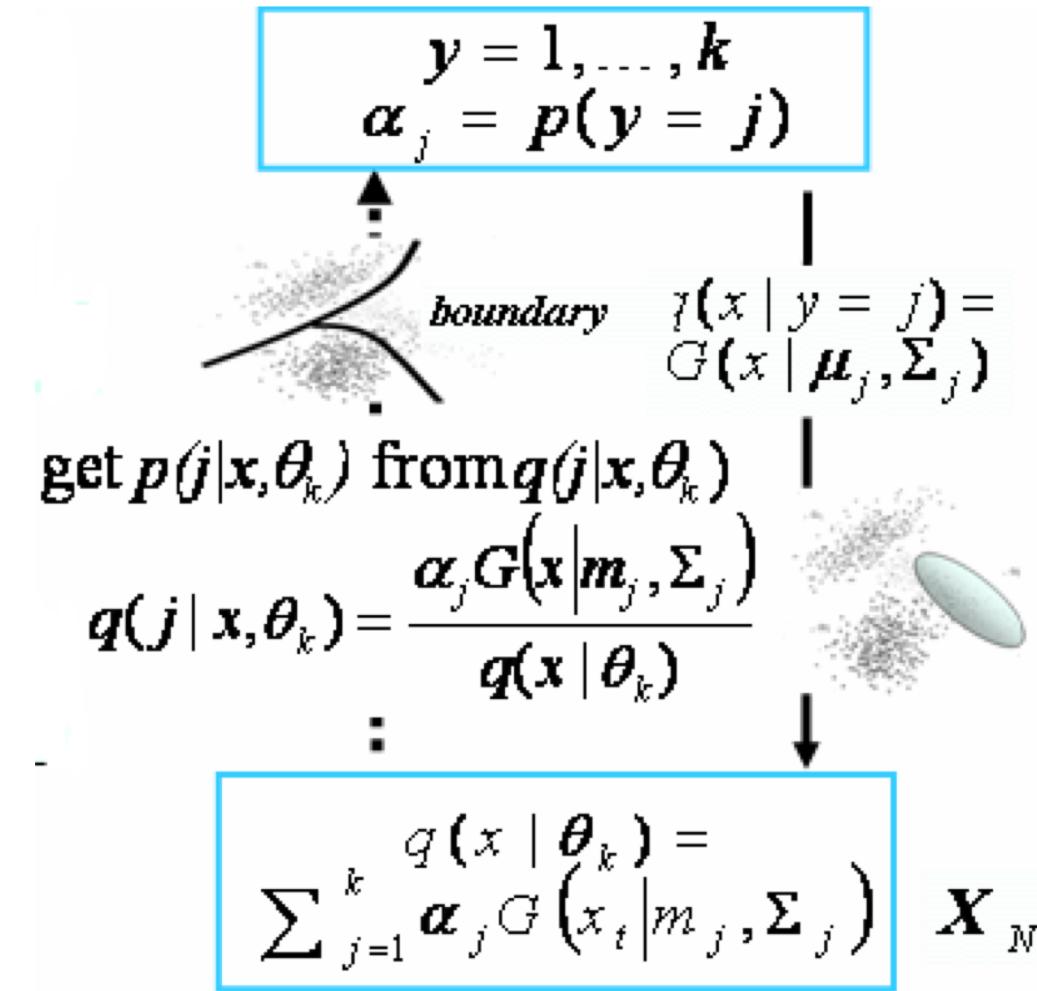
Outline

- **Fundamentals and challenges of learning**
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- **Bi-directional learning**
 - Inbound learning theory
 - Outbound learning theory
 - Bi-directional architectures
- **Early bi-directional deep learning**

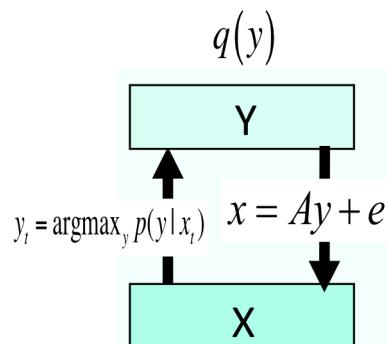
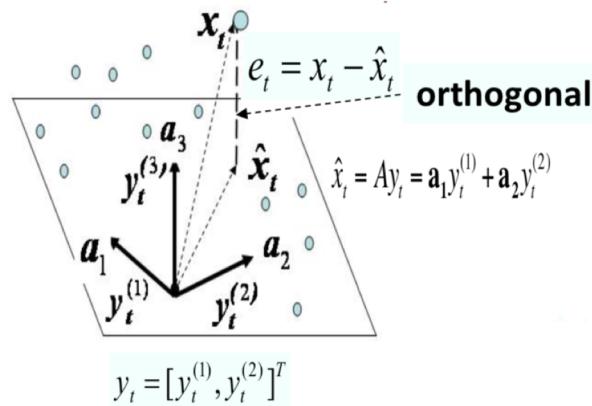
Three levels of inverse problems



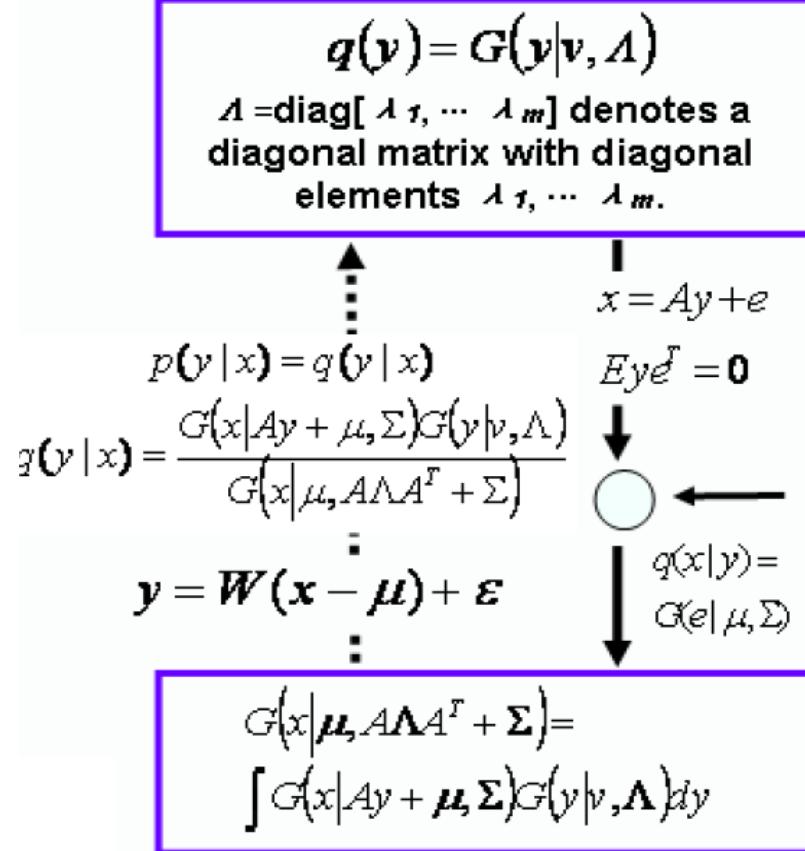
Learning Gaussian Mixture Models



Learning Factor Analysis



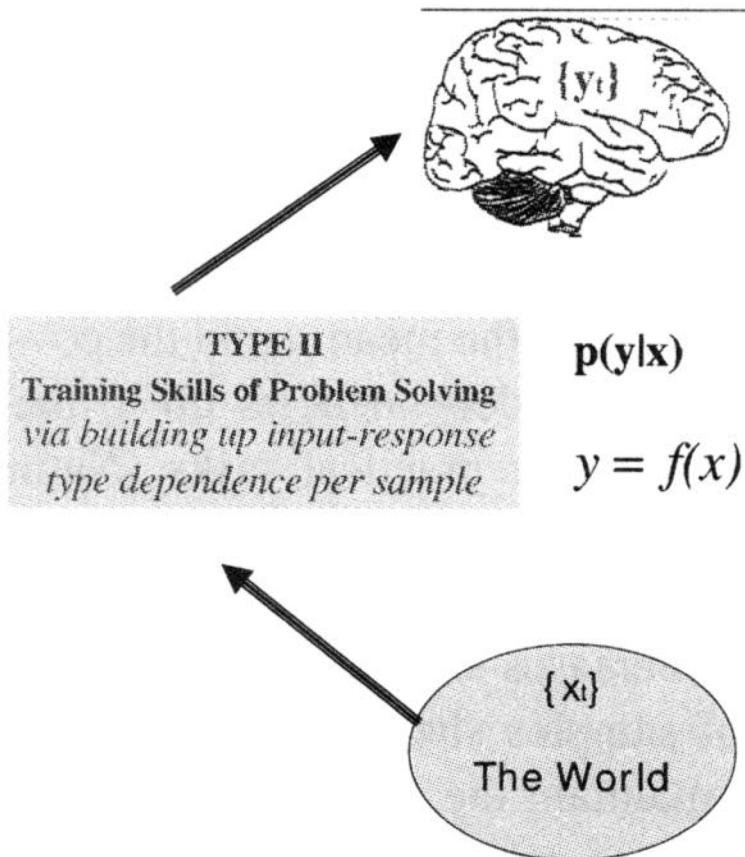
$$\max_{\theta} \sum_t \ln q(x_t | \theta)$$



Outline

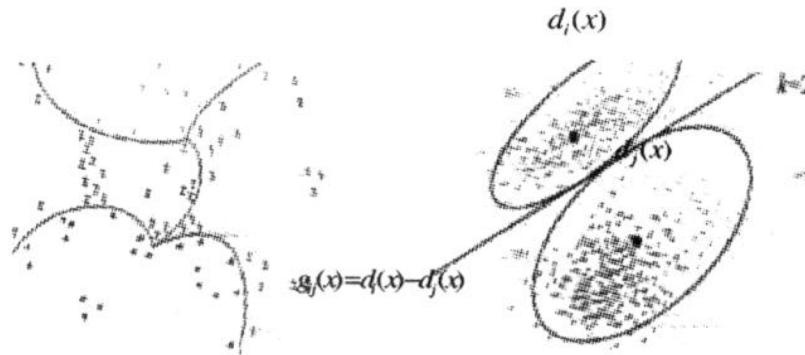
- Fundamentals and challenges of learning
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- Bi-directional learning
 - **Inbound learning theory**
 - Outbound learning theory
 - Bi-directional architectures
- Early bi-directional deep learning

Inbound architecture

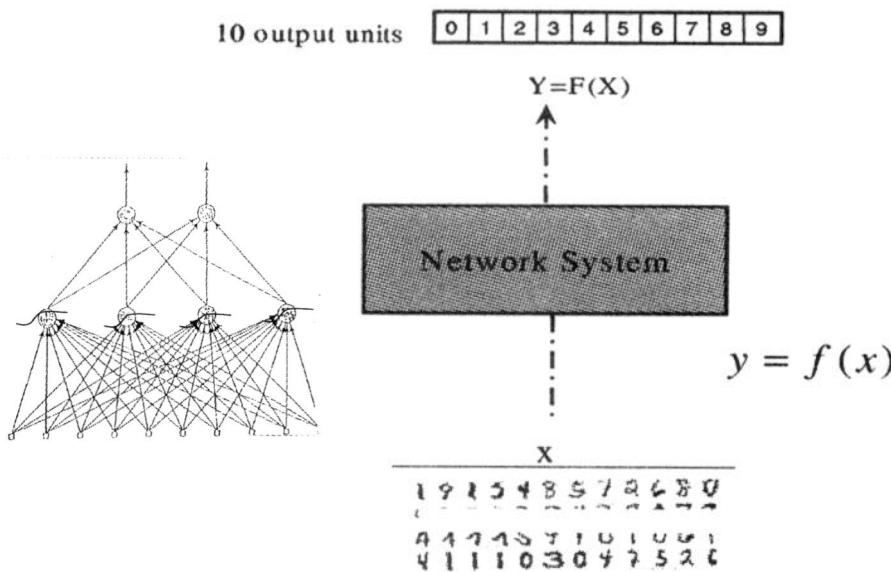


Pairwise structures

$\{(x_t, y_t)\}$



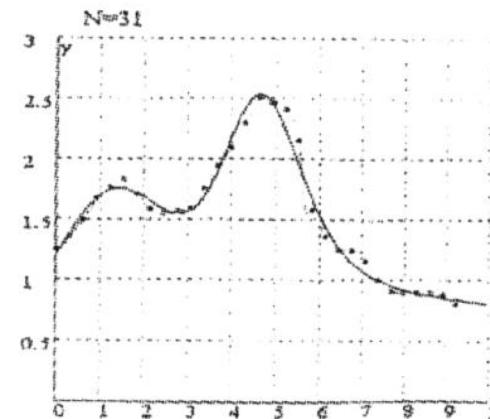
(a) classification



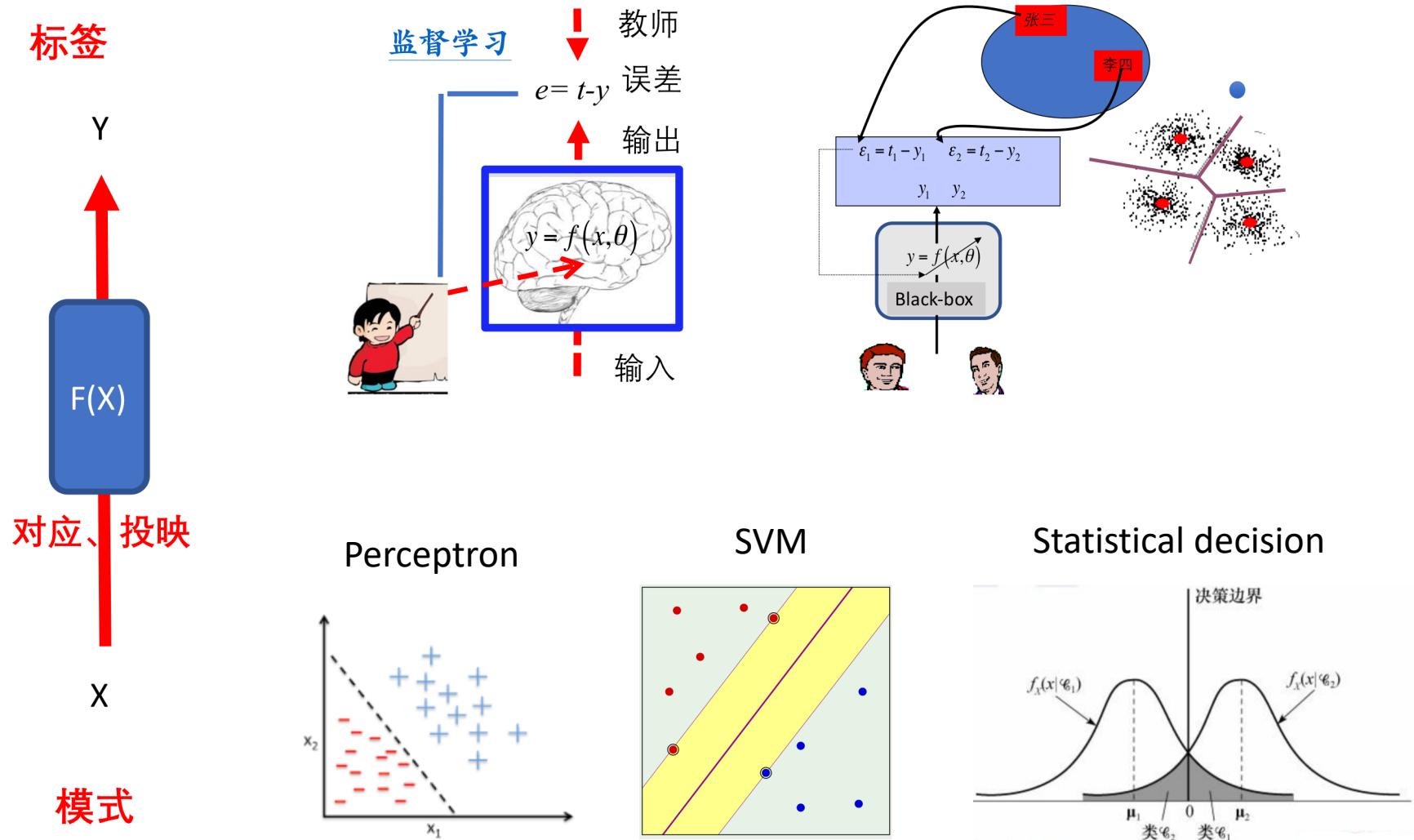
(b) regression



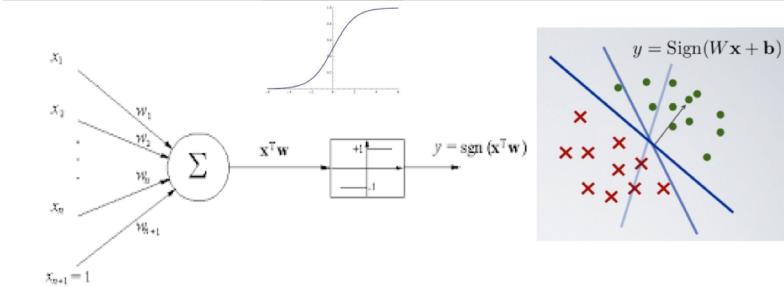
$$E(y|x) = \int y p(y|x) dy$$



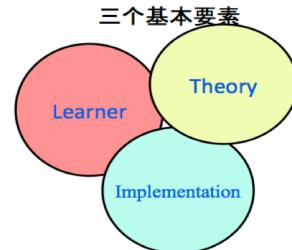
最佳拟合、最大区分、最小错误



From perceptron to deep learning



Frank Rosenblatt
(1928-1971)



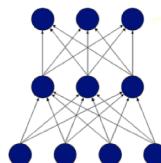
Rosenblatt's Perceptron Learning(1957)

Minsky and Papert Perceptrons (1969)
(unable to learn an XOR function ???)

到Rumelhart, McClelland & Hinton 的BP-learning(1986)

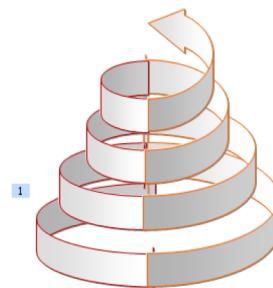
再到Hinton 的 Deep learning (2006),

突破的主要原因是计算技术和信息技术的发展积累所
产生的高速计算能力和大数据的获得和处理能力



Structure	Type of Decision Regions	Exclusive-OR Problem	Classes with Mixed Regions	Most General Region Shapes
Single-layer	Half plane bounded by hyperplane	(A) (B) (B) (A)	B (A)	
Two-layers	Convex open or closed regions	(A) (B) (B) (A)	B (A)	
Three-layers	Arbitrary (Complexity limited by number of nodes)	(A) (B) (B) (A)	B (A)	

完成了螺旋上升之肯定 – 否定 –
再肯定 – 再否定 – 出现突破





David H. Hubel



Torsten N. Wiesel

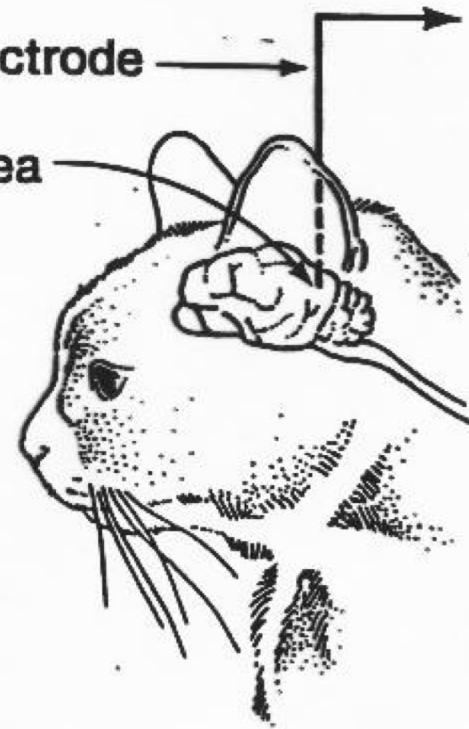
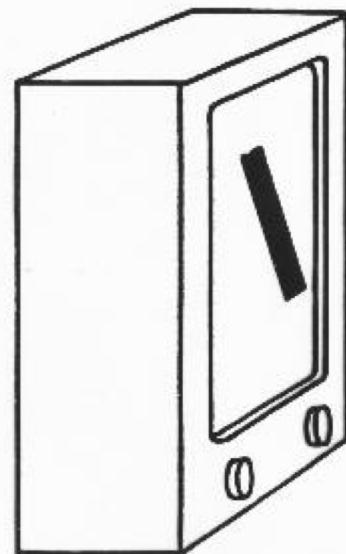
Wiesel and Hubel 特征检测理论

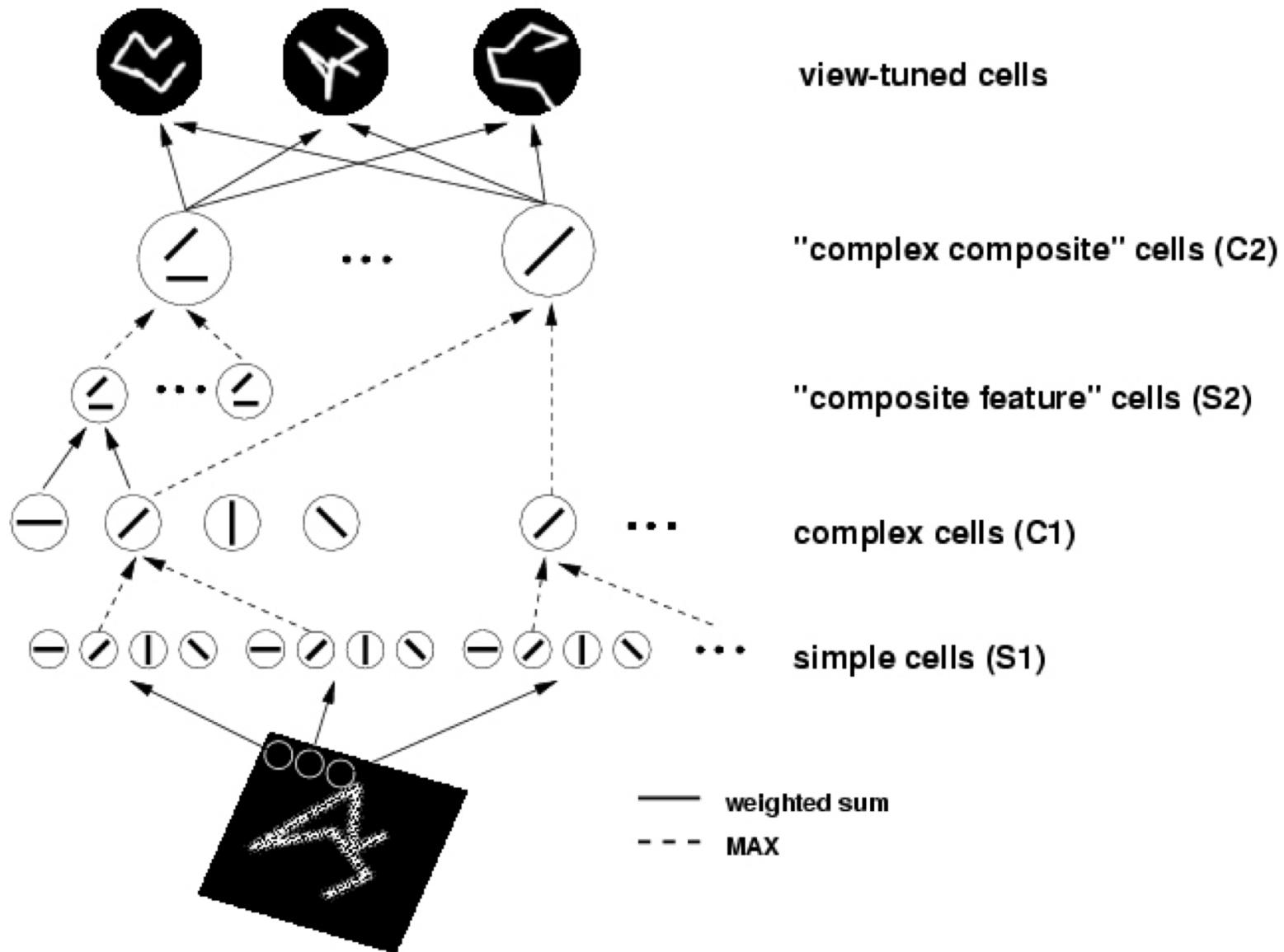
Electrical signal
from brain

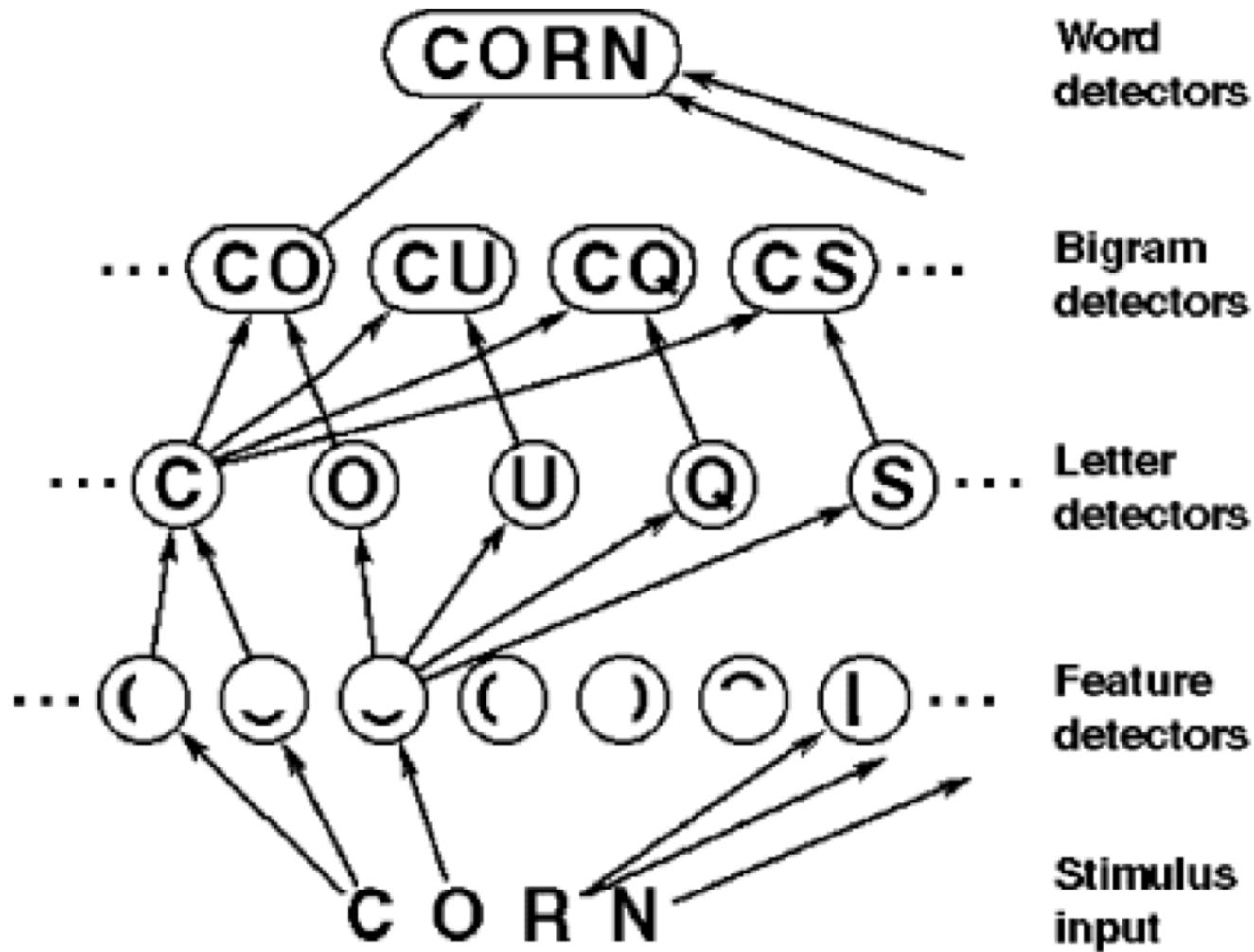
Recording electrode

Visual area
of brain

Stimulus



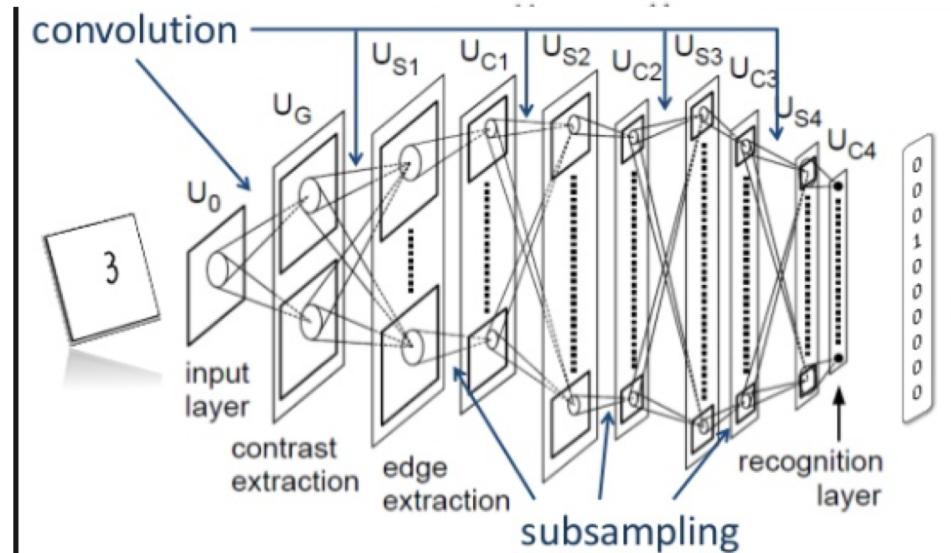




Neocognitron

神经认知机

Fukushima (1980). Hierarchical multilayered neural network

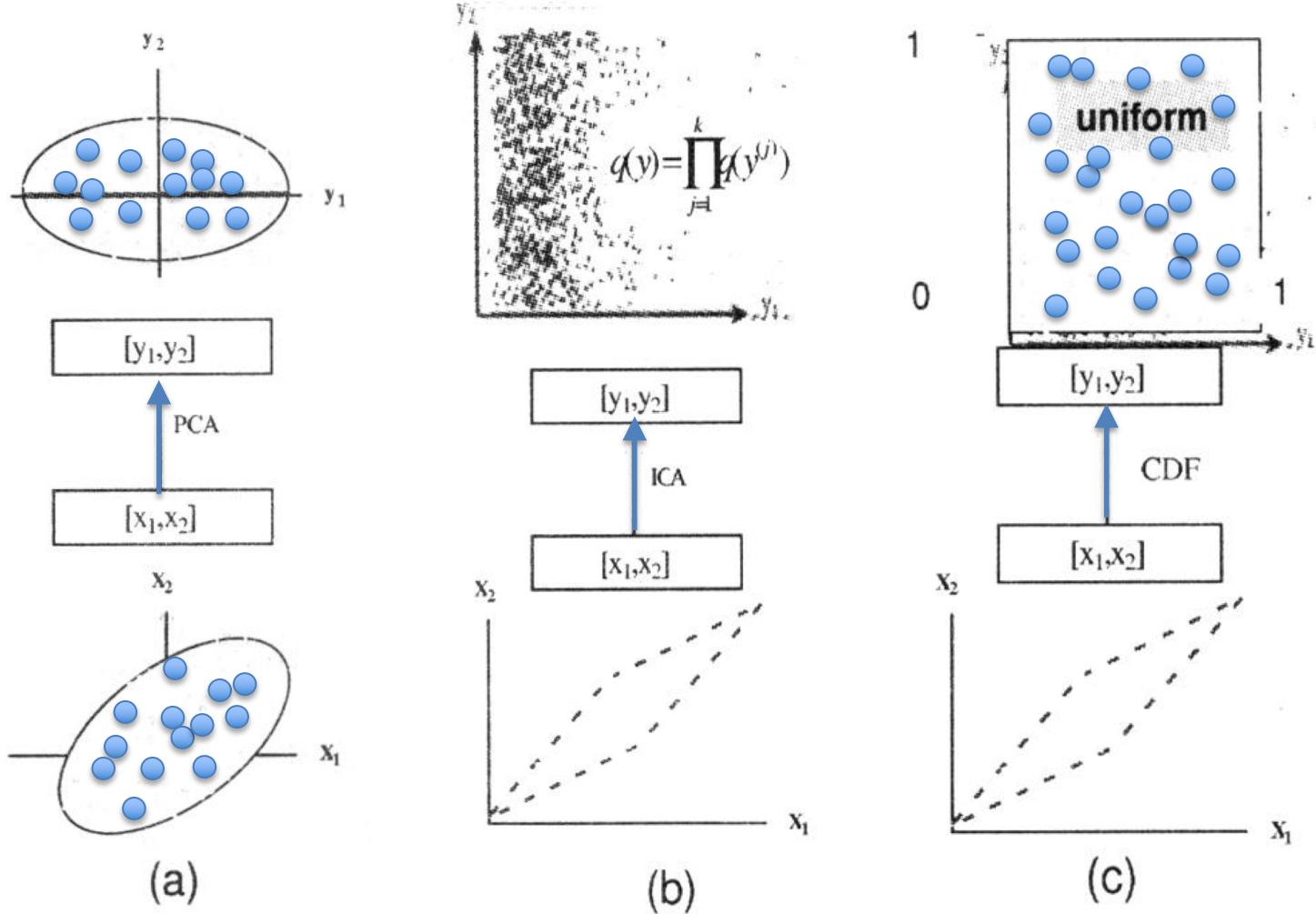


S-cells work as feature-extracting cells. They resemble simple cells of the primary visual cortex in their response.

C-cells, which resembles complex cells in the visual cortex, are inserted in the network to allow for positional errors in the features of the stimulus. The input connections of C-cells, which come from S-cells of the preceding layer, are fixed and invariable. Each C-cell receives excitatory input connections from a group of S-cells that extract the same feature, but from slightly different positions. The C-cell responds if at least one of these S-cells yield an output.

Transformation structures

$$p(y) = \int p(y|x)p(x)\mu(dx) \quad q(y) \Rightarrow q(y) = \prod_{j=1}^k q(y^{(j)}) \text{ by } \min \int p(y) \ln \frac{p(y)}{q(y|\theta)} dy$$

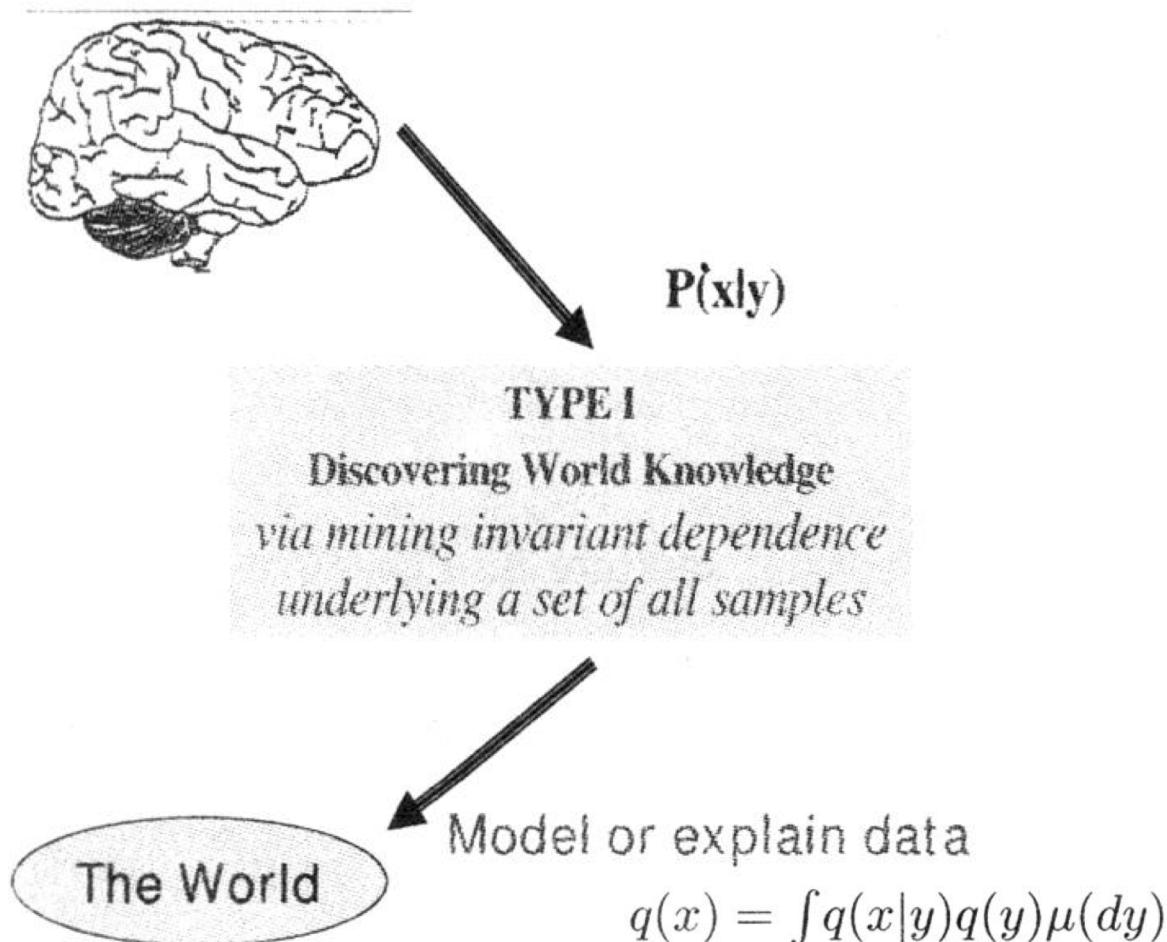


最佳抽象（信息之最佳保持）：最大传递、最小互信息、最小冗余

Outline

- **Fundamentals and challenges of learning**
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- **Bi-directional learning**
 - Inbound learning theory
 - **Outbound learning theory**
 - Bi-directional architectures
- Early bi-directional deep learning

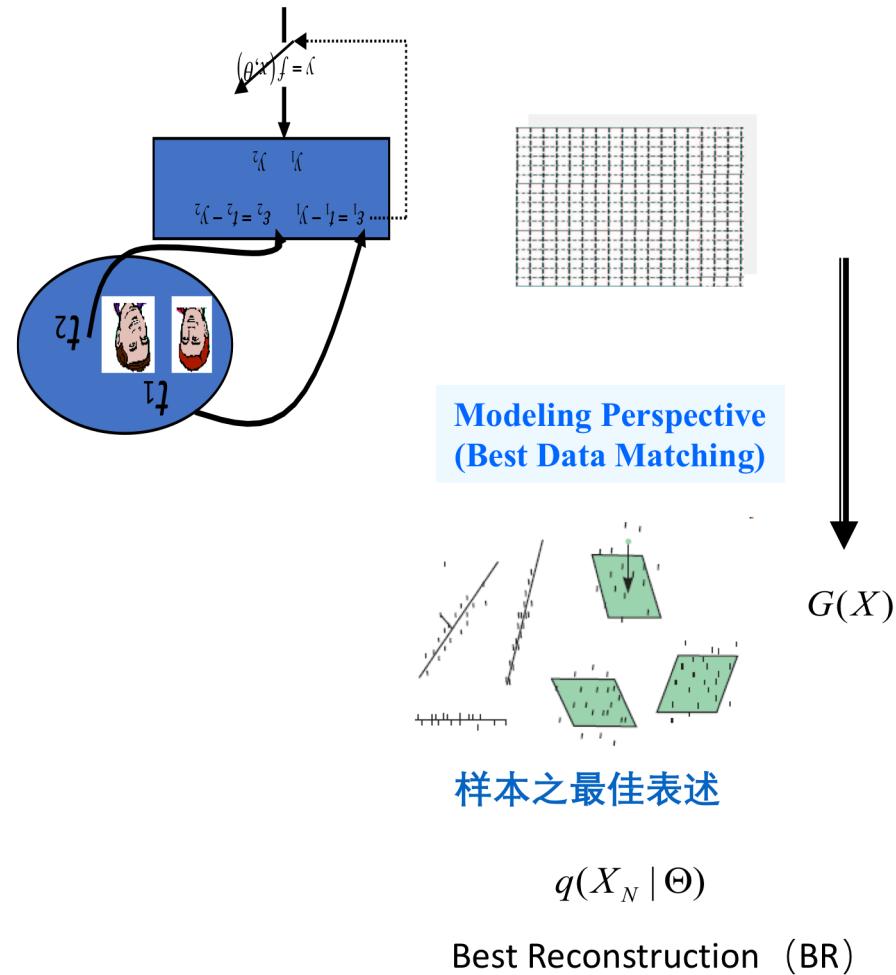
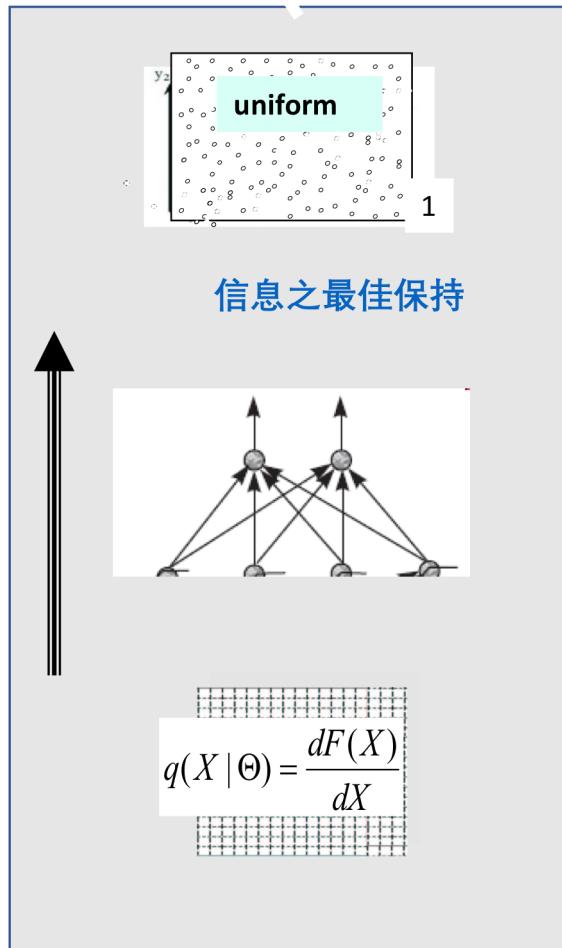
Outbound architecture



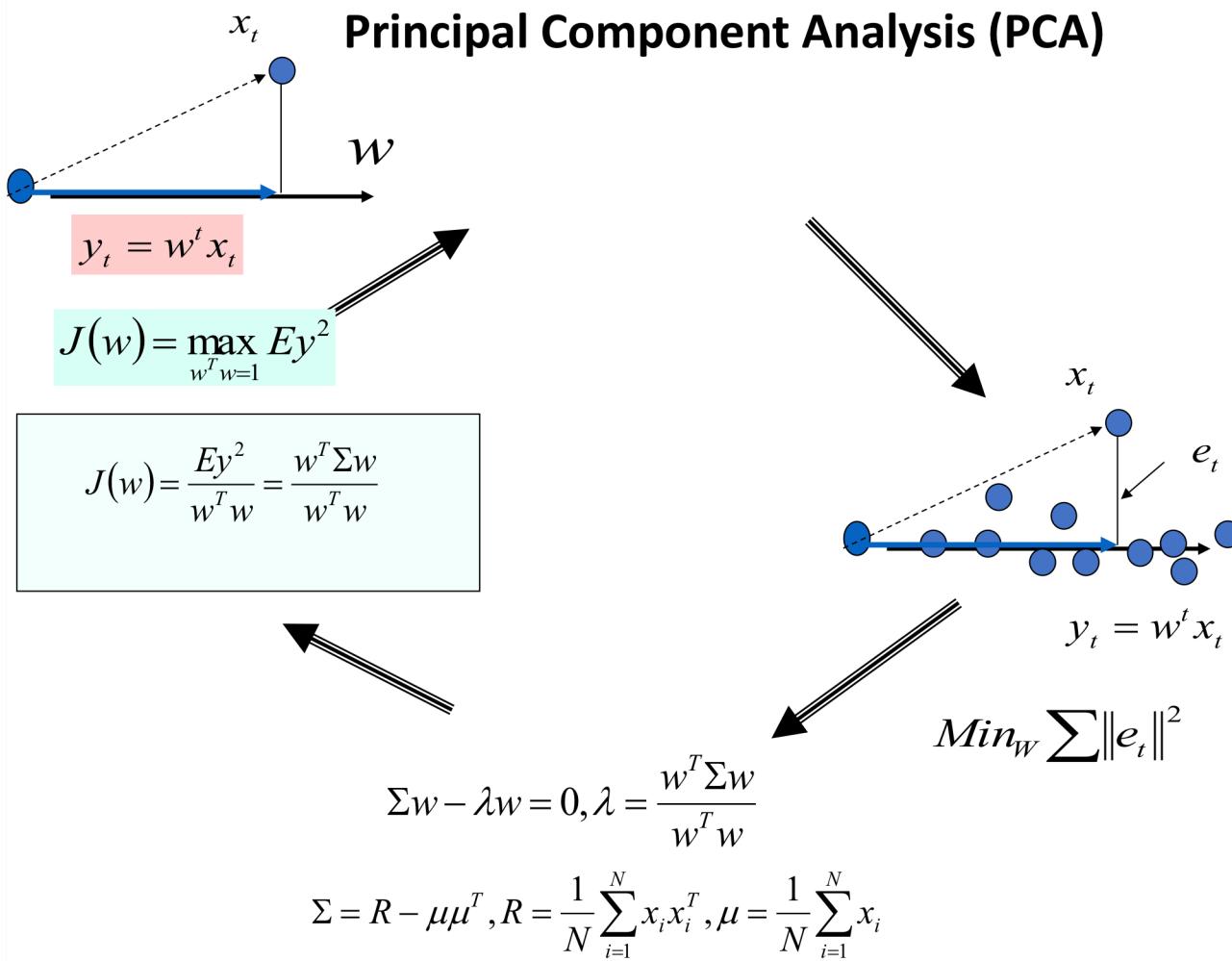
It describes dependence among variables $x^{(1)}, \dots, x^{(d)}$ by attempting to reconstruct observations of x from certain inner factors y via $q(x|y)$ in an appropriate parametric form.

Best reconstruction

最佳重建

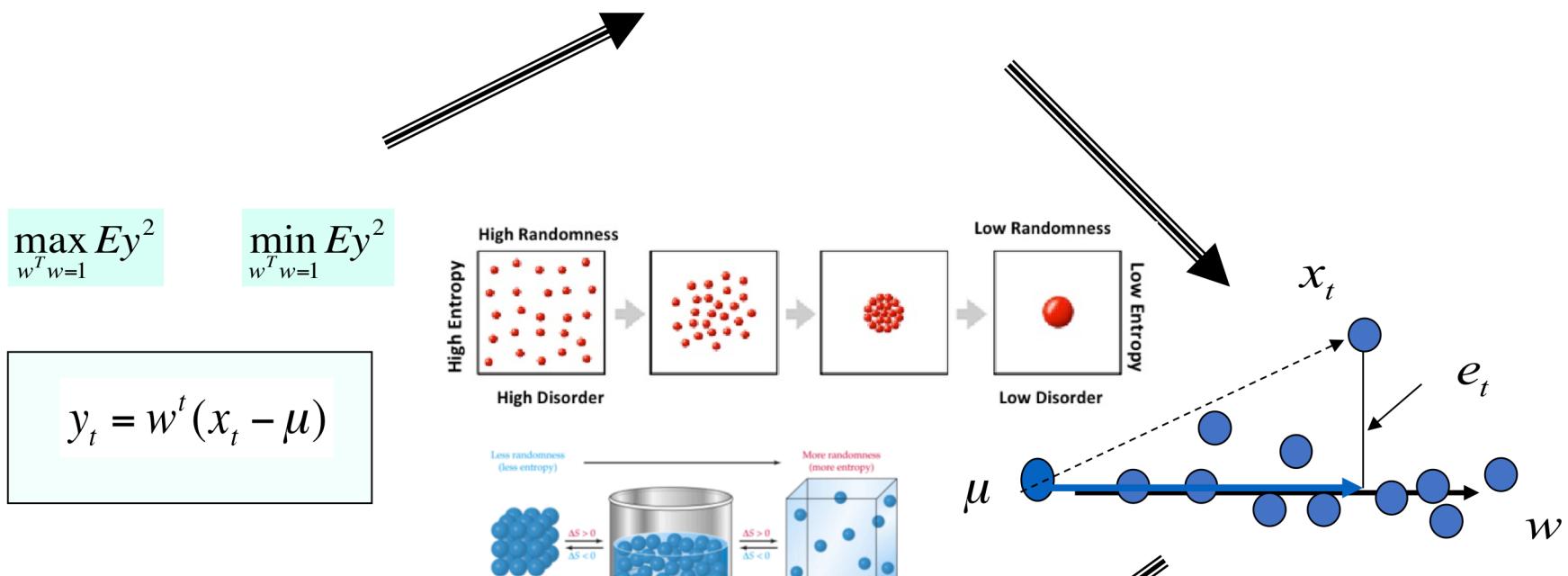


PCA reconstruction



BA theory : maximum versus minimum information transfer

best inverse



Entropy

$$-\int p(x) \ln p(x) dx$$

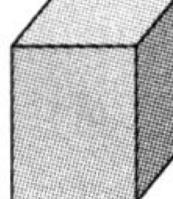
$$-\int G(\mathbf{x} | \mu, \sigma^2) \ln G(\mathbf{x} | \mu, \sigma^2) d\mathbf{x} = 0.5 \ln(2\pi e) + 0.5\sigma^2$$

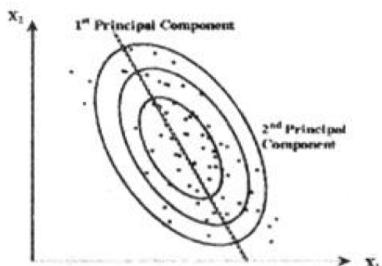
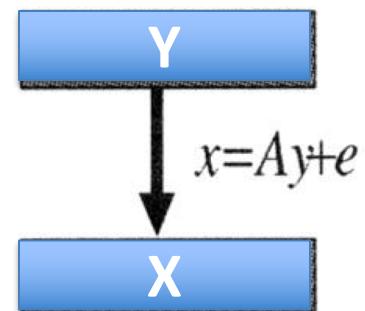
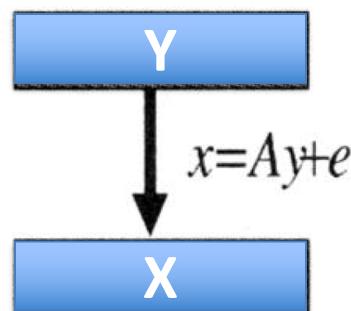
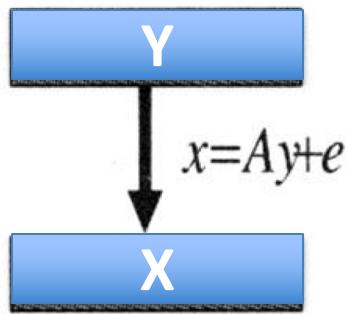
BR theory : minimum error fitting

Linear latent structures

$$q(y) = G(y|0, I)$$


$$q(y) = \prod_{j=1}^k q(y^{(j)})$$


$$q(y^{(j)}) = q_j^{y^{(j)}} (1 - q_j)^{1-y^{(j)}}$$

$$y_j = 1 \text{ or } 0$$



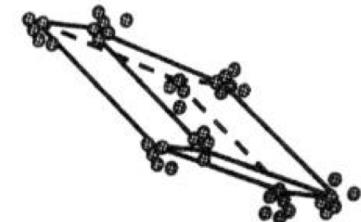
$$x = y + e$$

$y = [y_1 \ y_2 \ \dots \ y_n]^T$

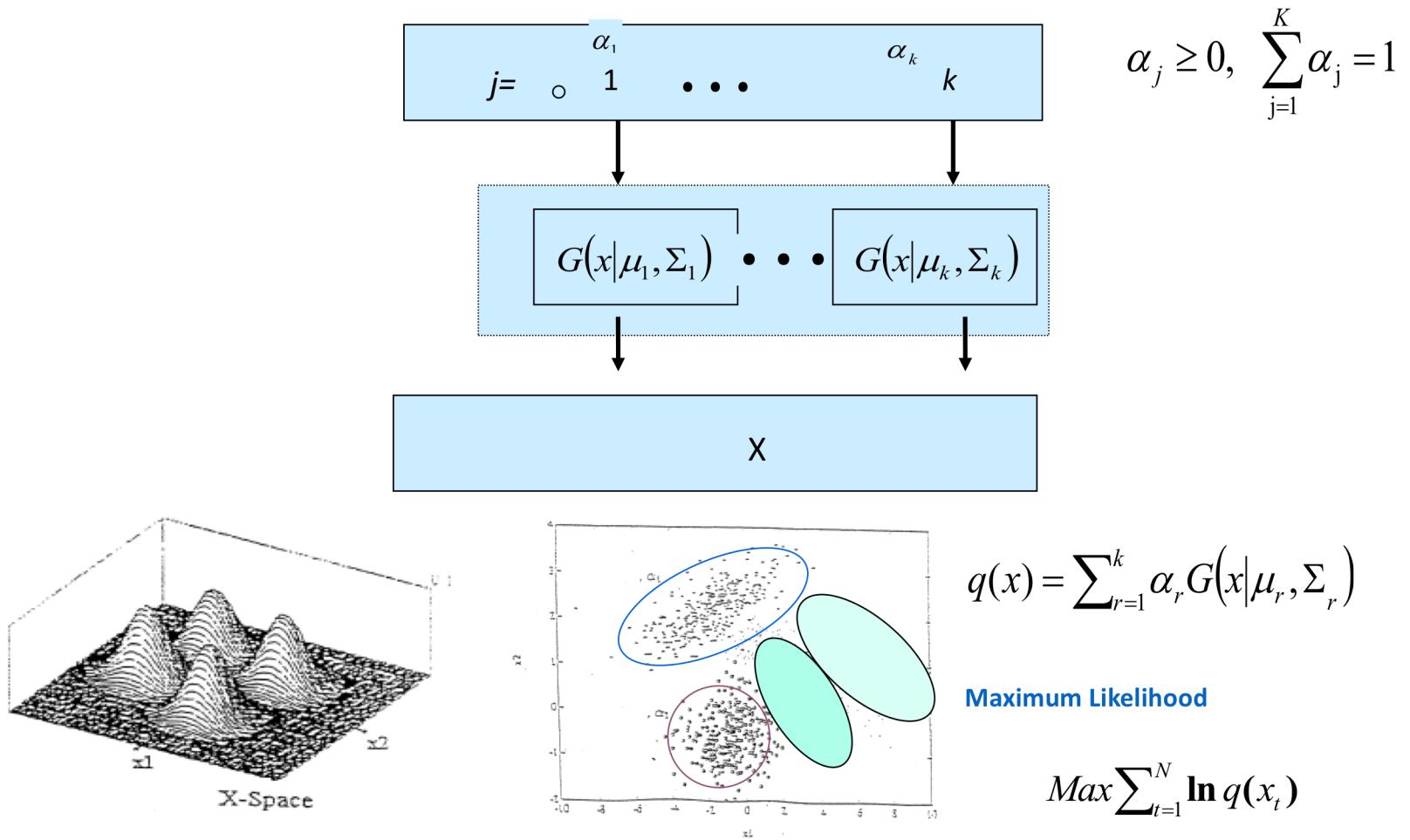
$e = [e_1 \ e_2 \ \dots \ e_n]^T$

Gaussian noise





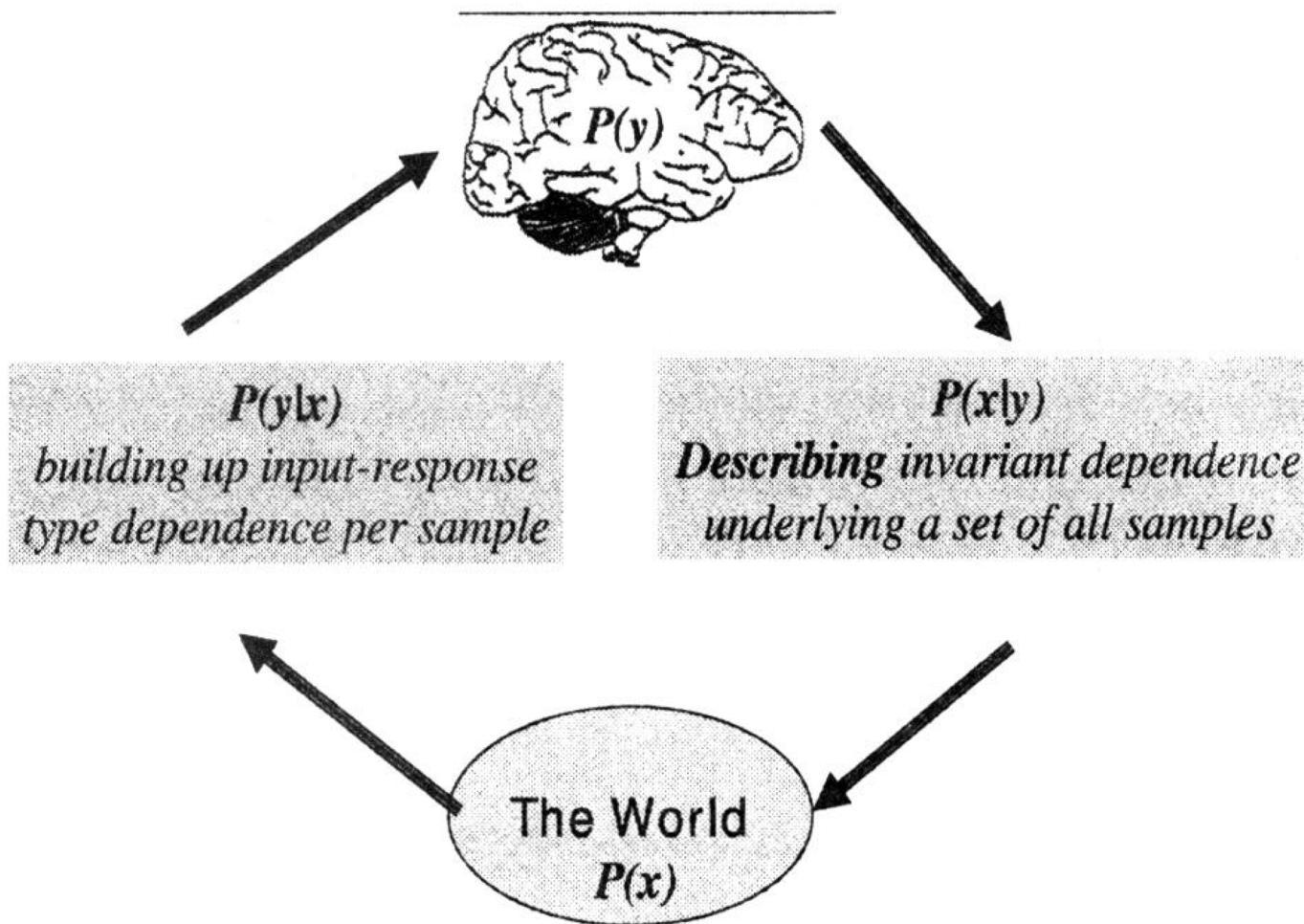
Mixture structures



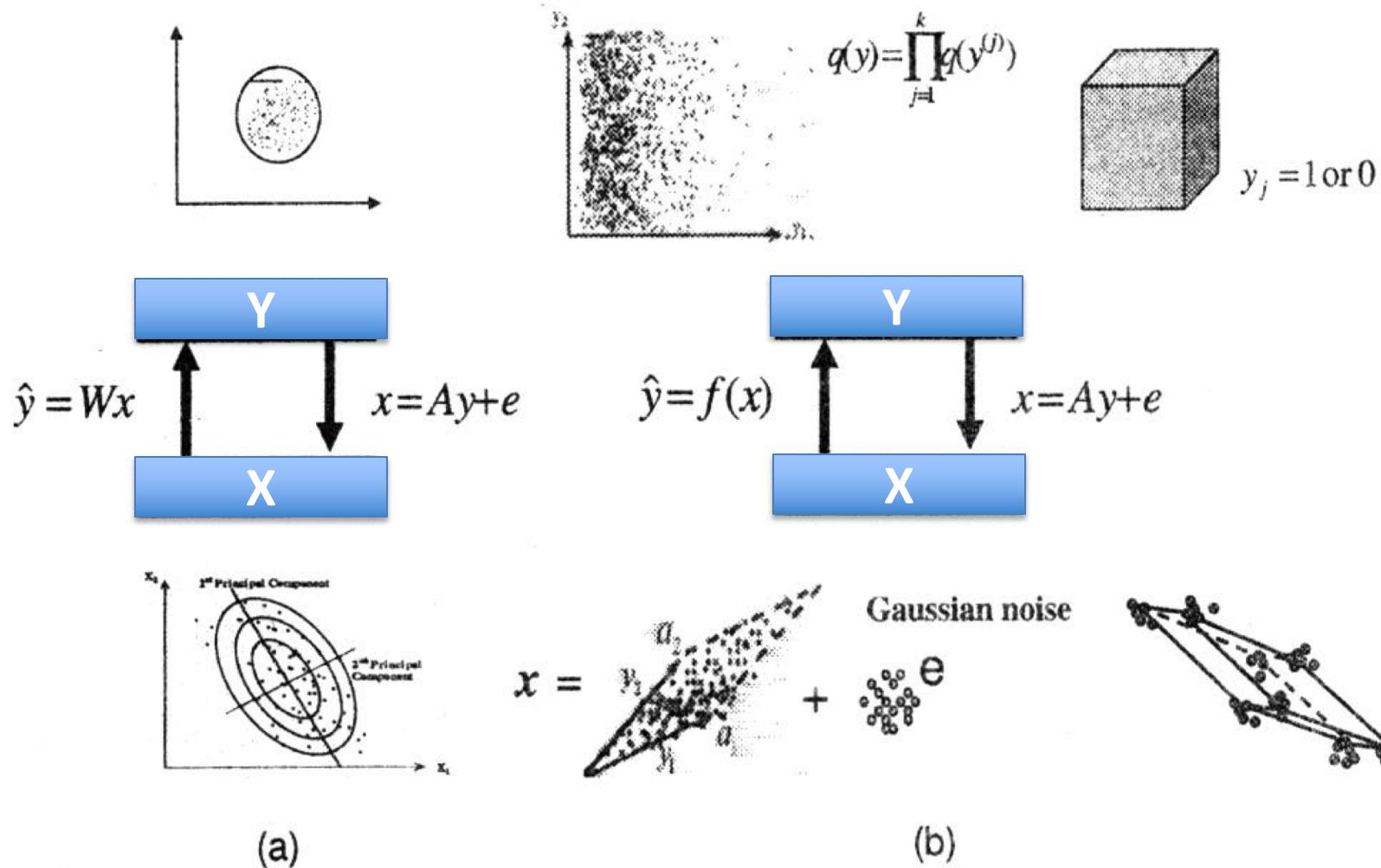
Outline

- **Fundamentals and challenges of learning**
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- **Bi-directional learning**
 - Inbound learning theory
 - Outbound learning theory
 - **Bi-directional architectures**
- **Early bi-directional deep learning**

Bi-directional architectures



Bi-directional latent structures

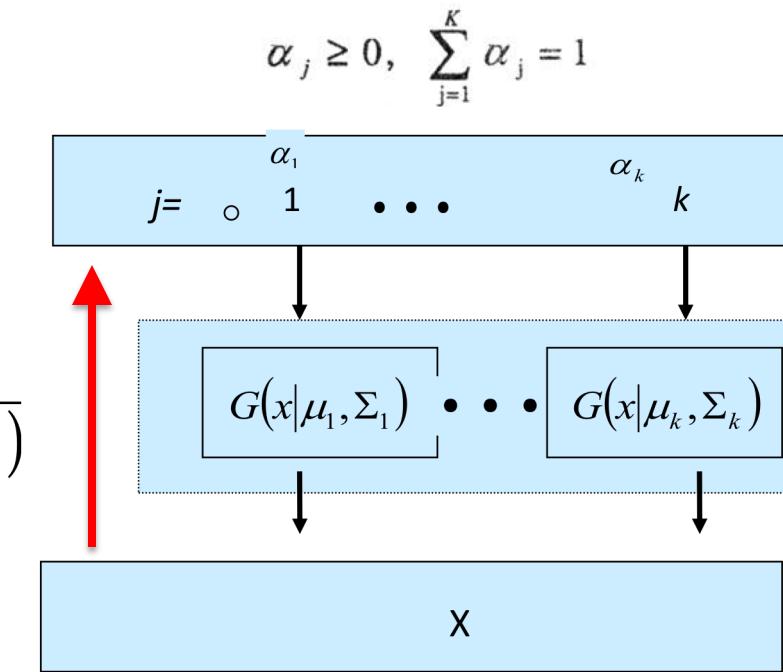
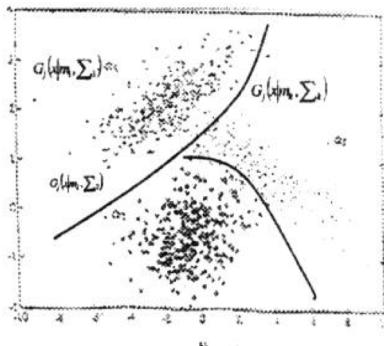


GMM in bi-directional architecture

Bayesian Inversion

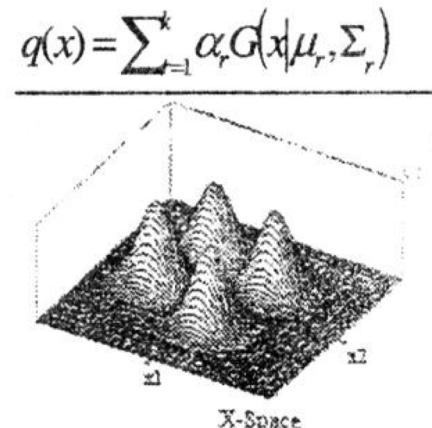
$$p(j | x_t) = \frac{\alpha_j G(x_t | m_j, \Sigma_j)}{\sum_r \alpha_r G(x_t | m_r, \Sigma_r)}$$

$$p_{\ell,t} = \begin{cases} p(\ell | x_t), & \text{no label} \\ \ell(x_t), & \text{labeled} \end{cases}$$



$$p(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x_i)$$

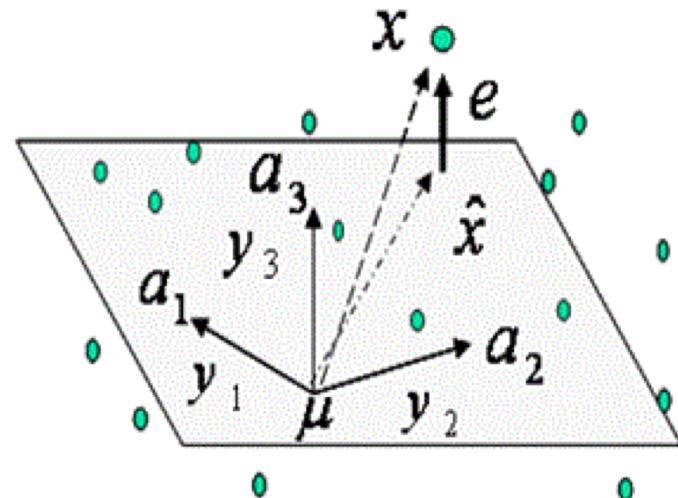
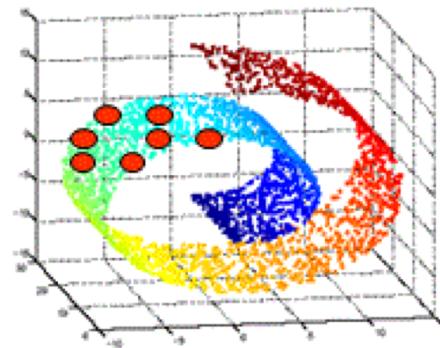
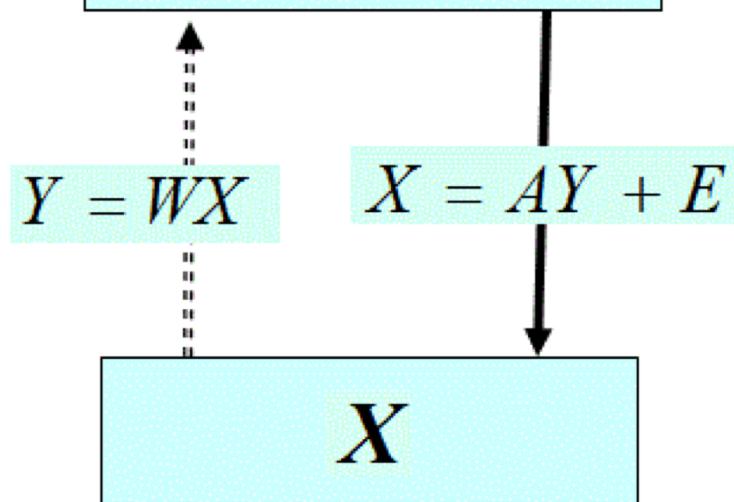
↑↑↑↑↑↑↑↑↑↑



Manifold factor analysis

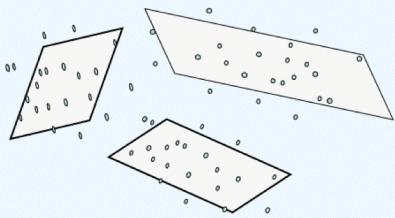
$$q(Y | \Lambda) = \frac{\exp\{-0.5 \text{Tr}[\Lambda^{-1} Y Y^T]\}}{Z(L, \Lambda)}$$

$$q(Y | \theta_y) = q(Y | \Lambda)$$



Given Laplacian \mathbf{L} of the nearest neighbor graph \mathbf{G} of X

Local factor analysis (LFA)

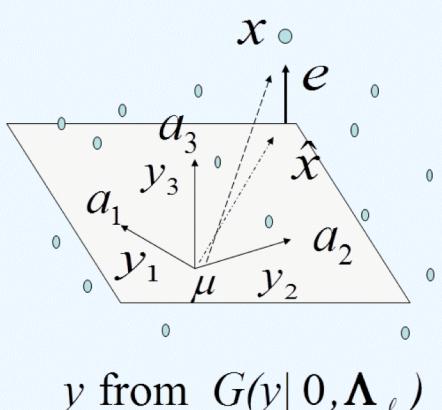


$$x = A_\ell y + \mu_\ell + e,$$

$$E[ey^T] = 0,$$

e from $G(e|0, \Sigma_\ell)$,

$$\ell = 1 \dots k.$$



(a)

$$q(\theta_\ell) = q(\alpha)q(A_\ell)q(\mu_\ell)q(\Lambda_\ell)q(\Sigma_\ell)$$

(1) Dirichlet distribution $q(\alpha) \propto \prod_\ell \alpha_\ell^{\gamma_\ell}$, $\gamma_\ell \geq -1$, particularly, it becomes Jeffrey prior if $\gamma_\ell = -0.5$.

(2) Jeffreys prior

$$q(\Lambda_\ell) \propto |\Lambda_\ell|^{-0.5}, \quad q(\Sigma_\ell) \propto |\Sigma_\ell|^{-0.5}.$$

(3) Gaussian - Jeffreys joint prior

For $A_\ell = [a_\ell^{(ij)}]$,

$$q(A_\ell, \{\eta_{\ell,ij}^A\}) \propto \prod_{ij} \{G(a_\ell^{(ij)}|0, \eta_{\ell,ij}^A) / \sqrt{\eta_{\ell,ij}^A}\}.$$

For $\mu_\ell = [\mu_\ell^{(1)}, \dots, \mu_\ell^{(d)}]^T$,

$$q(\mu_\ell, \{\eta_{\ell,i}^\mu\}) \propto \prod_i \{G(\mu_\ell^{(i)}|0, \eta_{\ell,i}^\mu) / \sqrt{\eta_{\ell,i}^\mu}\}.$$

(b)

Learning LFA

$$p(y, \ell | x) = p(y | x, \ell)p(\ell | x)$$

$$p(y | x_t, \ell) =$$

$$G(y | W_\ell(x_t - \mu_\ell), \Pi_\ell^{y-1} + P_\ell)$$

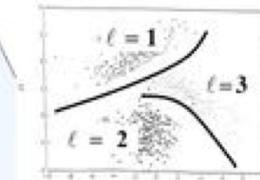
$$\begin{aligned} y_{t,\ell} &= \operatorname{argmin}_y \varepsilon(y, x_t, \theta_\ell) \\ &= \Pi_\ell^{y-1} A_\ell^T \Sigma_\ell^{-1} (x_t - \mu_\ell), \end{aligned}$$

$$q(\ell | x, \theta_\ell) = \frac{\exp[-\varepsilon_t(\theta_\ell) - \frac{1}{2} \ln \{|\Pi_\ell^y| / (2\pi)^{m_\ell}\}]}{\sum_j \exp[-\varepsilon_t(\theta_j) - \frac{1}{2} \ln \{|\Pi_j^y| / (2\pi)^{m_j}\}]}$$

$$p(\ell | x_t, \theta_\ell) = \begin{cases} q(\ell | x_t, \theta_\ell), & \ell \in C_\kappa(x_t), \\ 0, & \ell \notin C_\kappa(x_t). \end{cases}$$

$C_\kappa(x_t) = \{\ell : \text{for the first } \kappa$
 largest values of $L(x_t, \theta_\ell)$ by Eq.(6)\}.

$$p_{\ell,t} = p(\ell | x_t, \theta_\ell^{\text{old}}) [1 + \Delta_{\ell,t}(\theta_\ell^{\text{old}})] \text{ by Eq.(7).}$$



$$q(y, \ell) = q(y | \ell)q(\ell)$$

$$q(\ell) = \alpha_\ell = e^{c_\ell} / \sum_j e^{c_\ell}$$

$$q(y | \ell) = G(y | 0, \Lambda_\ell)$$

$$q(x | y, \ell) = G(x | A_\ell y + \mu_\ell, \Sigma_\ell)$$

$$\varepsilon(y, x_t, \theta_\ell) =$$

$$-\ln[\alpha_\ell G(y | 0, \Lambda_\ell) G(x | A_\ell y + \mu_\ell, \Sigma_\ell)],$$

$$\Pi_\ell^y = A_\ell^T \Sigma_\ell^{-1} A_\ell + \Lambda_\ell^{-1},$$

$$\varepsilon_t(\theta_\ell) = \varepsilon(y_{t,\ell}, x_t, \theta_\ell),$$

$$\delta_{t,\ell} = y_{t,\ell} - W_\ell(x_t - \mu_\ell).$$

$$\theta_\ell = \{A_\ell, \mu_\ell, \Lambda_\ell, \Sigma_\ell, W_\ell, c_\ell\},$$

$$\theta_\ell^{\text{new}} - \theta_\ell^{\text{old}} \propto \sum_{t \in T_\ell} G_t(\theta_\ell^{\text{old}}).$$

$$G_t(\theta_\ell) = p_{\ell,t} \nabla_{\theta_\ell} \varepsilon_t(\theta_\ell) + p(\ell | x_t, \theta_\ell^{\text{old}}) \nabla_{\theta_\ell} r_t(\theta_\ell),$$

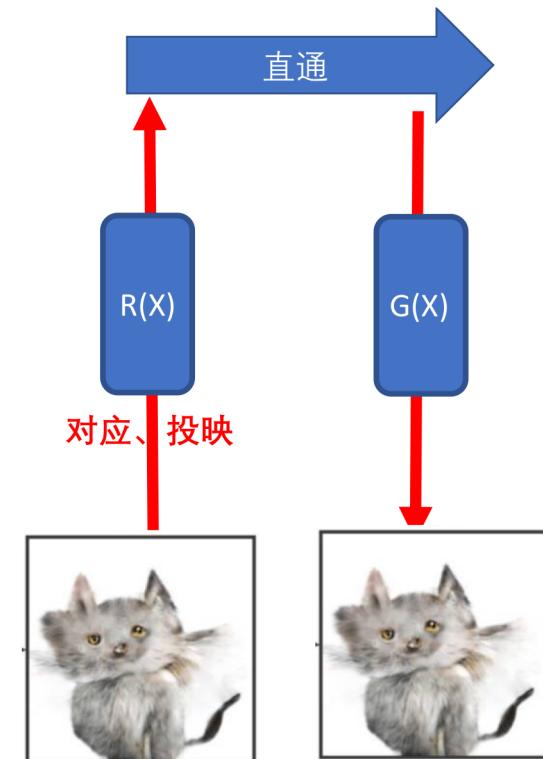
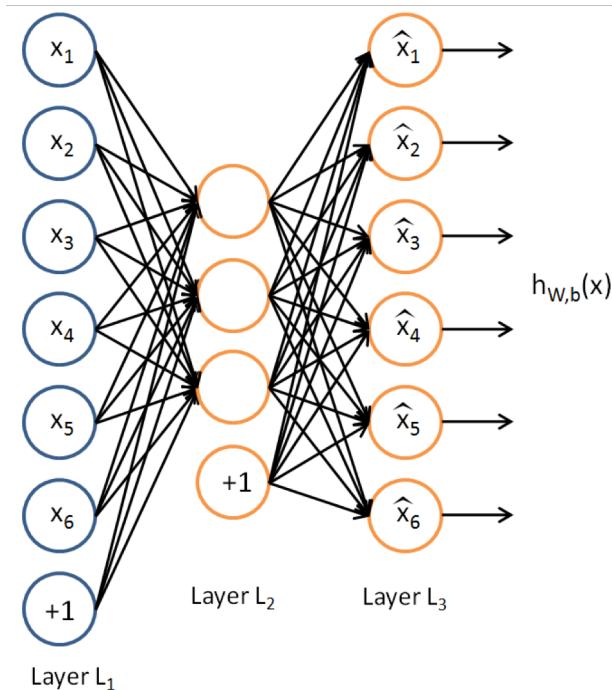
$$r_t(\theta_\ell) = \frac{1}{2} \delta_{t,\ell}^T \Pi_\ell^{y-1} \delta_{t,\ell} + \frac{1}{2} h^2 \text{Tr}[\Sigma_\ell^{-1}] - \ln q(\theta_\ell).$$

Outline

- Fundamentals and challenges of learning
 - Two types of intelligent abilities and learning
 - Three levels of inverse problems
- Bi-directional learning
 - Inbound learning theory
 - Outbound learning theory
 - Bi-directional architectures
- **Early bi-directional deep learning**

Autoencoder (AE)

- An **autoencoder** neural network is an unsupervised learning algorithm that applies backpropagation, setting the target values to be equal to the inputs.



The autoencoder tries to learn a function $h_{W,b}(x) \approx x$. In other words, it is trying to learn an approximation to the identity function, so as to output \hat{x} that is similar to x . The identity function seems a particularly trivial function to be trying to learn; but by placing constraints on the network, such as by limiting the number of hidden units, we can discover interesting structure about the data.

General architecture

Sparse autoencoder

Imposing sparsity on hidden units

Informally, we will think of a neuron as being “active” (or as “firing”) if its output value is close to 1, or as being “inactive” if its output value is close to 0. We would like to constrain the neurons to be inactive most of the time.³

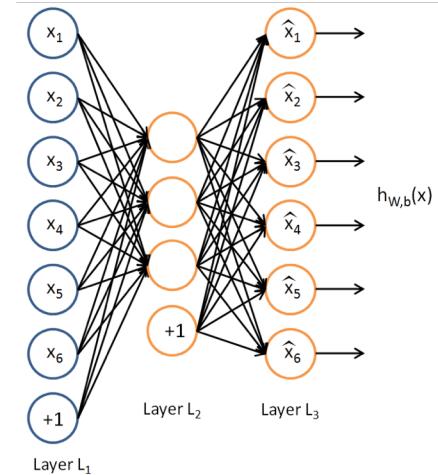
Recall that $a_j^{(2)}$ denotes the activation of hidden unit j in the autoencoder. However, this notation doesn’t make explicit what was the input x that led to that activation. Thus, we will write $a_j^{(2)}(x)$ to denote the activation of this hidden unit when the network is given a specific input x . Further, let

$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m [a_j^{(2)}(x^{(i)})]$$

be the average activation of hidden unit j (averaged over the training set). We would like to (approximately) enforce the constraint

$$\hat{\rho}_j = \rho,$$

where ρ is a **sparsity parameter**, typically a small value close to zero (say $\rho = 0.05$). In other words, we would like the average activation of each hidden neuron j to be close to 0.05 (say). To satisfy this constraint, the hidden unit’s activations must mostly be near 0.



Penalize the objective with sparsity constraints

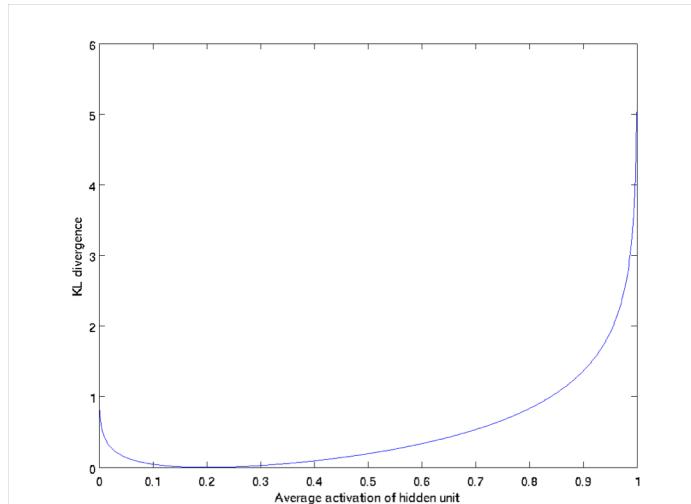
Using KL divergence to represent the sparsity constraint:

$$\sum_{j=1}^{s_2} \text{KL}(\rho || \hat{\rho}_j) = \sum_{j=1}^{s_2} \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}.$$

Average activation

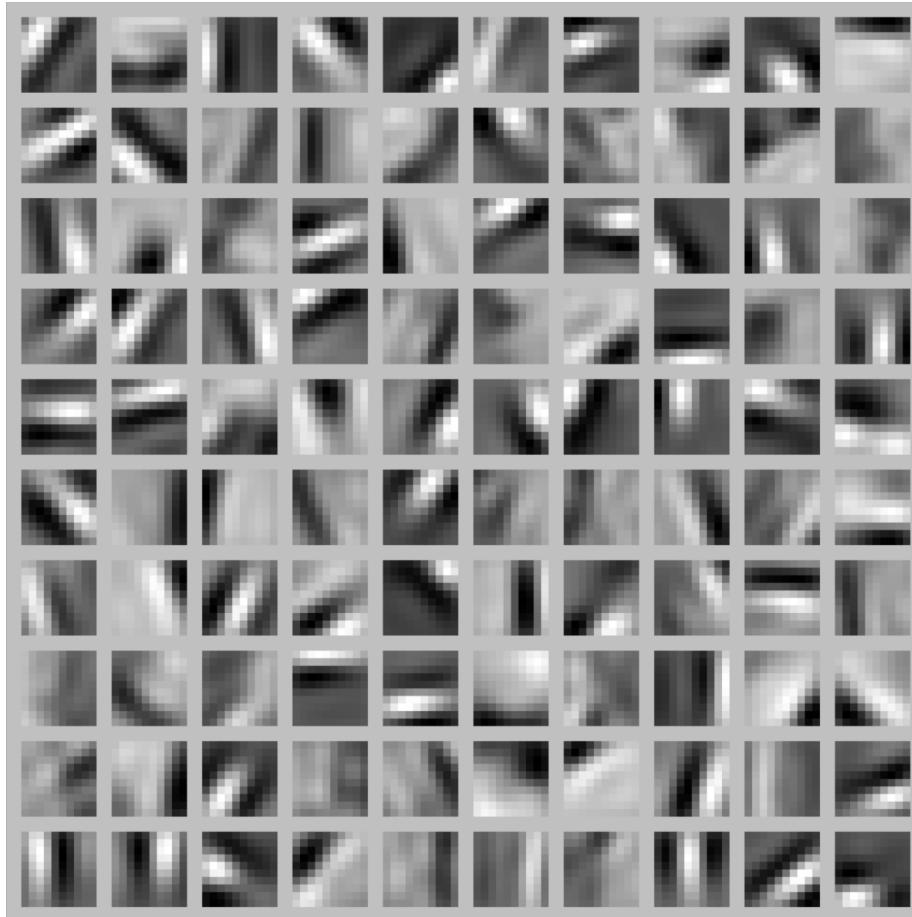
$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m [a_j^{(2)}(x^{(i)})]$$

This penalty function has the property that $\text{KL}(\rho || \hat{\rho}_j) = 0$ if $\hat{\rho}_j = \rho$, and otherwise it increases monotonically as $\hat{\rho}_j$ diverges from ρ . For example, in the figure below, we have set $\rho = 0.2$, and plotted $\text{KL}(\rho || \hat{\rho}_j)$ for a range of values of $\hat{\rho}_j$:



Visualize the weights: input->hidden

Input is image



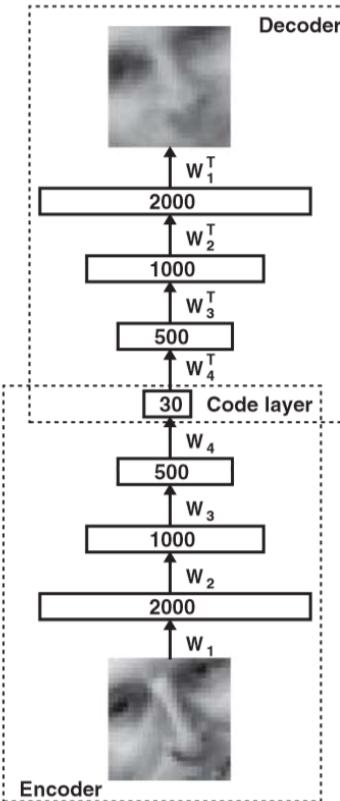
$$x_j = \frac{W_{ij}^{(1)}}{\sqrt{\sum_{j=1}^{100} (W_{ij}^{(1)})^2}}.$$

Each square in the figure above shows the (norm bounded) input image x that maximally activates one of 100 hidden units. We see that the different hidden units have learned to detect edges at different positions and orientations in the image.

早期双向深度学习：auto-encoder

Stacked RBMs

Hinton, et al (2006)

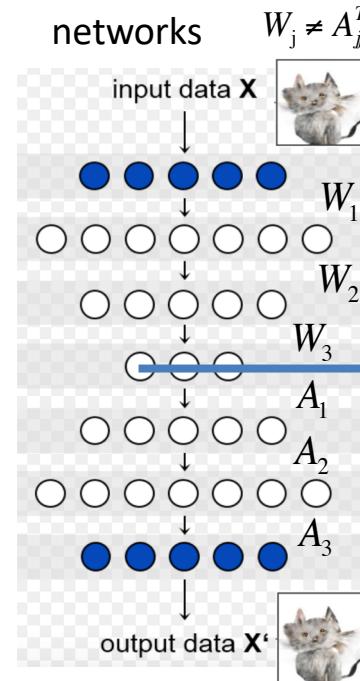


强制了参数对称性

$$W_j = A_j^T$$

Early auto-associative networks

$$W_j \neq A_j^T$$



Bourlard, H., Kamp, Y.
Auto-association by multilayer perceptrons
Biological cybernetics 59(4-5), 291-294 (1988)

Xu, Proc. IJCNN 91, Singapore, pp. 2362-2373
Neural Networks, vol. 6, pp. 627-648, 1993

Least Mean Square Error Reconstruction Principle
for Self-Organizing Neural-Nets

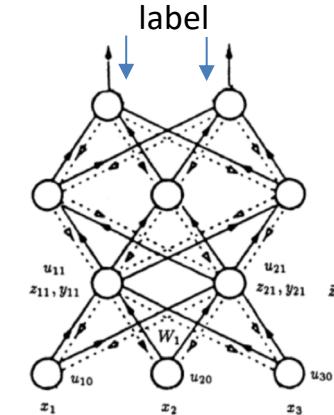
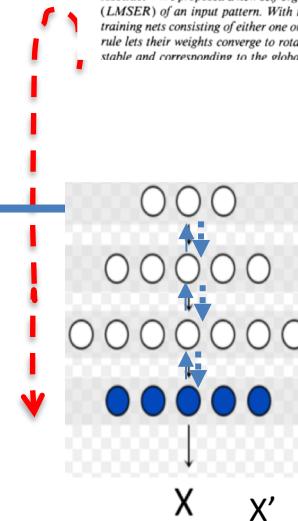
LEI XU

Peking University and Harvard University

(Received 29 July 1991; revised and accepted 16 October 1992)

Abstract—We proposed a new self-organizing net based on the principle of Least Mean Square Error Reconstruction (LMSER) of an input pattern. With this principle, a local learning rule called LMSER is naturally obtained for training nets consisting of either one or several layers. We proved that for one layer with n_i linear units, the LMSER rule lets their weights converge to rotations of the data's first n_i principal components. These converged points are stable and corresponding to the global minimum in the Mean Square Error (MSE) landscape which have many

多层LMSER自组织学习



强制了上下两半部的参数对称性
和神经元强度的对称性 $W_j = A_j^T$

Reducing the Dimensionality of Data with Neural Networks

Science

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

28 JULY 2006 VOL 313 SCIENCE www.sciencemag.org

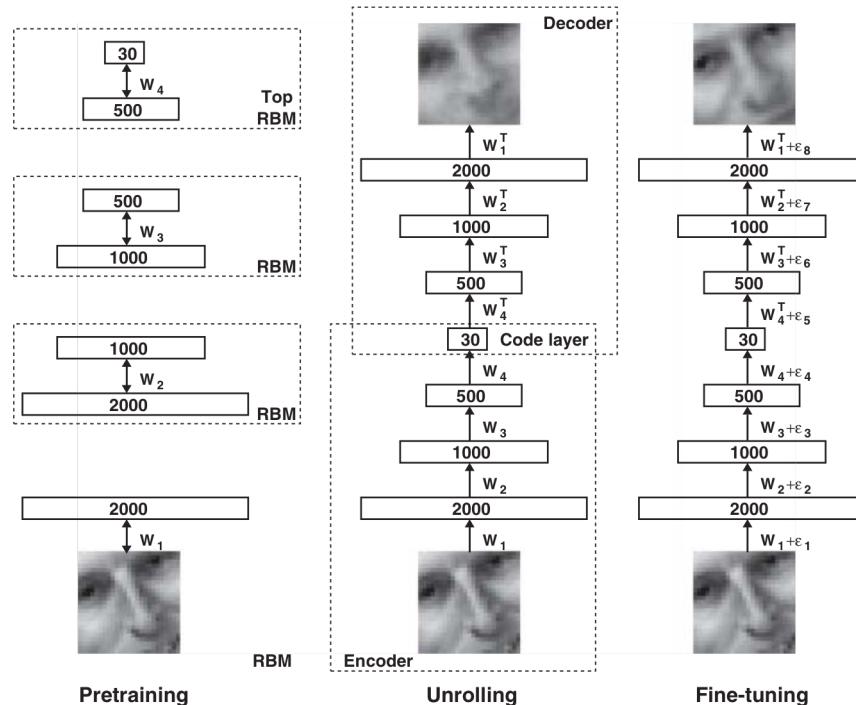


Fig. 2. (A) Top to bottom: Random samples of curves from the test data set; reconstructions produced by the six-dimensional deep autoencoder; reconstructions by “logistic PCA” (8) using six components; reconstructions by logistic PCA and standard PCA using 18 components. The average squared error per image for the last four rows is 1.44, 7.64, 2.45, 5.90. (B) Top to bottom: A random test image from each class; reconstructions by the 30-dimensional autoencoder; reconstructions by 30-dimensional logistic PCA and standard PCA. The average squared errors for the last three rows are 3.00, 8.01, and 13.87. (C) Top to bottom: Random samples from the test data set; reconstructions by the 30-dimensional autoencoder; reconstructions by 30-dimensional PCA. The average squared errors are 126 and 135.

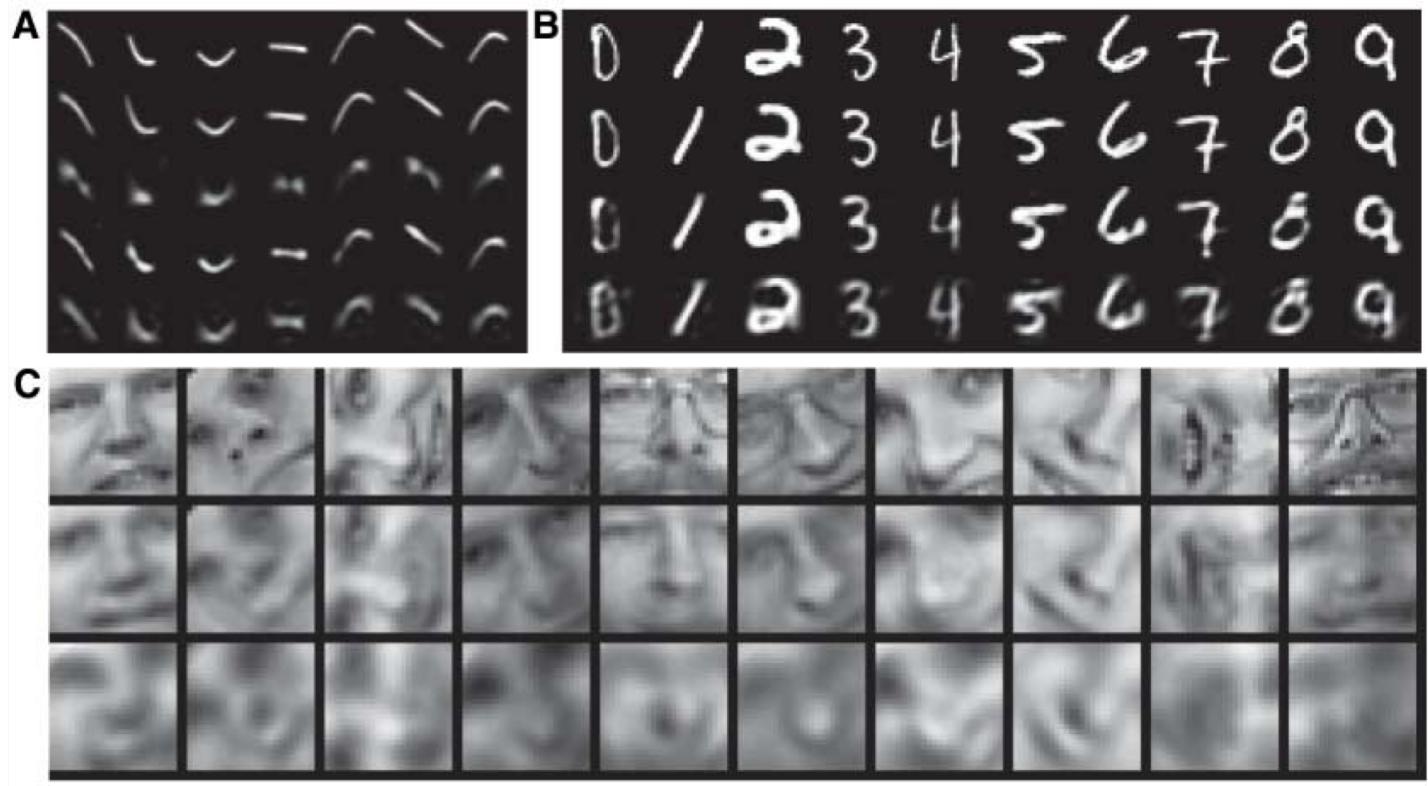
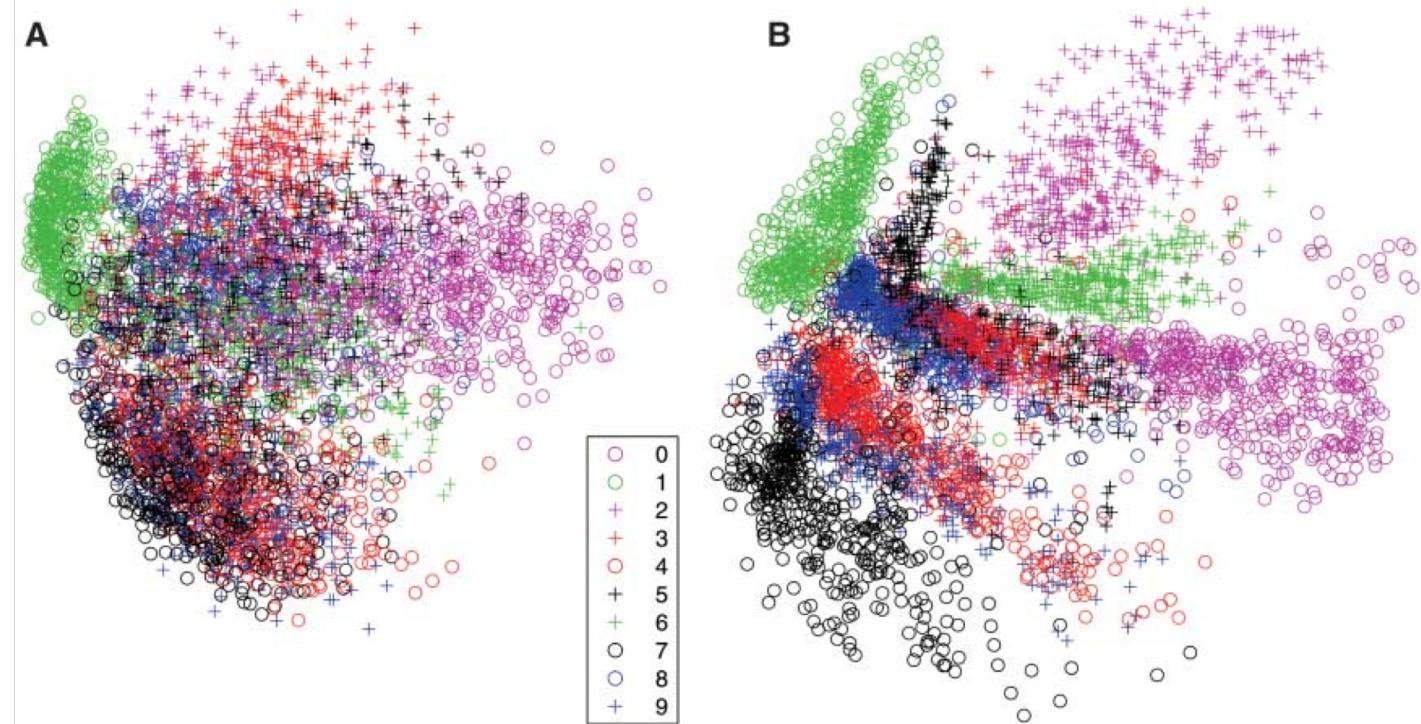


Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).



Thank you!