

深度学习-神经网络基础知识

史春奇

2017 年

目录

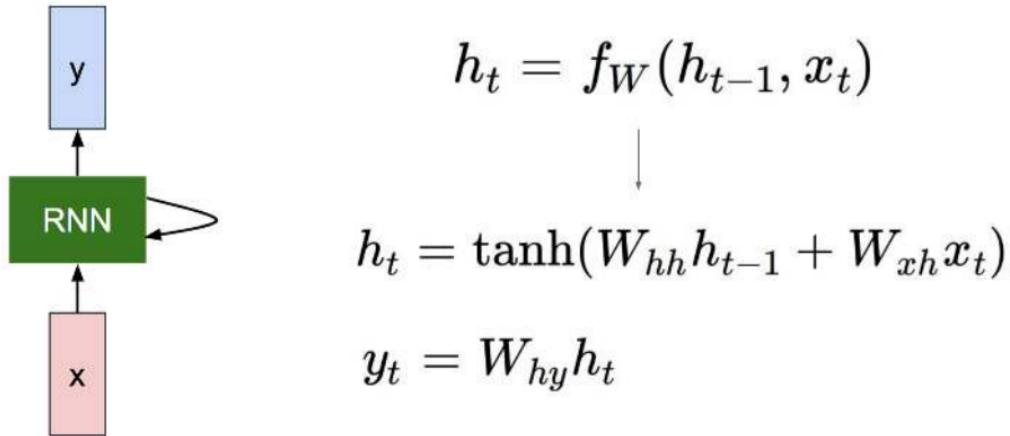
1 递归神经网络 RNN	4
2 玻尔兹曼机 BM	69
3 自动编码器 AE	157
4 生成对抗网络 GAN	188

第二部分（下）

- DBN 和 DBM 的特性
 - RBM 如何训练？
 - DBN 如何 Tune ？
 - DBM 和 DBN 相比？
- 自动编码器的理解
 - AE 和 PCA 如何联系？
 - 为什么要有 Sparse AE？
 - VAE 的优点？
- 生成对抗网络的理解
 - GAN 和 VAE 的不同？
 - GAN ZOO 的关系？

1 递归神经网络 RNN

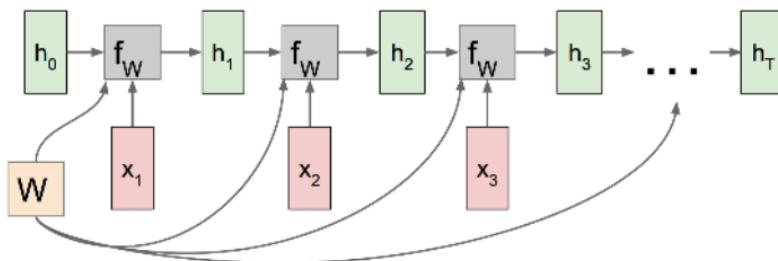
- 基本的 RNN 结构



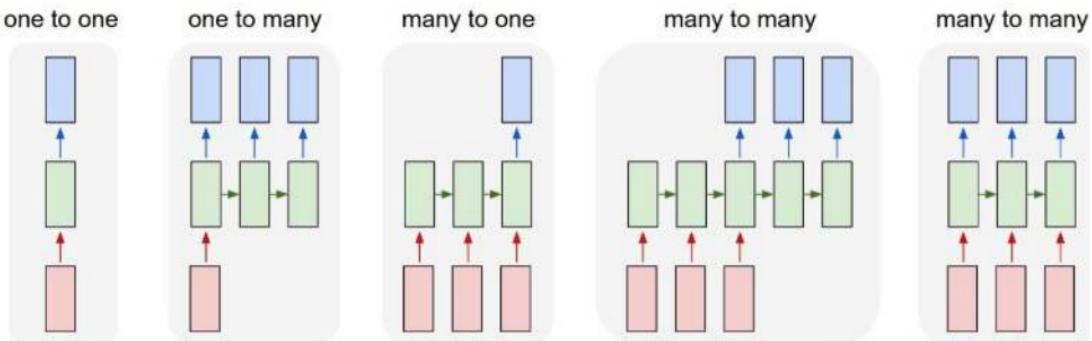
- 基本的 RNN 结构，按时间展开

RNN: Computational Graph

Re-use the same weight matrix at every time-step



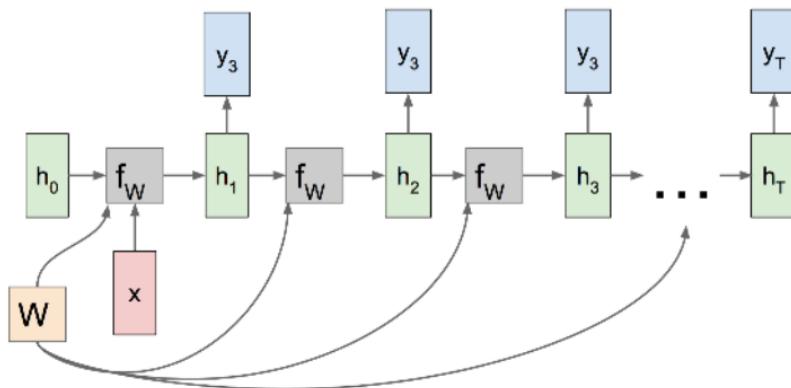
- 1 到多：图像到句子



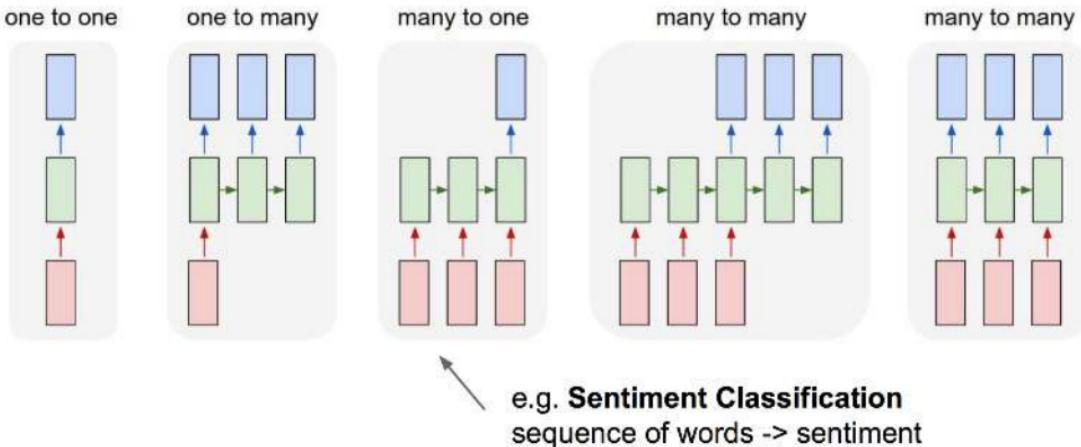
e.g. **Image Captioning**
image -> sequence of words

- 1 到多，按时间展开

RNN: Computational Graph: One to Many

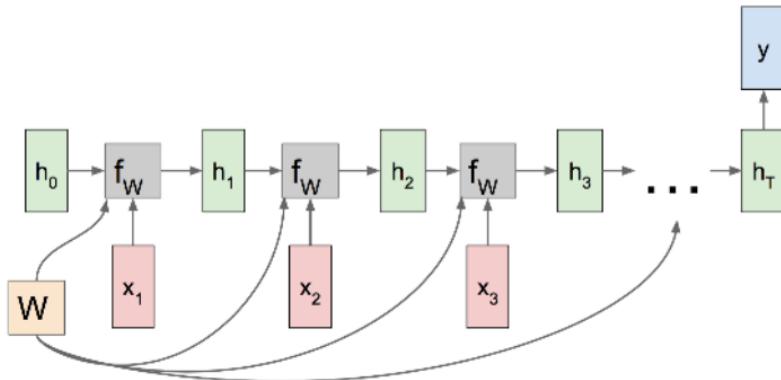


- 多到 1: 句子到情感

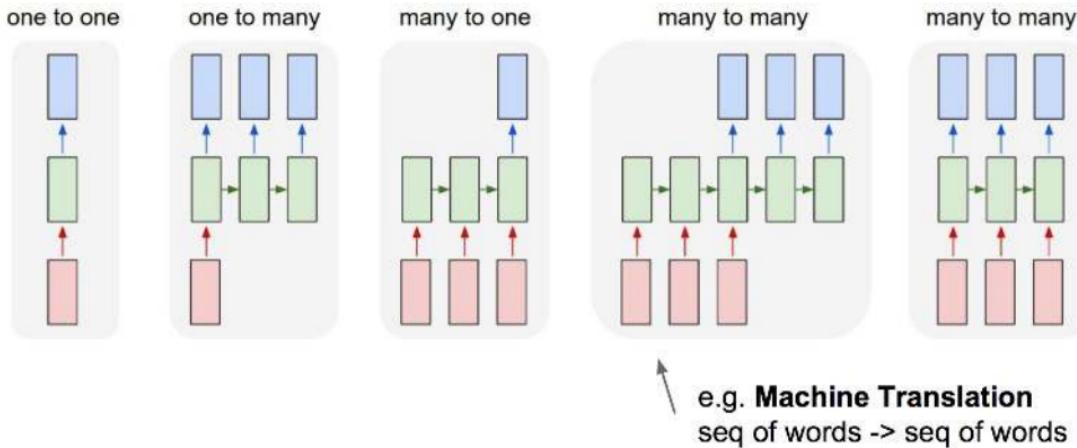


- 多到 1，按时间展开

RNN: Computational Graph: Many to One

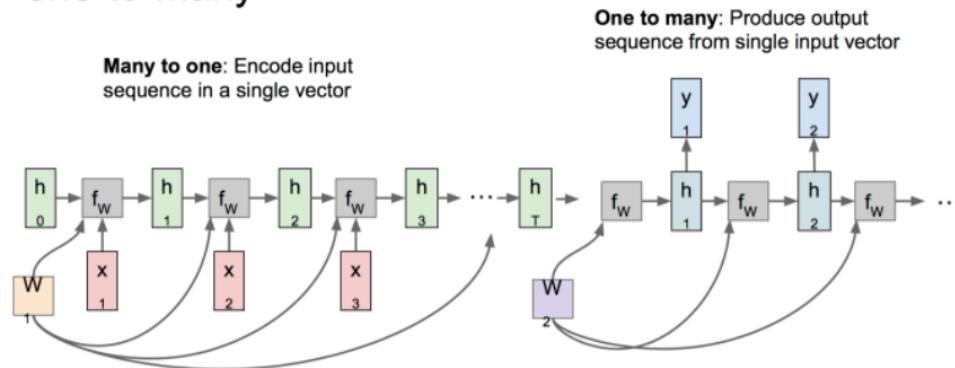


- 多到多：seq2seq 机器翻译

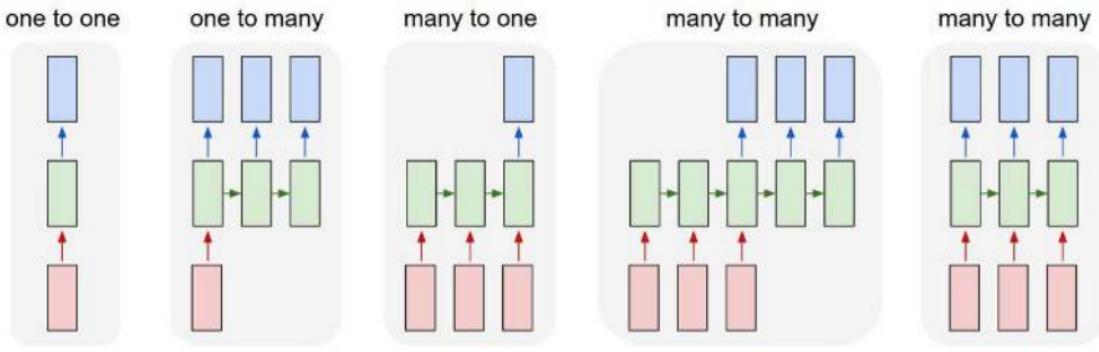


- 多到多 (seq2seq), 按时间展开

Sequence to Sequence: Many-to-one + one-to-many

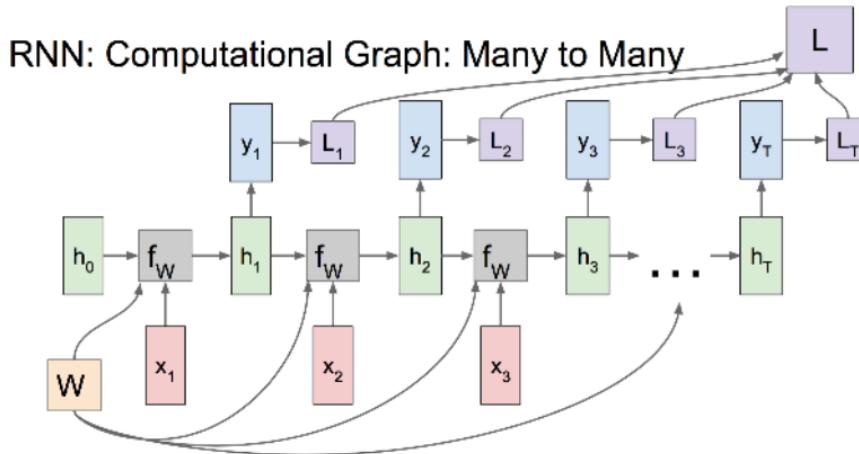


- 多到多：视频基于帧的分类

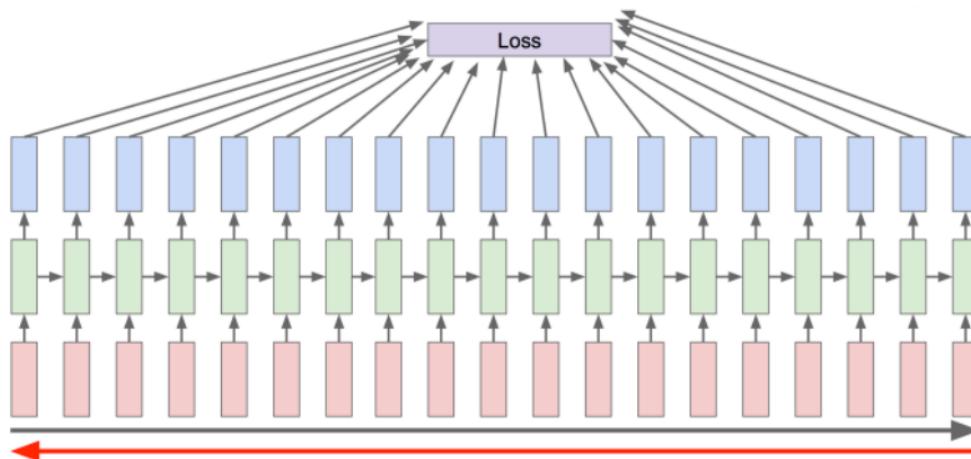


e.g. **Video classification on frame level**

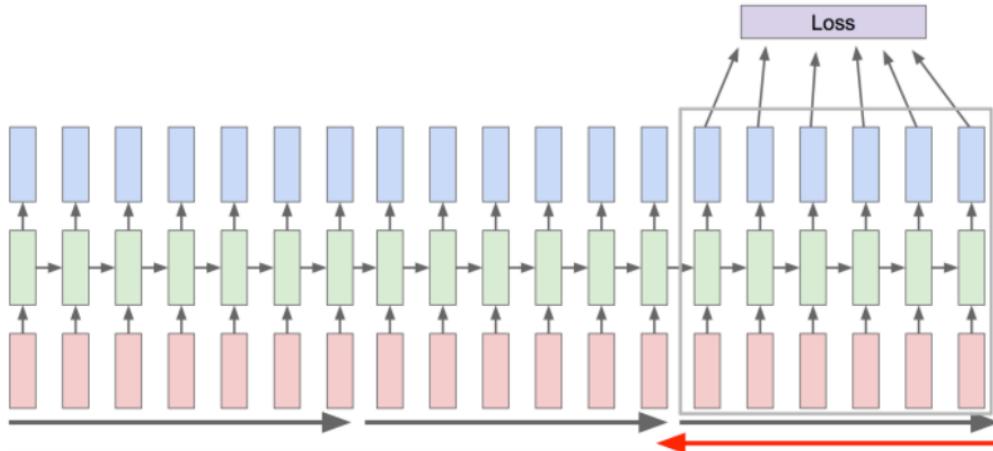
- 多到多，按时间展开



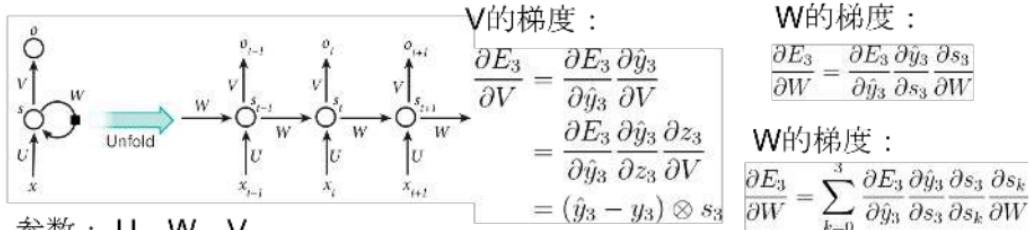
- BPTT: backpropagation through time



- Truncated BPTT: 划分成少量步骤



- Truncated BPTT



参数： U, W, V

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

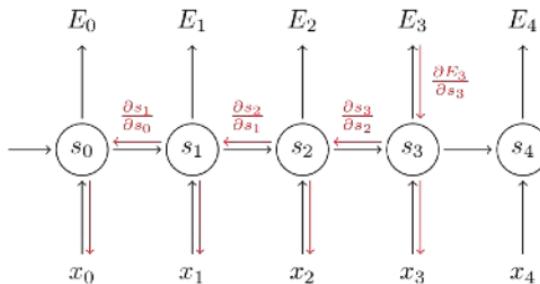
$$o_t = \text{softmax}(Vs_t)$$

Log 损失：

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

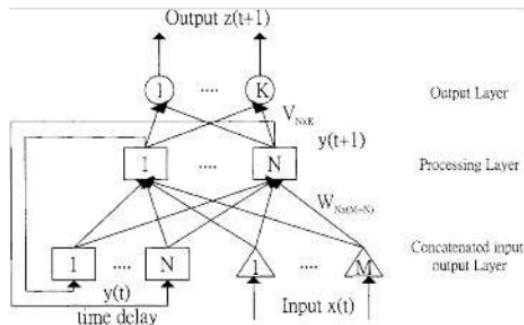
$$= -\sum_t y_t \log \hat{y}_t$$



- Truncated BPTT

```
def bptt(self, x, y):  
    T = len(y)  
    # Perform forward propagation  
    o, s = self.forward_propagation(x)  
    # We accumulate the gradients in these variables  
    dLdU = np.zeros(self.U.shape)  
    dLdV = np.zeros(self.V.shape)  
    dLdW = np.zeros(self.W.shape)  
    delta_o = o  
    delta_o[np.arange(len(y)), y] -= 1.  
    # For each output backwards...  
    for t in np.arange(T)[::-1]:  
        dLdV += np.outer(delta_o[t], s[t].T)  
        # Initial delta calculation: dL/dz  
        delta_t = self.V.T.dot(delta_o[t]) * (1 - (s[t]**2))  
        # Backpropagation through time (for at most self.bptt_truncate steps)  
        for bptt_step in np.arange(max(0, t-self.bptt_truncate), t+1)[::-1]:  
            # print "Backpropagation step t=%d bptt step=%d" % (t, bptt_step)  
            # Add to gradients at each previous step  
            dLdW += np.outer(delta_t, s[bptt_step-1])  
            dLdU[:,x[bptt_step]] += delta_t  
            # Update delta for next step dL/dz at t-1  
            delta_t = self.W.T.dot(delta_t) * (1 - s[bptt_step-1]**2)  
    return [dLdU, dLdV, dLdW]
```

- RTRL real time recurrent learning



$$p_{ij}^k(t) = \frac{\partial y_k(t)}{\partial w_{ij}}$$

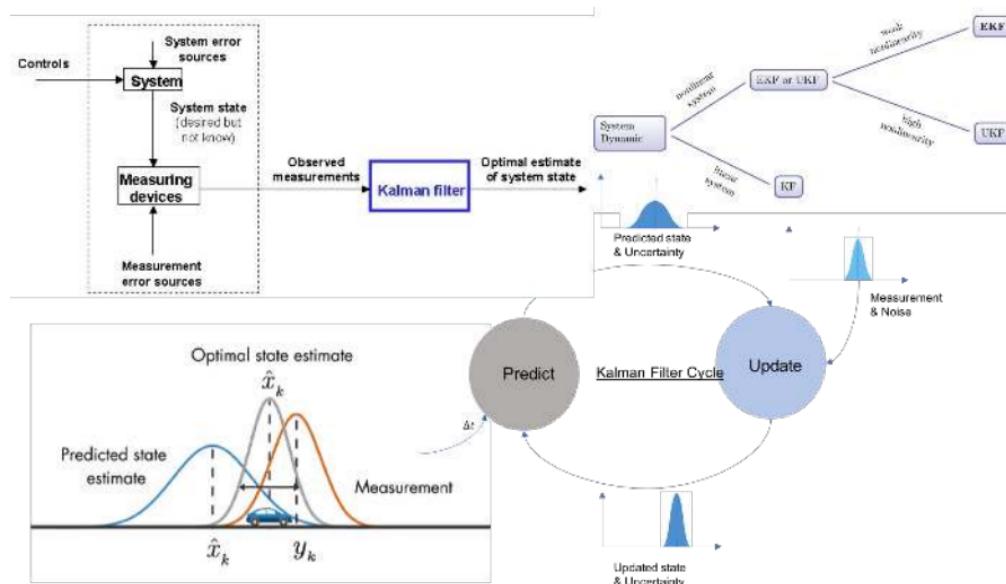
$$\frac{\partial J(t)}{\partial w_{ij}} = \sum_{k \in U} e_k(t) p_{ij}^k(t)$$

$$p_{ij}^k(t+1) = f'_k(s_k(t+1)) \left[\sum_{l \in U} w_{kl} p_{lj}^l(t) + \delta_{ik} x_j(t) \right]$$

$$p_{ij}^k(t_0) = \frac{\partial y_k(t_0)}{\partial w_{ij}} = 0 \quad \delta_{ik} = 1 \text{ if } i = k \text{ and } 0 \text{ otherwise}$$

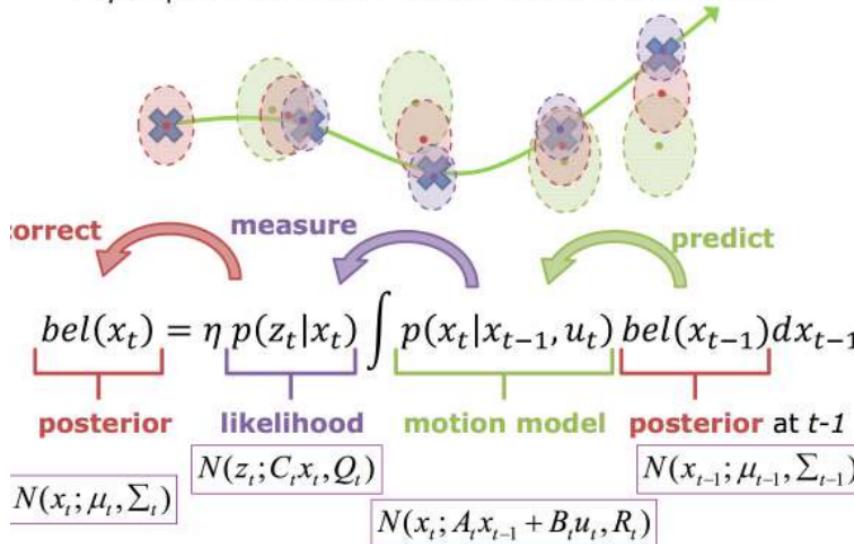
和 BPTT 的等价，反向和前向的区别！

- ELF extended Kalman filter



- Kalman filter

- Bayes update rule + linear Gaussian models → Kalman filter



- Kalman filter

1. **Kalman_filter(μ_{t-1} , Σ_{t-1} , u_t , z_t):**

2. Prediction:

Requirements:

$$\bar{\mu}_t = A_t \mu_{t-1} + B_t u_t$$

- Initial belief is Gaussian

$$\bar{\Sigma}_t = A_t \Sigma_{t-1} A_t^T + R_t$$

- Linear Gaussian state-space model

5. Correction:

$$K_t = \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1}$$

$$\mu_t = \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t)$$

$$\Sigma_t = (I - K_t C_t) \bar{\Sigma}_t$$

9. Return μ_t , Σ_t

- Extended Kalman filter

1. **Extended_Kalman_filter**($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$):

2. Prediction:

$$3. \bar{\mu}_t = g(u_t, \mu_{t-1})$$

$$\bar{\mu}_t = A_t \mu_{t-1} + B_t u_t$$

$$4. \bar{\Sigma}_t = G_t \Sigma_{t-1} G_t^T + R_t$$

$$\bar{\Sigma}_t = A_t \Sigma_{t-1} A_t^T + R_t$$

5. Correction:

$$6. K_t = \bar{\Sigma}_t H_t^T (H_t \bar{\Sigma}_t H_t^T + Q_t)^{-1}$$

$$K_t = \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1}$$

$$7. \mu_t = \bar{\mu}_t + K_t (z_t - h(\bar{\mu}_t))$$

$$\mu_t = \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t)$$

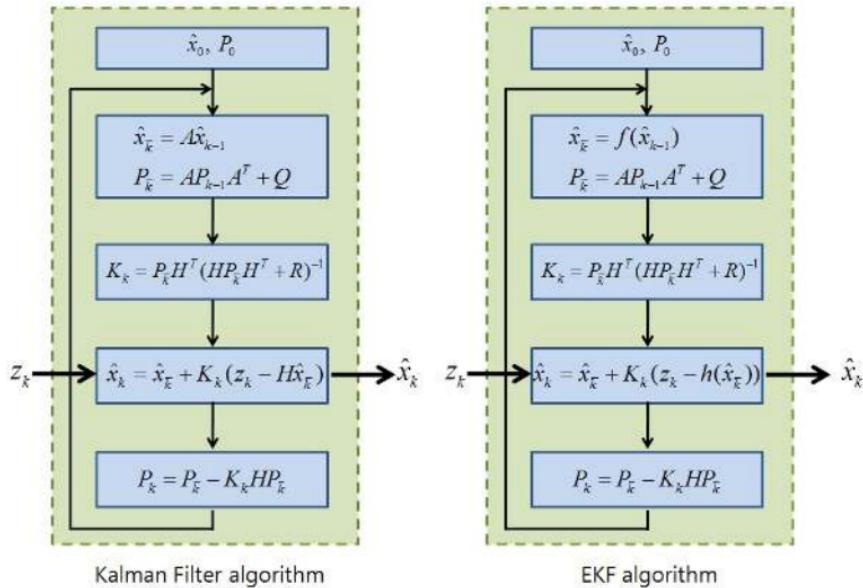
$$8. \Sigma_t = (I - K_t H_t) \bar{\Sigma}_t$$

$$\Sigma_t = (I - K_t C_t) \bar{\Sigma}_t$$

9. Return μ_t, Σ_t

$$H_t = \frac{\partial h(\bar{\mu}_t)}{\partial x_t} \quad G_t = \frac{\partial g(u_t, \mu_{t-1})}{\partial x_{t-1}}$$

- Extended Kalman filter



- Extended Kalman filter with Noise

1. **Extended_Kalman_filter**($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$):

2. Prediction:

$$\bar{\mu}_t = g(u_t, \mu_{t-1})$$

$$\bar{\Sigma}_t = G_t \Sigma_{t-1} G_t^T + \boxed{W_t R_t W_t^T}$$

$$G_t = \frac{\partial g(u_t, \mu_{t-1}, 0)}{\partial x_{t-1}}$$

5. Correction:

$$6. K_t = \bar{\Sigma}_t H_t^T (H_t \bar{\Sigma}_t H_t^T + \boxed{V_t Q_t V_t^T})^{-1}$$

$$W_t = \frac{\partial g(u_t, \mu_{t-1}, \varepsilon_t)}{\partial \varepsilon_t}$$

$$7. \mu_t = \bar{\mu}_t + K_t (z_t - h(\bar{\mu}_t))$$

$$H_t = \frac{\partial h(\bar{\mu}_t, 0)}{\partial x_t}$$

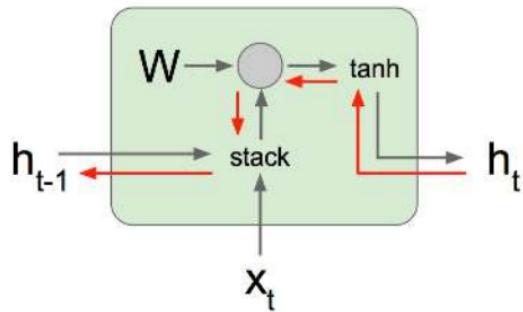
$$8. \Sigma_t = (I - K_t H_t) \bar{\Sigma}_t$$

9. Return μ_t, Σ_t

$$V_t = \frac{\partial h(\bar{\mu}_t, \delta)}{\partial \delta_t}$$

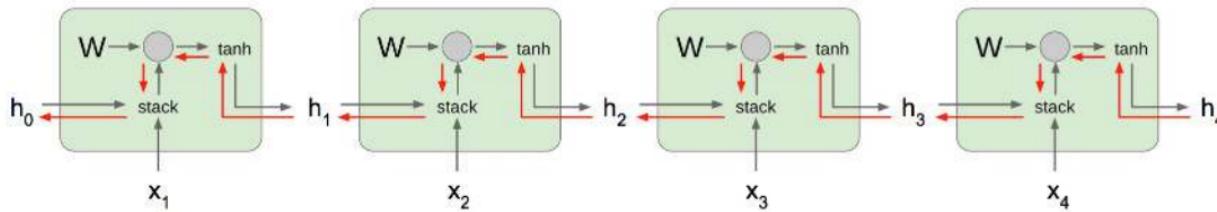
- Vanilla RNN 矩阵表示

Backpropagation from h_t
to h_{t-1} multiplies by W
(actually W_{hh}^T)



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

- Vanilla RNN 的梯度消失、爆炸

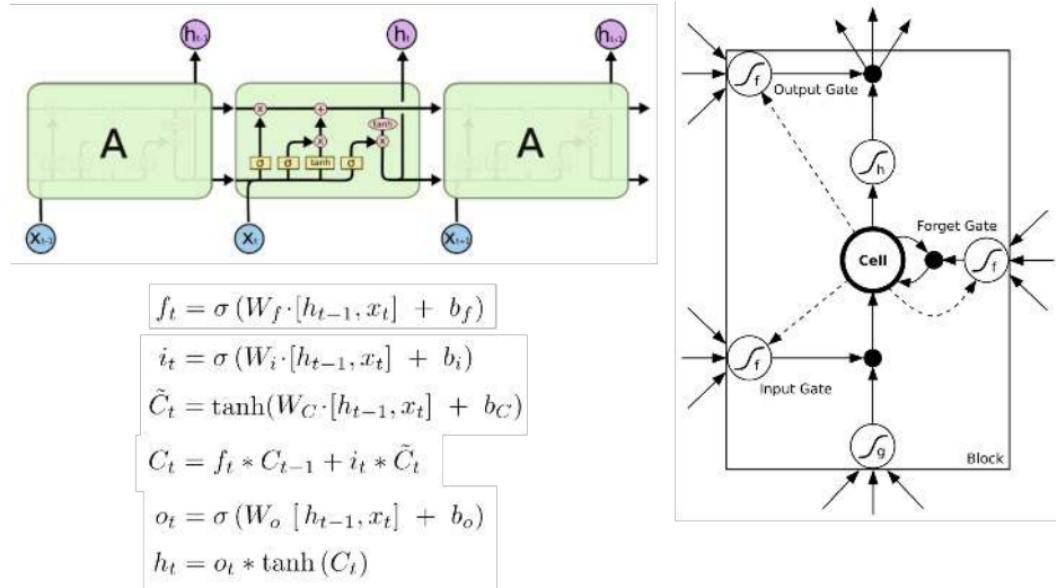


Computing gradient of h_0 involves many factors of W (and repeated \tanh)

Largest singular value > 1 :
Exploding gradients

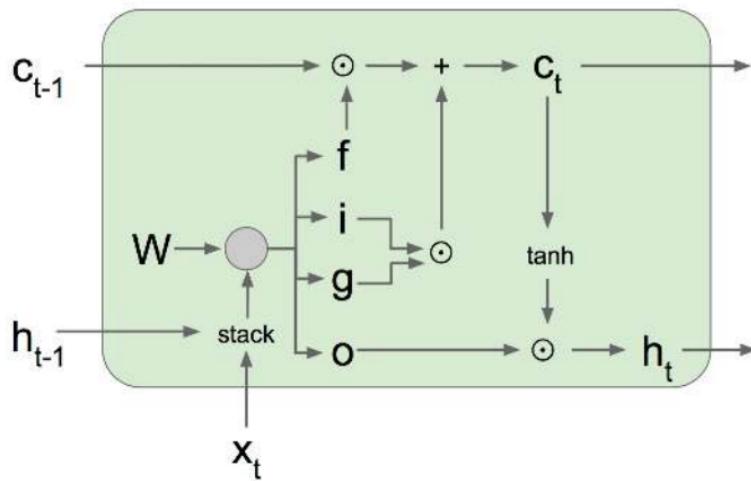
Largest singular value < 1 :
Vanishing gradients

- Long Short Term Memory



- LSTM 矩阵表示

[Hochreiter et al., 1997]



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

- LSTM 矩阵表示

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

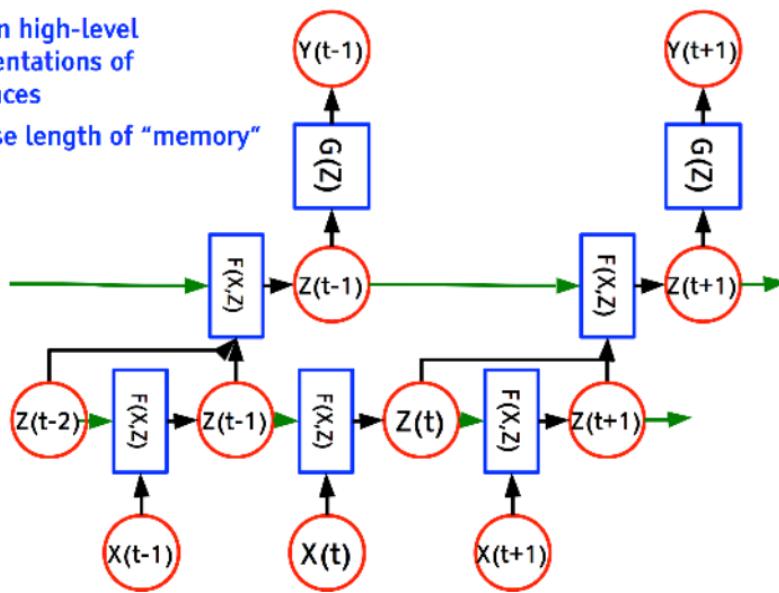
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

- LSTM 记忆单元

- To learn high-level representations of sequences
- Increase length of "memory"



- LSTM 记忆单元

- Some state variable are constrained to evolve slowly (matrix near unity)
 - Slow state used as context for faster state.
- "Learning Longer Memory in Recurrent Neural Networks",
 - Tomas Mikolov, Armand Joulin, Sumit Chopra, Michael Mathieu, Marc'Aurelio Ranzato [arxiv:1412.7753]

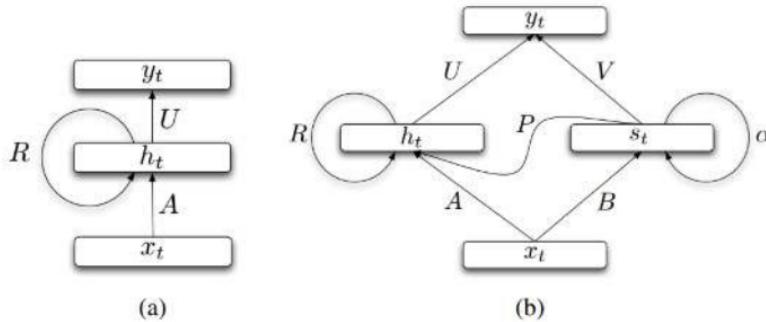
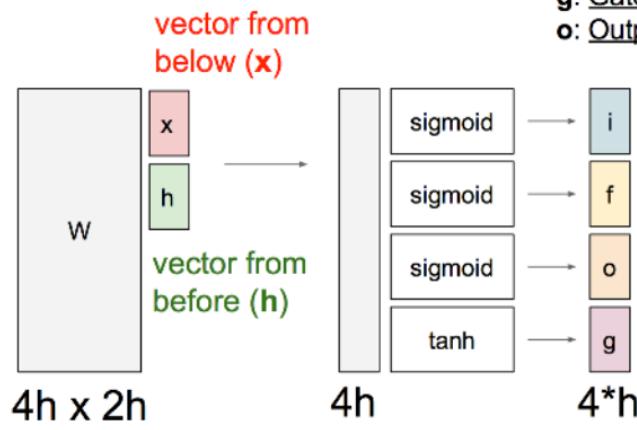


Figure 1: (a) Simple recurrent network. (b) Recurrent network with context features.

- LSTM 如何学习

[Hochreiter et al., 1997]



- f: Forget gate, Whether to erase cell
- i: Input gate, whether to write to cell
- g: Gate gate (?), How much to write to cell
- o: Output gate, How much to reveal cell

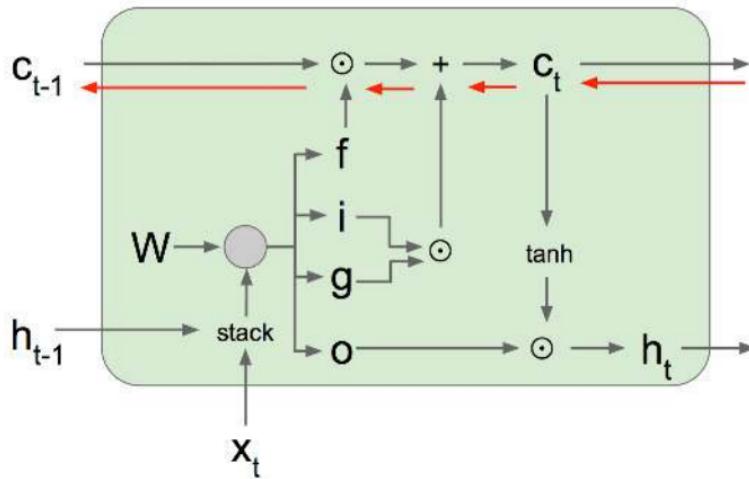
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

- LSTM 梯度反传

[Hochreiter et al., 1997]



Backpropagation from c_t to c_{t-1} only elementwise multiplication by f , no matrix multiply by W

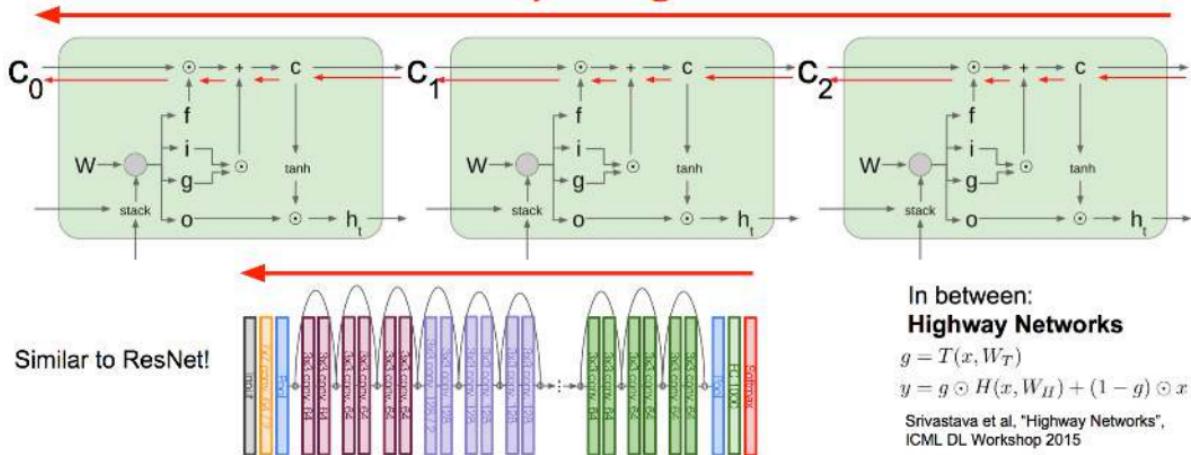
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

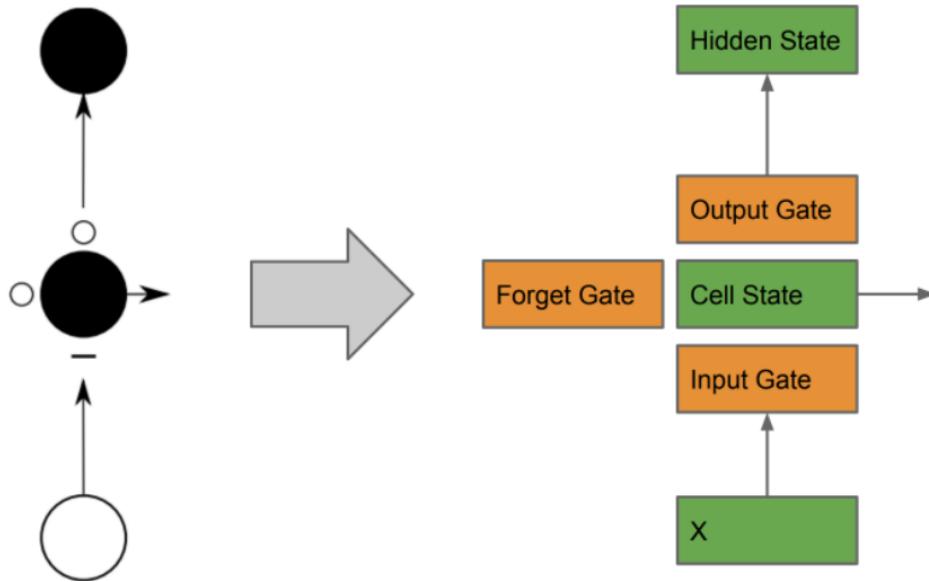
$$h_t = o \odot \tanh(c_t)$$

- LSTM 梯度反传

Uninterrupted gradient flow!

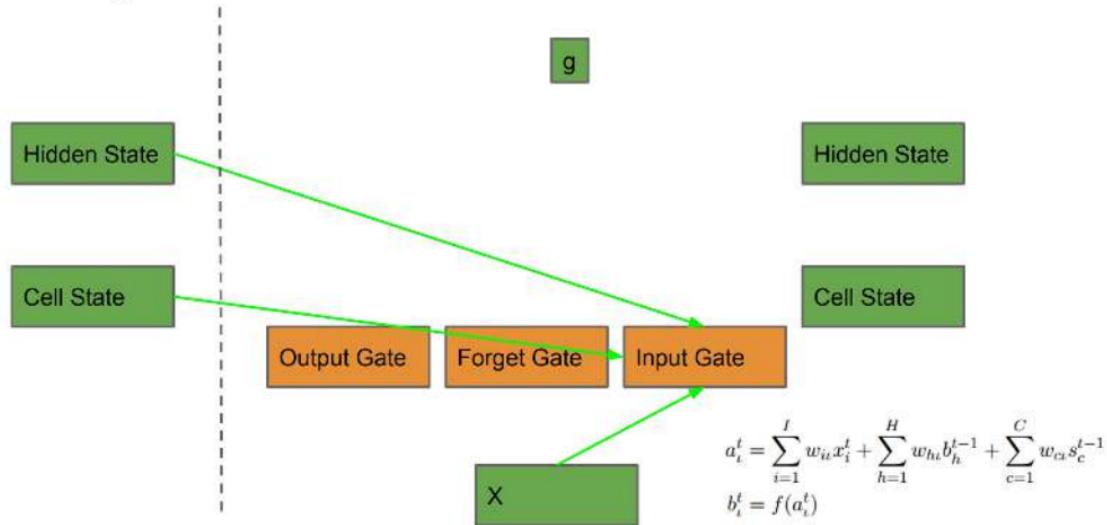


- LSTM 正向传播



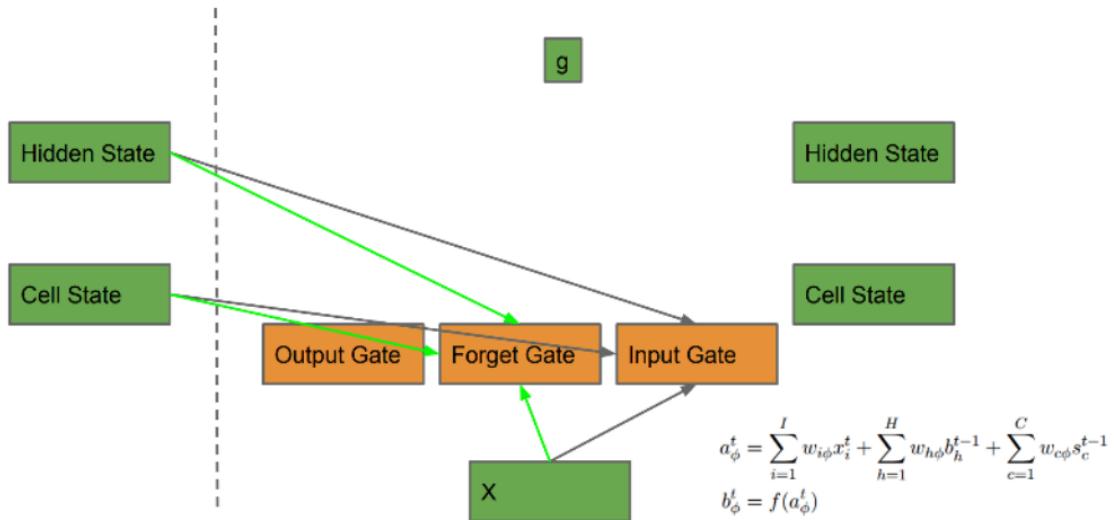
- LSTM 正向传播：输入门

1. Input Gate



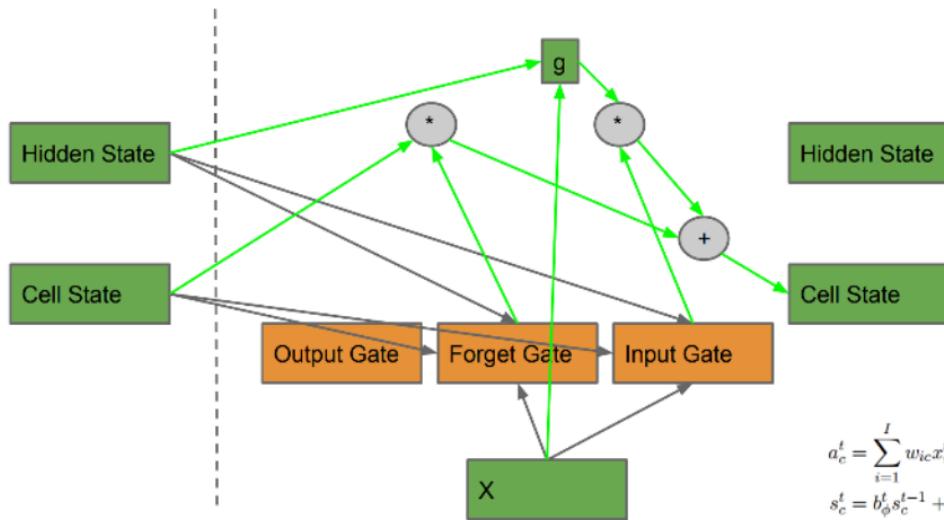
- LSTM 正向传播：忘记门

2. Forget Gate



- LSTM 正向传播：记忆单元 Cell 状态

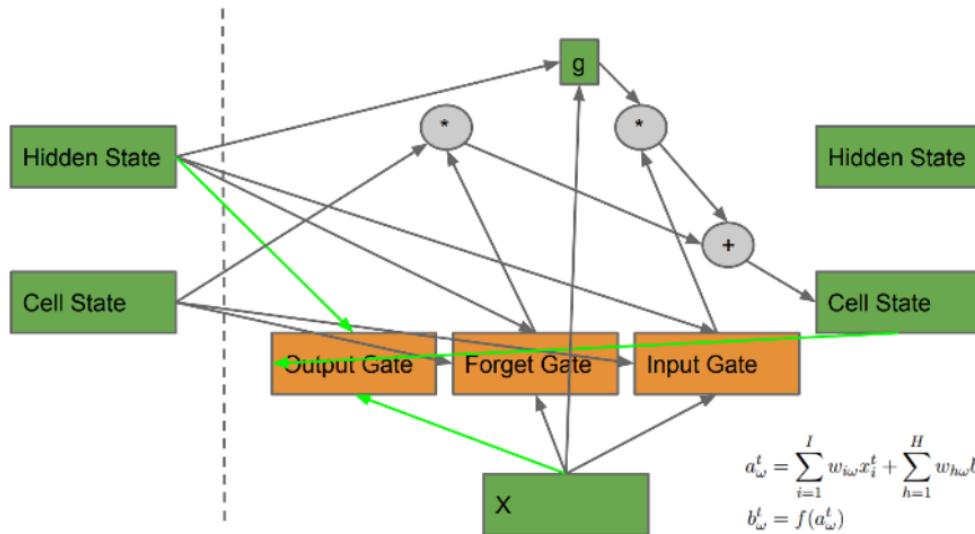
3. Cell State



$$a_c^t = \sum_{i=1}^I w_{ic} x_i^t + \sum_{h=1}^H w_{hc} b_h^{t-1}$$
$$s_c^t = b_\phi^t s_c^{t-1} + b_i^t g(a_c^t)$$

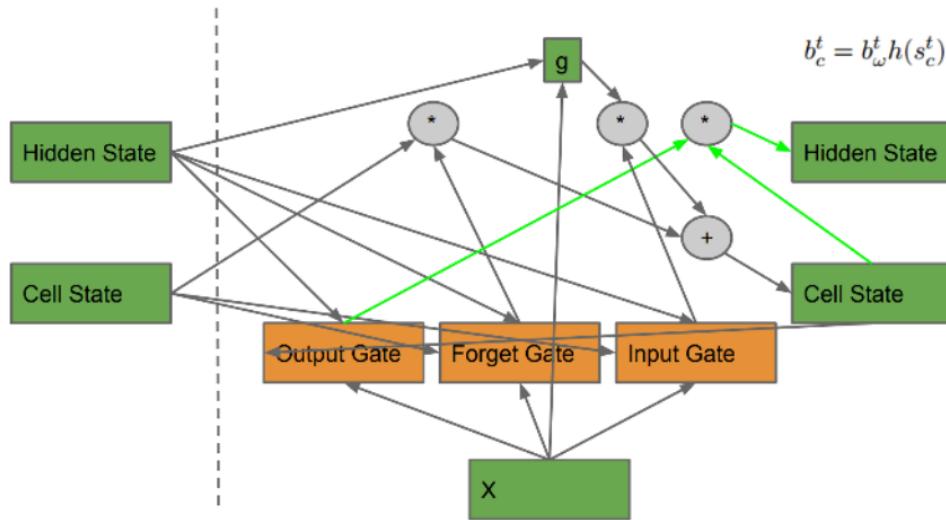
- LSTM 正向传播：输出门

4. Output Gate



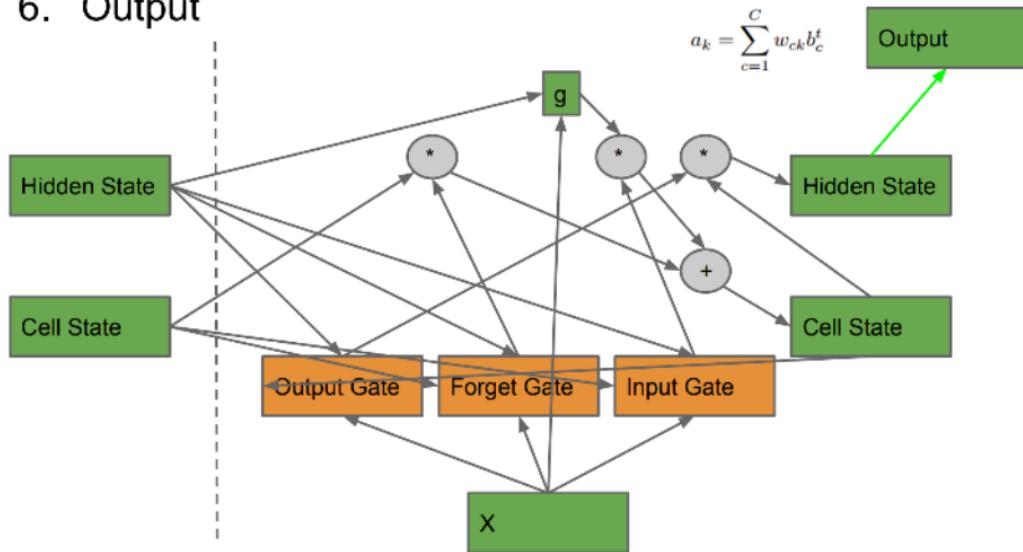
- LSTM 正向传播: 隐藏状态

5. Hidden State



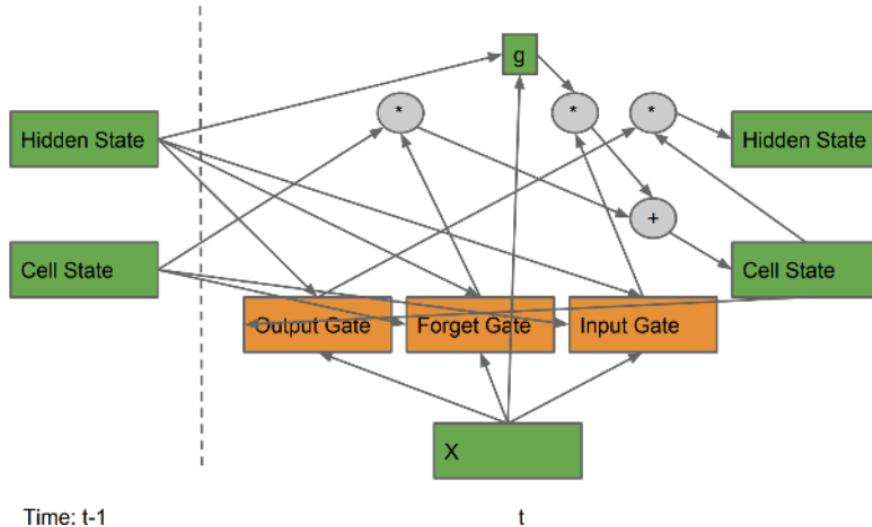
- LSTM 正向传播: 最后输出

6. Output



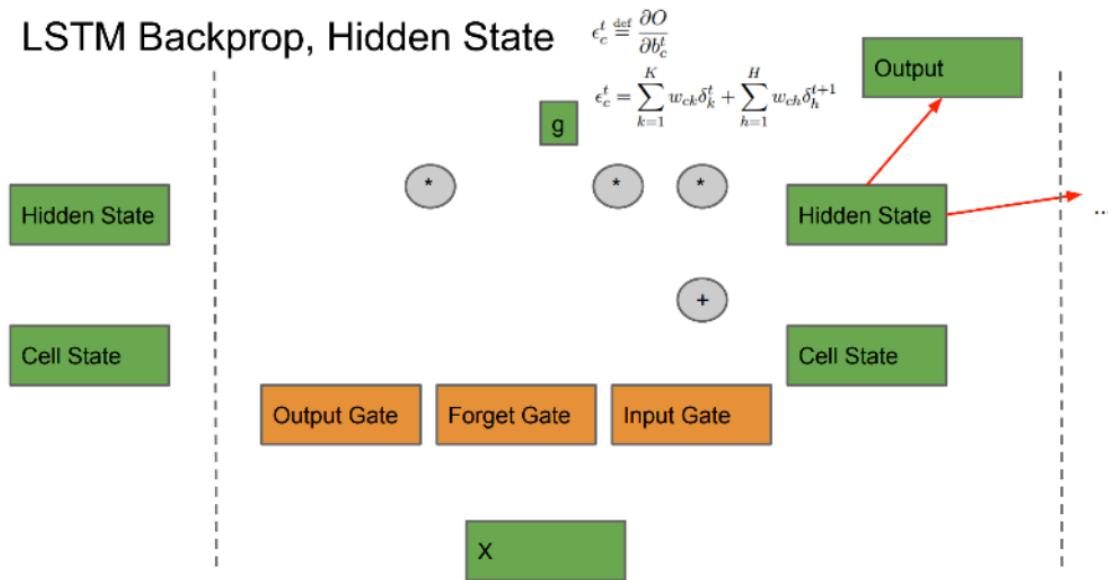
- LSTM 正向传播

LSTM Forward Propagation



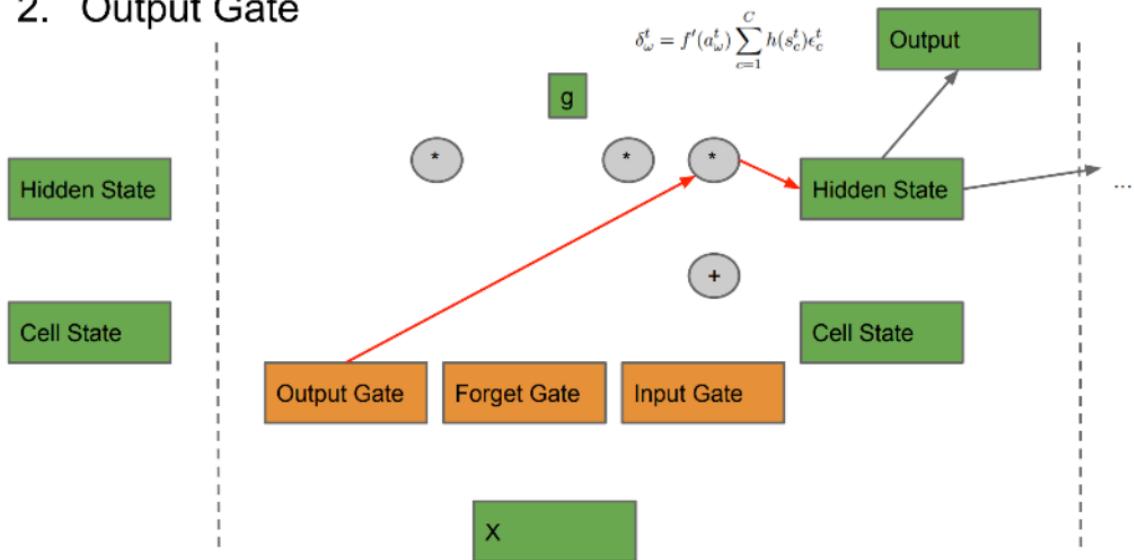
- LSTM 反向传播：输出和隐藏状态

LSTM Backprop, Hidden State



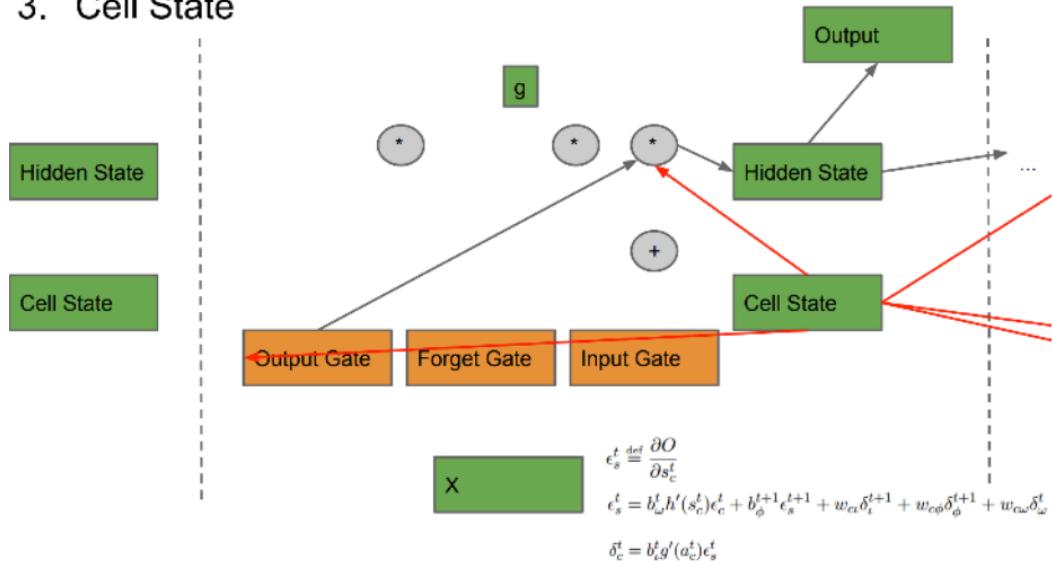
- LSTM 反向传播：输出门

2. Output Gate



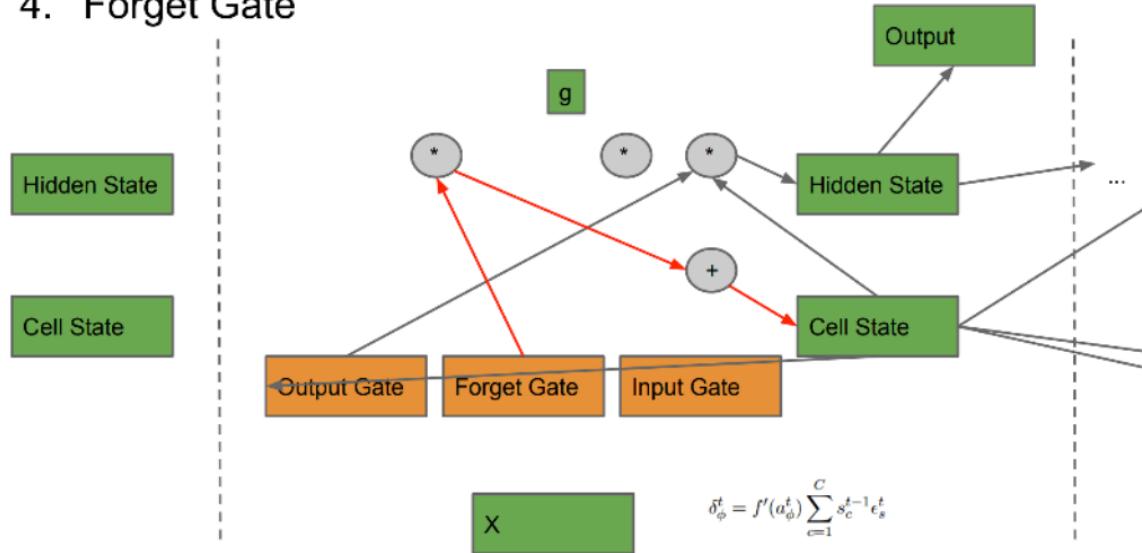
- LSTM 反向传播：记忆 Cell 状态

3. Cell State



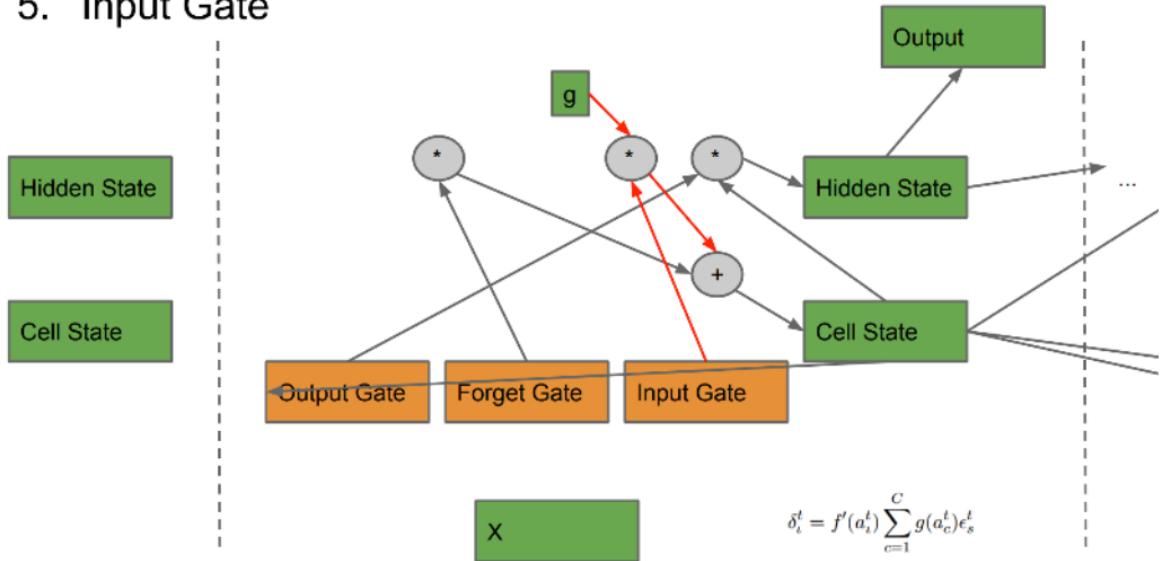
- LSTM 反向传播：忘记门

4. Forget Gate



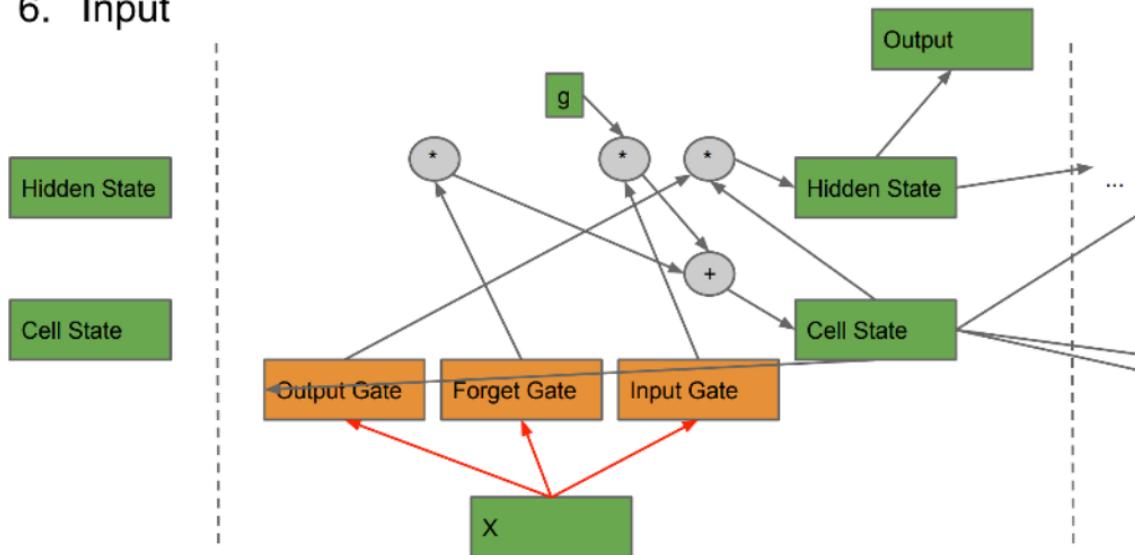
- LSTM 反向传播：输入门

5. Input Gate



- LSTM 反向传播：输入

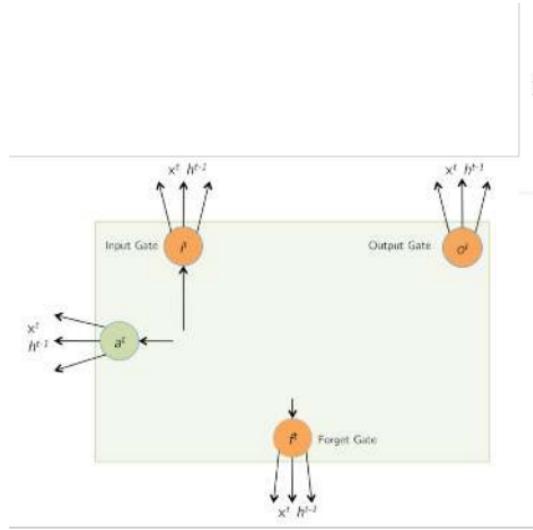
6. Input



- LSTM 反向传播详细

$\begin{aligned} a^t &= \tanh(W_c x^t + U_c h^{t-1}) = \tanh(\hat{a}^t) \\ i^t &= \sigma(W_i x^t + U_i h^{t-1}) = \sigma(\hat{i}^t) \\ f^t &= \sigma(W_f x^t + U_f h^{t-1}) = \sigma(\hat{f}^t) \\ o^t &= \sigma(W_o x^t + U_o h^{t-1}) = \sigma(\hat{o}^t) \end{aligned}$	$\begin{aligned} \frac{\partial E}{\partial c_i^t} &= \frac{\partial E}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial c_i^t} \\ &= \delta h_i^t \cdot o_i^t \cdot (1 - \tanh^2(c_i^t)) \\ \therefore \delta c^t &= \delta h^t \odot o^t \odot (1 - \tanh^2(c^t)) \end{aligned}$	$\begin{aligned} \frac{\partial E}{\partial o_i^t} &= \frac{\partial E}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial o_i^t} \\ &= \delta h_i^t \cdot \tanh(c_i^t) \\ \therefore \delta o^t &= \delta h^t \odot \tanh(c^t) \end{aligned}$
$\begin{aligned} c^t &= i^t \odot a^t + f^t \odot c^{t-1} \\ h^t &= o^t \odot \tanh(c^t) \end{aligned}$	$\begin{aligned} \frac{\partial E}{\partial i_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial i_i^t} \\ &= \delta c_i^t \cdot a_i^t \\ \therefore \delta i^t &= \delta c^t \odot a^t \end{aligned}$	$\begin{aligned} \frac{\partial E}{\partial f_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial f_i^t} \\ &= \delta c_i^t \cdot c_i^{t-1} \\ \therefore \delta f^t &= \delta c^t \odot c^{t-1} \end{aligned}$
	$\begin{aligned} \frac{\partial E}{\partial a_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial a_i^t} \\ &= \delta c_i^t \cdot i_i^t \\ \therefore \delta a^t &= \delta c^t \odot i^t \end{aligned}$	$\begin{aligned} \frac{\partial E}{\partial c_i^{t-1}} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial c_i^{t-1}} \\ &= \delta c_i^t \cdot f_i^t \\ \therefore \delta c^{t-1} &= \delta c^t \odot f^t \end{aligned}$

- LSTM 的 BPTT



Forward Pass: $z^t = \begin{bmatrix} \hat{a}^t \\ \hat{i}^t \\ \hat{f}^t \\ \hat{o}^t \end{bmatrix} = W \times I^t$

Given $\delta a^t, \delta i^t, \delta f^t, \delta o^t$, find δz^t

$$\begin{aligned}\delta \hat{a}^t &= \delta a^t \odot (1 - \tanh^2(\hat{a}^t)) \\ \delta \hat{i}^t &= \delta i^t \odot i^t \odot (1 - i^t) \\ \delta \hat{f}^t &= \delta f^t \odot f^t \odot (1 - f^t) \\ \delta \hat{o}^t &= \delta o^t \odot o^t \odot (1 - o^t) \\ \delta z^t &= [\delta \hat{a}^t, \delta \hat{i}^t, \delta \hat{f}^t, \delta \hat{o}^t]^T\end{aligned}$$

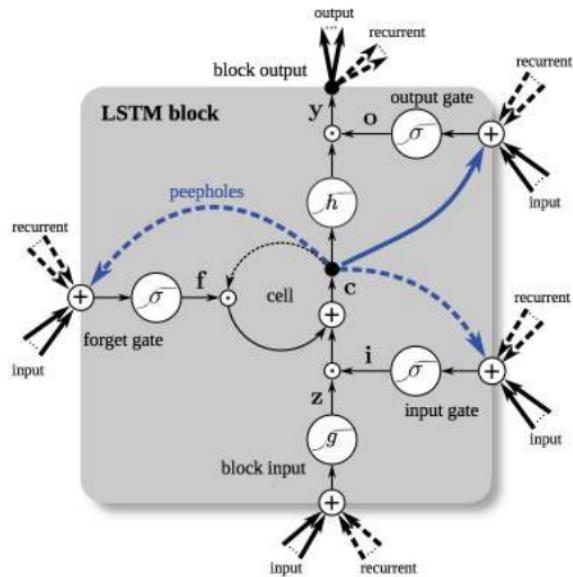
Forward Pass: $z^t = W \times I^t$
Given δz^t , find $\delta W^t, \delta h^{t-1}$

$$\begin{aligned}\delta I^t &= W^T \times \delta z^t \\ \text{As } I^t &= \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix}, \\ \delta h^{t-1} &\text{ can be retrieved from } \delta I^t \\ \delta W^t &= \delta z^t \times (I^t)^T\end{aligned}$$

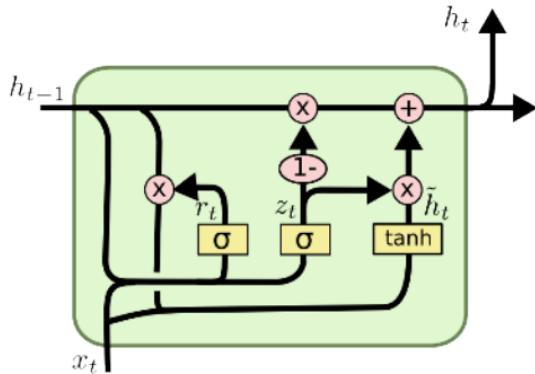
$x = [x^1, x^2, \dots, x^T]$

$\delta W = \sum_{t=1}^T \delta W^t$

- LSTM Peepholes



- GRU



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

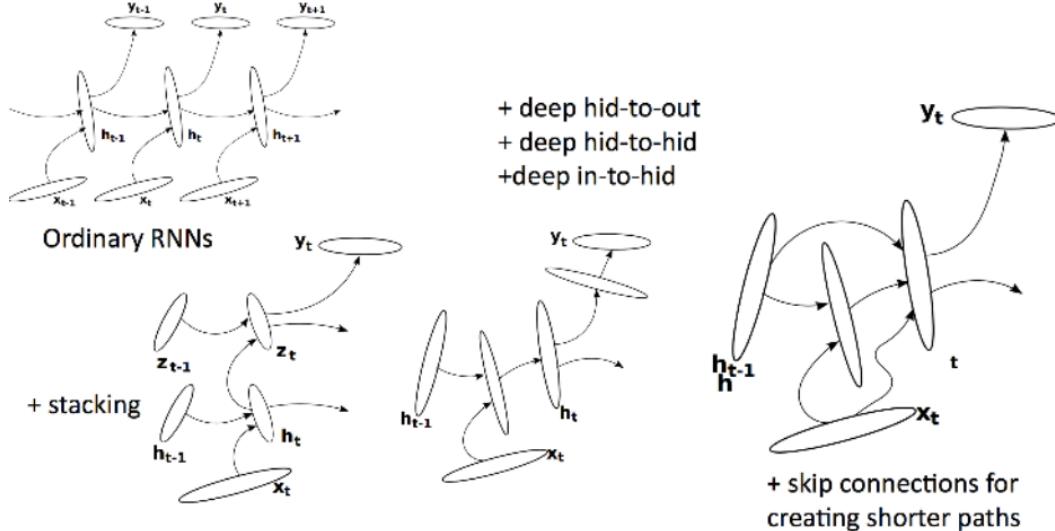
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- GRU
 1. 复位门 Reset Gate: 如何结合以前隐藏状态和当前输入
 2. 更新门 Update Gate: 如何保留内部状态

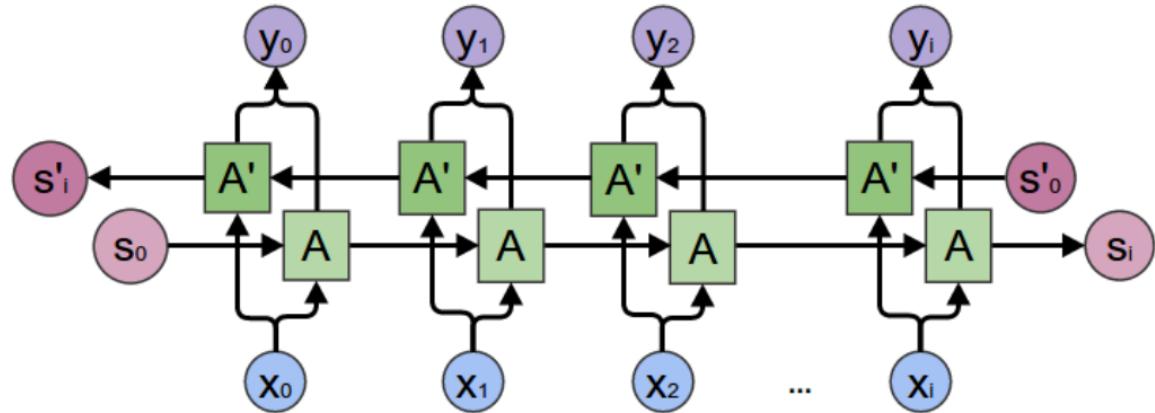
- LSTM vs GRU
 - 1. 效果差不多
 - 2. GRU 参数更少，需要的数据量少
 - 3. LSTM 表达能力强，需要的数据量大

- 深度 RNN

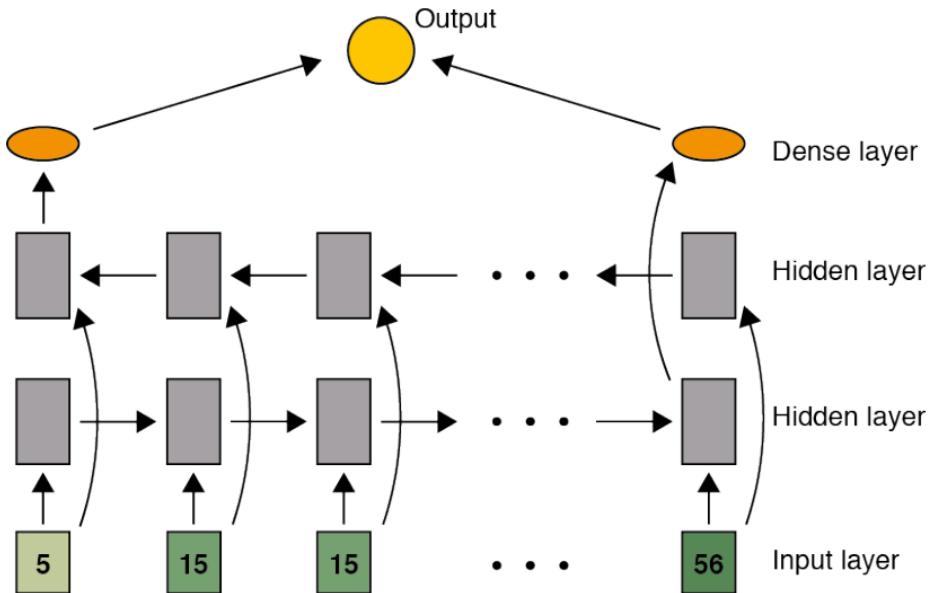
- ICLR 2014, *How to construct deep recurrent neural networks*



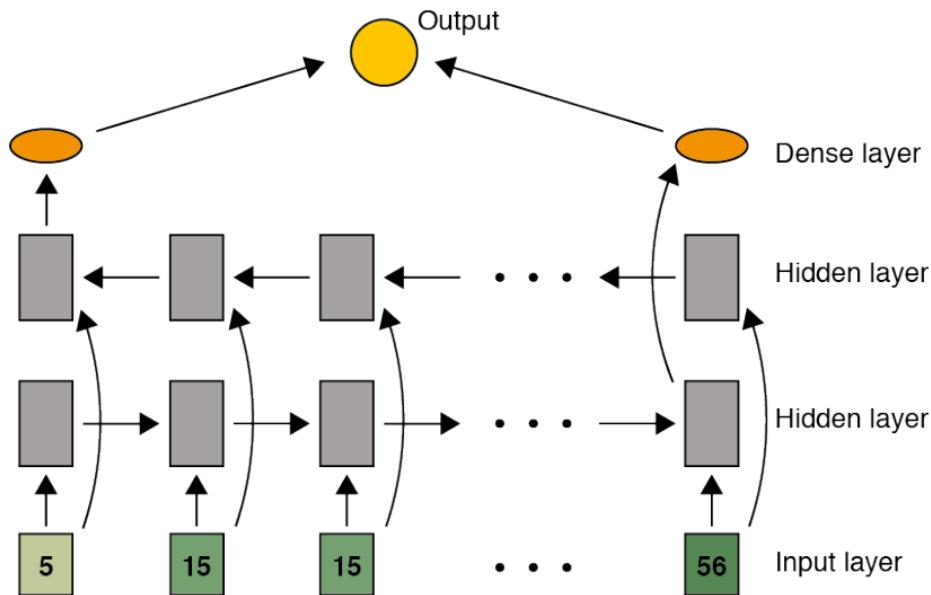
- Bidirectional RNN



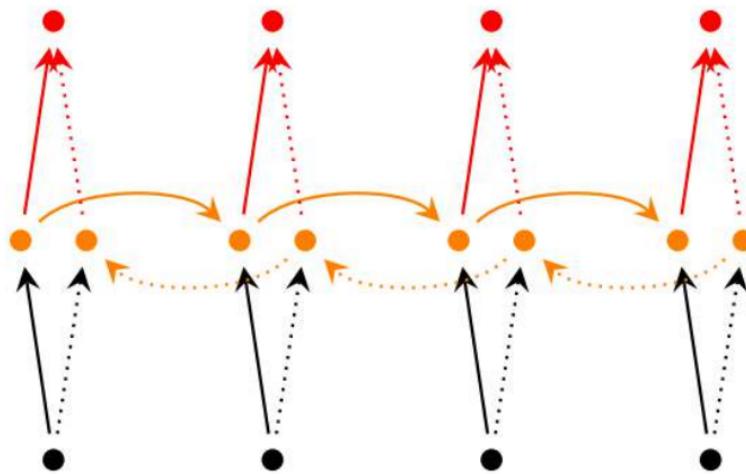
- Bidirectional RNN



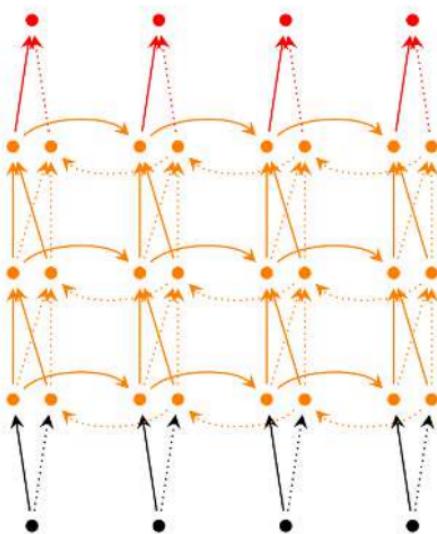
- Bidirectional RNN



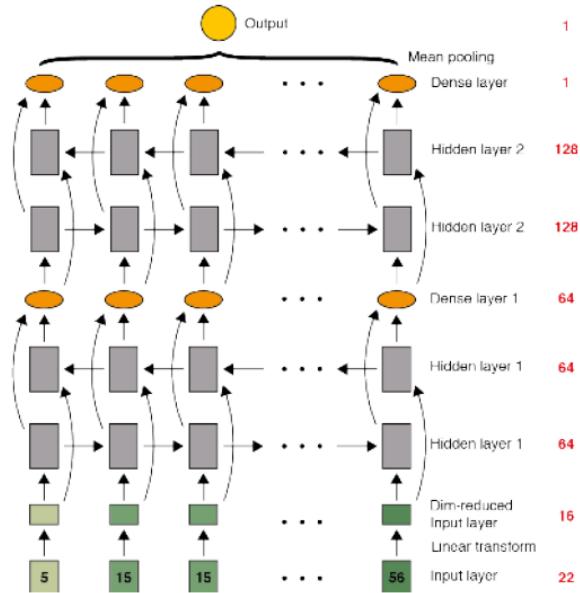
- Bidirectional RNN



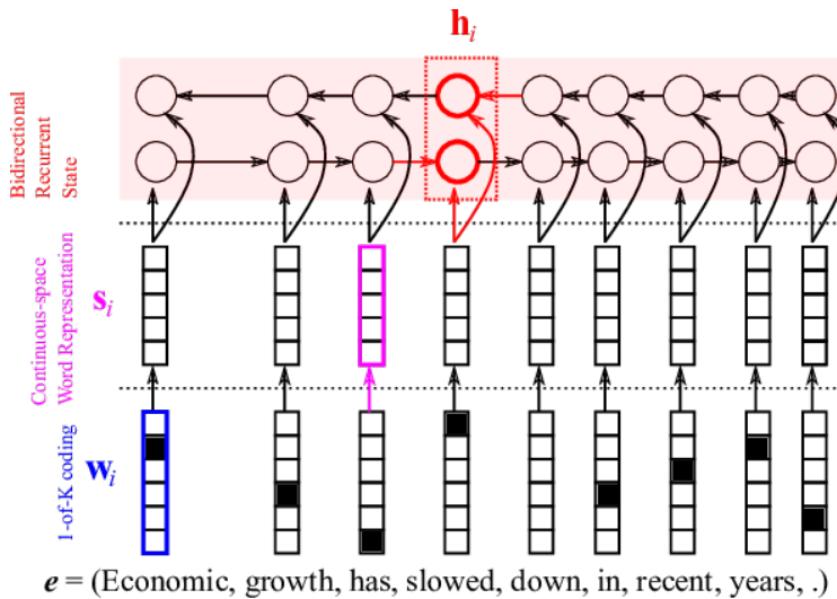
- Deep (Bidirectional) RNNs



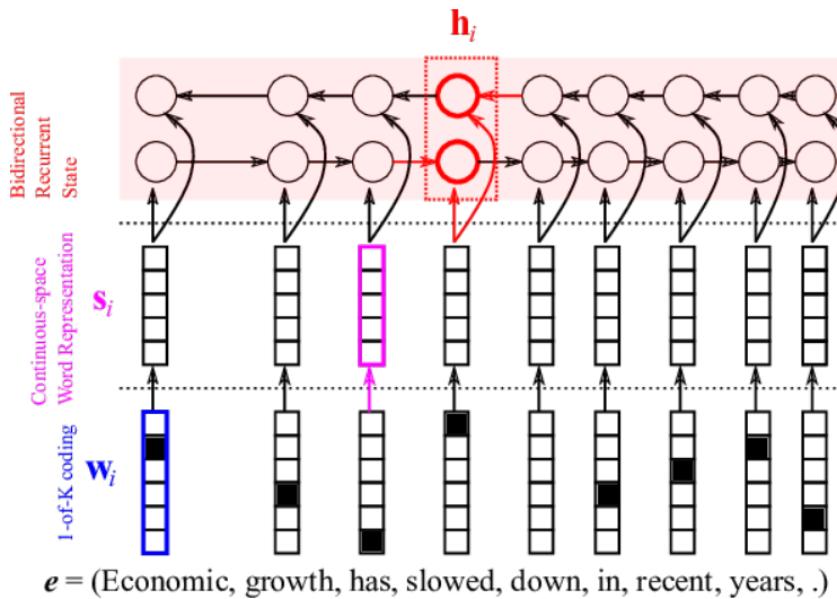
- Deep (Bidirectional) RNNs



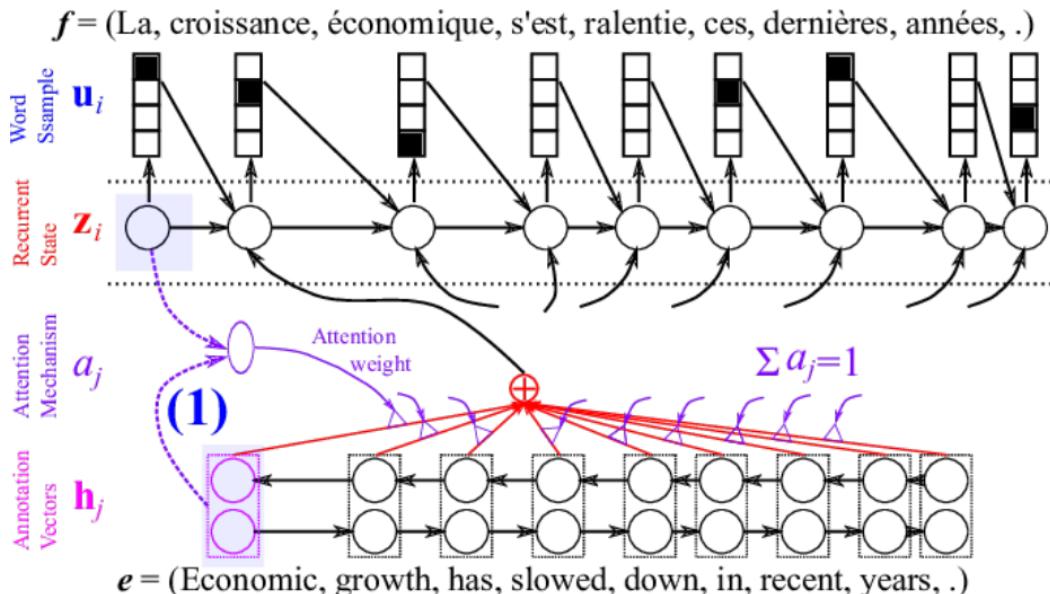
- Bidirectional RNNs 应用



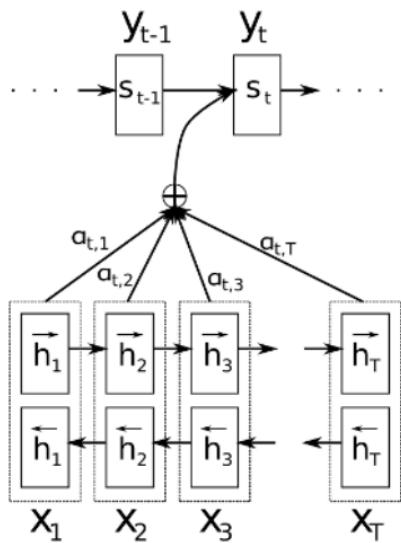
- Bidirectional RNNs 应用



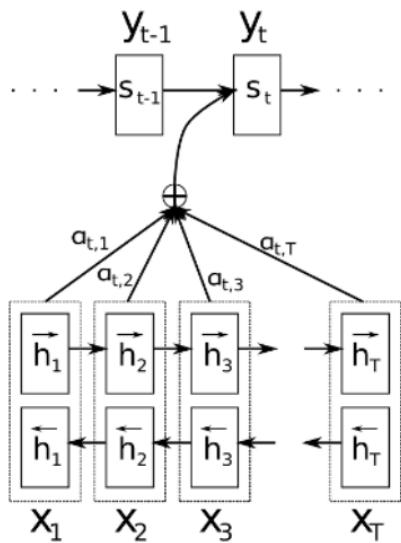
- Bidirectional RNN 和 Attention 机制



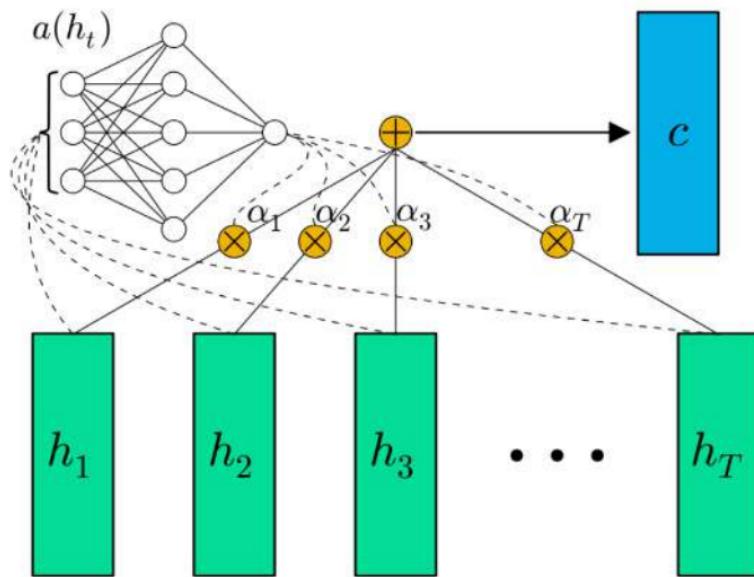
- Attention 机制



- Attention 机制



- Attention 机制：权重学习

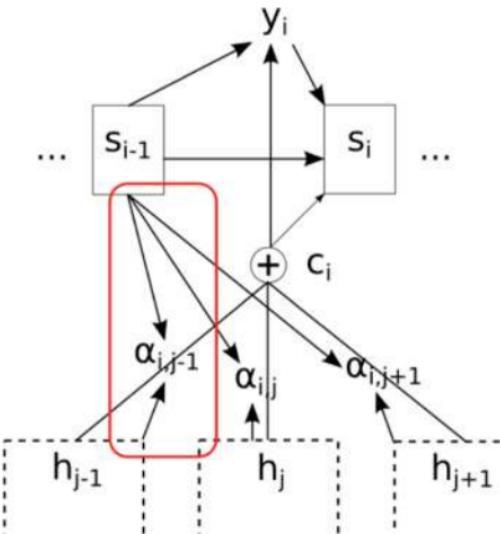


- Attention 机制

$$e_{ij} = v^T \tanh(Ws_{i-1} + Vh_j) \quad (1)$$

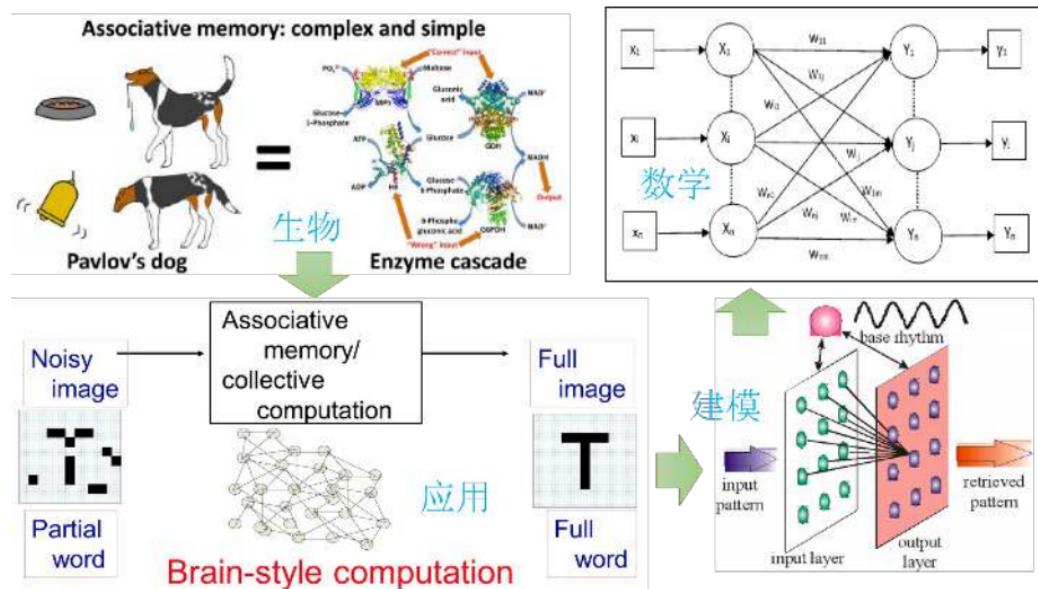
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^L \exp(e_{ik})} \quad (2)$$

- nonlinearity (\tanh) is crucial!
- simplest model possible
- Vh_j is precomputed => quadratic complexity with low constant

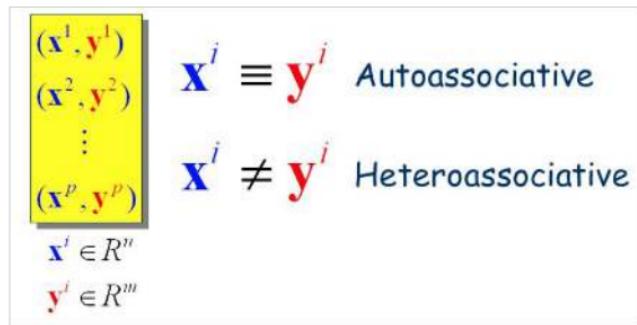
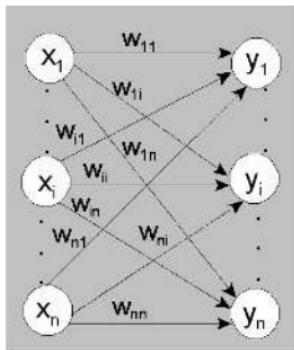


2 玻尔兹曼机 BM

- 联想记忆 Associative Memory



- 能量函数 Energy Memory



Hebb学习：

$$y = \sum_j w_j x_j$$

$$\Delta w_i = \eta x_i y$$

梯度下降：

$$w := w - \eta \nabla E(w)$$

能量函数：

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} x_i y_j$$

- Hopfield 网络能量函数

能量函数：

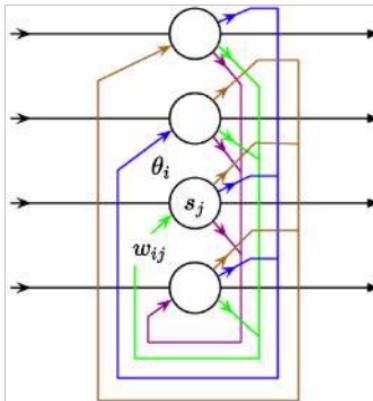
$$E = -\frac{1}{2} \sum_{i,j} w_{ij} x_i y_j$$

神经单元设置：

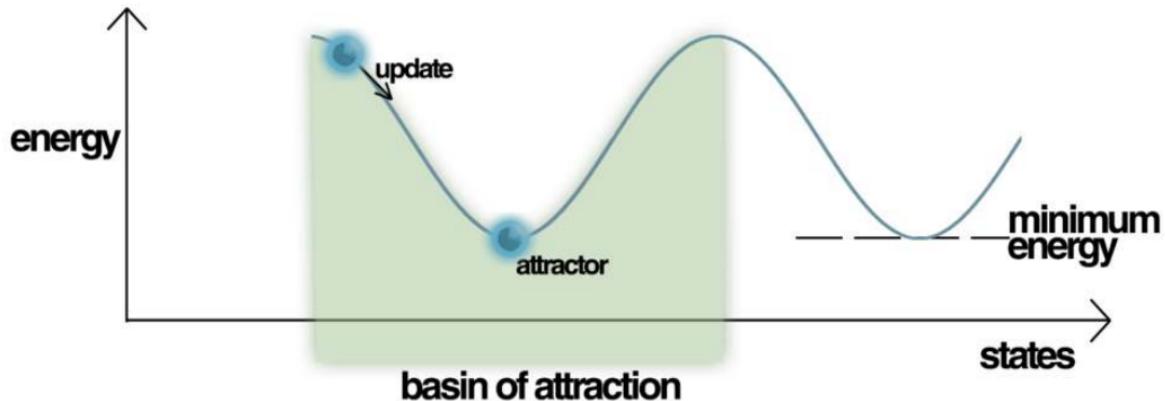
$$s_i \leftarrow \begin{cases} +1 & \text{if } \sum_j w_{ij} s_j \geq \theta_i \\ -1 & \text{otherwise.} \end{cases}$$

Hopfield 网络能量函数：

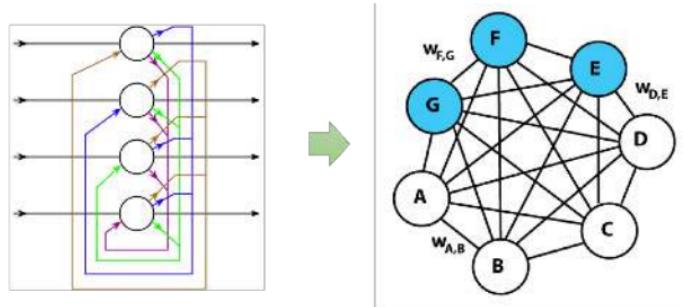
$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$



- Hopfield 网络能量



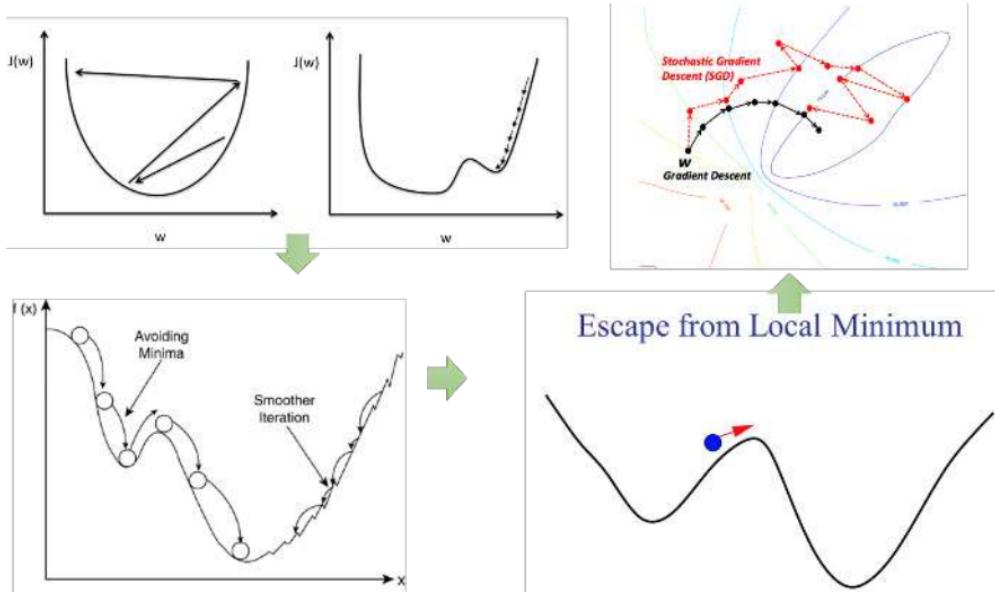
- 玻尔兹曼机 Boltzmann Machine



$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i \quad \Rightarrow \quad E = - \left(\sum_{i < j} w_{ij} s_i s_j + \sum_i \theta_i s_i \right)$$



- 引入随机、跳出局部最小值



- 玻尔兹曼分布

能级状态：

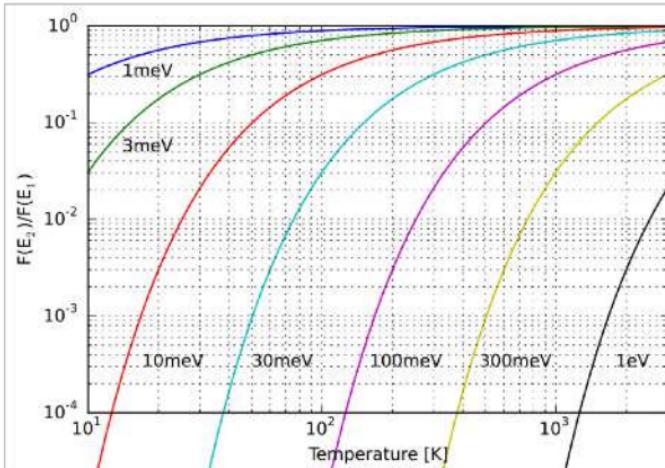
$$F(\text{state}) \propto e^{-\frac{E}{kT}}$$

粒子数量占比：

$$\frac{N_i}{N} = \frac{e^{-\varepsilon_i/kT}}{\sum_{j=1}^M e^{-\varepsilon_j/kT}}$$

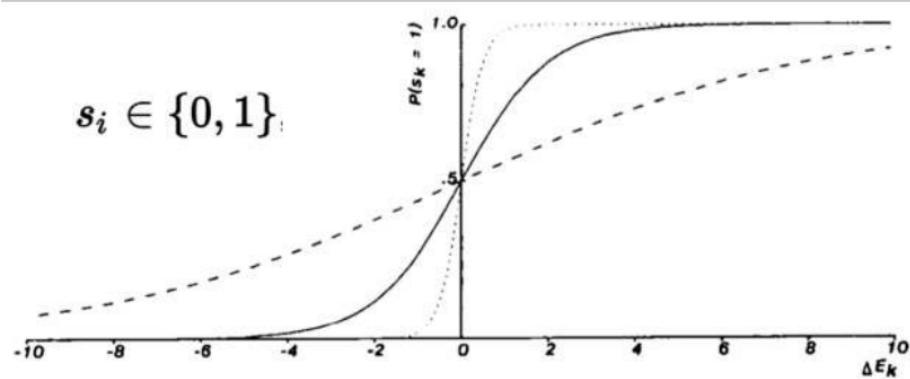
能级概率：

$$p_i = \frac{e^{-\varepsilon_i/kT}}{\sum_{j=1}^M e^{-\varepsilon_j/kT}}$$



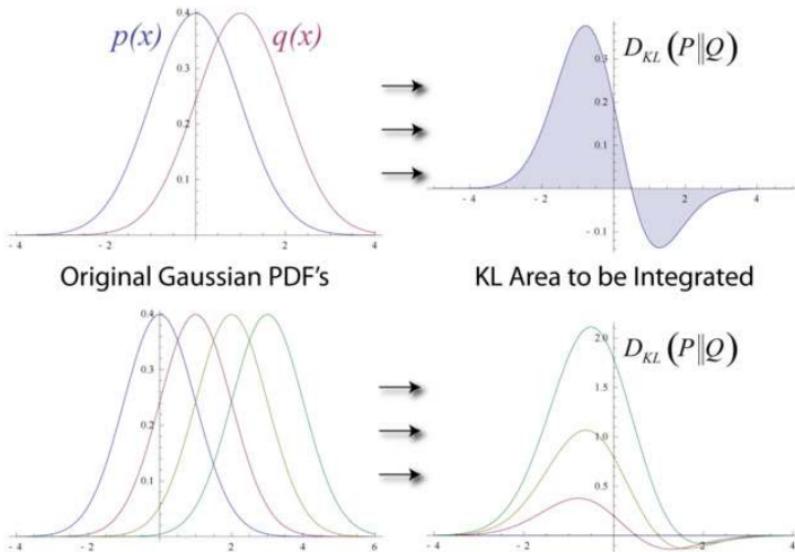
- 玻尔兹曼机的分布

$$p(s_i=1) = \frac{1}{1+e^{-\Delta E_i/T}}$$



$$s_i \in \{0, 1\}$$

- 分布间距离：KL 散度



- KL 散度

连续定义：

$$D_{\text{KL}}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx,$$

离散定义：

$$D_{\text{KL}}(P\|Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}.$$

熵的解释：交互熵 - 熵

$$\begin{aligned} D_{\text{KL}}(P\|Q) &= -\sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) \\ &= H(P, Q) - H(P) \end{aligned}$$

- KL 散度称为 KL 距离：度量大于 0

$$\begin{aligned} D(p||q) &= \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \\ &= - \sum_{x \in \mathcal{X}} p(x) \log \frac{q(x)}{p(x)} \\ &= -\mathbb{E} \left[\log \frac{q}{p} \right] \\ &\geq -\log \left(\mathbb{E} \left[\frac{q}{p} \right] \right) \\ &= -\log \left(\sum_{x \in \mathcal{X}} p(x) \frac{q(x)}{p(x)} \right) \\ &= 0, \end{aligned}$$

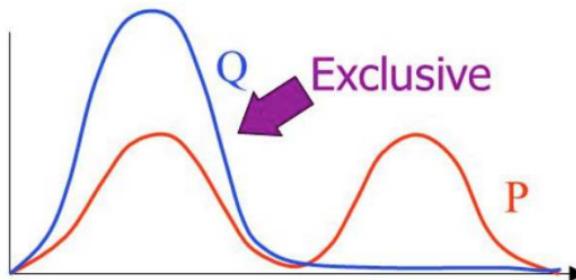
- KL 散度性质 1:

$$\begin{aligned} D(p||q) &= \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \\ &= - \sum_{x \in \mathcal{X}} p(x) \log \frac{q(x)}{p(x)} \\ &= -\mathbb{E} \left[\log \frac{q}{p} \right] \\ &\geq -\log \left(\mathbb{E} \left[\frac{q}{p} \right] \right) \\ &= -\log \left(\sum_{x \in \mathcal{X}} p(x) \frac{q(x)}{p(x)} \right) \\ &= 0, \end{aligned}$$

- KL 距离性质 2: 非对称度量

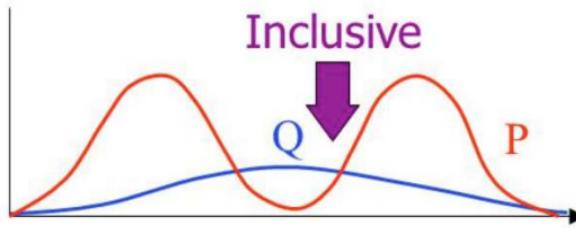
Minimising
 $\text{KL}(Q||P)$

$$= \sum_H Q(H) \ln \frac{Q(H)}{P(H|V)}$$



Minimising
 $\text{KL}(P||Q)$

$$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$$



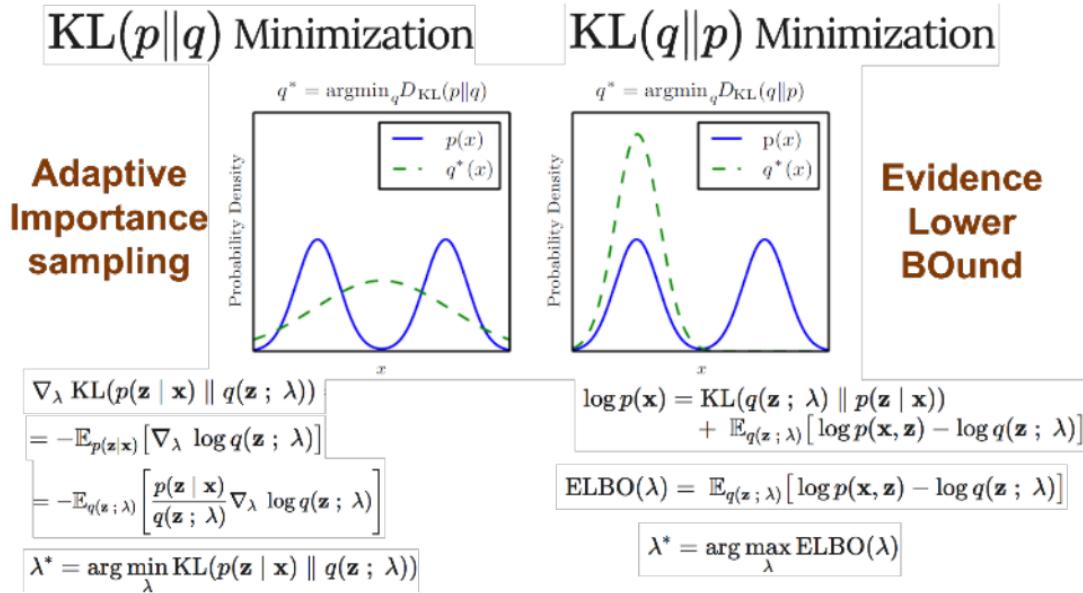
- KL 距离性质 3: 最大似然函数等价最小 KL 距离

带参数分布 **最小KL距离**

$$\begin{aligned} D(\hat{p}||p_{\theta}) &= \sum_{x \in \mathcal{X}} \hat{p}(x) \log \frac{\hat{p}(x)}{p_{\theta}(x)} \\ &= -H(\hat{p}) - \sum_{x \in \mathcal{X}} \hat{p}(x) \log p_{\theta}(x) \\ &= -H(\hat{p}) - \frac{1}{N} \sum_{x \in \mathcal{X}} \sum_{n=1}^N \delta(x - x_n) \log p_{\theta}(x) \\ &= -H(\hat{p}) - \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(x_n) \end{aligned}$$

最大Log似然度 带参数分布

- KL 距离性质 4: ELBO vs Importance Sampling



- BM 的最小 KL 距离 (最大似然估计)

$$\begin{aligned} \text{KL}(\mathbb{P}_{\text{target}} || \mathbb{P}_\theta) &\equiv \sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \log \frac{\mathbb{P}_{\text{target}}(\tilde{\mathbf{x}})}{\mathbb{P}_\theta(\tilde{\mathbf{x}})} \\ &= \sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \log \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) - \sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \log \mathbb{P}_\theta(\tilde{\mathbf{x}}) \end{aligned}$$

$$\begin{aligned} \nabla \log \mathbb{P}_\theta(\mathbf{x}) &= \nabla \log \frac{\exp(-E_\theta(\mathbf{x}))}{\sum_{\hat{\mathbf{x}}} \exp(-E_\theta(\hat{\mathbf{x}}))} \\ &= -\nabla E_\theta(\mathbf{x}) - \nabla \log \sum_{\hat{\mathbf{x}}} \exp(-E_\theta(\hat{\mathbf{x}})) \\ &= -\nabla E_\theta(\mathbf{x}) + \frac{\sum_{\hat{\mathbf{x}}} \exp(-E_\theta(\hat{\mathbf{x}})) \nabla E_\theta(\hat{\mathbf{x}})}{\sum_{\hat{\mathbf{x}}} \exp(-E_\theta(\hat{\mathbf{x}}))} \\ &= -\nabla E_\theta(\mathbf{x}) + \sum_{\hat{\mathbf{x}}} \mathbb{P}_\theta(\hat{\mathbf{x}}) \nabla E_\theta(\hat{\mathbf{x}}), \end{aligned}$$

$$\begin{aligned} \nabla f(\theta) &= \sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \nabla \log \mathbb{P}_\theta(\tilde{\mathbf{x}}) \\ &= -\sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \nabla E_\theta(\tilde{\mathbf{x}}) + \\ &\quad \sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \sum_{\hat{\mathbf{x}}} \mathbb{P}_\theta(\hat{\mathbf{x}}) \nabla E_\theta(\hat{\mathbf{x}}) \\ &= -\sum_{\tilde{\mathbf{x}}} \mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) \nabla E_\theta(\tilde{\mathbf{x}}) + \sum_{\hat{\mathbf{x}}} \mathbb{P}_\theta(\hat{\mathbf{x}}) \nabla E_\theta(\hat{\mathbf{x}}) \\ &= -\sum_{\tilde{\mathbf{x}}} (\mathbb{P}_{\text{target}}(\tilde{\mathbf{x}}) - \mathbb{P}_\theta(\tilde{\mathbf{x}})) \nabla E_\theta(\tilde{\mathbf{x}}). \end{aligned}$$

$$\nabla f(\theta) = [-\mathbb{E}_{\text{target}} [\nabla E_\theta(\mathbf{X})] + \mathbb{E}_\theta [\nabla E_\theta(\mathbf{X})]]$$

- BM 的参数更新

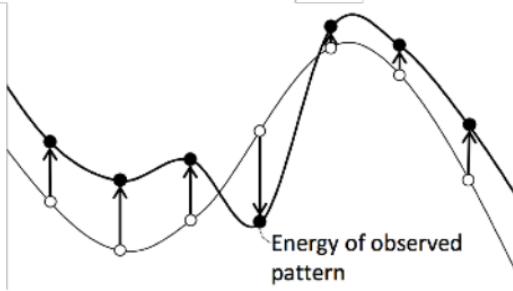
$$\nabla f(\theta) = [-\mathbb{E}_{\text{target}}[\nabla E_\theta(\mathbf{X})] + \mathbb{E}_\theta[\nabla E_\theta(\mathbf{X})]]$$



$$\begin{aligned}\frac{\partial}{\partial b_i} f(\theta) &= \mathbb{E}_{\text{target}}[X_i] - \mathbb{E}_\theta[X_i] \\ \frac{\partial}{\partial w_{i,j}} f(\theta) &= \mathbb{E}_{\text{target}}[X_i X_j] - \mathbb{E}_\theta[X_i X_j]\end{aligned}$$

$$\begin{aligned}E_\theta(\mathbf{x}) &= -\sum_{i=1}^N b_i x_i - \sum_{i=1}^{N-1} \sum_{j=i+1}^N w_{i,j} x_i x_j \\ &= -\mathbf{b}^\top \mathbf{x} - \mathbf{x}^\top \mathbf{W} \mathbf{x}.\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial b_i} E_\theta(\mathbf{x}) &= -x_i \\ \frac{\partial}{\partial w_{i,j}} E_\theta(\mathbf{x}) &= -x_i x_j\end{aligned}$$



- BM 的参数更新的直观理解

$$\frac{\partial \log p(\mathbf{v})}{\partial w_{ij}} = \left\langle s_i s_j \right\rangle_{\mathbf{v}} - \left\langle s_i s_j \right\rangle_{model}$$

Derivative of log probability of one training vector, \mathbf{v} under the model.

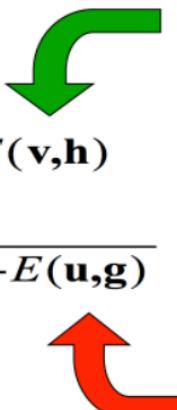
Expected value of product of states at thermal equilibrium when \mathbf{v} is clamped on the visible units

Expected value of product of states at thermal equilibrium with no clamping

$$\Delta w_{ij} \propto \left\langle s_i s_j \right\rangle_{data} - \left\langle s_i s_j \right\rangle_{model}$$

1. 正向: 在可见端固定 \mathbf{x} , 等稳定后, 计算期望值
2. 反向: 松开可见端 \mathbf{x} , 等稳定后, 计算期望值

- BM 的最大似然的直观理解

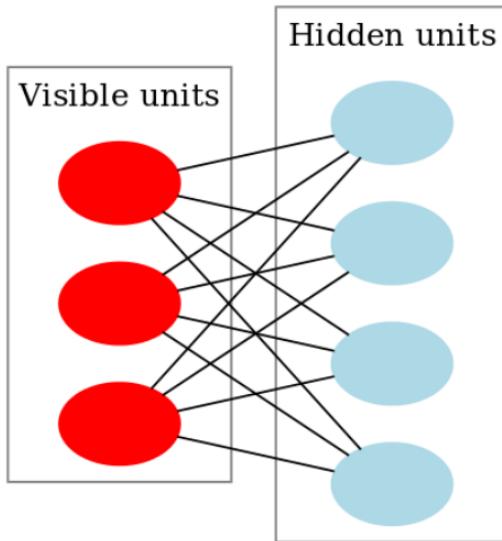
$$p(\mathbf{v}) = \frac{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{u}} \sum_{\mathbf{g}} e^{-E(\mathbf{u}, \mathbf{g})}}$$


The positive phase finds hidden configurations that work well with \mathbf{v} and lowers their energies.

The negative phase finds the joint configurations that are the best competitors and raises their energies.

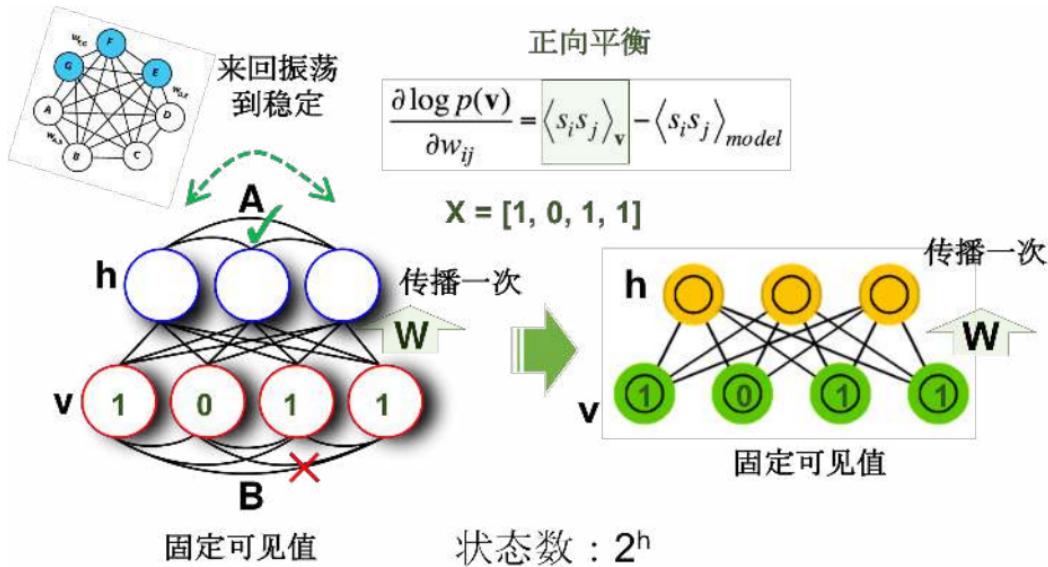
1. 正向: 找到隐藏层参数, 修正降低能量
2. 反向: 找到联合参数, 修正提高能量

- 受限玻尔兹曼机 Restricted Boltzmann Machines

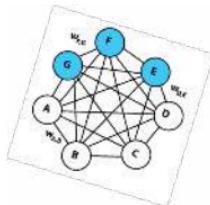


BM 过于复杂，难以训练！

- RBM 的好处：正向平衡计算一次

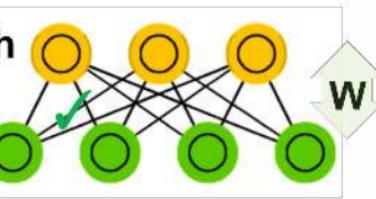
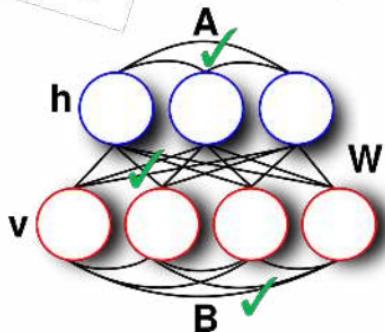


- RBM 的好处：反向平衡下上来回振荡



反向平衡

$$\frac{\partial \log p(\mathbf{v})}{\partial w_{ij}} = \langle s_i s_j \rangle_{\mathbf{v}} - \langle s_i s_j \rangle_{model}$$



上下来回振荡到稳定

全部来回振荡到稳定 状态数： 2^{h+v}

- RBM 的好处，条件独立性

RBM 能量函数 (分离了 V、H 内部乘法):

$$E(v, h) = -\sum_i a_i v_i - \sum_j b_j h_j - \sum_i \sum_j v_i w_{i,j} h_j$$

矩阵形式:

$$E(v, h) = -a^T v - b^T h - v^T W h$$

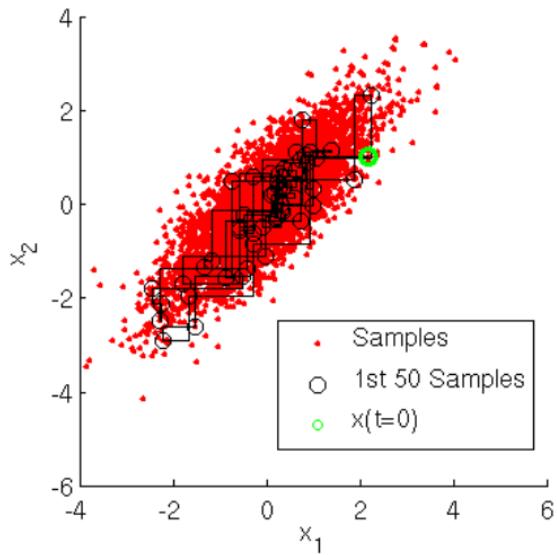
可见层条件独立性:

$$P(v|h) = \prod_{i=1}^m P(v_i|h)$$

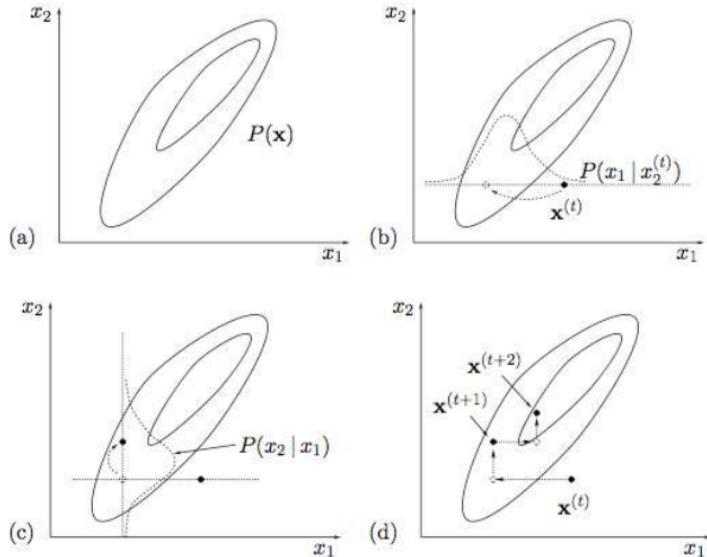
隐藏层条件独立性:

$$P(h|v) = \prod_{j=1}^n P(h_j|v)$$

- 吉布斯采样求解：Gibbs Sampling



- 吉布斯采样求解：Gibbs Sampling



- 蒙特卡洛马尔科夫链：MCMC

Markov Chain Monte Carlo (MCMC)

Von Neumann



he contributed to the development of the Monte Carlo method, which allowed solutions to complicated problems to be approximated using random numbers.

Ulam



Metropolis



Paper 1953:
They applied MCMC for
a chemical problem

Paper 1949: Using Markov
chain for Monte Carlo
approximation

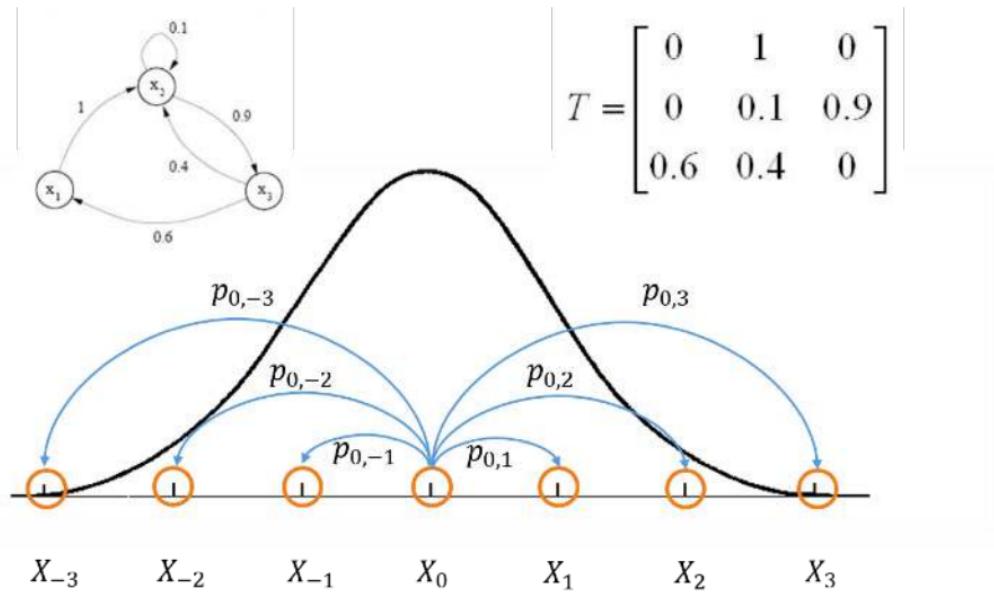


Rossenborg

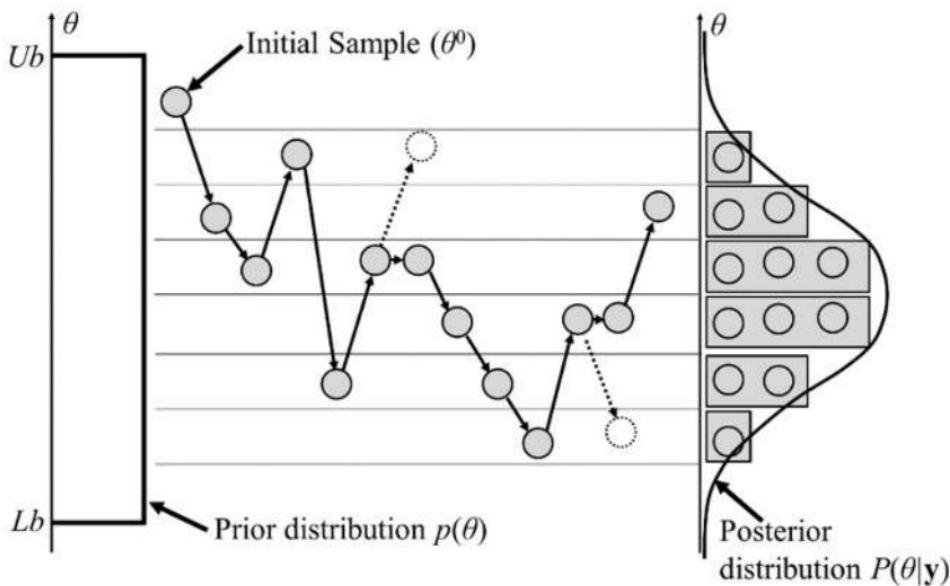
Teller



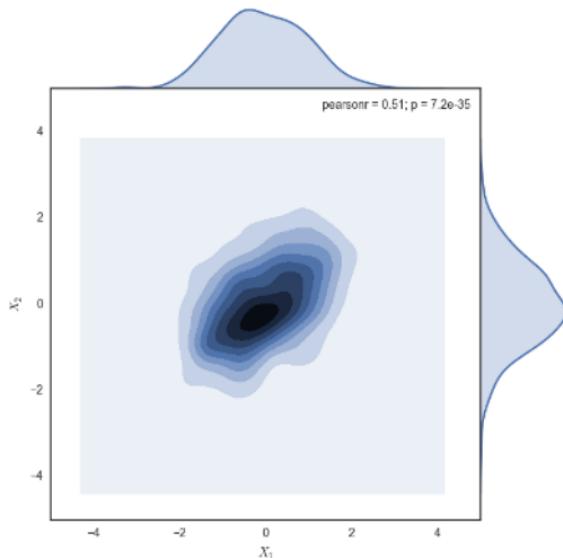
- MCMC 迁移



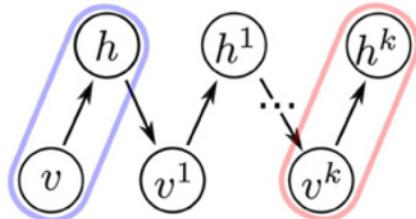
- MCMC 分布抽样



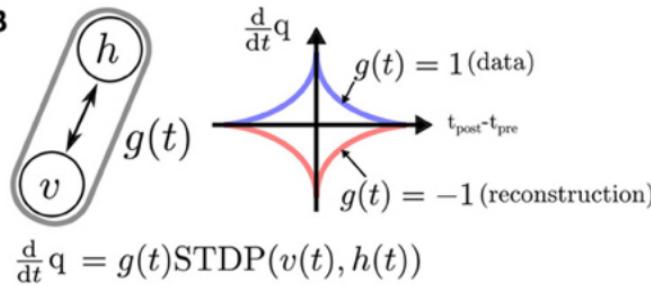
- Gibbs 抽样收敛：两个独立变量核心分布



- RBM Gibbs 抽样

A

$$\Delta w \propto \langle vh \rangle_{\text{data}} - \langle v^k h^k \rangle_{\text{recon}}$$

B

- 正向反向平衡的重新解读

$$\begin{aligned}\frac{\partial E(\mathbf{X}; \Theta)}{\partial \Theta} &= \frac{\partial \log Z(\Theta)}{\partial \Theta} - \frac{1}{K} \sum_{i=1}^K \frac{\partial \log f(x_i; \Theta)}{\partial \Theta} \\ &= \boxed{\frac{\partial \log Z(\Theta)}{\partial \Theta}} - \left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{\mathbf{x}}\end{aligned}$$

反向平衡

$\mathbf{X}^0 \equiv \mathbf{X}$ 输入

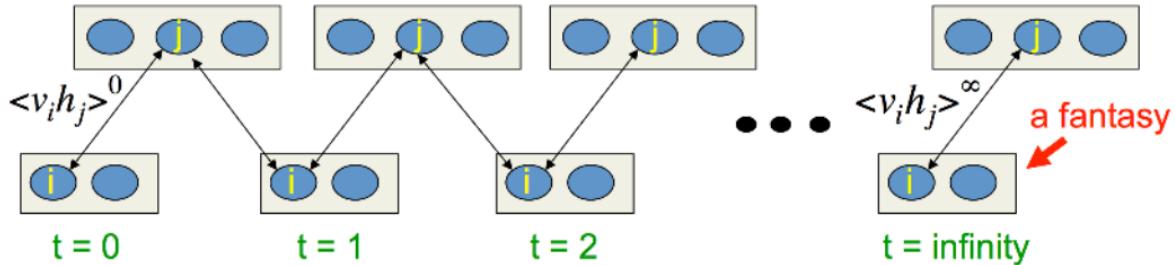
$$\frac{\partial E(\mathbf{X}; \Theta)}{\partial \Theta} = \left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{\mathbf{x}^\infty} - \left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{\mathbf{x}^0}$$

重建

$$\Theta_{t+1} = \Theta_t + \eta \left(\left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{\mathbf{x}^0} - \left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{\mathbf{x}^1} \right)$$

$$\begin{aligned}\frac{\partial \log Z(\Theta)}{\partial \Theta} &= \frac{1}{Z(\Theta)} \frac{\partial Z(\Theta)}{\partial \Theta} \\ &= \frac{1}{Z(\Theta)} \frac{\partial}{\partial \Theta} \int f(x; \Theta) dx \\ &= \frac{1}{Z(\Theta)} \int \frac{\partial f(x; \Theta)}{\partial \Theta} dx \\ &= \frac{1}{Z(\Theta)} \int f(x; \Theta) \frac{\partial \log f(x; \Theta)}{\partial \Theta} dx \\ &= \int p(x; \Theta) \frac{\partial \log f(x; \Theta)}{\partial \Theta} dx \\ &= \left\langle \frac{\partial \log f(x; \Theta)}{\partial \Theta} \right\rangle_{p(x; \Theta)}\end{aligned}$$

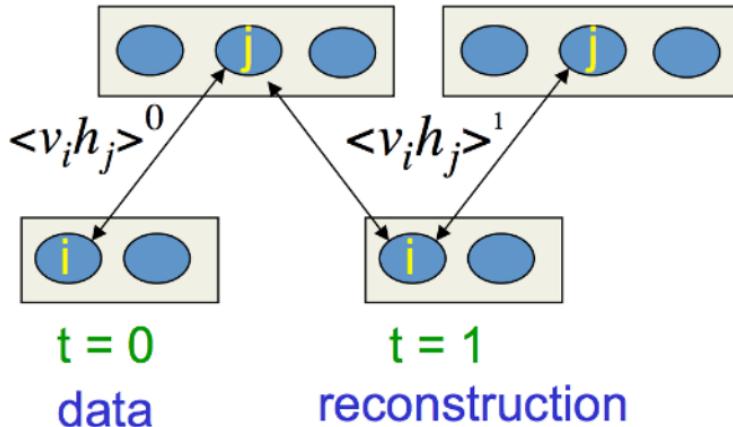
- RBM Gibbs 抽样：可见层一端视角



Start with a training vector on the visible units. Then alternate between updating all the hidden units in parallel and updating all the visible units in parallel.

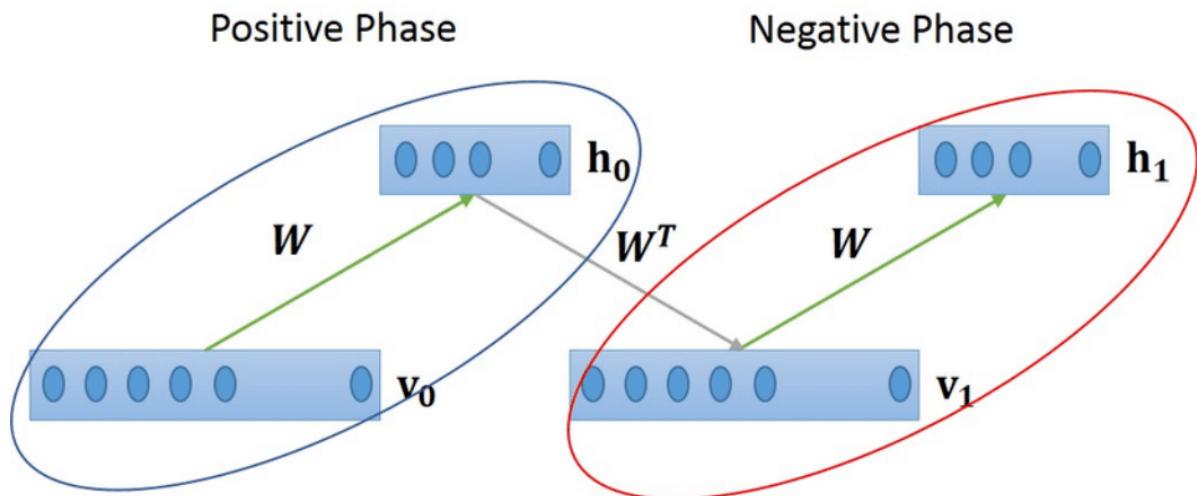
$$\Delta w_{ij} = \varepsilon (\langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^\infty)$$

- RBM Gibbs 1 步抽样：重建视角

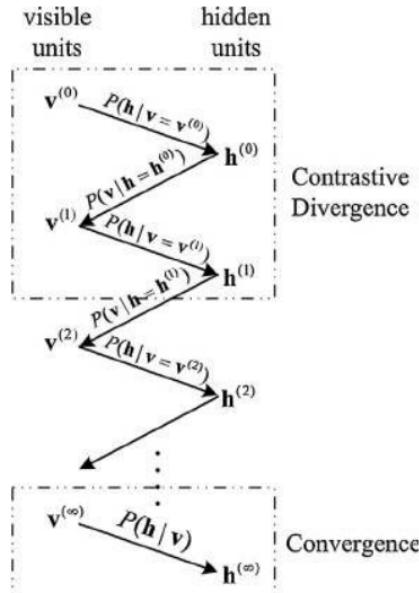


$$\Delta w_{ij} = \varepsilon (\langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1)$$

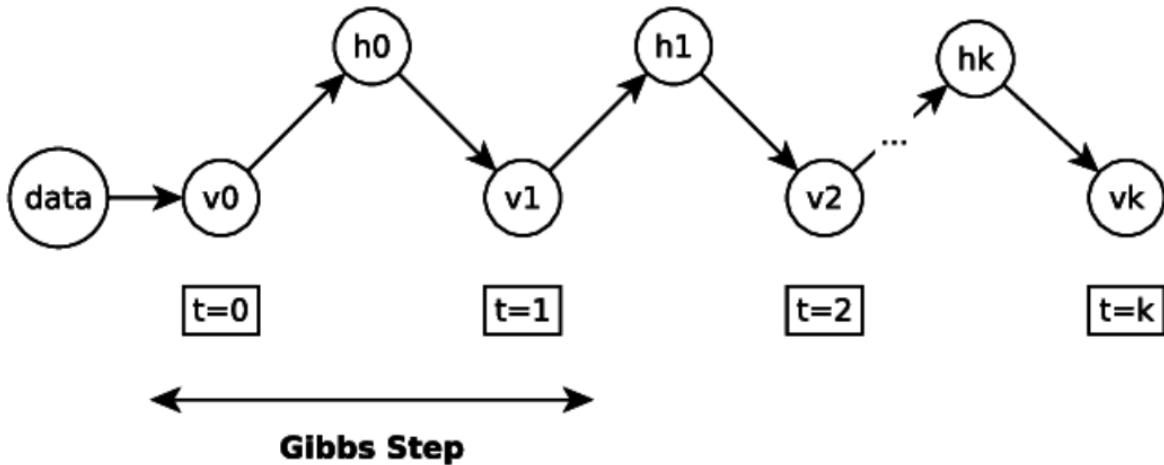
- Contrastive Divergence: Gibbs 1 步抽样



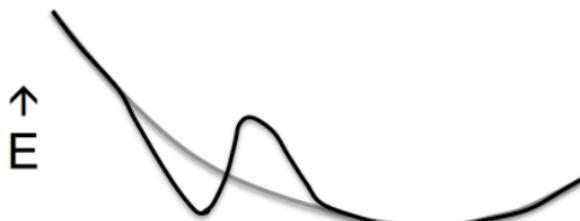
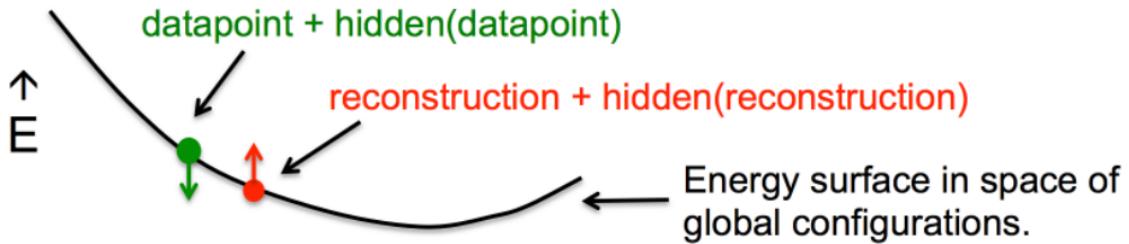
- Contrastive Divergence: Gibbs 1 步抽样



- CD-K: Gibbs K 步抽样



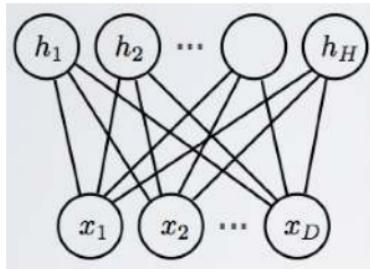
- CD-K: 全局优化观



Change the weights to pull the energy down at the datapoint.

Change the weights to pull the energy up at the reconstruction.

- CD 算法具体更新



$$\frac{\partial -\log p(\mathbf{x}^{(t)})}{\partial \theta} = \underbrace{\mathbb{E}_{\mathbf{h}} \left[\frac{\partial E(\mathbf{x}^{(t)}, \mathbf{h})}{\partial \theta} \mid \mathbf{x}^{(t)} \right]}_{\text{positive phase}} - \underbrace{\mathbb{E}_{\mathbf{x}, \mathbf{h}} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right]}_{\text{negative phase}}$$

$$E(\mathbf{v}, \mathbf{h}; \theta) = - \sum_{ij} W_{ij} v_i h_j - \sum_i b_i v_i - \sum_j a_j h_j$$

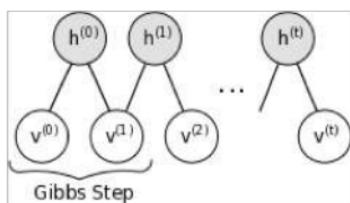
$$P(\mathbf{h}|\mathbf{v}) = \prod_i P(h_j|\mathbf{v}) \quad P(h_j=1|\mathbf{v}) = \frac{1}{1 + \exp(-\sum_i W_{ij} v_i - a_j)}$$

$$P(\mathbf{v}|\mathbf{h}) = \prod_i P(v_i|\mathbf{h}) \quad P(v_i=1|\mathbf{h}) = \frac{1}{1 + \exp(-\sum_j W_{ij} h_j - b_i)}$$

网络权重 $\mathbf{W} \Leftarrow \mathbf{W} + \alpha (\mathbf{h}(\mathbf{x}^{(t)}) \mathbf{x}^{(t)\top} - \mathbf{h}(\tilde{\mathbf{x}}) \tilde{\mathbf{x}}^\top)$

隐藏层偏差 $\mathbf{b} \Leftarrow \mathbf{b} + \alpha (\mathbf{h}(\mathbf{x}^{(t)}) - \mathbf{h}(\tilde{\mathbf{x}}))$

输入层偏差 $\mathbf{c} \Leftarrow \mathbf{c} + \alpha (\mathbf{x}^{(t)} - \tilde{\mathbf{x}})$



Sigmoid 函数

- Product of Experts (PoE) 框架

似然函数

$$P(x|\{\theta_j\}) = \frac{1}{Z} \prod_{j=1}^M f_j(x|\theta_j)$$

$$Z = \int dx \prod_{j=1}^M f_j(x|\theta_j)$$

Log似然函数

$$L(\{\theta_j\}| \{x_n\}) = \sum_{n=1}^N \sum_{j=1}^M \log f_j(x_n|\theta_j) - N \log Z$$

梯度

$$\nabla_j L = \sum_{n=1}^N \nabla_j \log f_j(x_n) - N \langle \nabla_j \log f_j(x) \rangle_{P(x)}$$

梯度下降

$$\theta_j \rightarrow \theta_j + \eta \nabla_j L$$

Restricted Boltzmann Machines (RBM)

$$P(x, h) = \frac{1}{Z} \exp \left(\sum_i \alpha_i x_i + \sum_j \beta_j h_j + \sum_{ij} W_{ij} x_i h_j \right)$$

$$P(x) = \frac{1}{Z} \prod_i \exp(\alpha_i x_i) \prod_j \left(1 + \exp(\beta_j + \sum_i W_{ij} x_i) \right)$$

Independent Components Analysis (ICA)

$$P(x|\{w_j\}) = |\det(W)| \prod_{j=1}^M p_j \left(\sum_i w_{ji} x_i \right)$$

- PoE vs MoE (Mixture of Experts)

Mixture Model

$$p(\vec{d} | \theta_1, \dots, \theta_n) = \sum_m \alpha_m f_m(\vec{d} | \theta_m)$$

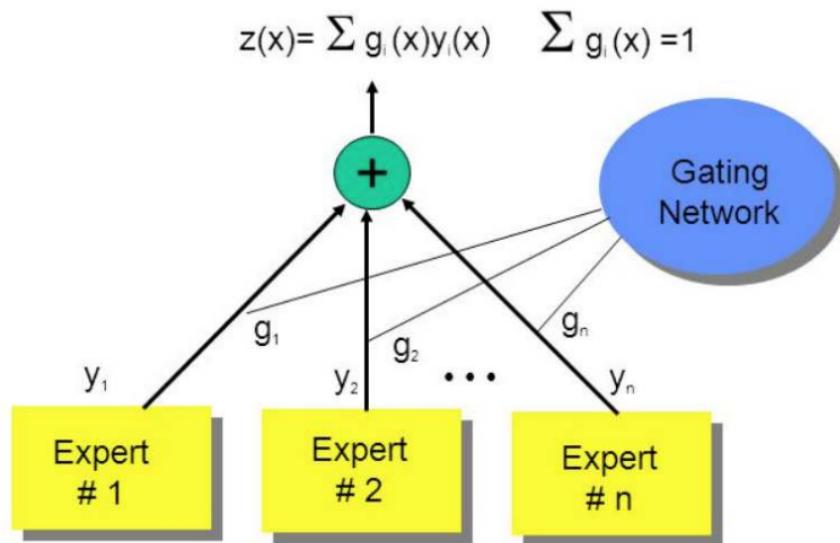
Use EM to learn parameters

Product of Experts

$$p(\vec{d} | \theta_1, \dots, \theta_n) = \frac{\prod_m f_m(\vec{d} | \theta_m)}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}$$

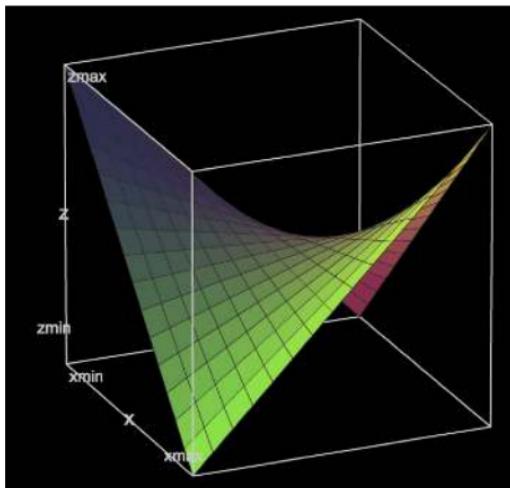
Use Contrastive Divergence to learn parameters.

- MoE (Mixture of Experts)

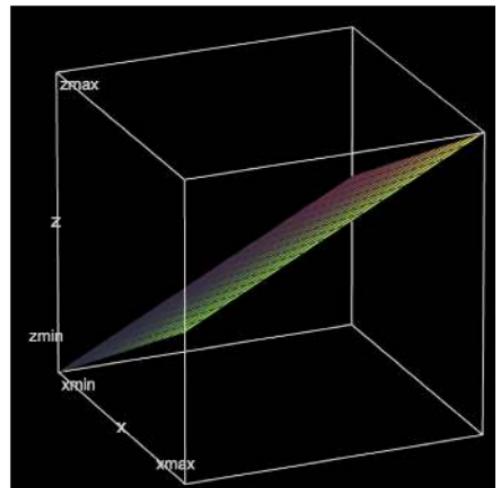


- PoE vs MoE

$$Z = X * Y$$

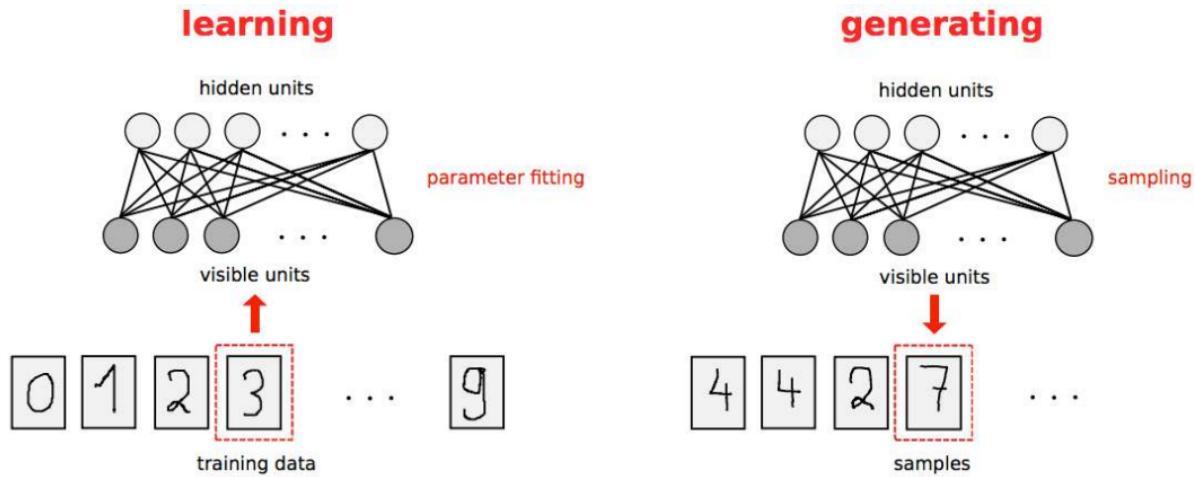


$$Z = X + Y$$



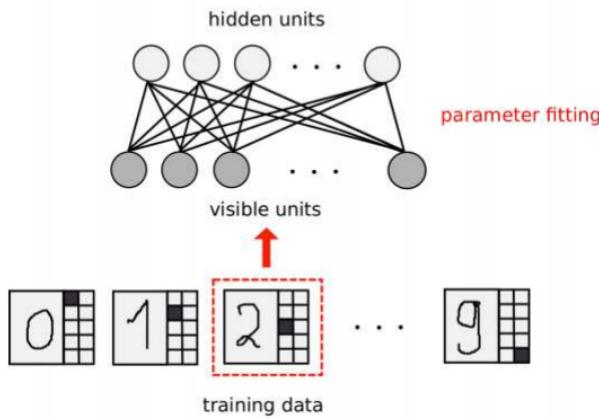
乘法比加法更为尖锐，加法类似投票

- RMB 做无监督学习

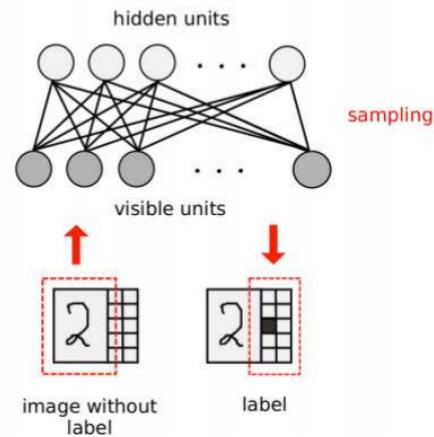


- RMB 做有监督学习

learning with labels

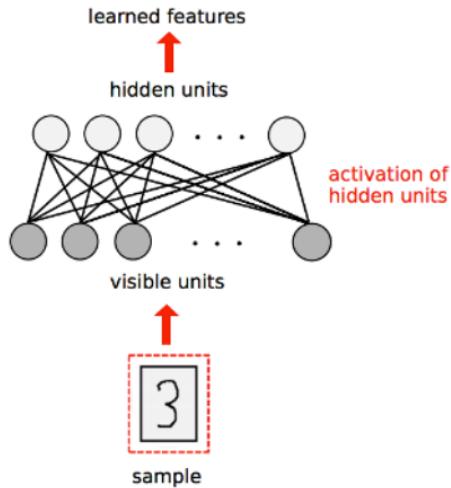


classification

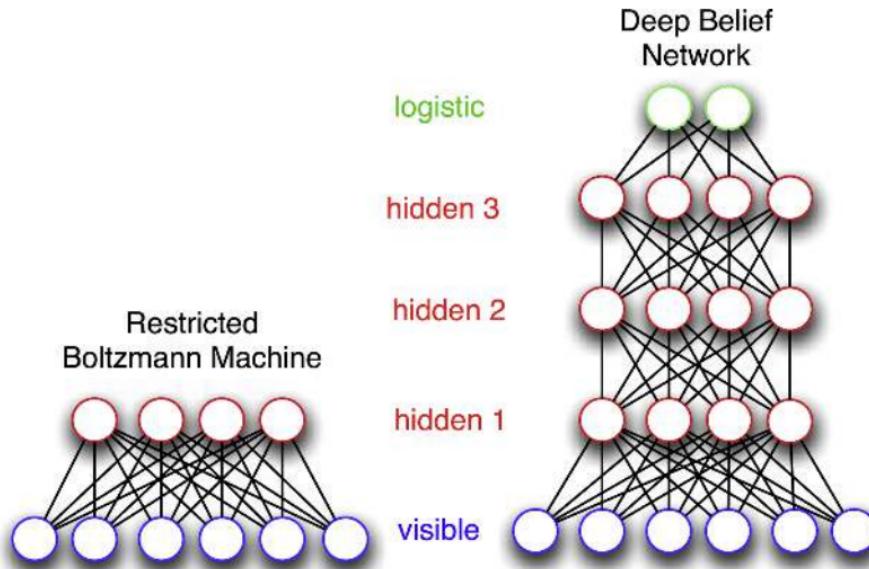


- RMB 做特征提取

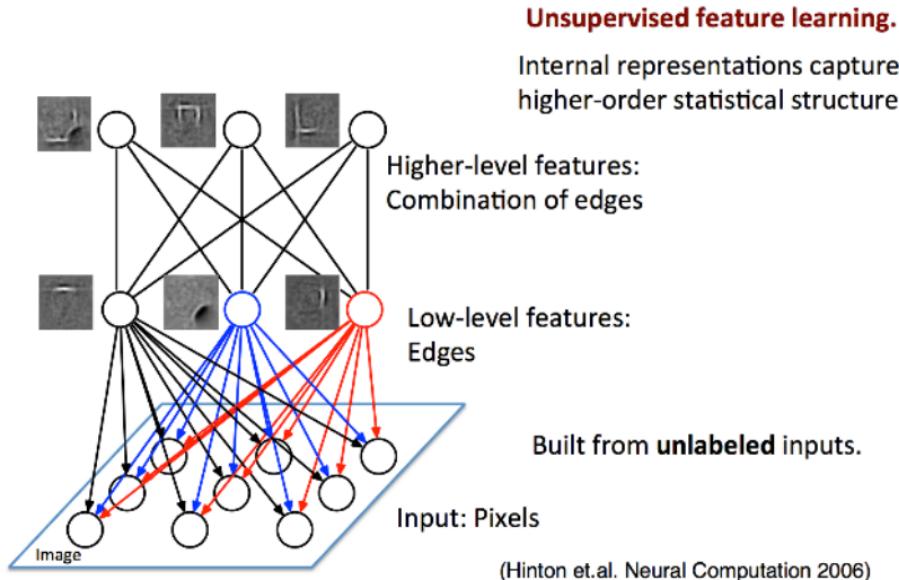
feature mapping



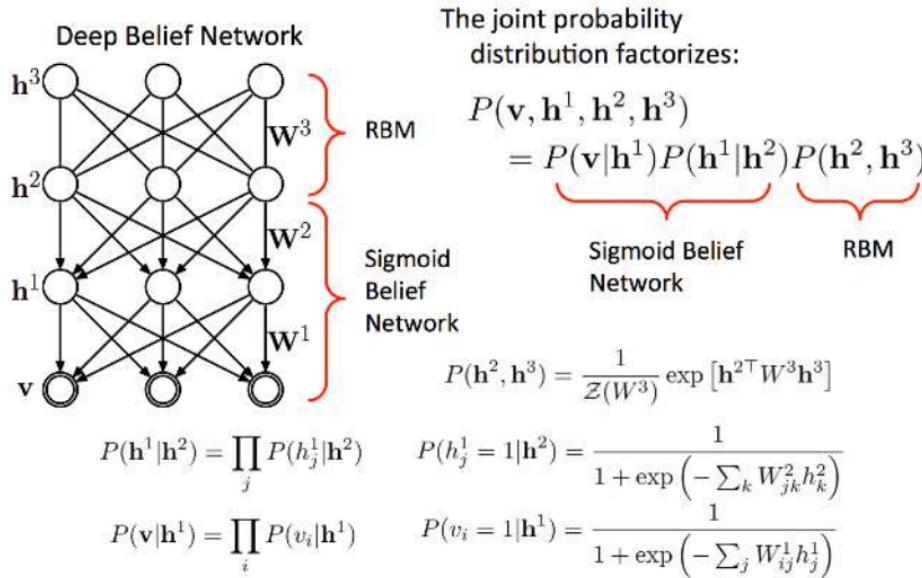
- Deep Belief Network (DBN)



- DBN 分层特征表示

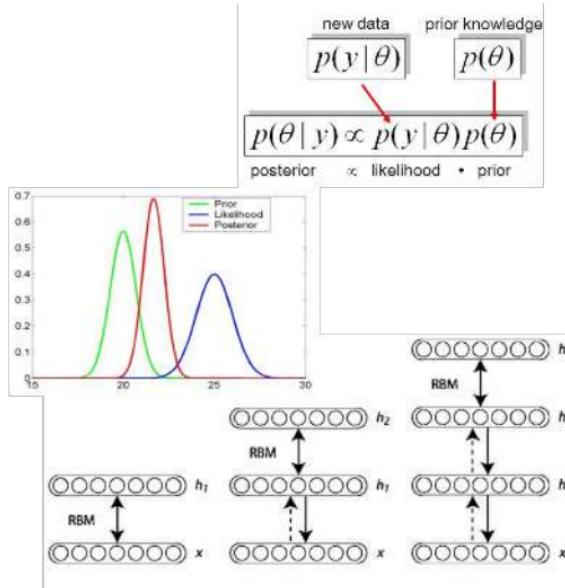
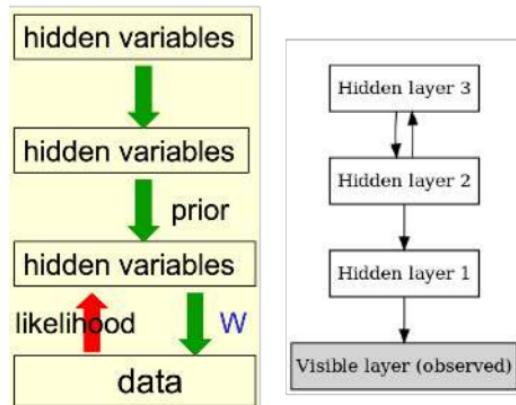


- DBN = Sigmoid Belief Network + RBM



- DBN 增强第一层隐藏层表达能力

第一层隐藏层，需要先验支撑
来更好的表达数据



- DBN Layer-Wise Training

$$P(v_i = 1 | \mathbf{h}) = \frac{1}{1 + \exp(-\sum_j W_{ij} h_j - b_i)}$$

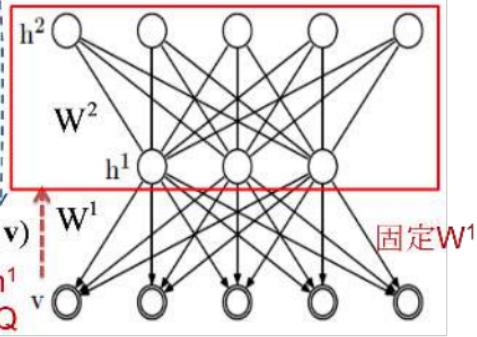
$$P(h_j = 1 | \mathbf{v}) = \frac{1}{1 + \exp(-\sum_i W_{ij} v_i - a_j)}$$

$$P(\mathbf{h} | \mathbf{v}) = \prod_i P(h_i | \mathbf{v})$$

两种近似 $Q(\mathbf{h} | \mathbf{v})$

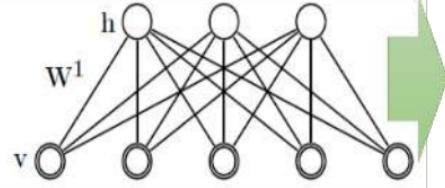
1. 直接抽样计算
2. 通过 sigmoid 输出计算均值

按RBM训练倒数第二层



按RBM训练最低层

近似

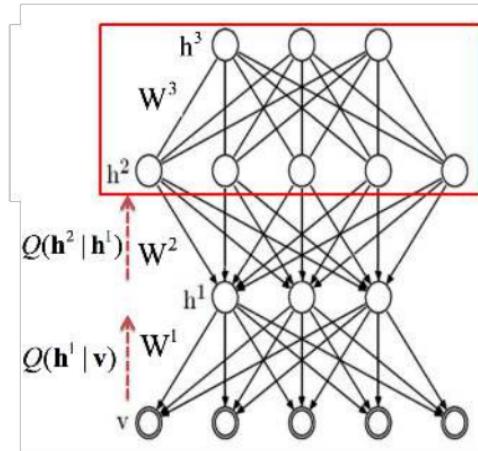
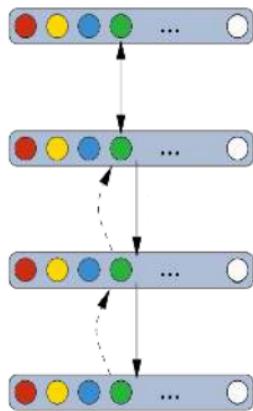


2. 均值计算
1. 抽样 h^1 获取 Q

- Greedy-Layer Wise Pre-Training

无监督训练深度网络Pre-training

1. 训练本身不需要监督数据
2. 逐层训练避免梯度消失



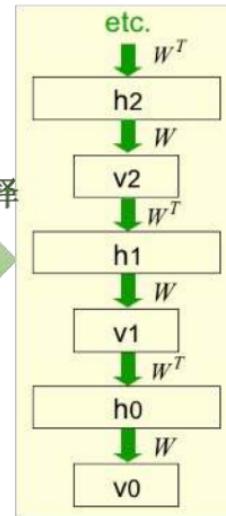
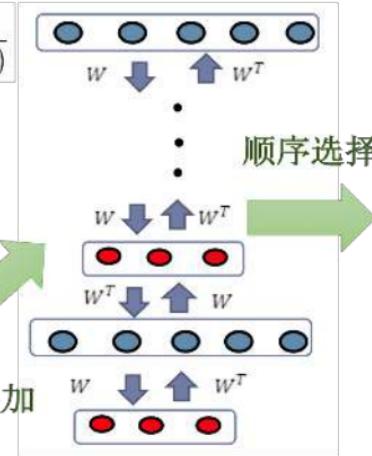
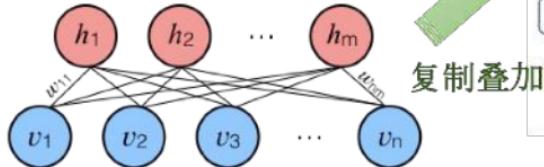
- RBM 等价无限 Sigmoid Belief Network

$$P(v_i = 1 | \mathbf{h}) = \frac{1}{1 + \exp(-\sum_j W_{ij} h_j - b_i)}$$

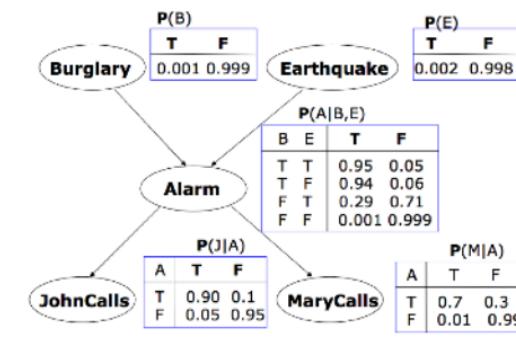
$$P(h_j = 1 | \mathbf{v}) = \frac{1}{1 + \exp(-\sum_i W_{ij} v_i - a_j)}$$

$$P(v_i | \mathbf{h}) = \text{sigm}(\sum_j w_{ij} h_j + b_i)$$

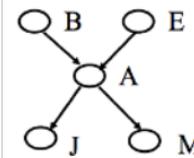
$$P(h_i | \mathbf{v}) = \text{sigm}(\sum_i w_{ij} v_i + c_j)$$



- 贝叶斯网络 (Bayes Belief Network)



$$P(X_1, X_2, \dots, X_n) = \prod_{i=1, \dots, n} P(X_i | pa(X_i))$$



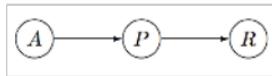
$$P(B=T, E=T, A=T, J=T, M=F) =$$

按信念Belief传播
的方式计算条件
概率的乘积

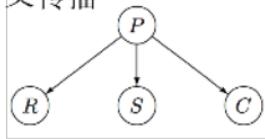
$$\begin{aligned}
 &= P(J=T | B=T, E=T, A=T, M=F) P(B=T, E=T, A=T, M=F) \\
 &= P(J=T | A=T) P(B=T, E=T, A=T, M=F) \\
 &\quad P(M=F | B=T, E=T, A=T) P(B=T, E=T, A=T) \\
 &\quad P(M=F | A=T) P(B=T, E=T, A=T) \\
 &\quad P(A=T | B=T, E=T) P(B=T, E=T) \\
 &\quad P(B=T) P(E=T) \\
 &= P(J=T | A=T) P(M=F | A=T) P(A=T | B=T, E=T) P(B=T) P(E=T)
 \end{aligned}$$

- Sigmoid Belief Network

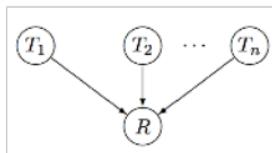
顺序传播



分叉传播



聚合传播



$$p(\mathbf{U}) = p(A_1 = a_1, A_2 = a_2, \dots, A_n = a_n) = \prod_i p(a_i | pa(A_i)).$$

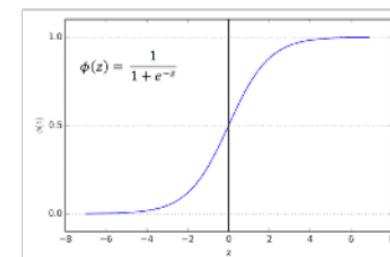
Binomial 二项分布

$$L(\boldsymbol{\mu} | Y) = \prod_{i=1}^n (1_{y_i=1}(\mu_i) + 1_{y_i=0}(1 - \mu_i))$$

$$\ln\left(\frac{P}{1-P}\right) = \alpha + \beta x \quad \frac{P}{1-P} = e^{\alpha+\beta x}$$

Sigmoid 激活函数

Sigmoid 信念网络

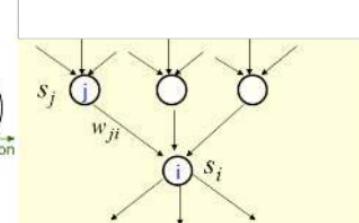
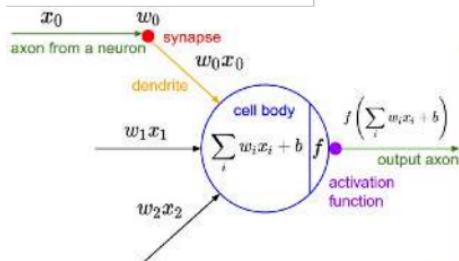


- Sigmoid Belief Network 计算

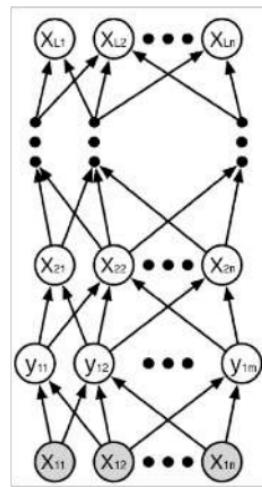
$$\mathbf{P}(X_1, X_2, \dots, X_n) = \prod_{i=1,..,n} \mathbf{P}(X_i | pa(X_i))$$

$$p(y_{k+1} = 1 | x_{k+1}) = \sigma(W^\top x_{k+1} + c)$$

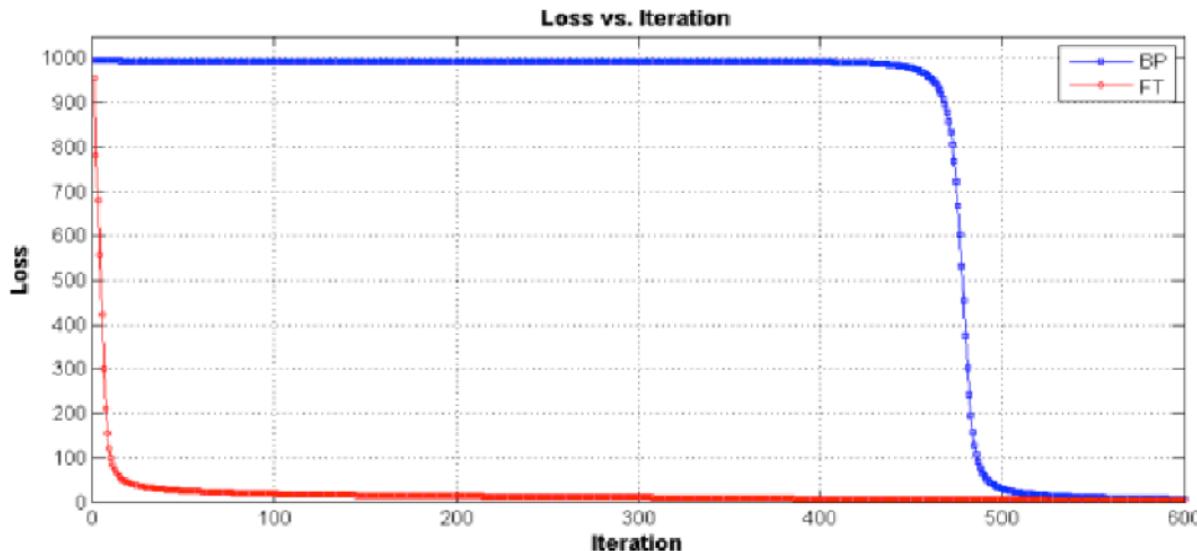
$$p(x_k = 1 | y_k) = \sigma(Wy_k + b)$$



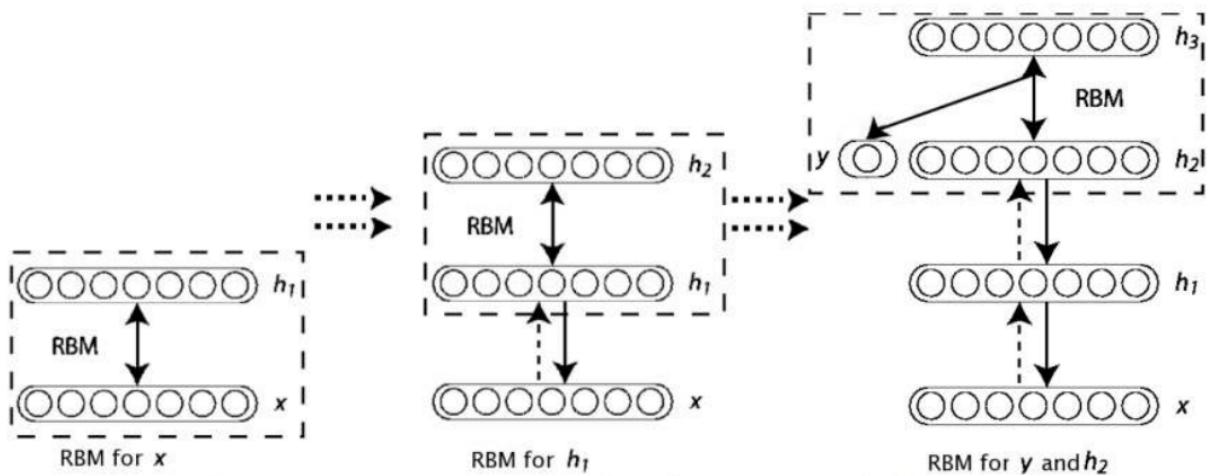
$$p_i \equiv p(s_i = 1) = \frac{1}{1 + \exp(-\sum_j s_j w_{ji})}$$



- Pre-Training 加速了收敛

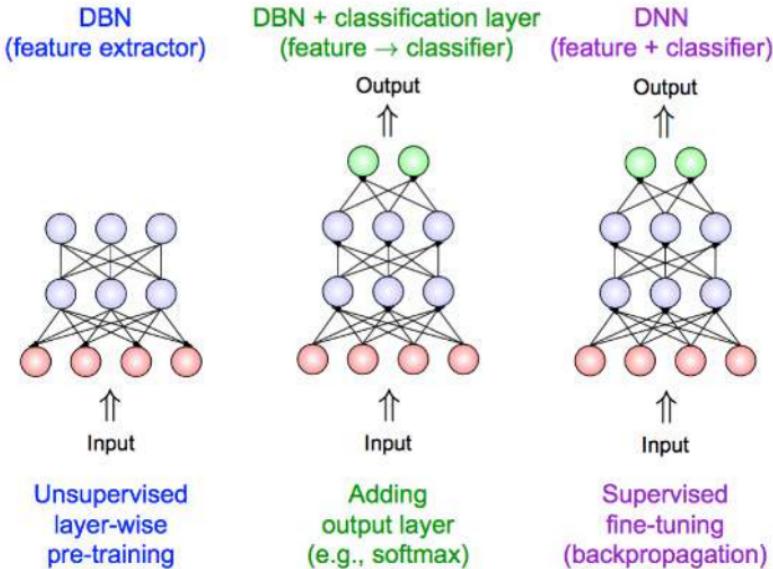


- DBN 做分类学习

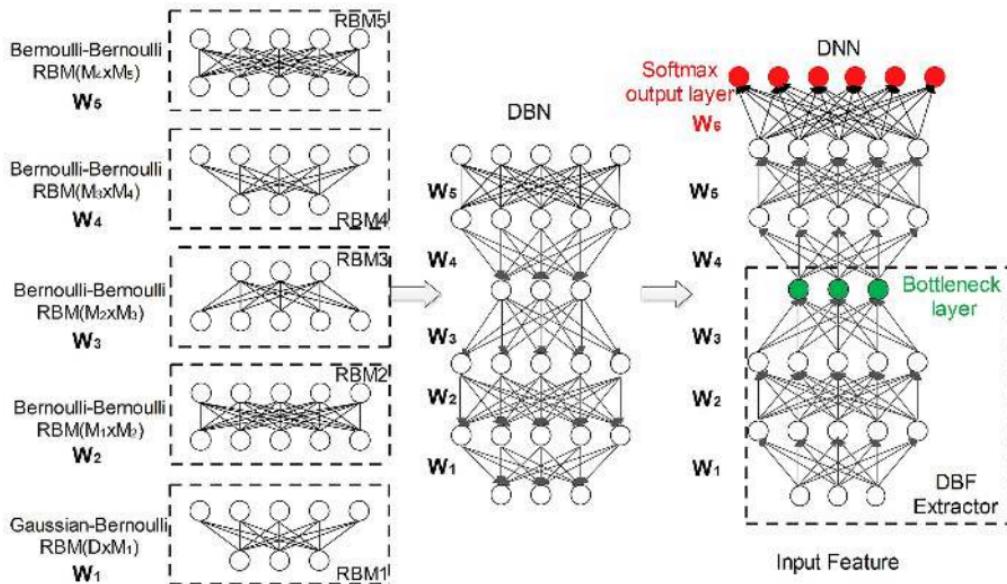


Stacking Restricted Boltzmann Machines (RBM) \rightarrow Deep Belief Network (DBN)

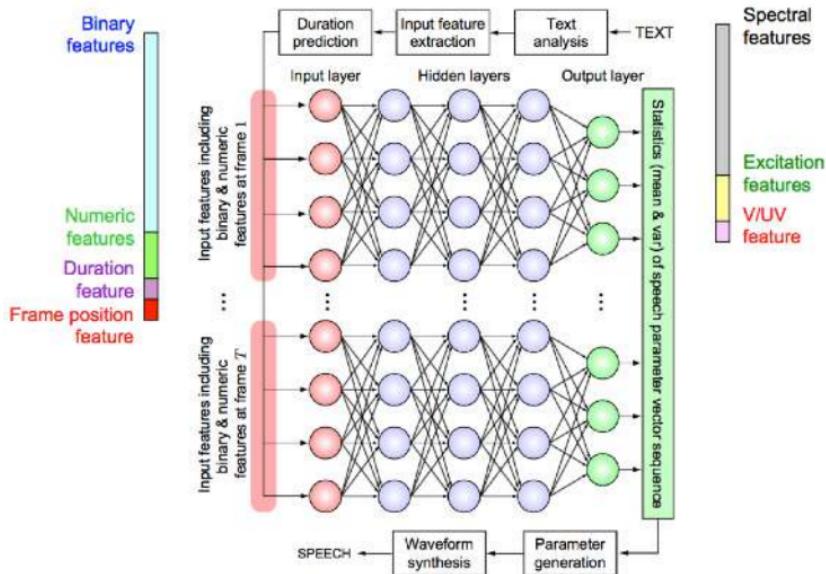
- DBN 做特征学习



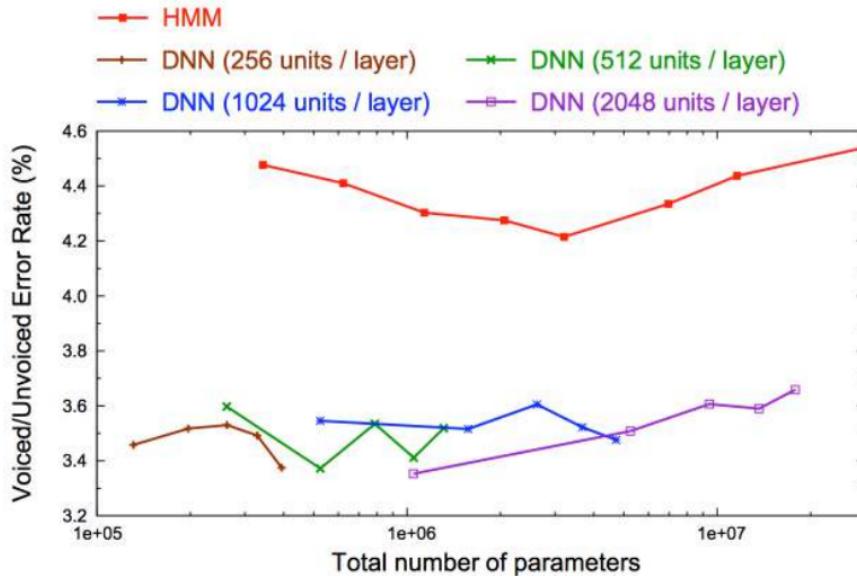
- DNN Deep Neural Network



- DNN 在语言识别

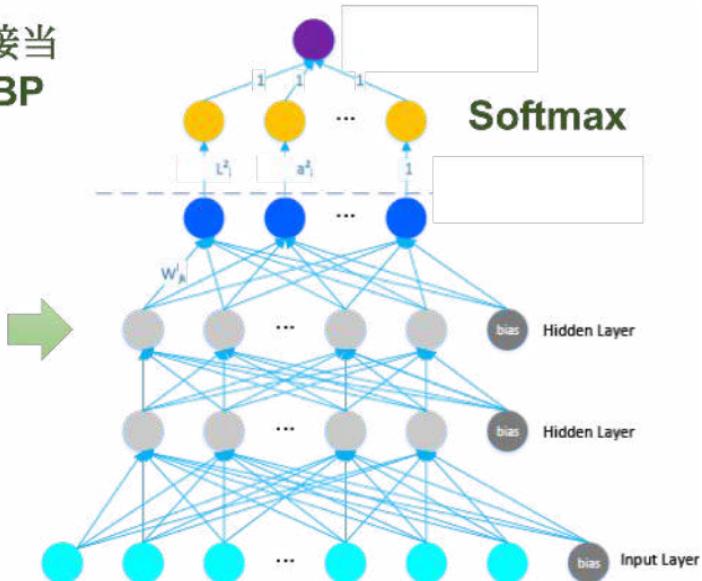
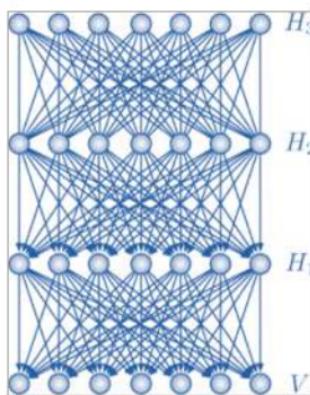


- DNN 在语言识别大获成功

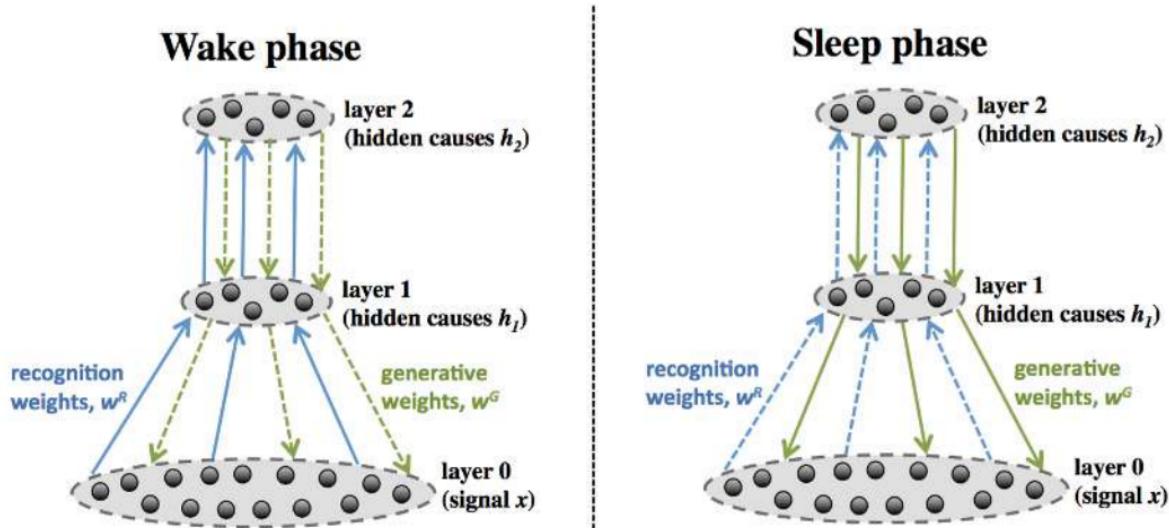


- Fine-Tuning 有监督学习：BP 算法

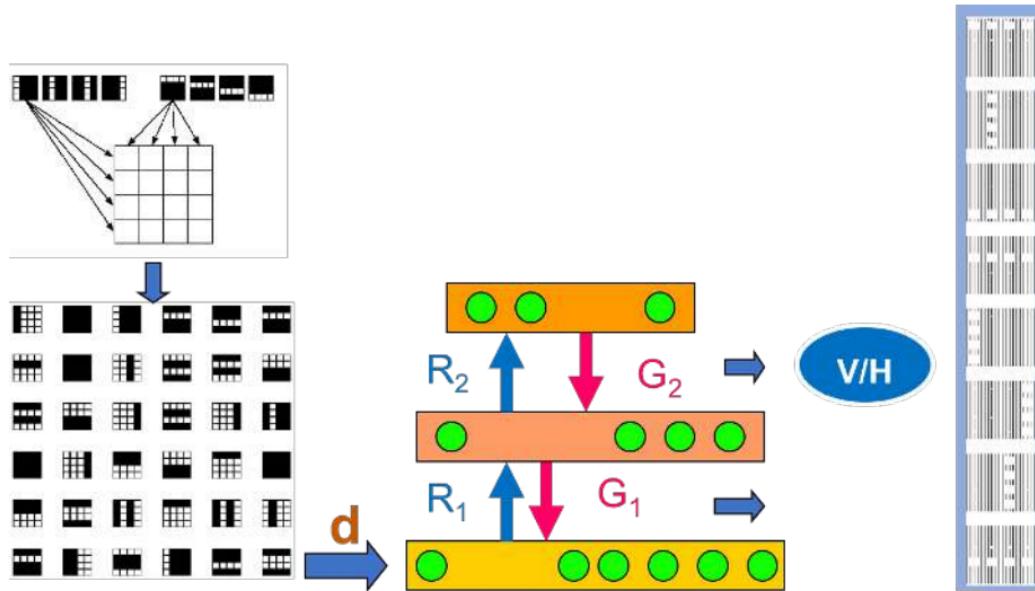
忘记DBN的结构，直接当成DNN (MLP) 利用BP
进行Fine-Tuning



- Fine-Tuning 无监督学习: wake-sleep



- wake-sleep 算法



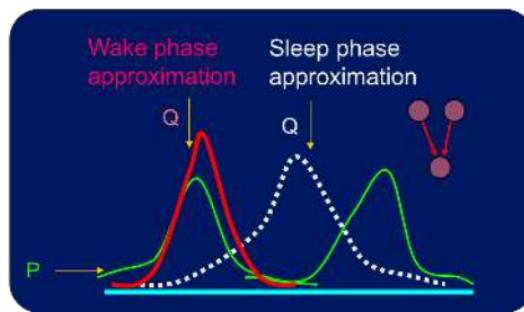
- wake-sleep 算法: $\text{KL}(P, Q)$ vs $\text{KL}(Q, P)$

$$\log P[d; \mathcal{G}] \geq \boxed{\log P[d; \mathcal{G}] - \text{KL}(Q[h|d; \mathcal{R}], P[h|d; \mathcal{G}])}$$

$\equiv -\mathcal{F}(d; \mathcal{R}, \mathcal{G})$

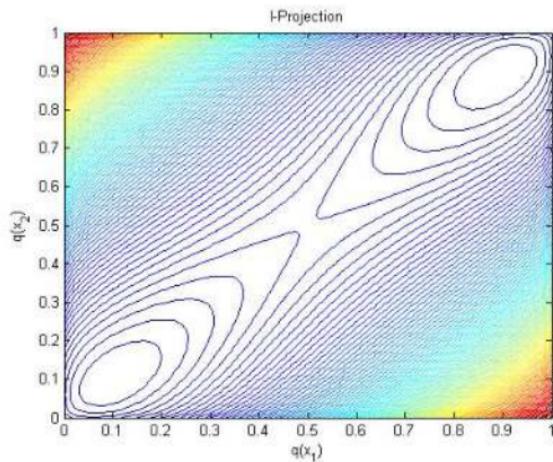
Free energy

$$\mathcal{F}(d; \mathcal{R}, \mathcal{G}) = -\log P[d; \mathcal{G}] + \text{KL}(Q[h|d; \mathcal{R}], P[h|d; \mathcal{G}])$$

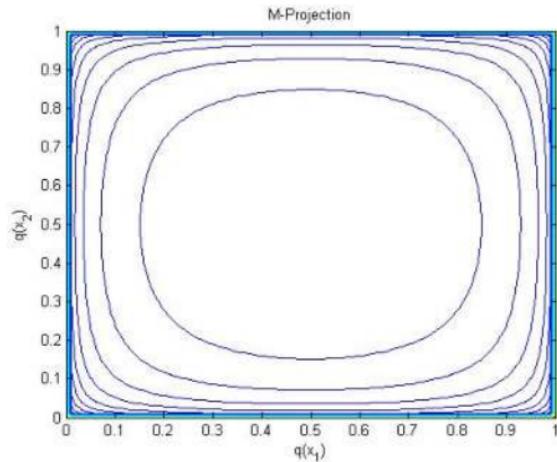


- $\text{KL}(P, Q)$: I-Projection vs $\text{KL}(Q, P)$: M-Projection

I-Projection



M-Projection



- wake-sleep 算法：认知和生成

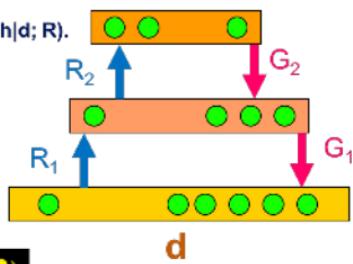
Wake phase

$$\Delta \mathcal{G} \propto -\nabla_{\mathcal{G}} \mathcal{F}(d; \mathcal{R}, \mathcal{G})$$

$$\begin{aligned} \Delta g_j^y &\propto (y_j^\circ - \hat{y}_j(x^\circ)), \quad \Delta G_{ij}^{xy} \propto x_i^\circ (y_j^\circ - \hat{y}_j(x^\circ)) \\ \hat{y}_j(x^\circ) &= \sigma(g_j^y + \sum_i G_{ij}^{xy} x_i^\circ) \end{aligned}$$

Remind:: $\Delta \mathcal{G} \propto \nabla_{\mathcal{G}} \log P[d; \mathcal{G}] = \sum_h P[h|d; \mathcal{G}] \nabla_{\mathcal{G}} \log P[h; \mathcal{G}]$

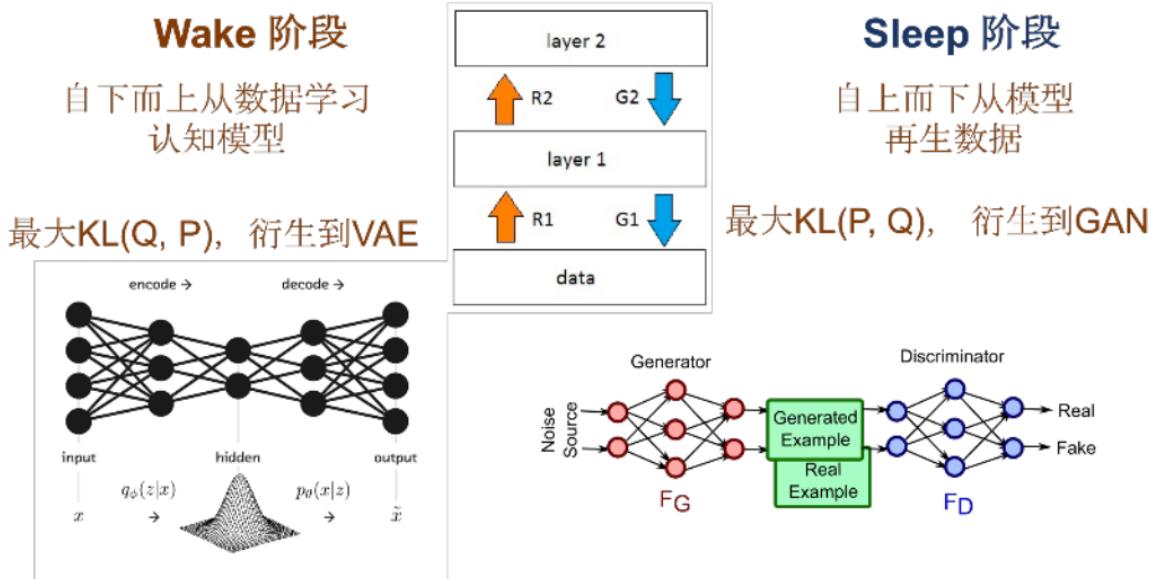
Replaced by



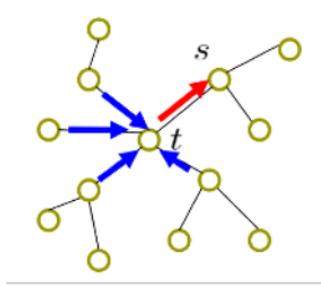
Sleep Phase

$$\begin{aligned} \Delta r_j^x &\propto (x_j^\bullet - \hat{x}_j(x^\bullet)), \quad \Delta R_{ij}^{yx} \propto y_i^\bullet (x_j^\bullet - \hat{x}_j(y^\bullet)) \\ \hat{x}_j(y^\bullet) &= \sigma(r_j^x + \sum_i R_{ij}^{yx} y_i^\bullet) \end{aligned}$$

- wake-sleep 最新衍生：VAE(认知) 和 GAN(生成)

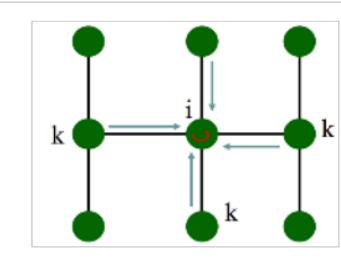
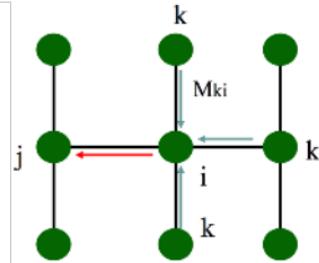


- Fine-Tuning 无监督学习: mean-field approximation



单向网络

$$p(x_1, \dots, x_m) = \frac{1}{Z} \prod_{s \in V} \psi_s(x_s) \prod_{(s,t) \in E} \psi_{st}(x_s, x_t)$$



多源单向网络

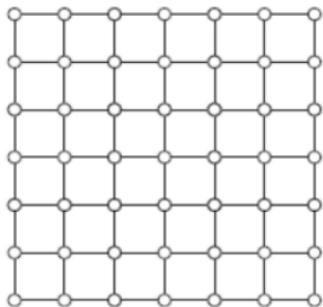
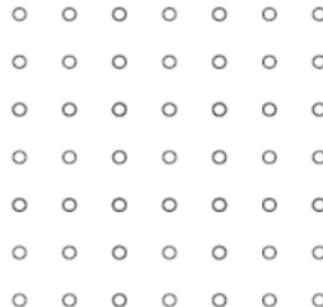
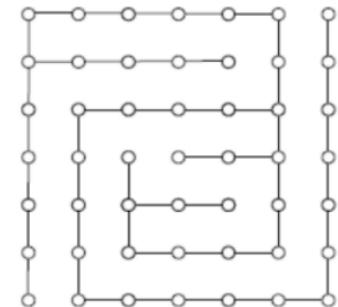
$$M_{i \rightarrow j}(x_j) \propto \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_k M_{k \rightarrow i}(x_k)$$

↑
external evidence
↑
Compatibilities (interactions)

双向网络

$$b_i(x_i) \propto \psi_i(x_i) \prod_k M_k(x_k)$$

- 什么是 mean-field approximation? 近似 $p(x)$

original G (Naïve) MF H_0 structured MF H_s 

$$p(x) \propto \prod_c \phi_c(x_c)$$

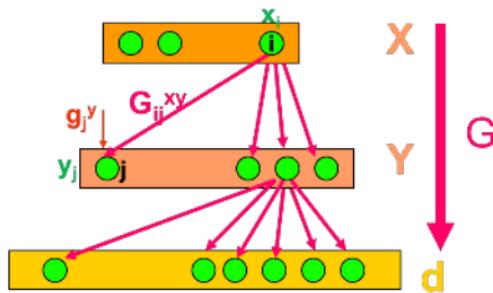
$$q(x) \propto \prod_i q_i(x_i)$$

$$q(x) \propto q_A(x_A) q_B(x_B)$$

找个可以分解的 $q(x)$ 使得 $\text{KL}(p(x)||q(x))$ 最小。

- 为什么可以使用 mean-field approximation?

$$p(y_j = 1) = \frac{1}{1 + \exp(-\sum_i x_i G_{ij}^{xy} - g_j^y)}$$

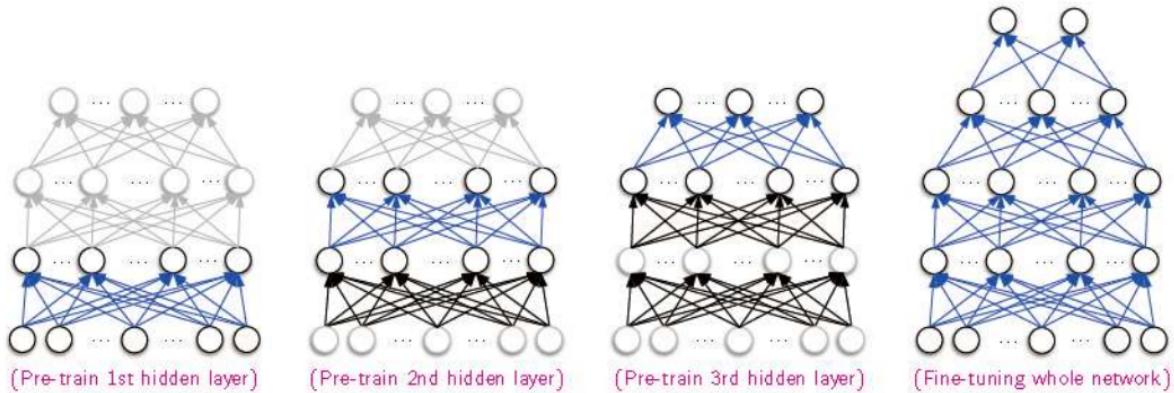


条件概率都是可以
由独立子单元乘积

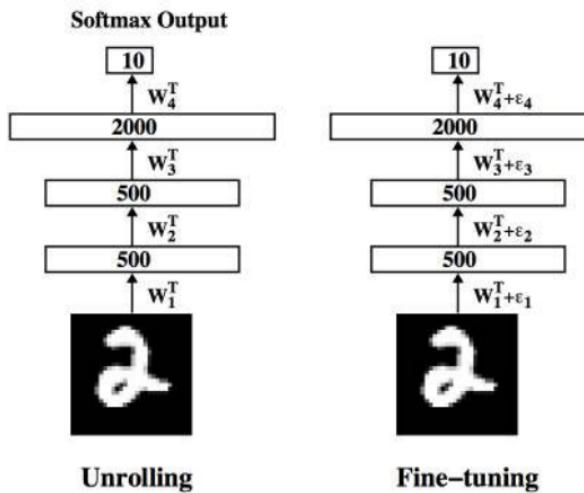
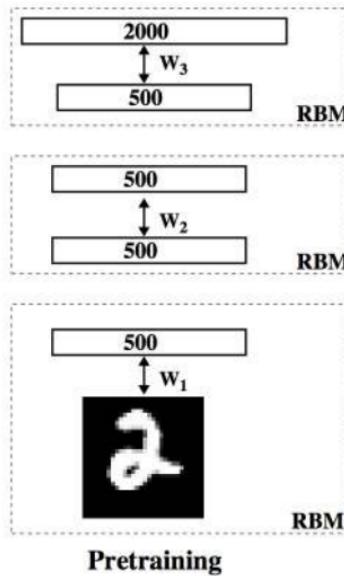
$$P[d|y; \mathcal{G}] = \prod_k P[d_k|y; \mathcal{G}]$$

$$P[x; \mathcal{G}] = \prod_i P[x_i; \mathcal{G}], \quad P[y|x; \mathcal{G}] = \prod_j P[y_j|x; \mathcal{G}]$$

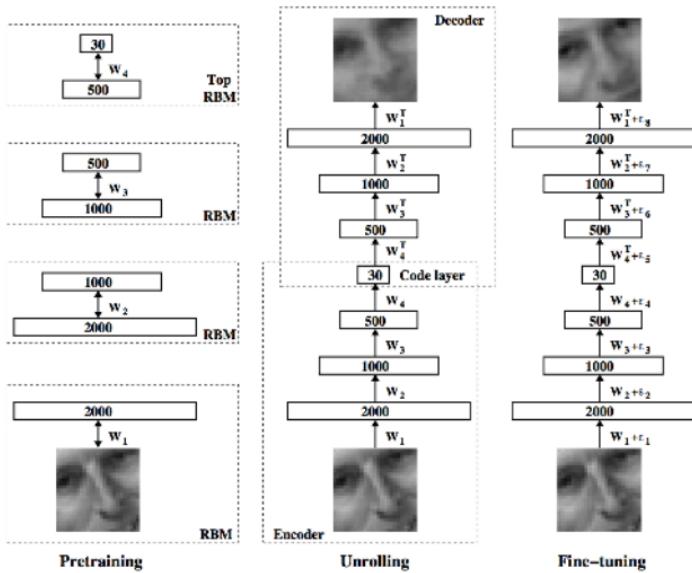
- DNN 有监督学习整个过程



- DNN 图像识别



- Deep Auto-Encoder



- DNN 研究带来了 ReLU 的诞生

Sigmoid → 权重共享

不同偏差：

$$-0.5, -1.5, -2.5, \dots - (N - 0.5)$$

概率求和：

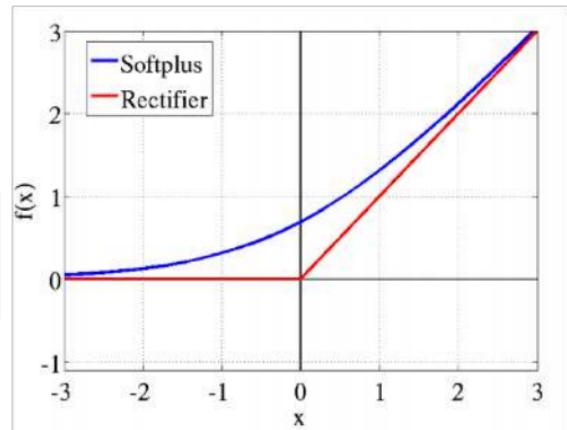
$$\sum_{i=1}^{\infty} \sigma(x - i + 0.5) \approx \log(1 + e^x)$$

softplus 函数和导数

$$f(x) = \ln[1 + \exp(x)]$$

$$f'(x) = \exp(x)/[\exp(x) + 1] = 1/[1 + \exp(-x)]$$

ReLU Rectified linear units

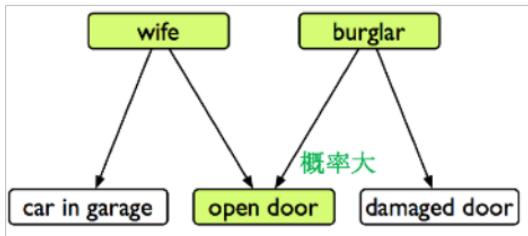


- DBN 潜在问题

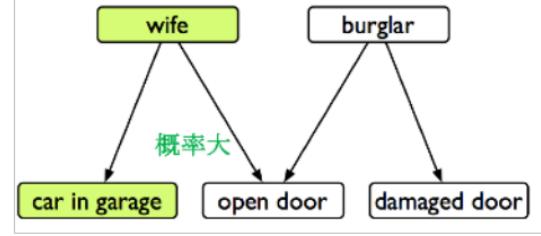
1. 由于是单向网络，存在“explaining away”的问题
 - 对于不确定性数据或者模糊数据潜在问题
 - 稳定性不强
2. 贪心逐层优化，很难全局最优
3. 只有前向推理，没有反向推理

- explaining away 现象

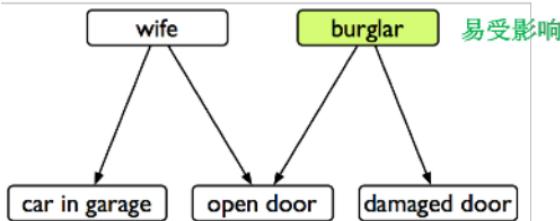
开门的两种可能性



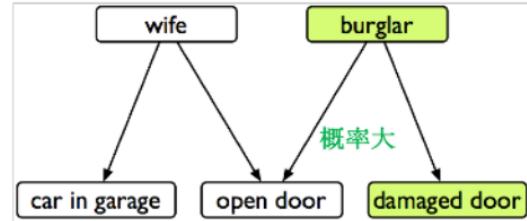
如果一个相邻条件独立现象成立



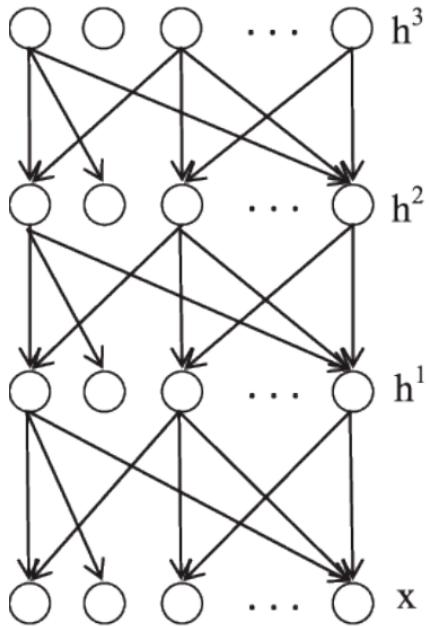
会对最初的理解改动



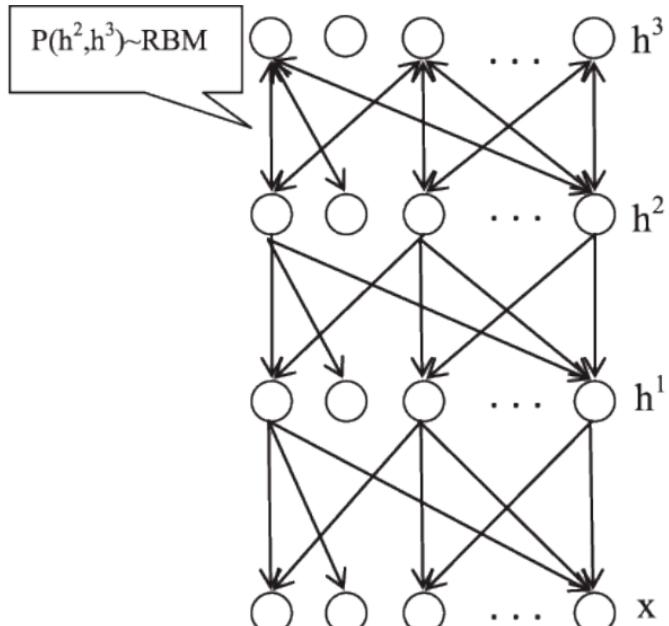
如果一个相邻条件独立现象成立



- Deep Boltzmann Machine (DBM)



Sigmoid Belief Networks

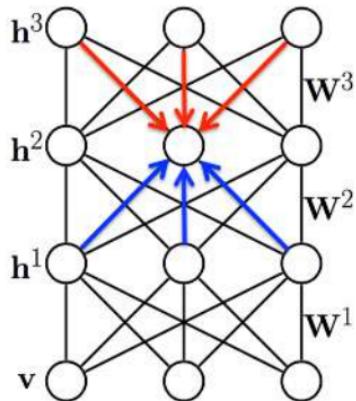


Deep Belief Networks

- DBM: 一种随机场 Random Field 无向图

$$P_{\theta}(\mathbf{v}) = \frac{P^*(\mathbf{v})}{\mathcal{Z}(\theta)} = \frac{1}{\mathcal{Z}(\theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3} \exp \left[\mathbf{v}^\top W^1 \mathbf{h}^1 + \underline{\mathbf{h}^{1\top} W^2 \mathbf{h}^2} + \underline{\mathbf{h}^{2\top} W^3 \mathbf{h}^3} \right]$$

Deep Boltzmann Machine


 $\theta = \{W^1, W^2, W^3\}$ model parameters

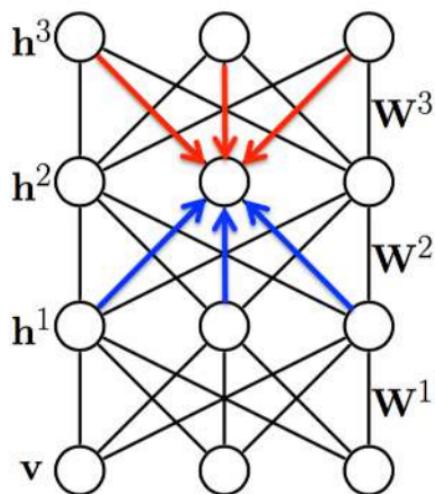
- Dependencies between hidden variables.
- All connections are undirected.
- Bottom-up and Top-down:

$$P(h_j^2 = 1 | \mathbf{h}^1, \mathbf{h}^3) = \sigma \left(\sum_k W_{kj}^3 h_k^3 + \sum_m W_{mj}^2 h_m^1 \right)$$

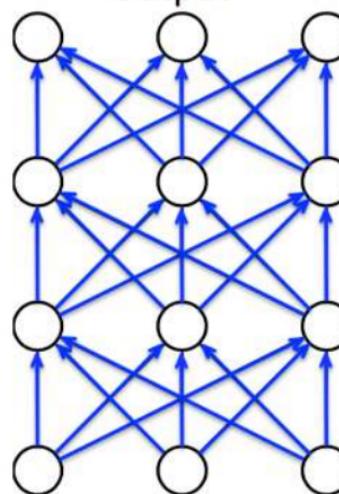
Top-down Bottom-up

- 神经网络、概率图模型 (贝叶斯网络、随机场)

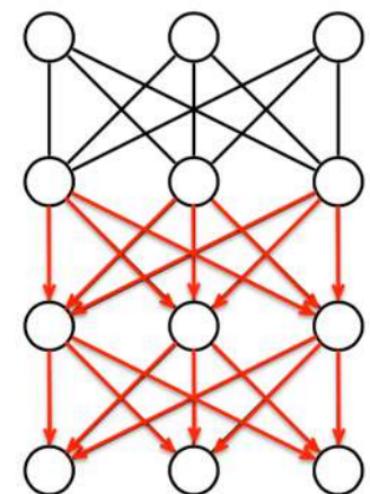
Deep Boltzmann Machine



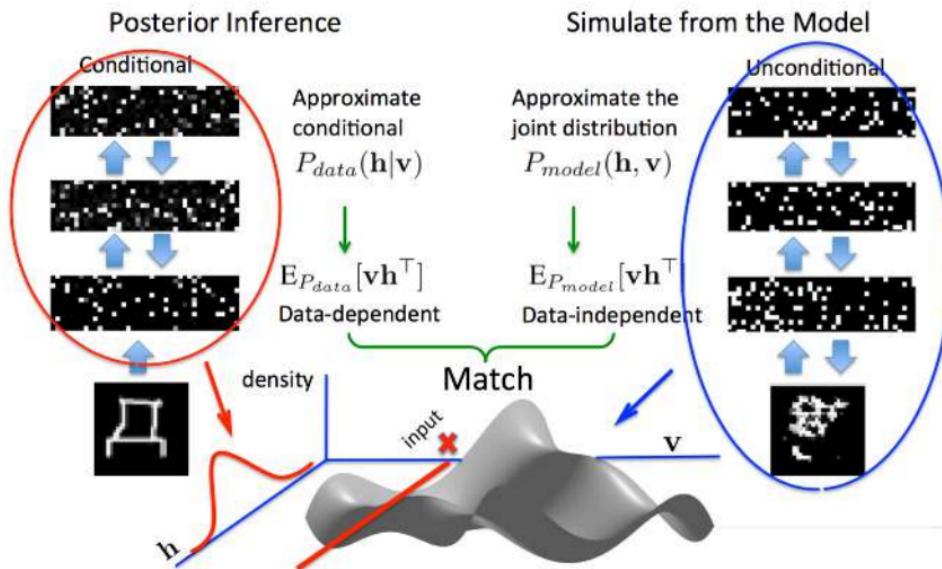
Neural Network Output



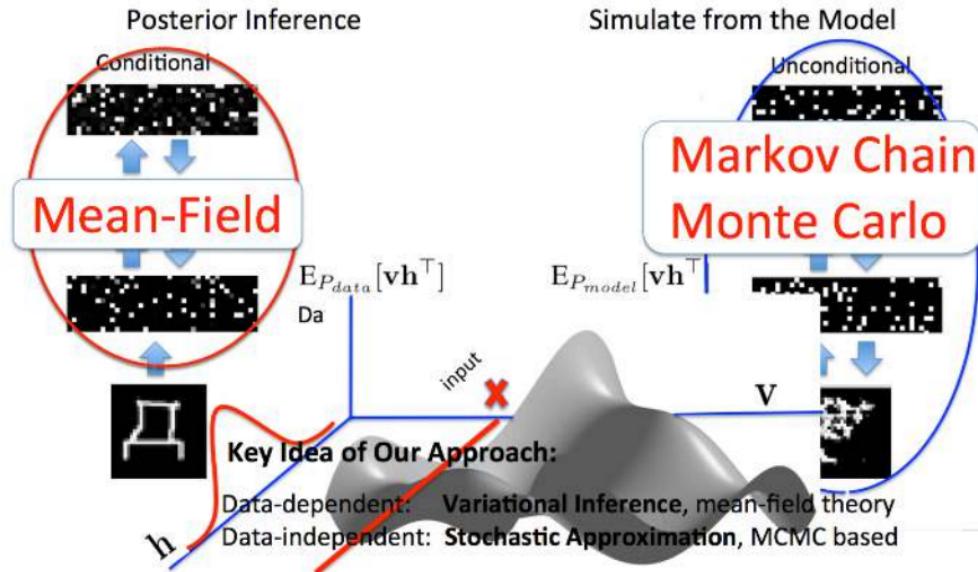
Deep Belief Network



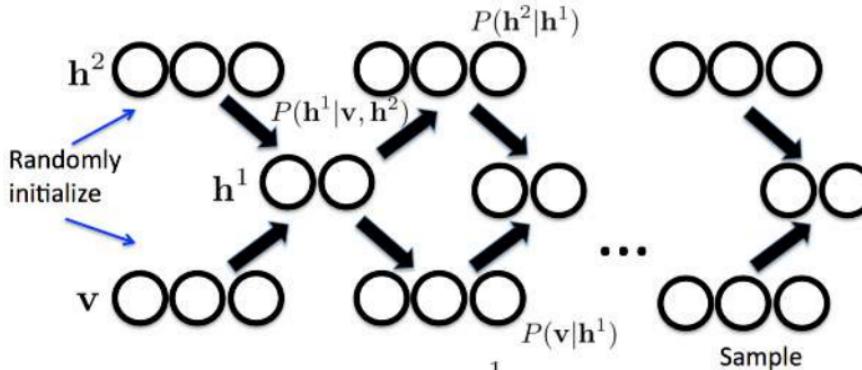
- 推理和数据重建都有



- 变分推理和 MCMC 方法



- 随机采样计算 CD 算法



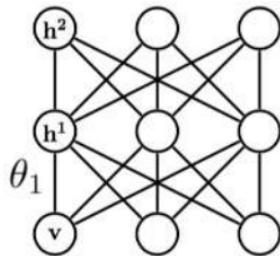
$$P(h_m^1 = 1 | v, h^2) = \frac{1}{1 + \exp(-\sum_i W_{im}^1 v_i - \sum_j W_{mj}^2 h_j^2)}$$

$$P(h_j^2 = 1 | h^1) = \frac{1}{1 + \exp(-\sum_m W_{mj}^2 h_m^1)}$$

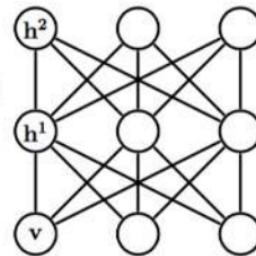
$$P(v_i = 1 | h^1) = \frac{1}{1 + \exp(-\sum_m W_{im}^1 h_m^1)}$$

- 顺序逐层更新

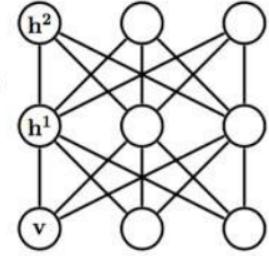
Time t=1

Update θ_1

t=2

Update θ_2

t=3



$$\mathbf{x}_1 \sim T_{\theta_1}(\mathbf{x}_1 \leftarrow \mathbf{x}_0)$$

$$\mathbf{x}_2 \sim T_{\theta_2}(\mathbf{x}_2 \leftarrow \mathbf{x}_1)$$

$$\mathbf{x}_3 \sim T_{\theta_3}(\mathbf{x}_3 \leftarrow \mathbf{x}_2)$$

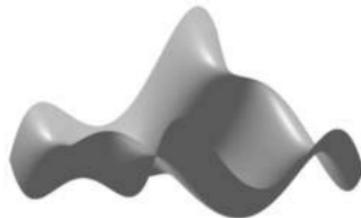
Update θ_t and \mathbf{x}_t sequentially, where $\mathbf{x} = \{v, h^1, h^2\}$

- CD 近似算法更新公式

Update rule decomposes:

$$\theta_{t+1} = \theta_t + \alpha_t \left(\underbrace{E_{P_{data}}[\mathbf{vh}^\top] - E_{P_{\theta_t}}[\mathbf{vh}^\top]}_{\text{True gradient}} \right) + \alpha_t \left(\underbrace{E_{P_{\theta_t}}[\mathbf{vh}^\top] - \frac{1}{M} \sum_{m=1}^M \mathbf{v}_t^{(m)} \mathbf{h}_t^{(m)\top}}_{\text{Noise term } \epsilon_t} \right)$$

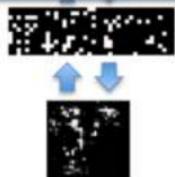
Almost sure convergence guarantees as learning rate $\alpha_t \rightarrow 0$



Problem: High-dimensional data:
the energy landscape is highly
multimodal

Key insight: The transition operator can be
any valid transition operator – Tempered
Transitions, Parallel/Simulated Tempering.

Markov Chain
Monte Carlo

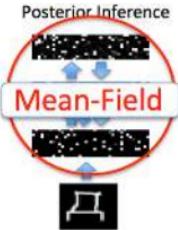


- 变分推理公式

Approximate intractable distribution $P_\theta(\mathbf{h}|\mathbf{v})$ with simpler, tractable distribution $Q_\mu(\mathbf{h}|\mathbf{v})$:

$$\text{KL}(Q||P) = \int Q(x) \log \frac{Q(x)}{P(x)} dx$$

$$\log P_\theta(\mathbf{v}) \geq \log P_\theta(\mathbf{v}) - \text{KL}(Q_\mu(\mathbf{h}|\mathbf{v}) || P_\theta(\mathbf{h}|\mathbf{v}))$$

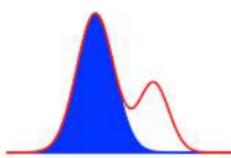


Variational Lower Bound

Mean-Field: Choose a fully factorized distribution:

$$Q_\mu(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^F q(h_j|\mathbf{v}) \text{ with } q(h_j = 1|\mathbf{v}) = \mu_j$$

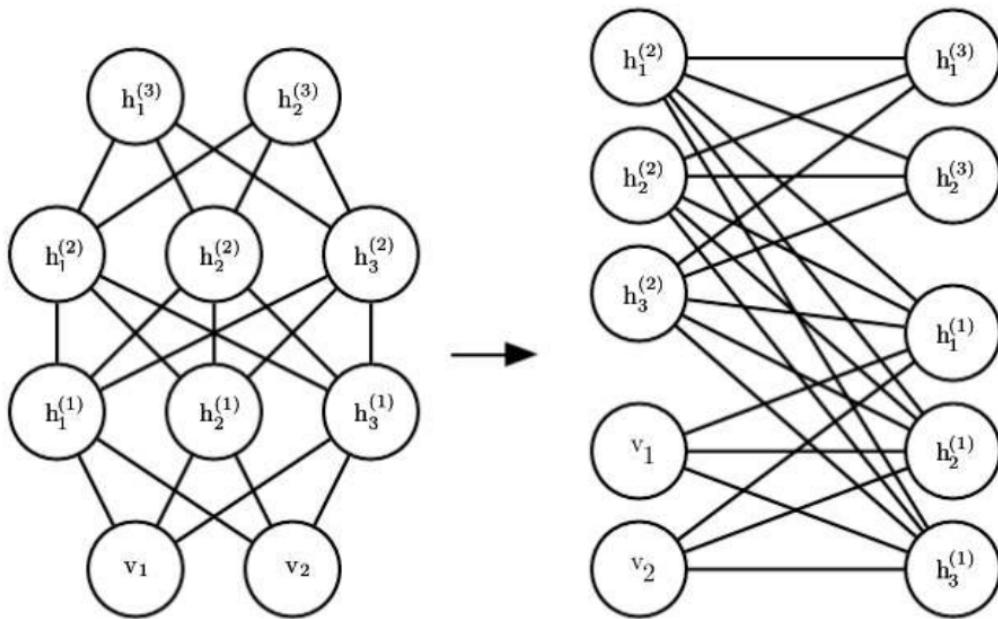
Variational Inference: Maximize the lower bound w.r.t. Variational parameters μ .



Nonlinear fixed-point equations:

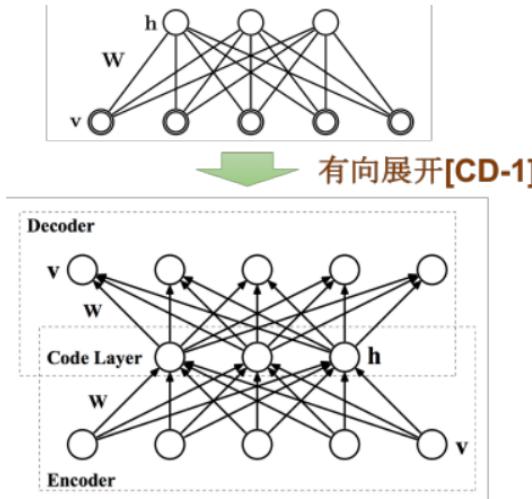
$$\begin{aligned}\mu_j^{(1)} &= \sigma \left(\sum_i W_{ij}^1 v_i + \sum_k W_{jk}^2 \mu_k^{(2)} \right) \\ \mu_k^{(2)} &= \sigma \left(\sum_j W_{jk}^2 \mu_j^{(1)} + \sum_m W_{km}^3 \mu_m^{(3)} \right) \\ \mu_m^{(3)} &= \sigma \left(\sum_k W_{km}^3 \mu_k^{(2)} \right)\end{aligned}$$

- DBM 的 RBM 视角

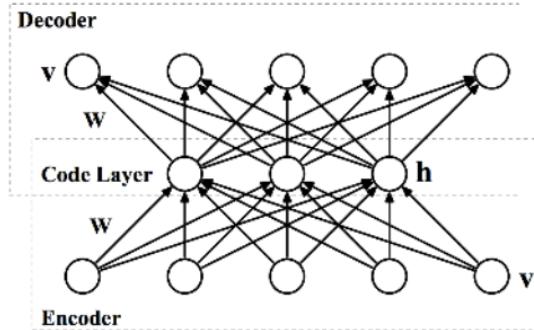


3 自动编码器 AE

- 从 RBM 到 AutoEncoder：一次重建



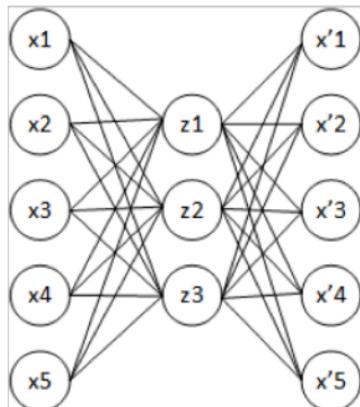
- 自动编码器 AutoEncoder



$$\text{Encoder: } h_j = \frac{1}{1 + \exp(-\sum_i v_i W_{ij})}, j = 1, \dots, K.$$

$$\text{Decoder: } \hat{v}_i = \frac{1}{1 + \exp(-\sum_j h_j W_{ij})}, i = 1, \dots, D.$$

- 重建角度下的 AE



$$z = f(Wx) \quad | \quad y = g(Vz)$$

理想重建：

$$J = \sum_{n=1}^N |x(n) - \hat{x}(n)|^2$$

$$\hat{x}(n) = Q^{-1}Qx(n)$$

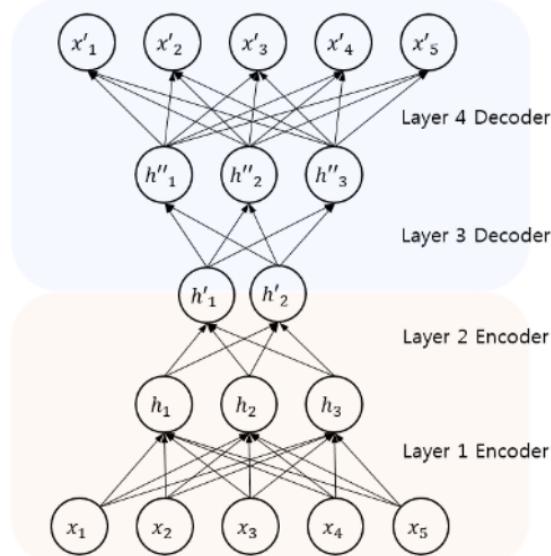
$$J = \sum_{n=1}^N |x(n) - Q^{-1}Qx(n)|^2$$

现实重建：(PCA 重建)

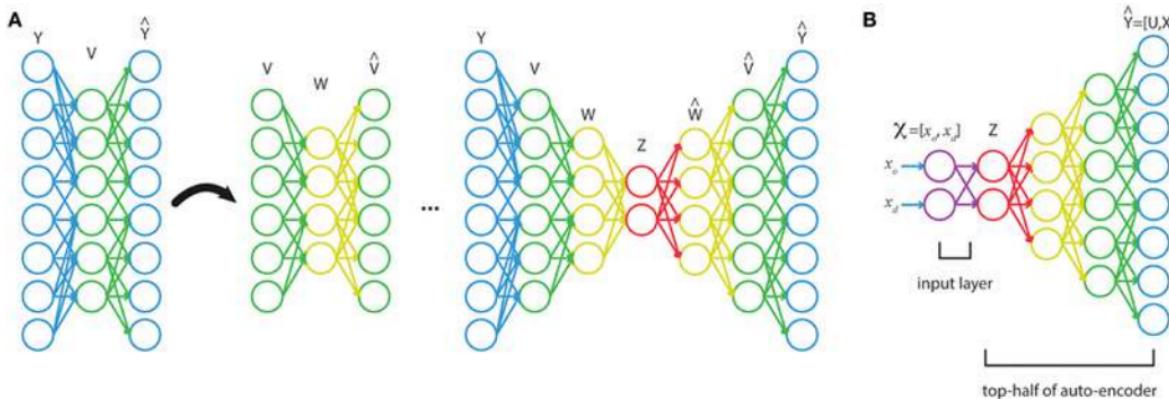
$$y = g(Vf(Wx)) = VWx$$

$$J = \sum_{n=1}^N |x(n) - VWx(n)|^2$$

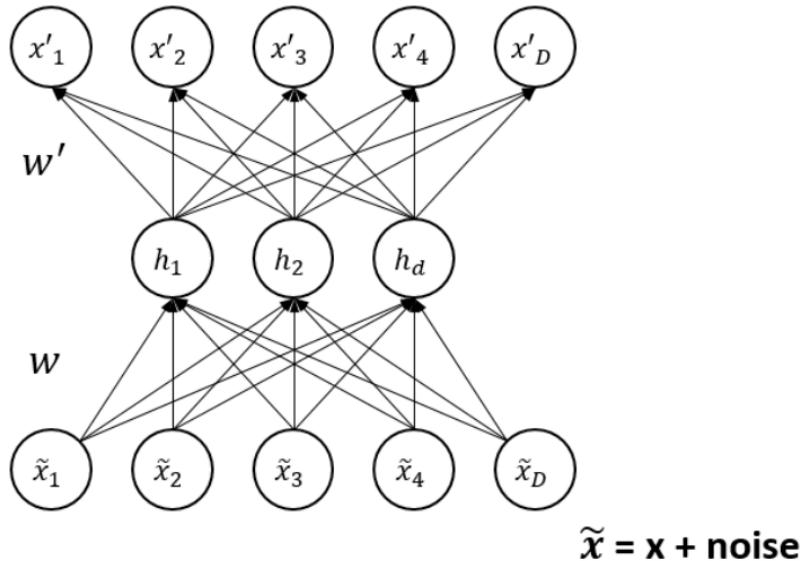
- 堆栈 (深度)AE: DBM 的有向展开



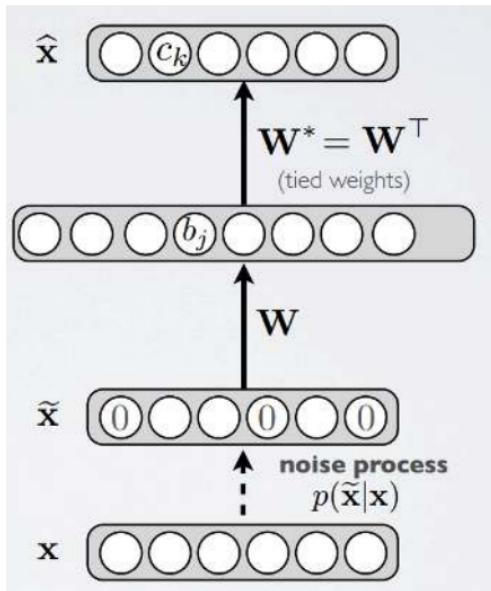
- Stacked AE 的对称逐层训练：也可以看成两层 DBN



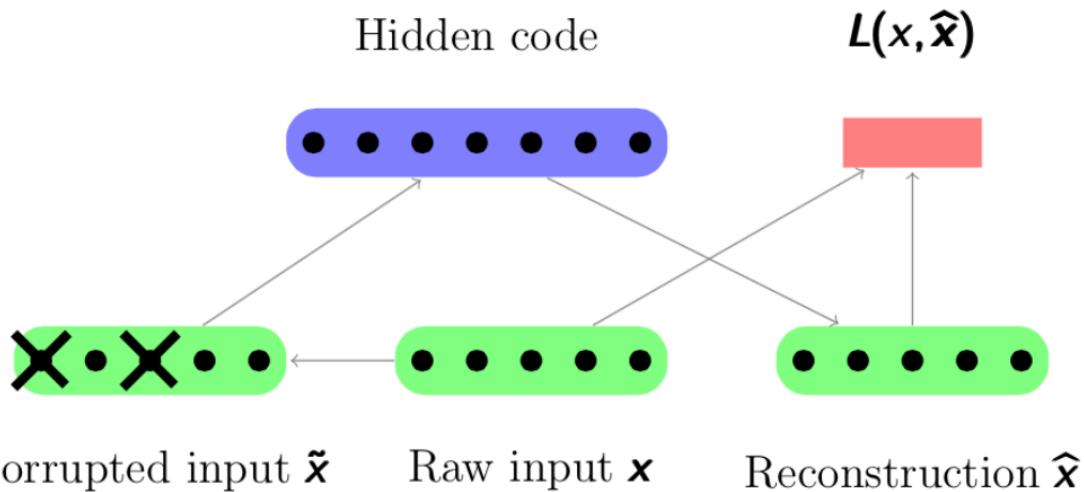
- 抗噪声 AE (Denoising AE): 如何满足带噪声情况?



- DAE: X 数据主动加噪声



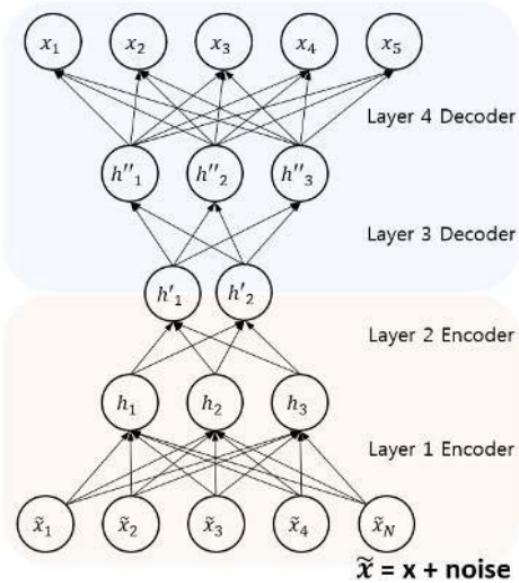
- DAE: X 数据破坏 Corrupted Data



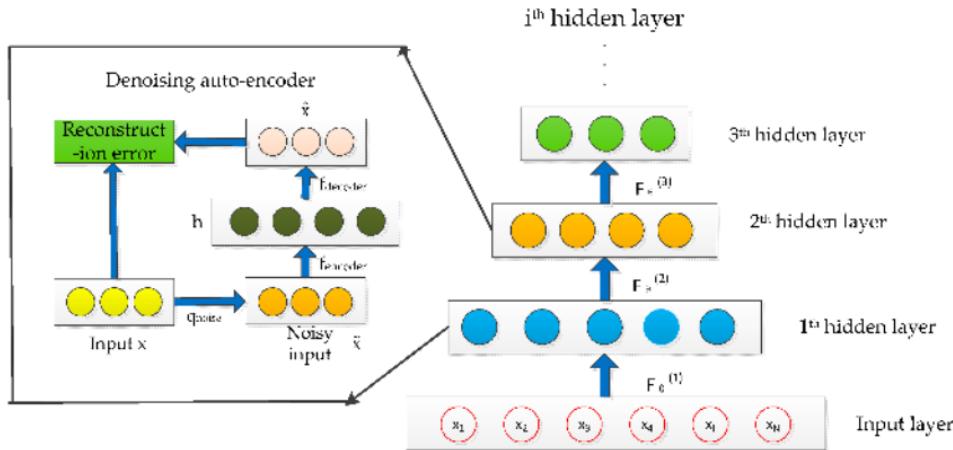
- 椒盐噪声 Salt-and-pepper noise



- 堆栈抗噪 AE (Stacked-Denoising AE)

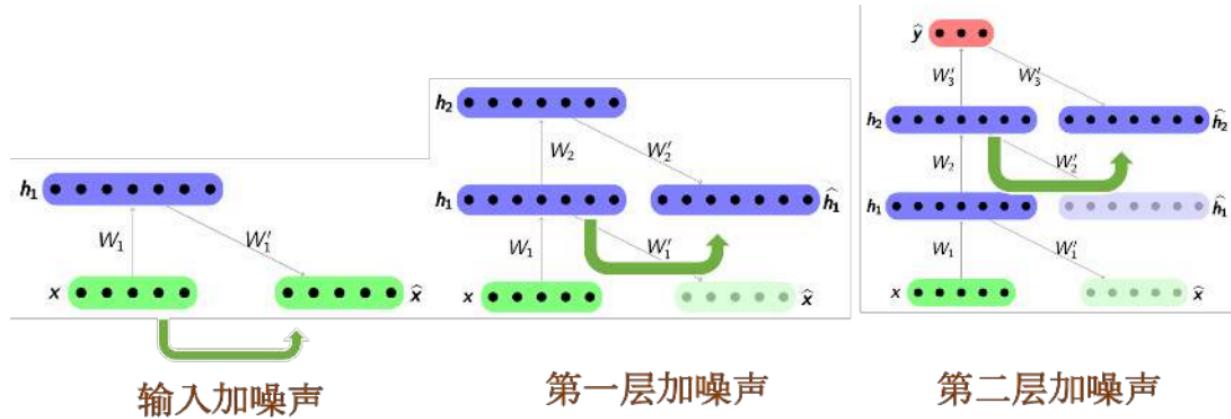


- 堆栈抗噪 AE (Stacked-Denoising AE)



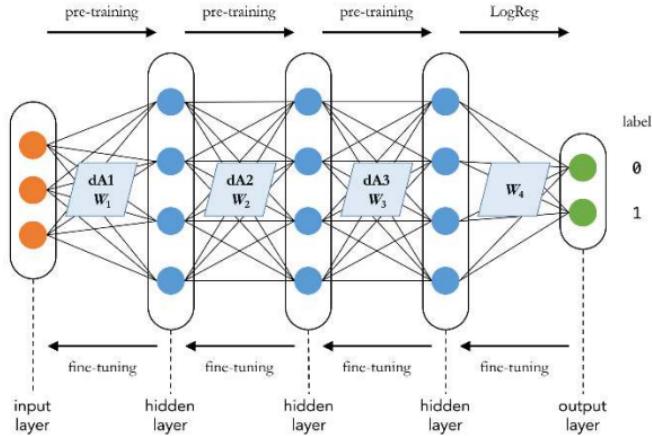
1) 逐层训练的时候每层加噪声

- 堆栈抗噪 AE (Stacked-Denoising AE): 部分类似加 Dropout



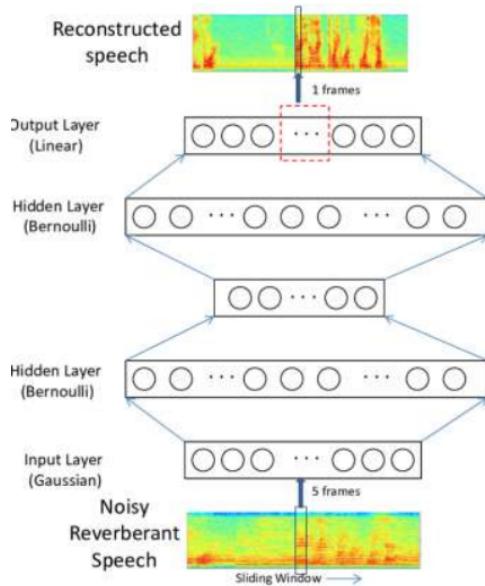
1) 逐层训练的时候每层加噪声

- 堆栈抗噪 AE (Stacked-Denoising AE): Fine-Tune

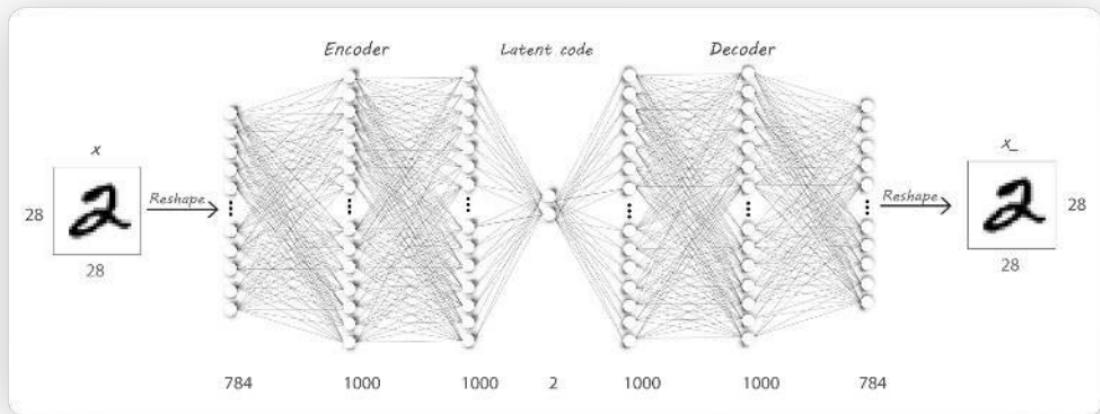


2) Fine-Tuning 的时候不加噪声

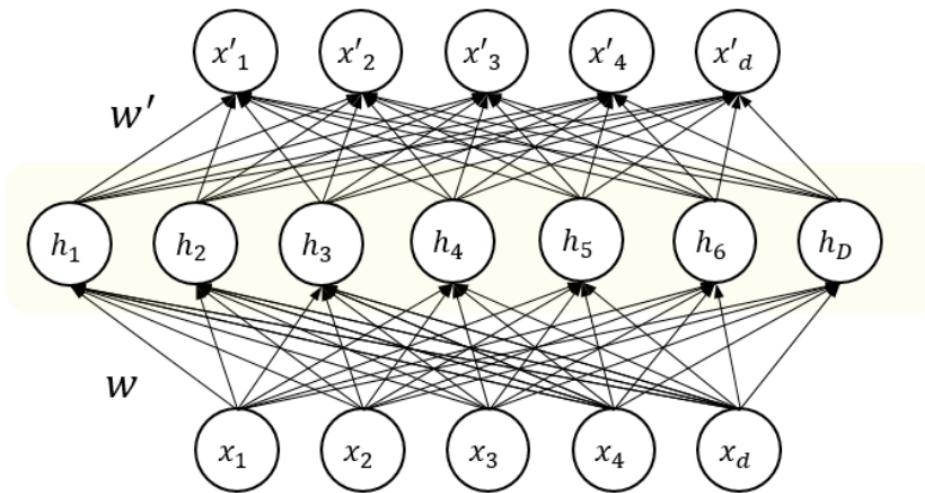
- 语音识别的噪声过滤



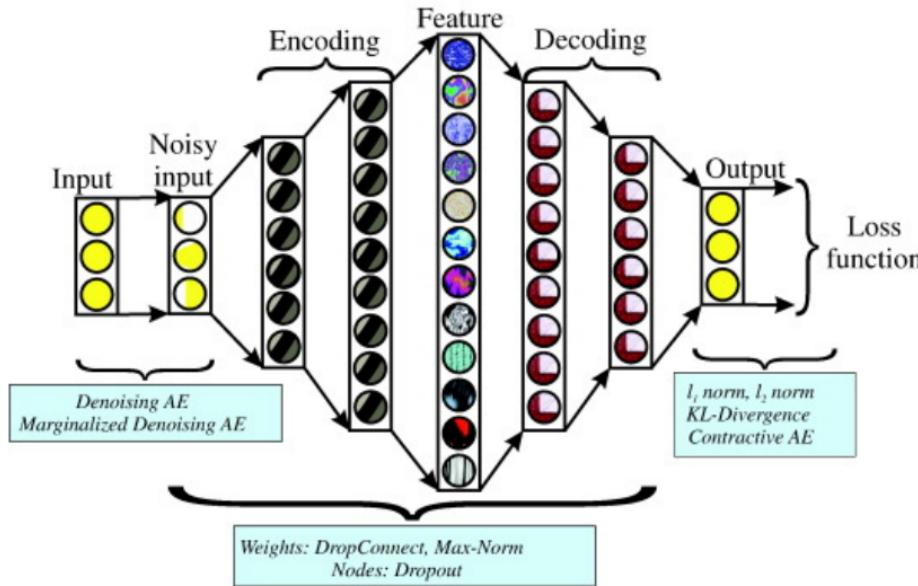
- 图像识别的噪声过滤



- 稀疏 AE (Sparse AE)



- 堆栈稀疏 AE (Stacked Sparse AE)



- Sparse AE 的目标: $\text{MSE} + \text{Reg}(\text{Weight}) + \text{Cost}(\text{Sparsity})$

平方距离 :

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

权重正则化项 :

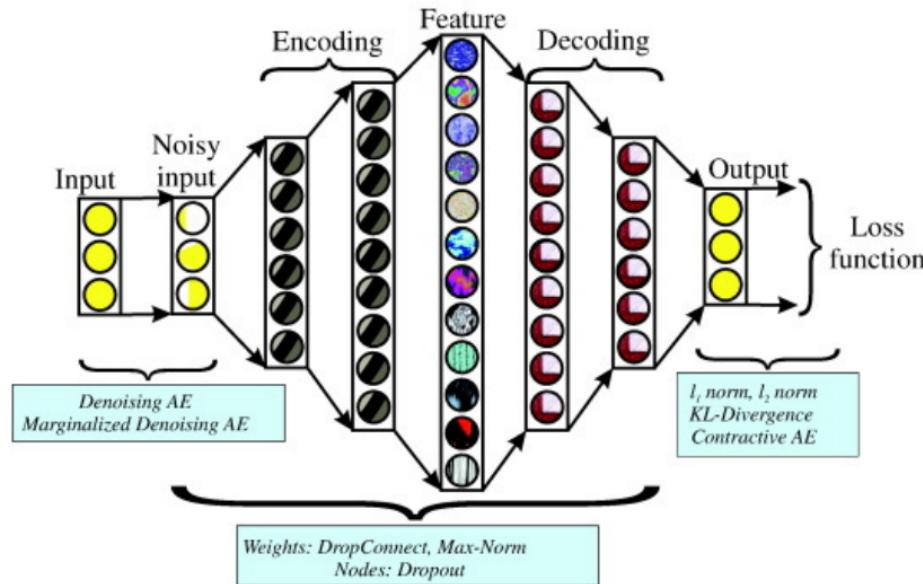
$$\frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{(l)} \right)^2$$

稀疏代价 :

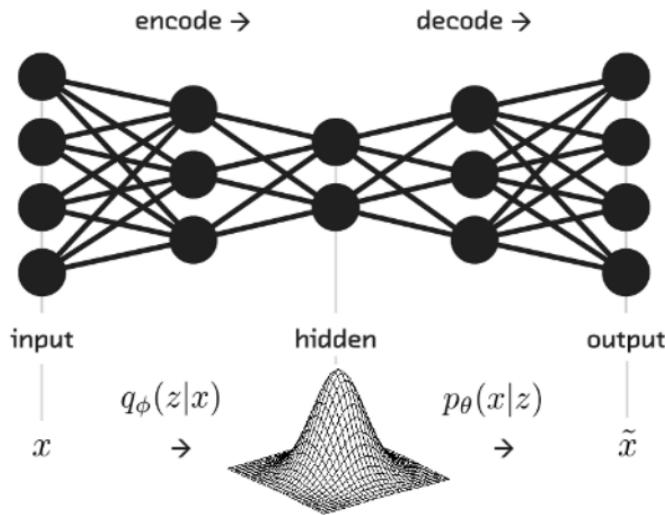
$$\sum_{j=1}^{s_2} \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}$$

其中概率估算 : $\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m \left[a_j^{(2)}(x^{(i)}) \right]$

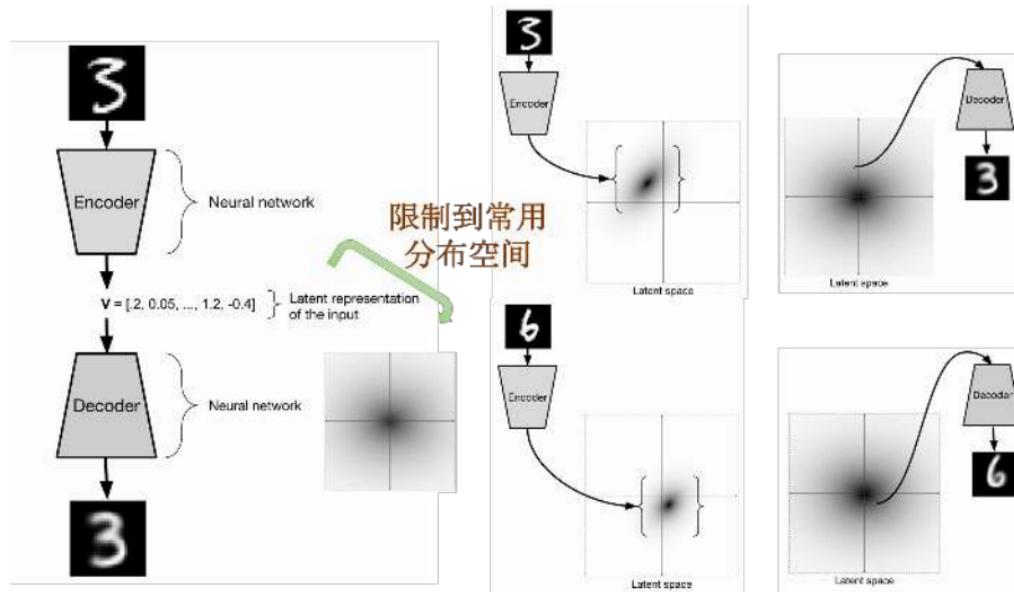
- Sparse AE 可以构建特征词典



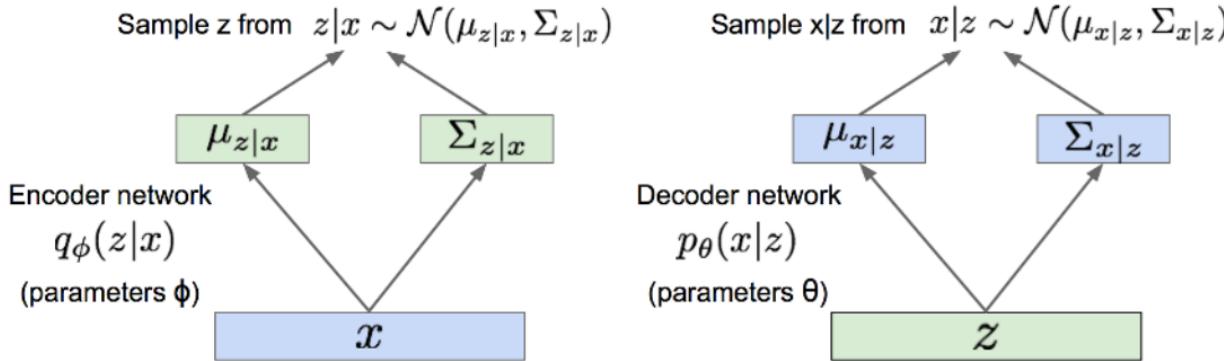
- 变分 AE (Variational AE)



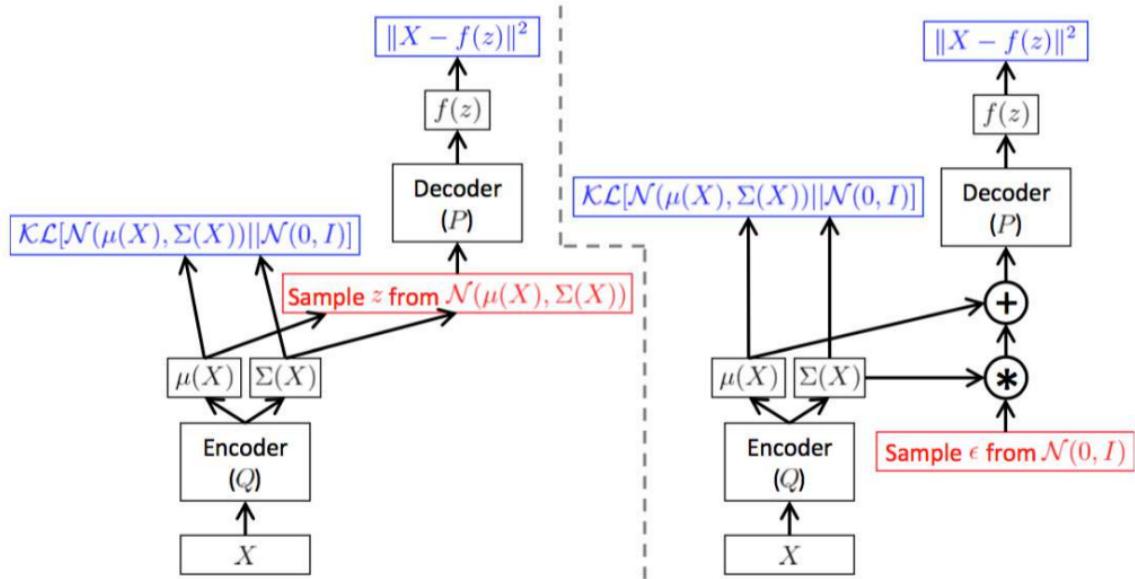
- VAE: 想法来源



- VAE 要求

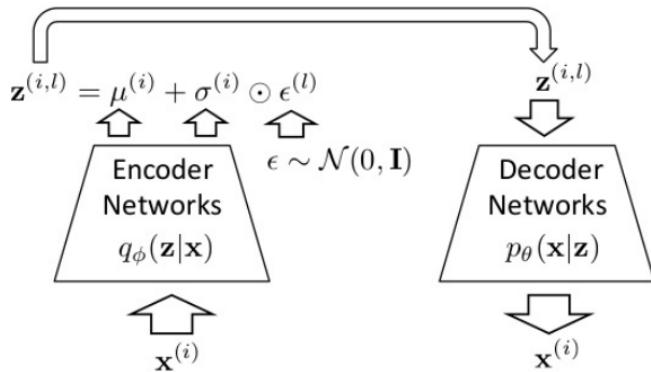


- VAE 构建需要的分布：基于标准正态分布



- 基于标准正态分布构建分布技巧：Reparameterization Trick

Reparameterization Trick



- ELBO 下限推理

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}
 \end{aligned}$$

Tractable lower bound which we can take gradient of and optimize! ($p_\theta(x|z)$ differentiable, KL term differentiable)

- ELBO 重建的含义

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

Reconstruct
the input data

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z) q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)}) q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}$$

$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$

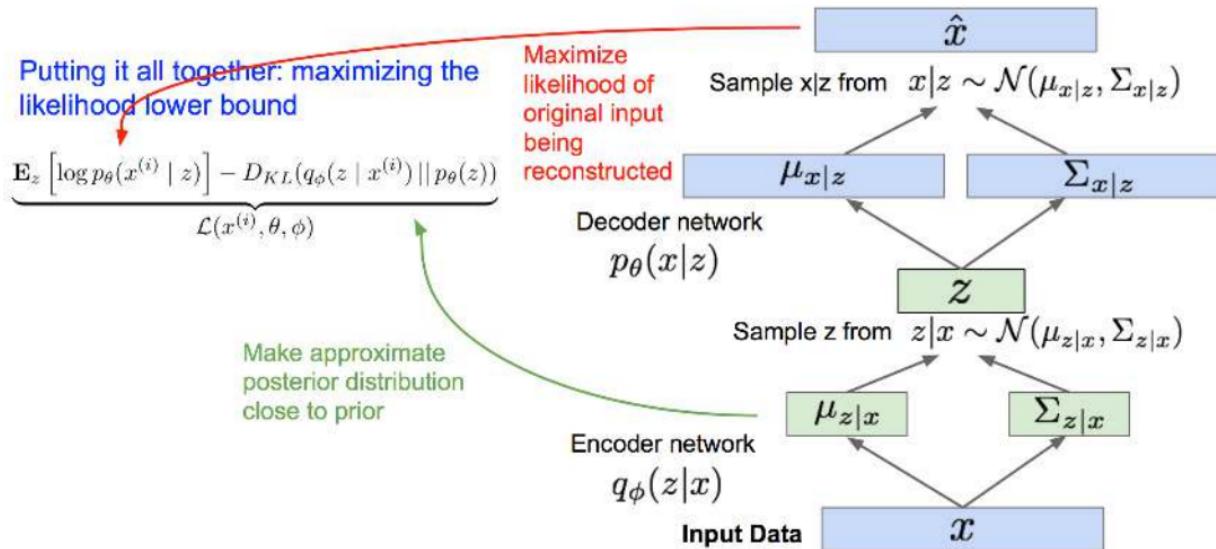
Variational lower bound ("ELBO")

Make approximate posterior distribution close to prior

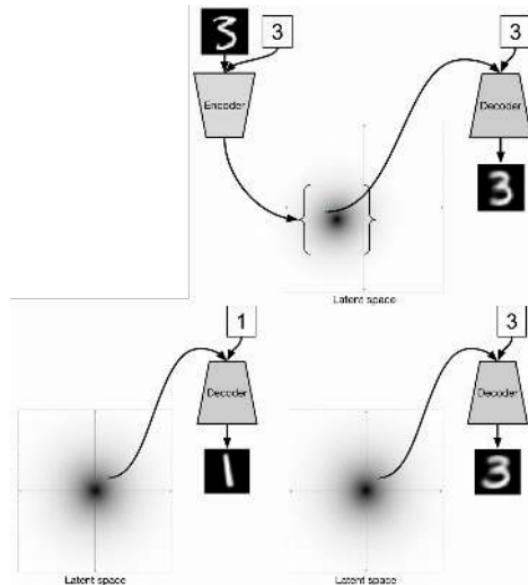
$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

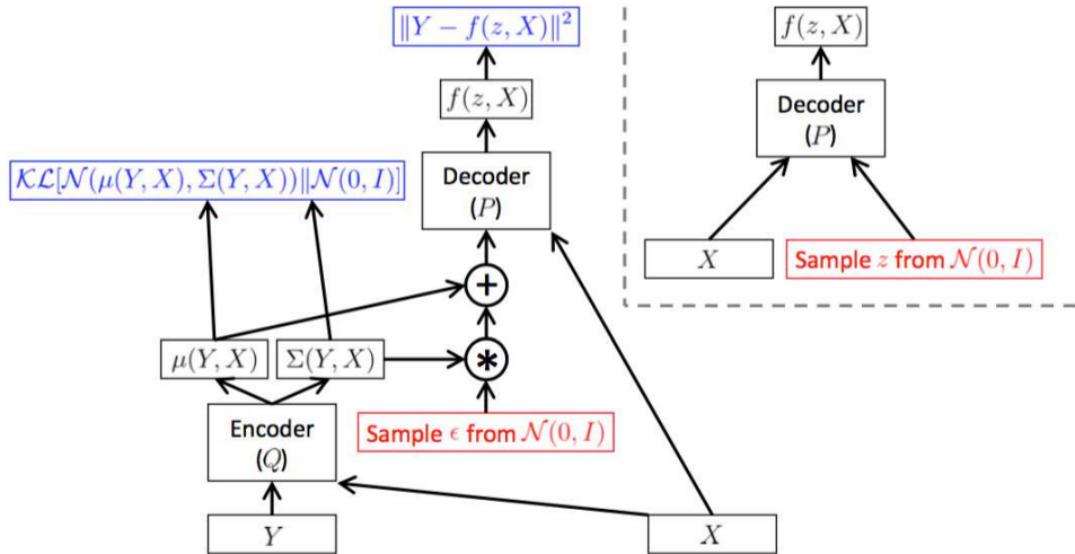
- ELBO: KL 分布距离作为正则化限制



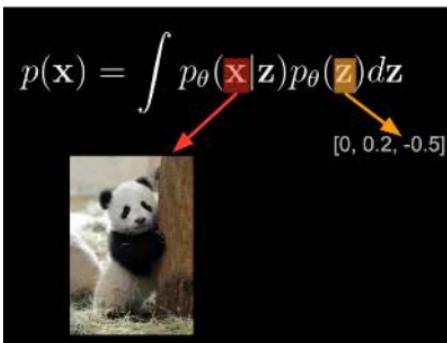
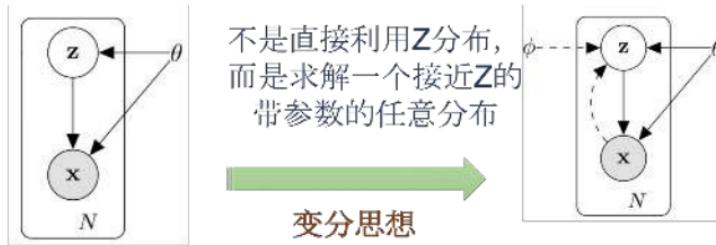
- 条件 VAE (Conditional VAE)



- 条件 VAE (Conditional VAE)



- VAE 的变分思想回顾



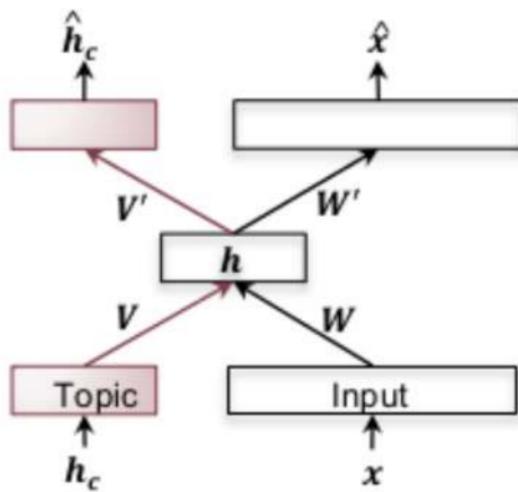
$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z}$$

$[0, 0.2, -0.5]$

$$\log p_{\theta}(\mathbf{x}_i) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}_i)}[\log p_{\theta}(\mathbf{x}_i|\mathbf{z})]$$

$$KL[q_{\phi}(\mathbf{z}|\mathbf{x}_i) || p_{\theta}(\mathbf{z})]$$

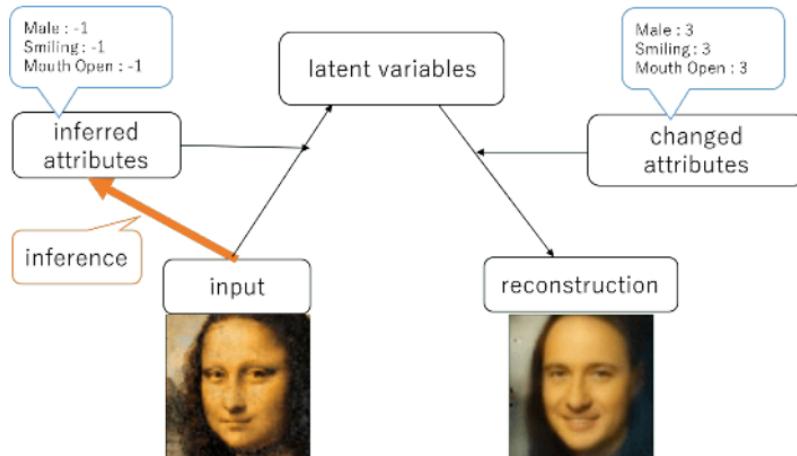
- 上下文 AE (Context-Sensitive AE)



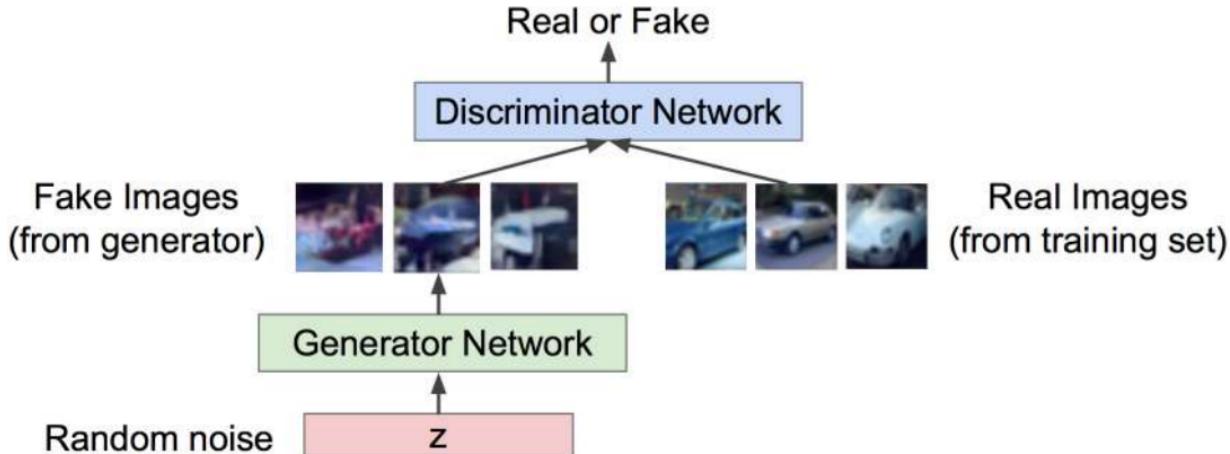
$$l(\mathbf{x}, \mathbf{h}_c) = ||\mathbf{x} - \hat{\mathbf{x}}||^2 + \lambda ||\mathbf{h}_c - \hat{\mathbf{h}}_c||^2$$

4 生成对抗网络 GAN

- VAE 更好的评价函数



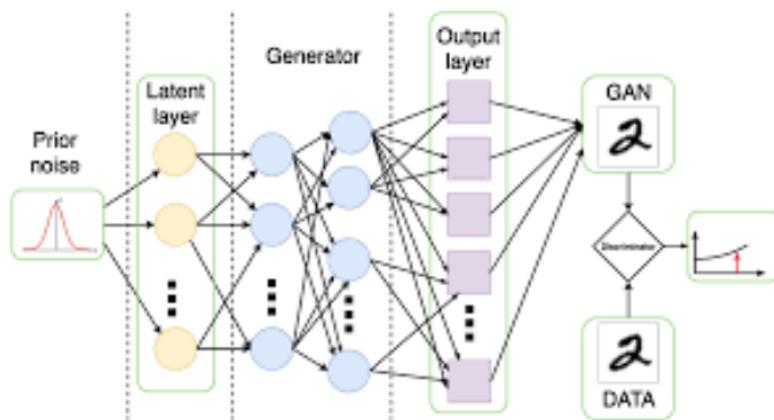
- 生成对抗网络 Generative Adversarial Network: 神经网络评价



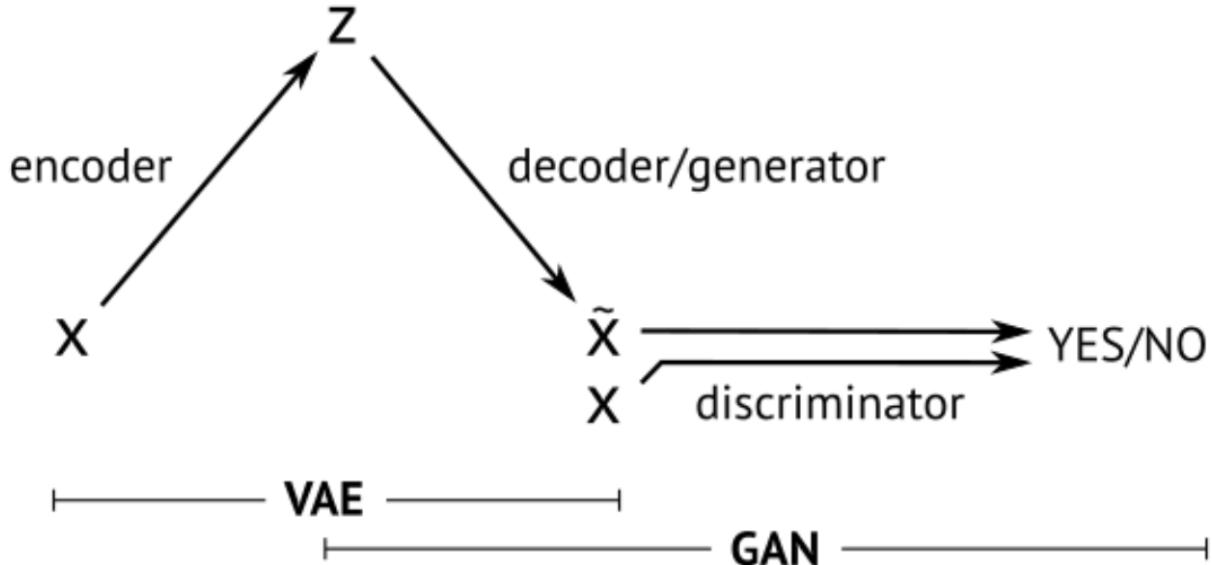
- 生成 + 判别的不停尝试



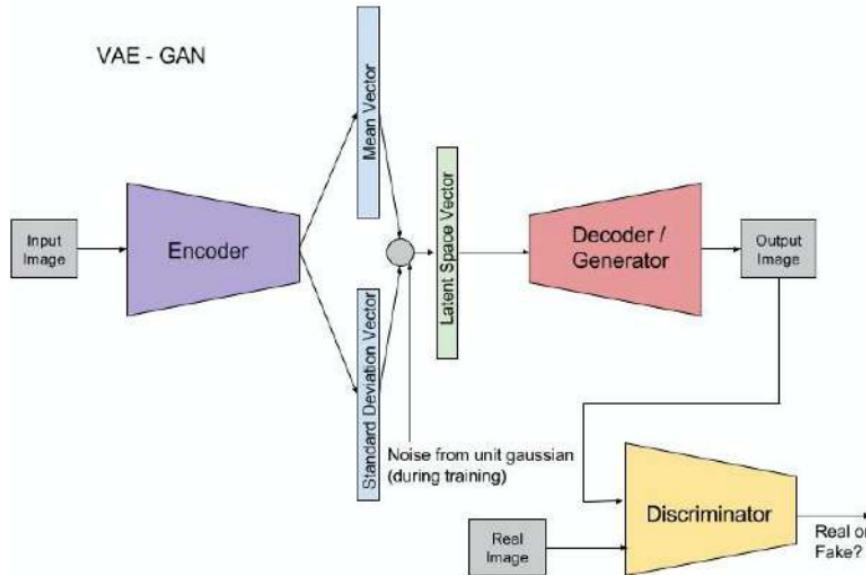
- 生成 + 判别的网络



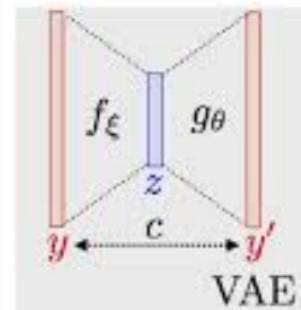
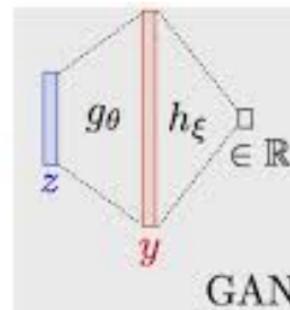
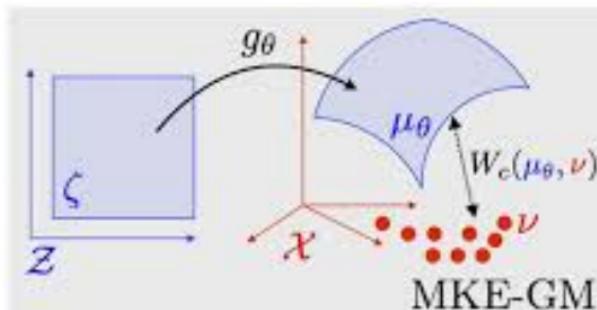
- GAN 和 VAE 的对比



- VAE 和 GAN



- VAE 和 GAN



- GAN 的优化目标

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. **Gradient ascent** on discriminator

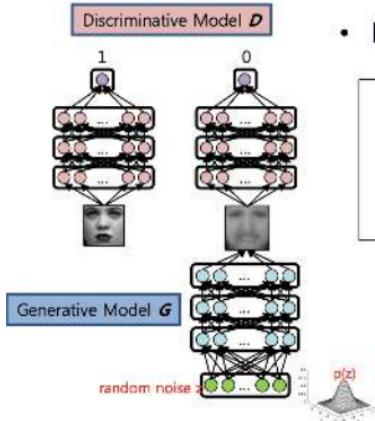
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

- GAN 的优化目标

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



- Fixed D, minimize $V(G)$:

$$\min_G V_D(G) = \min \left[\mathbb{E}_{z \sim p_{z(z)}} [\log(1 - D(G(z)))] \right]$$

Try to make $D(G(z)) = 1$

Stochastic Gradient

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right).$$

- GAN 的优化目标：求导

目标

$$\begin{aligned} J(G, D) &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D(G(\mathbf{z})))] \\ &= \int_{(\mathbf{x})} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_z(\mathbf{z}) \log(1 - D(G(\mathbf{z}))) dz \\ &= \int_{\mathbf{x}} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_G(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x} \end{aligned}$$

最优求导

$$D^* = \underset{D}{\operatorname{argmax}} J(G, D)$$

$$= \underset{D}{\operatorname{argmax}} \int_{\mathbf{x}} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_G(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x}$$

最优

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}$$

$$\begin{aligned} f(D) &= a \log(D) + b \log(1 - D). \\ \frac{\partial f}{\partial D} &\triangleq 0 \\ D^* &= \frac{a}{a+b} \end{aligned}$$

优化公式

- GAN 的优化目标: 最优情况

For G fixed, the optimal discriminator D is

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

Plug in the minimax objective

$$C(G) = \max_D V(G, D)$$

$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_g} [\log(1 - D_G^*(G(\mathbf{z})))]$$

$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

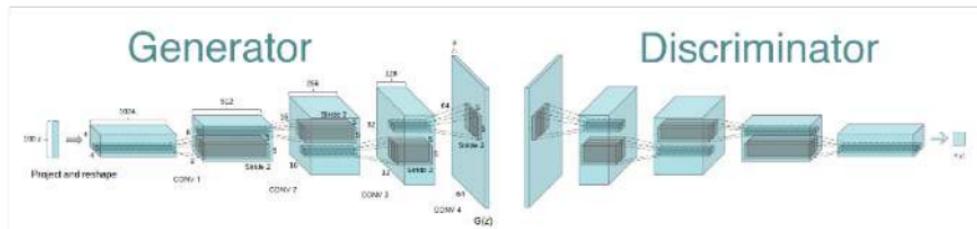
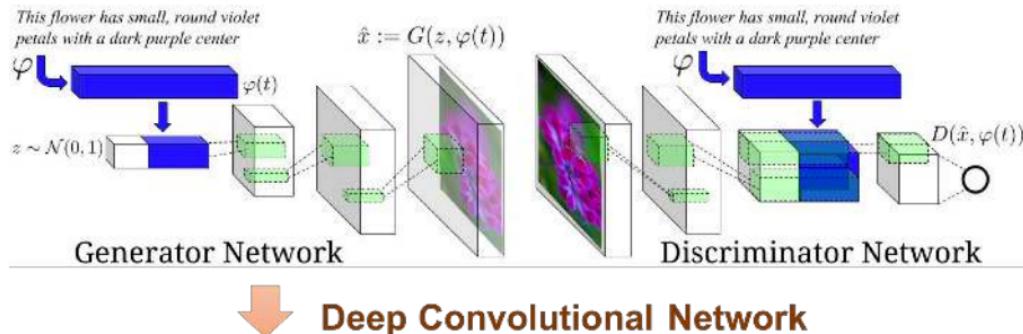
$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[\log \frac{p_g(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$C(G) = -\log(4) + KL \left(p_{\text{data}} \middle\| \frac{p_{\text{data}} + p_g}{2} \right) + KL \left(p_g \middle\| \frac{p_{\text{data}} + p_g}{2} \right)$$

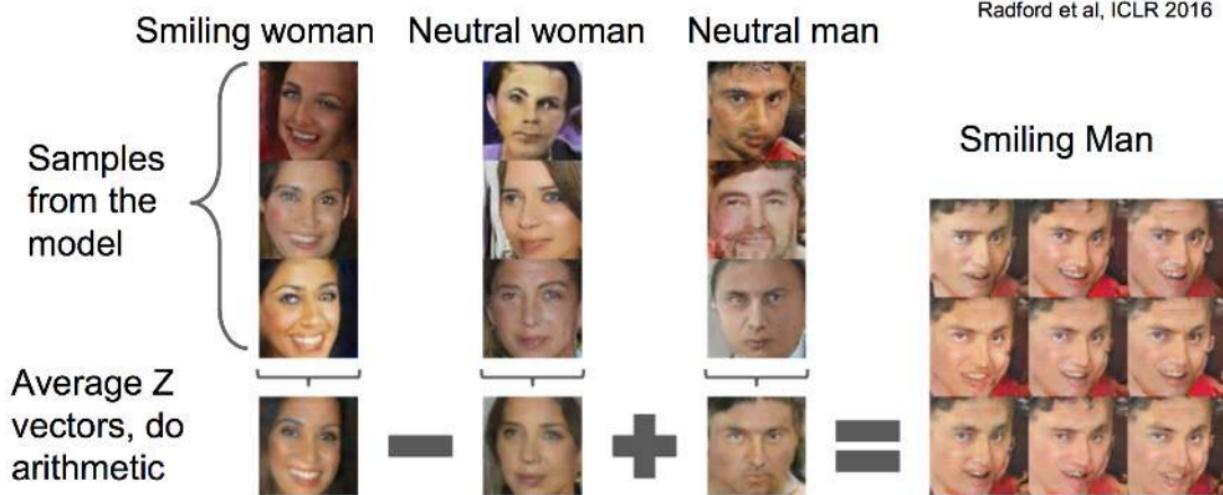
$$C(G) = -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$$

$C^* = -\log 4$ is the global minimum and the only solution is $p_g = p_{\text{data}}$

- DCGAN



- DCGAN 嵌入空间：发型



- DCGAN 嵌入空间：眼镜

Glasses man



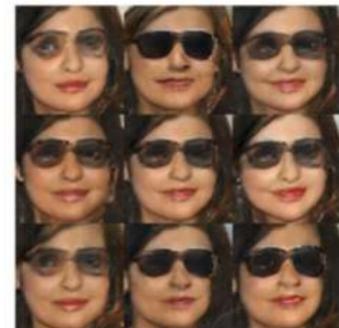
No glasses man



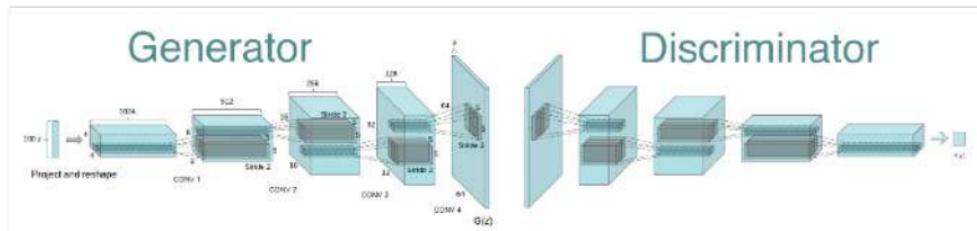
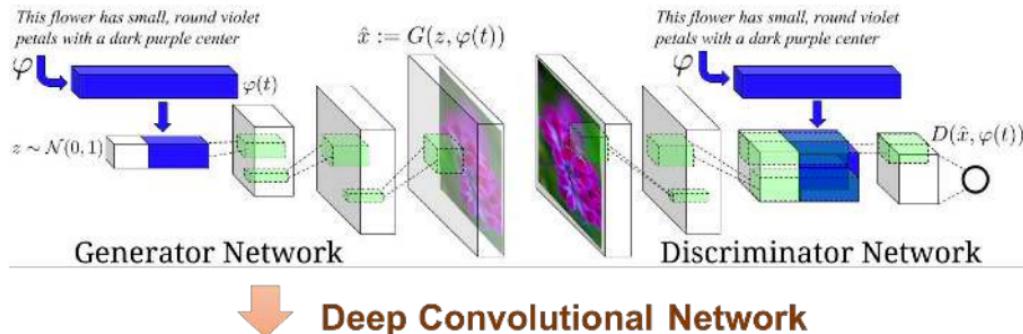
No glasses woman

Radford et al,
ICLR 2016

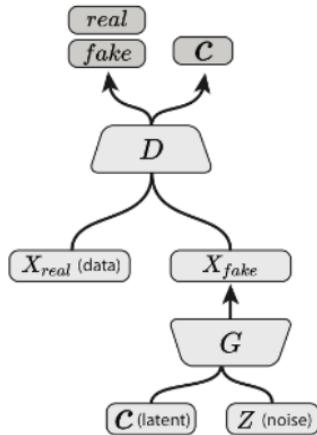
Woman with glasses



- DCGAN



- InfoGAN: Mutual information 衡量主要特征和内部变化

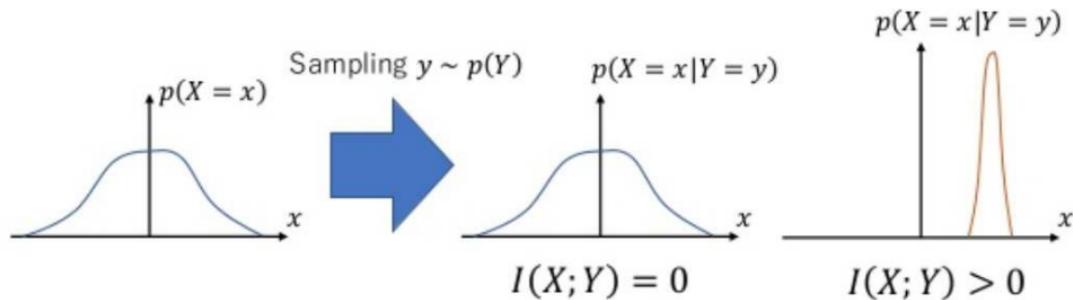


InfoGAN
(Chen, et al., 2016)

- Mutual Information

$$I(X; Y) = H(X) - H(X|Y), \text{ where}$$

- $H(X) = E_{x \sim p(x)}[-\ln p(X = x)]$:
Entropy of the prior distribution
- $H(X|Y) = E_{y \sim p(Y), x \sim p(x|Y=y)}[-\ln p(X = x|Y = y)]$:
Entropy of the posterior distribution



- Mutual Information to GAN

$$\begin{aligned}
 I[X, Y] &= H[Y] - \mathbb{E}_x H[Y|X=x] \\
 &= H[Y] + \mathbb{E}_x \mathbb{E}_{y|x} \log p(y|x) \\
 &= H[Y] + \mathbb{E}_x \mathbb{E}_{y|x} \log \frac{p(y|x)q(y|x)}{q(y|x)} \\
 &= H[Y] + \mathbb{E}_x \mathbb{E}_{y|x} \log q(y|x) + \mathbb{E}_x \mathbb{E}_{y|x} \log \frac{p(y|x)}{q(y|x)} \\
 &= H[Y] + \mathbb{E}_x \mathbb{E}_{y|x} \log q(y|x) + \mathbb{E}_x KL[p(y|x) \| q(y|x)] \\
 &\geq H[Y] + \mathbb{E}_x \mathbb{E}_{y|x} \log q(y|x)
 \end{aligned}$$

$$I[X, Y] = \max_q \{H[Y] + \mathbb{E}_{x,y} \log q(y|x)\}$$

$$I[X, Y] \geq H[Y] + \max_{\psi} \mathbb{E}_{x,y} \log q(y|x; \psi)$$

$$\begin{aligned}
 \ell_{GAN}(\theta) &= I[x, y] \\
 &\geq h(0.5) + \max_{\psi} \mathbb{E}_{x,y} \log q(y|x; \psi) \\
 &= h(0.5) + \max_{\psi} \left\{ \mathbb{E}_{x_{real}} \log q(1|x; \psi) + \mathbb{E}_{x_{fake}} \log q(0|x; \psi) \right\} \\
 \ell_{GAN}(\theta) + h(0.5) &\geq \max_{\psi} \left\{ \mathbb{E}_{x \sim P_{data}} \log D(x, \psi) + \mathbb{E}_{c,z} \log(1 - D(G(c, z, \theta), \psi)) \right\}
 \end{aligned}$$

- InfoGAN: 内部变化



(a) Azimuth (pose)

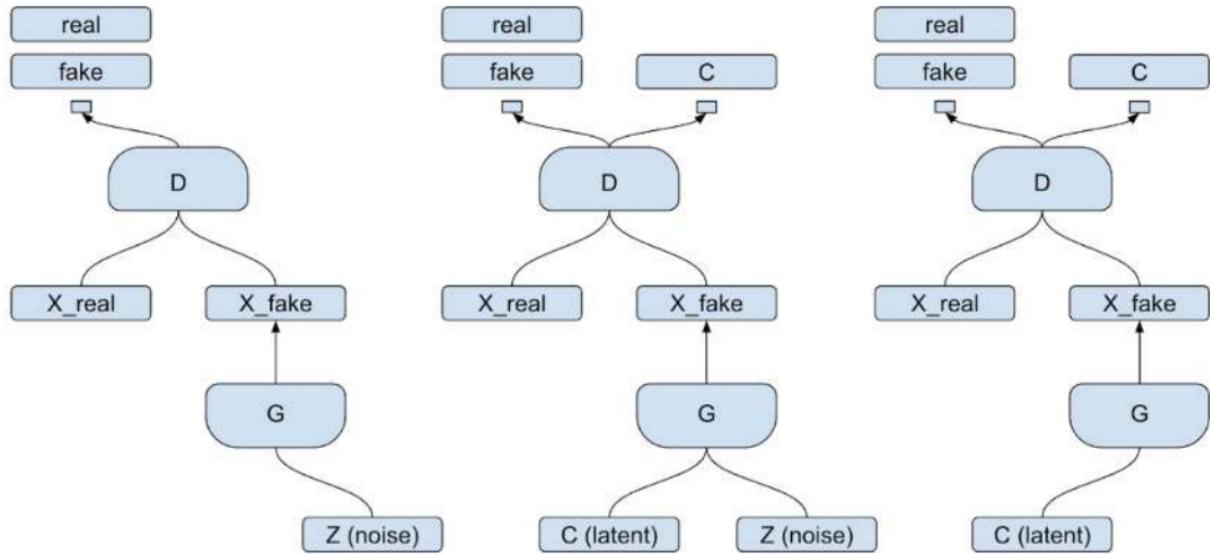
(b) Elevation



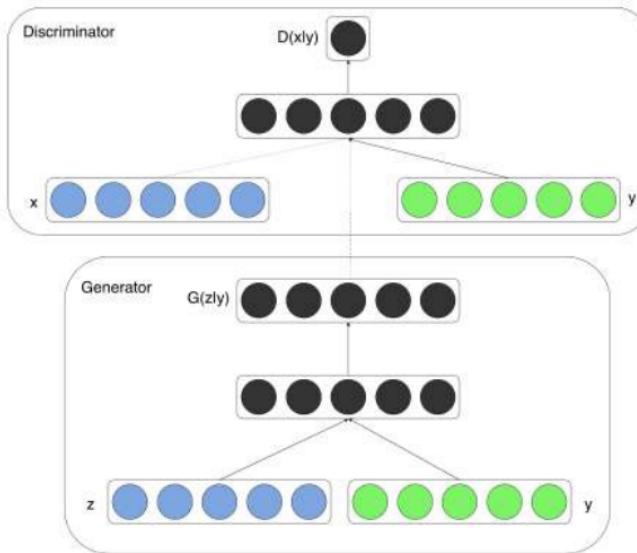
(c) Lighting

(d) Wide or Narrow

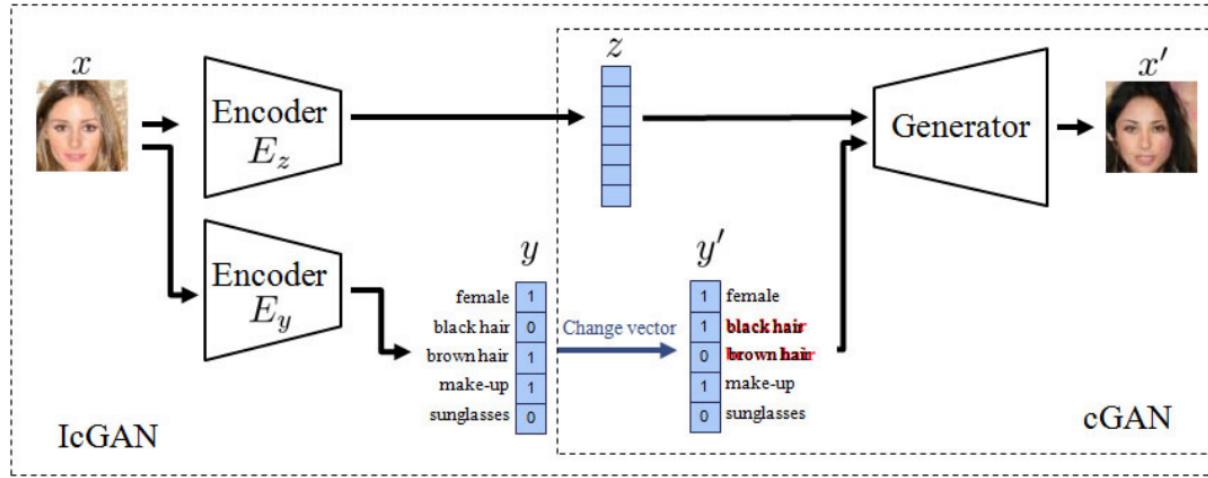
- InfoGAN: Vanilla GAN -> InfoGAN -> InfoGAN w



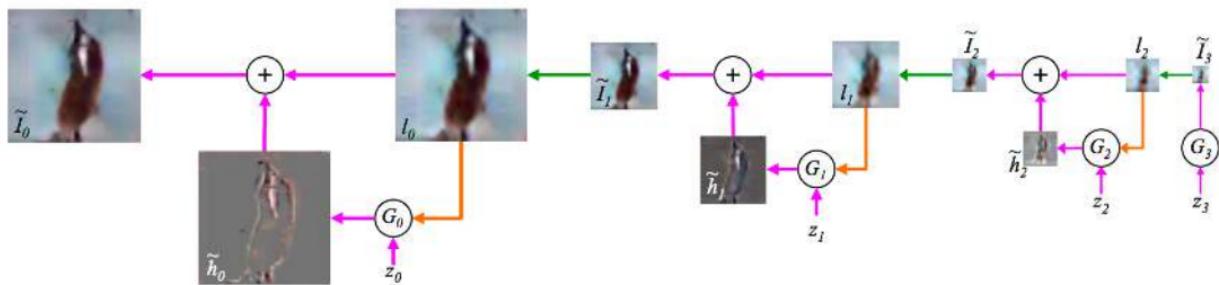
- Conditional GAN



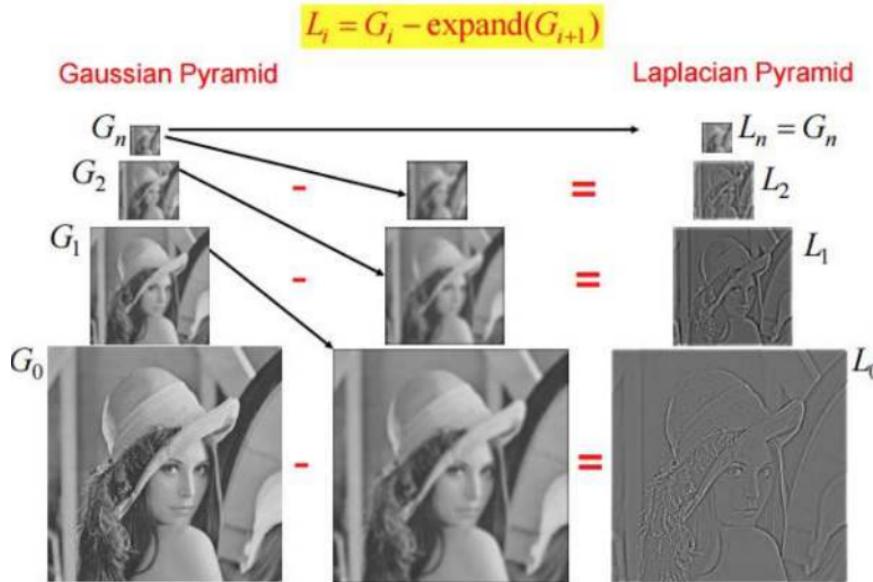
- Invertible Conditional GAN



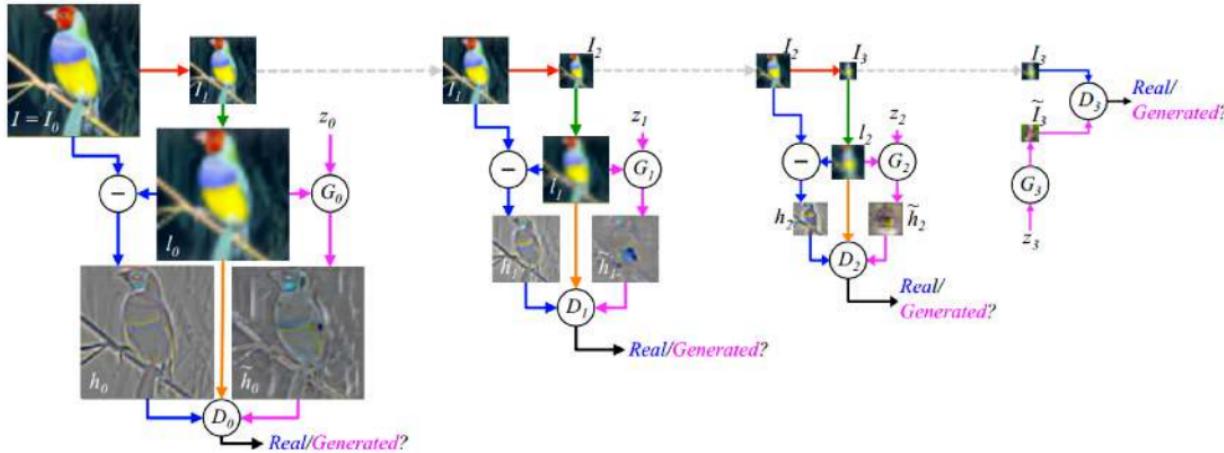
- LapGAN



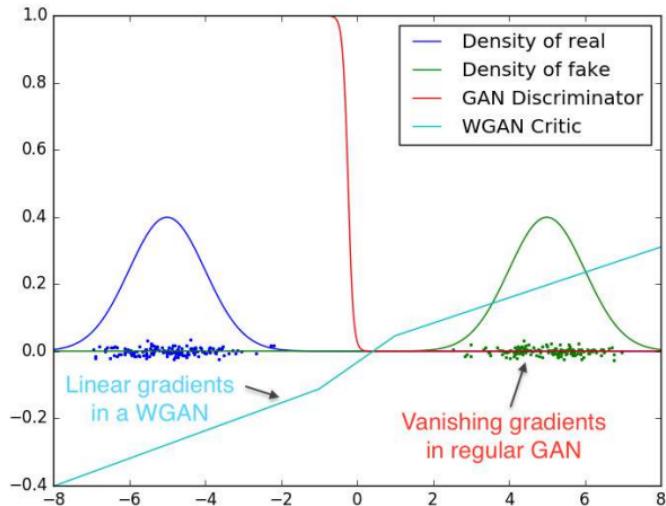
- Laplacian Pyramid



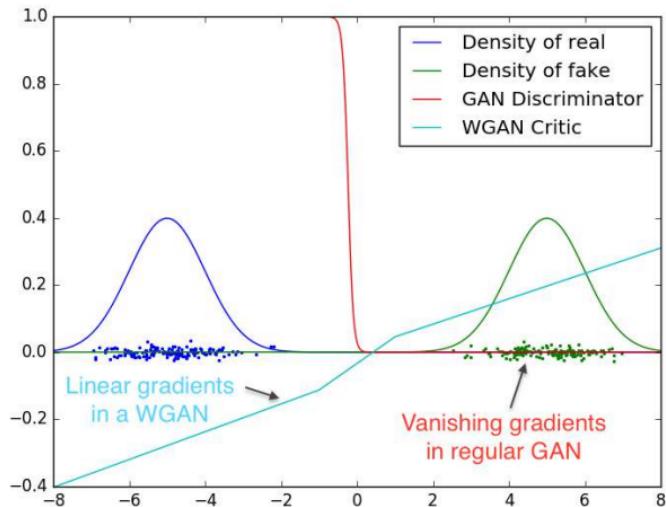
- LapGAN Train



- Wasserstein GAN



- Wasserstein GAN



- Wasserstein Distance

Jensen-Shannon divergence

$$JS(P_r, P_g) = \frac{1}{2}KL(P_r\|P_m) + \frac{1}{2}KL(P_g\|P_m)$$

Wasserstein distance

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

1-Lipschitz functions :

Kantorovich-Rubinstein duality

$$W(P_r, P_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_\theta}[f(x)]$$

感谢 Stanford, CMU, MIT 等网上公开课程和资料！

AI2ML

