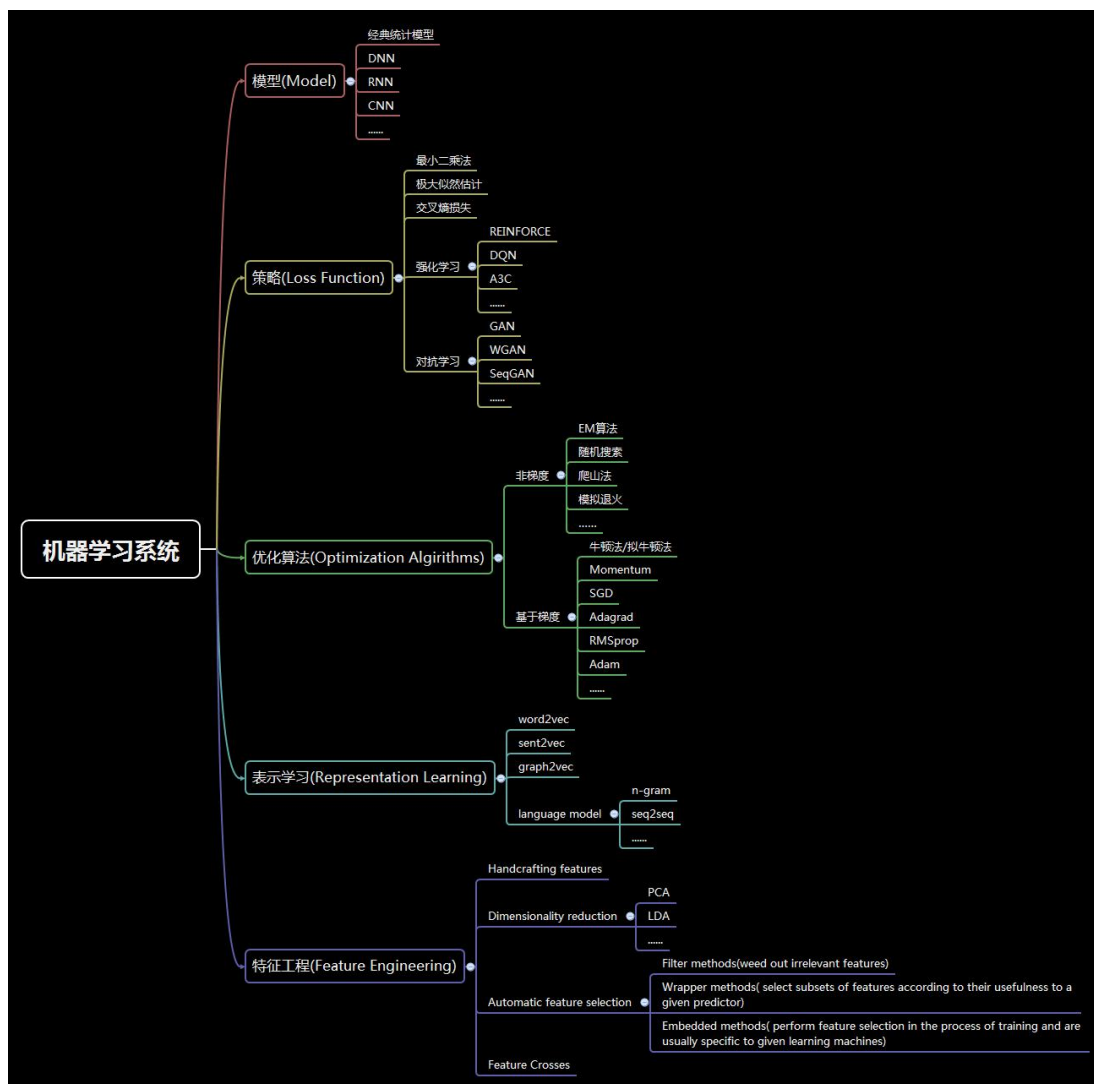# 机器学习与RMT结合探索

姜衡军(北京邮电大学计算机系)

# 大纲

- **机器学习系统概览**
- 我的NLP研究工作回顾
  - 端到端对话系统
  - 序列标注
  - 动态激活函数选择
  - 数据增强
- 机器学习与RMT结合探索

# 机器学习系统概览
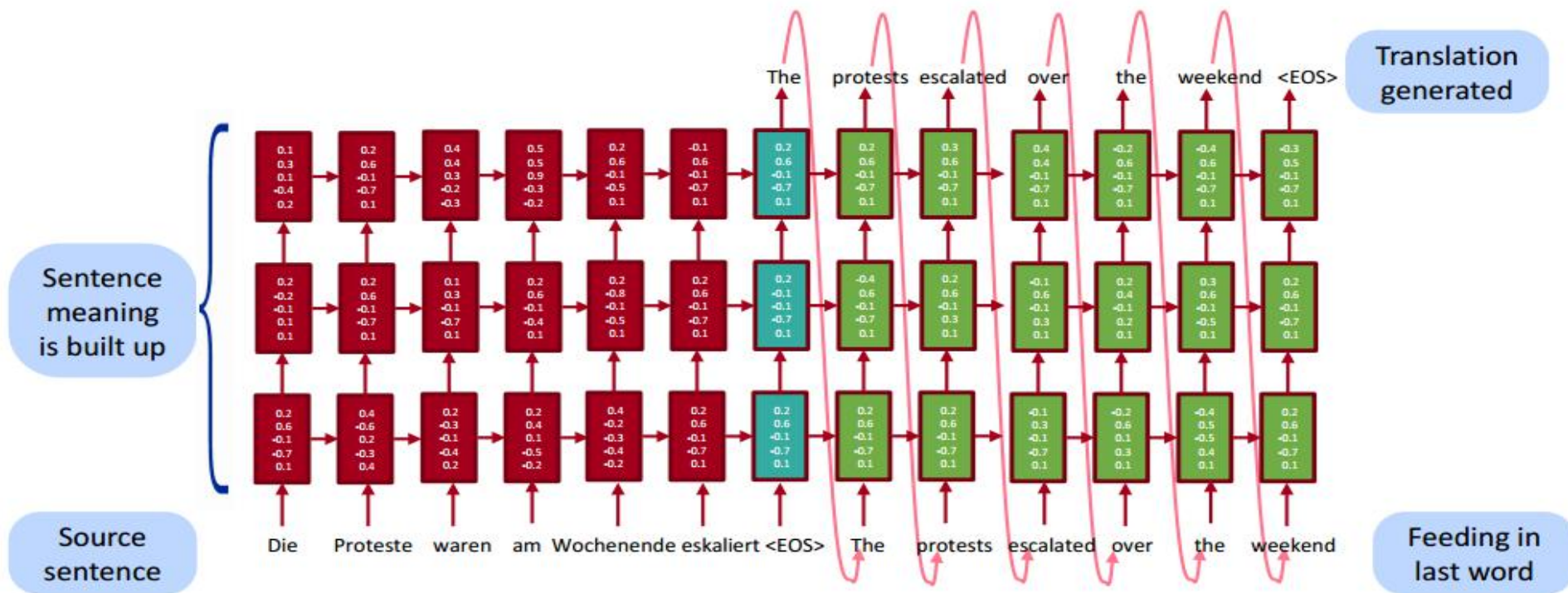


机器学习系统.PNG

# 大纲

- 机器学习系统概览
- **我的NLP研究工作回顾**
    - **端到端对话系统**
    - **序列标注**
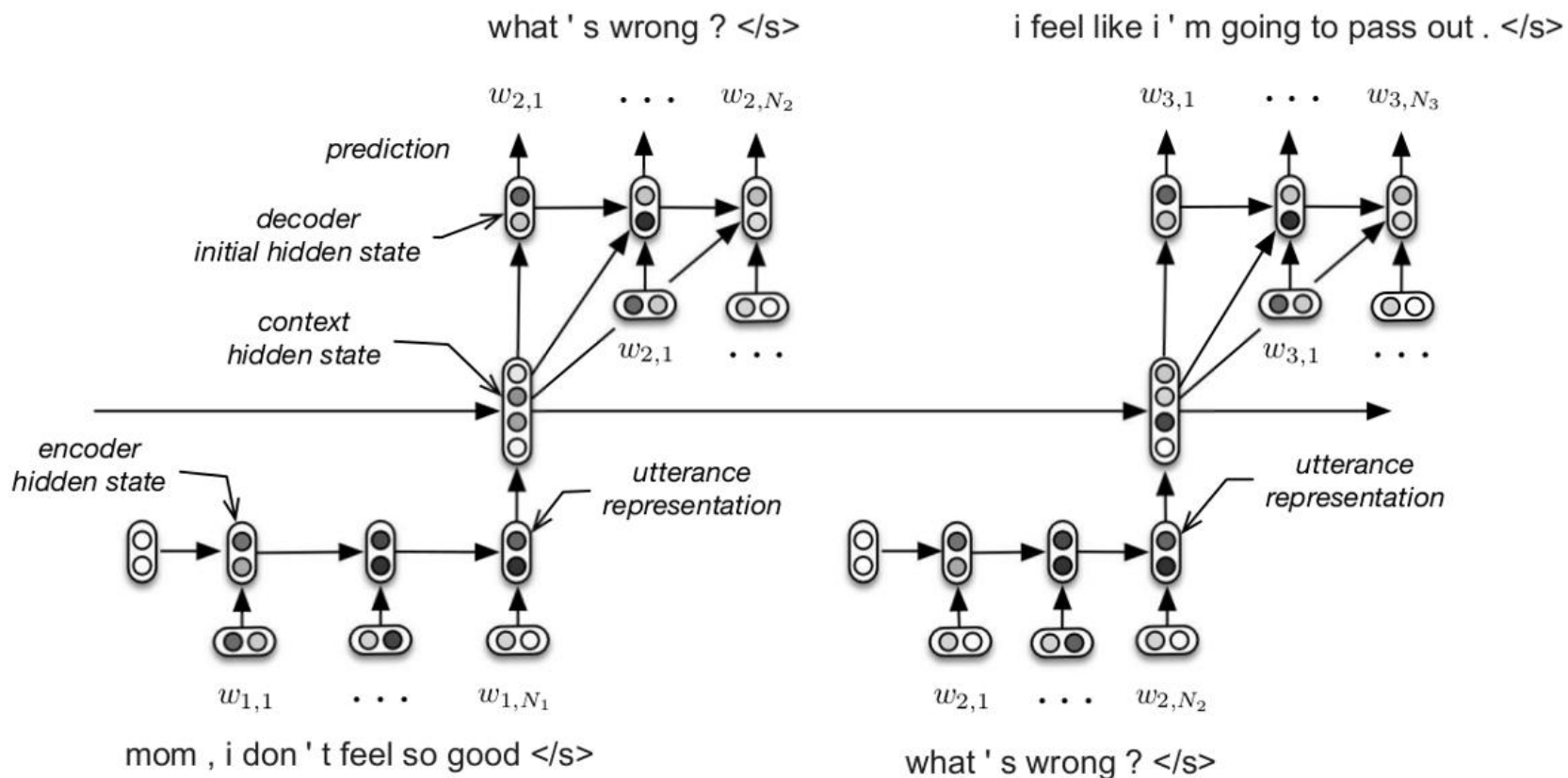    - **动态激活函数选择**
    - **数据增强**
- 机器学习与RMT结合探索

# 端到端对话系统

序列到序列模型：

# 方法与模型

层次化编解码模型：

what ' s wrong ? </s>

i feel like i ' m going to pass out . </s>

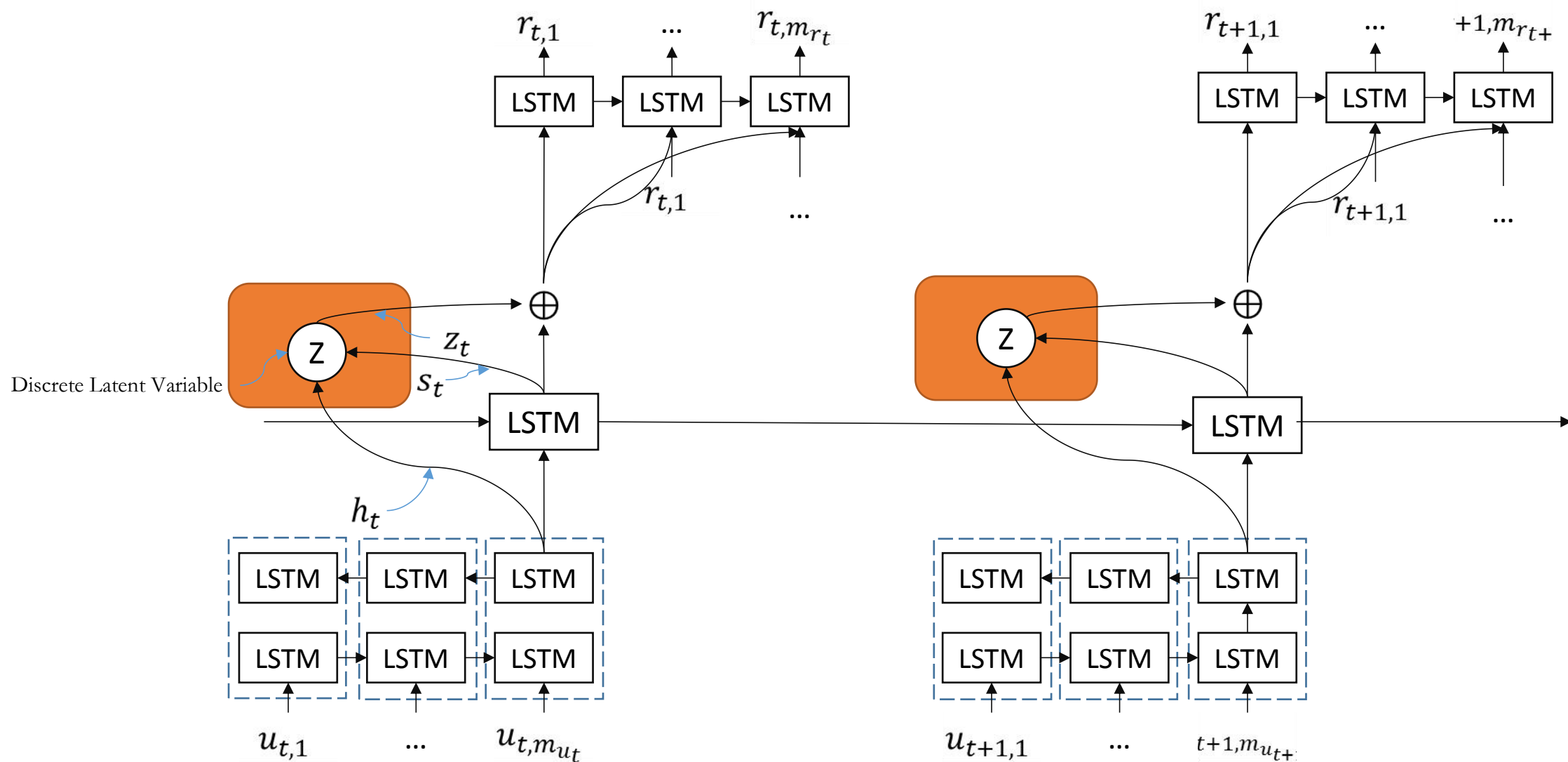# 方法与模型

引入基于**离散隐变量**的层次化编解码模型（DVHRED）：

基本原理：
　　在层次化编解码结构基础上引入离散隐变量，从而更好地对文本语义及用户意图进行建模

# 方法与模型

# 变分贝叶斯推断

潜在用户意图变量: z
近似后验概率分布: q(z|x)
真实后验概率分布: p(z|x)

$$\min KL(q \| p) = \int q(z|x) \log \frac{q(z|x)}{p(z|x)} dz$$

$$= -\int q(z|x) \log \frac{p(z|x)}{q(z|x)} dz$$

$$= -\int q(z|x) \log \frac{p(z,x)}{q(z|x)p(x)} dz$$

$$= \int q(z|x)[\log q(z|x) + \log p(x)]dz - \int q(z|x) \log p(z,x)dz$$

$$= \log p(x) + \int q(z|x) \log q(z|x)dz - \int q(z|x) \log p(z,x)dx$$

$$\Rightarrow \log p(x) = KL(q \| p) + L(q)$$

$$\Rightarrow \min KL(q \| p) == \max L(q)[\# ELOB(\text{Evidence Lower Bound})]$$

$$L(q) = \int q(z|x) \log p(z,x)dz - \int q(z|x) \log q(z|x)dz$$

$$= \int q(z|x) \log p(x|z)dz + \int q(z|x) \log p(z)dz - \int q(z|x) \log q(z|x)dz$$

$$= \int q(z|x) \log p(x|z)dz - \int q(z|x) \log \frac{q(z|x)}{p(z)}dz$$

$$= \int q(z|x) \log p(x|z)dz - KL[q(z|x) \| p(z)]$$

# 方法与模型

$$h_t = biLSTM(u_t)$$

$$s_t = LSTM(s_{t-1}, h_t)$$

$$z_t^1 \sim \pi_\theta(z_t \mid s_t, h_t)$$

# 方法与模型

由于上述模型缺少**领域知识**，因此只能与用户进行"闲聊"

  User: I want to order a chinese restaurant.

  Bot: sounds a good idea.

为了更好的完成任务型对话，需要在DVHRED框架上加入知识库（Knowledge Base, KB）(DVHRED+KB)

  User: I want to order a chinese restaurant.

  Bot: The good luck chinese food takeaway is in the south area.

# 方法与模型

# 方法与模型

$x_t$ : 10维的向量，[0, 0, 0, 0, 0, 1, 0, 1, 0, 1]

分别表示['phone_request'，'postcode_request', 'address_request', 'name_other', 'food_request', 'food_inform', 'pricerange_request', 'pricerange_inform', 'area_request', 'area_inform']
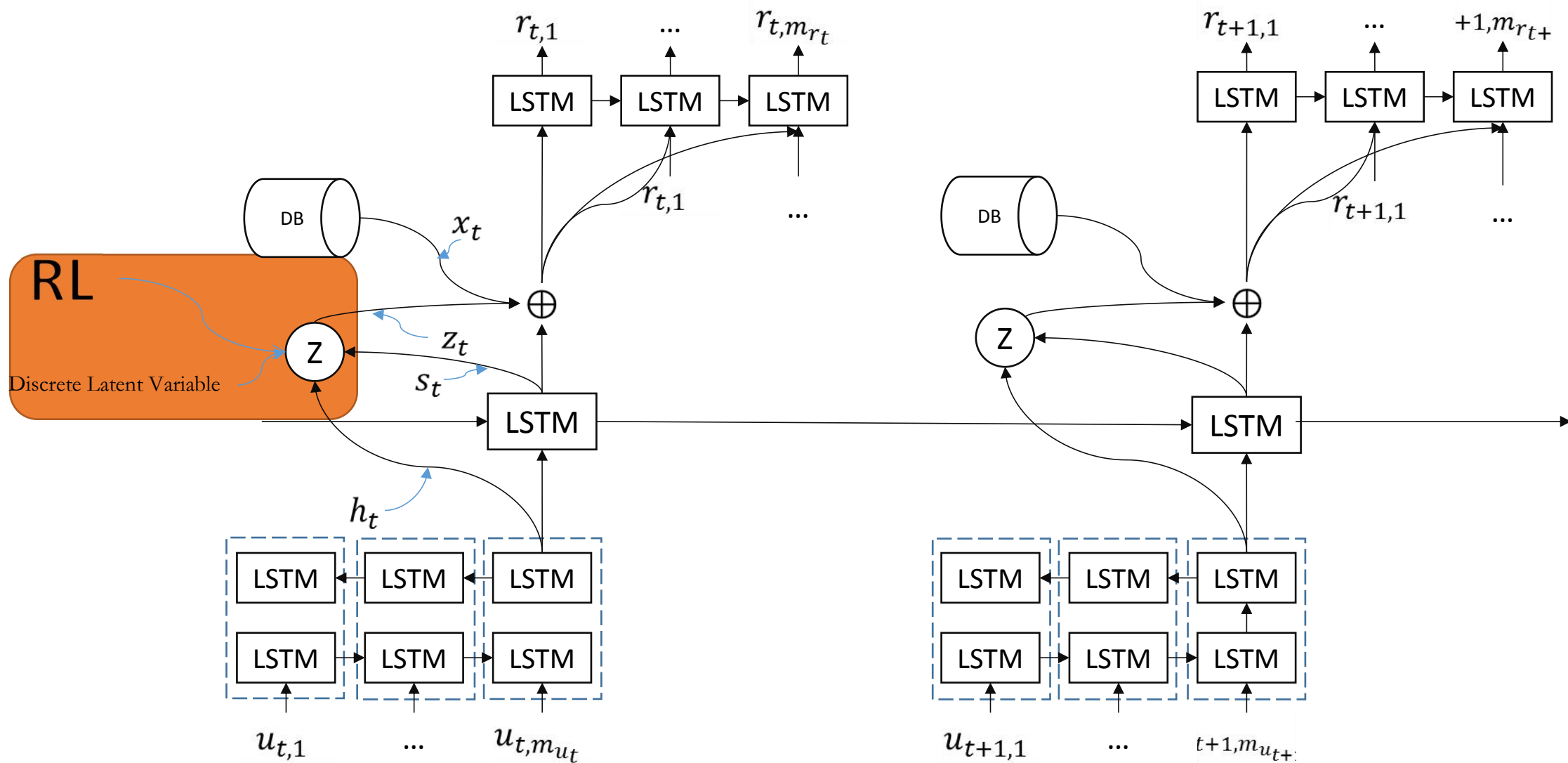
# 方法与模型

进一步，为了提升任务型对话的成功率，需要在选用户潜在意图时加入明确的激励

因此，我们在DVHRED+KB的基础上，通过强化学习，引入奖励函数

$$r_t = \eta \cdot \text{sBLEU}(m_t, \hat{m_t}) + \begin{cases} 1 & m_t \text{ improves} \\ -1 & m_t \text{ degrades} \\ 0 & \text{otherwise} \end{cases}$$

# 方法与模型

# REINFORCE Algorithm

$$\max E_\pi [\sum_k \gamma^k r_{t+k} \mid s_t = s] = \max \sum_{a_t} P_\theta(a_t \mid s_t; \theta) * Q(s_t, a_t)$$

$$l(\theta) = -\sum_{a_t} P_\theta(a_t \mid s_t; \theta) * Q(s_t, a_t)$$

$$\nabla_\theta l(\theta) = -\sum_{a_t} \nabla_\theta P_\theta(a_t \mid s_t; \theta) * Q(s_t, a_t)$$

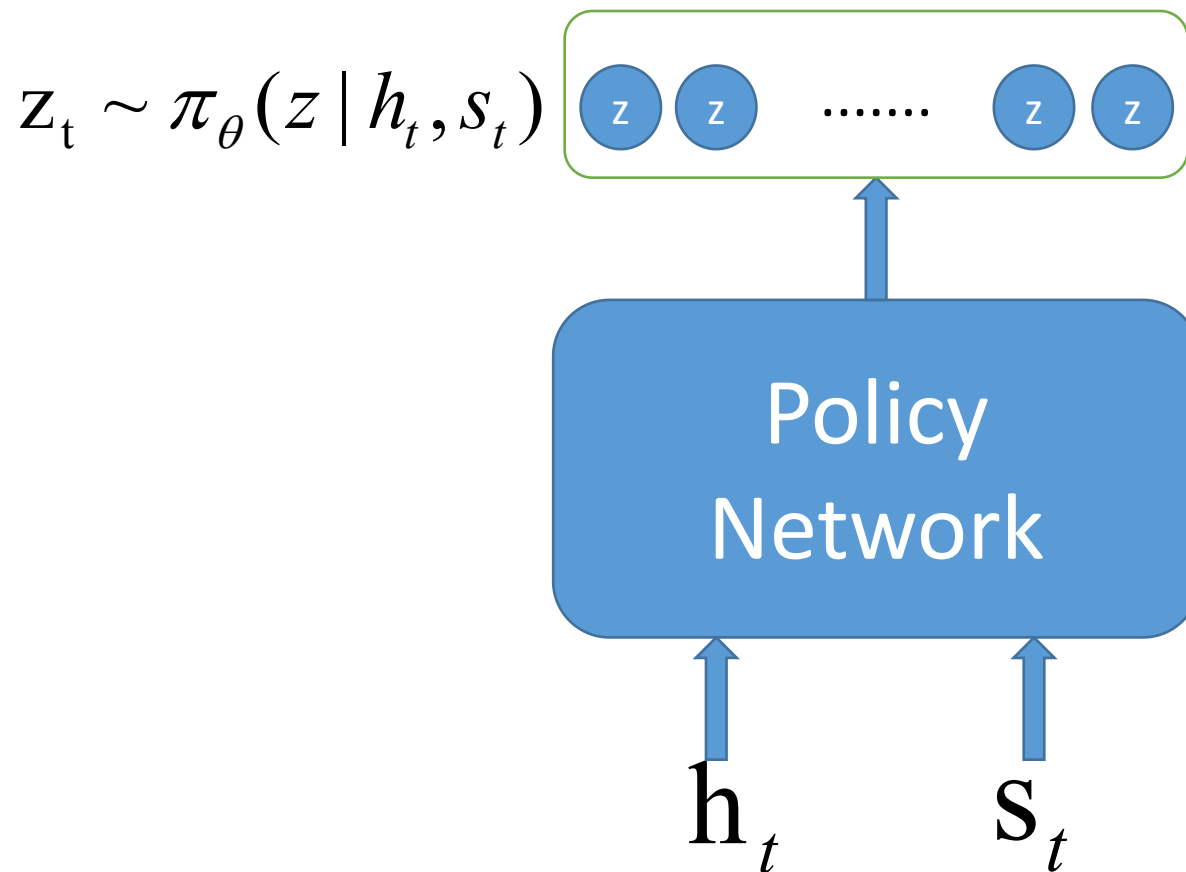$$= -\sum_{a_t} P_\theta(a_t \mid s_t; \theta) * \nabla_\theta \log(P_\theta(a_t \mid s_t; \theta)) * Q(s_t, a_t)$$

$$= -E_\pi [\nabla_\theta \log(P_\theta(a_t \mid s_t; \theta)) * Q(s_t, a_t) \mid s_t]$$

$$L(\theta) = -\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T} \log(P_\theta(a_t^n \mid s_t; \theta)) * Q(s_t, a_t^n) \quad (1) \qquad Q(s_t, a_t) = E [\sum_{k=0}^{T} \gamma^k r_{t+k} \mid s_t, a_t] \quad (2)$$

# 方法与模型

$$z_t \sim \pi_\theta(z \mid h_t, s_t)$$



Step 1.

    加载 pre-trained 模型

Step 2.

    微调策略网络

    For t in all_turn do

        For m in M do

            从策略网络中采样 $z_t^m \sim \pi_{\Theta_2}(z \mid s_t, h_t)$

        End for

    更新策略网络 $\pi_{\Theta_2}$ 参数：
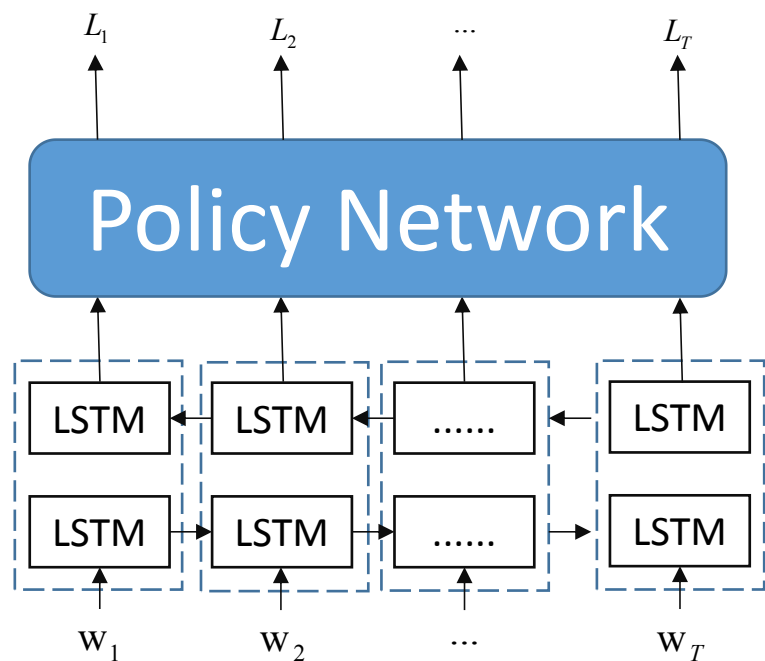
$$\frac{\partial J}{\partial \Theta_2} \approx \frac{1}{M}\sum_{m=1}^{M} R_t^m \frac{\partial \log \pi_{\Theta_2}(z_t^m \mid s_t, h_t)}{\partial \Theta_2}$$

    End for

# 序列标注



$$L(\theta) = -\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T} \log(P_\theta(L_t^n \mid s_t; \theta)) * Q(s_t, L_t^n) \quad (3)$$

$$Q(s_t, L_t) = E\left[\sum_{k=0}^{T} \gamma^k r_{t+k} \mid s_t, L_t\right] \quad (4)$$

# 动态激活函数选择



**Reward**

Softmax

$a \in \{'None', 'relu', 'relu6', 'rrelu', 'selu', 'sigmoid', 'elu', 'softplus', 'tanh', 'hardshrink', 'tanhshrink', 'hardtanh', 'leaky\_relu'\}$

Dense Layer

...

Dense Layer

Conv Layer

image input

$a_T$

$a_1$

$s_1$

$s_T$

$a_1$

$a_2$

...

$a_T$

Policy Network

LSTM    LSTM    ......    LSTM

LSTM    LSTM    ......    LSTM

# 结果

Policy Action Counter

- action_1_counter
- action_2_counter

| Label | action_1_counter | action_2_counter |
|---|---|---|
| None | 473 | 716 |
| ...ional.relu | 955 | 968 |
| ...onal.relu6 | 760 | 964 |
| ...onal.rrelu | 4359 | 3716 |
| ...ional.selu | 556 | 686 |
| ...l.sigmoid | 521 | 481 |
| ...tional.elu | 2301 | 1794 |
| ...l.softplus | 782 | 797 |
| ...onal.tanh | 644 | 727 |
| ...ardshrink | 454 | 728 |
| ...anhshrink | 732 | 844 |
| ...hardtanh | 817 | 981 |
| ...eaky_relu | 1983 | 1935 |

# 数据增强

# 数据增强

$$L(\theta_{enc}, \theta_{dec}, Z_{emb}) = \lambda_{auto}[L_{auto}(\theta_{enc}, \theta_{dec}, Z_{emb}, src) + L_{auto}(\theta_{enc}, \theta_{dec}, Z_{emb}, tgt)] +$$

$$\lambda_{cd}[L_{cd}(\theta_{enc}, \theta_{dec}, Z_{emb}, src, tgt) + L_{cd}(\theta_{enc}, \theta_{dec}, Z_{emb}, tgt, src)] +$$
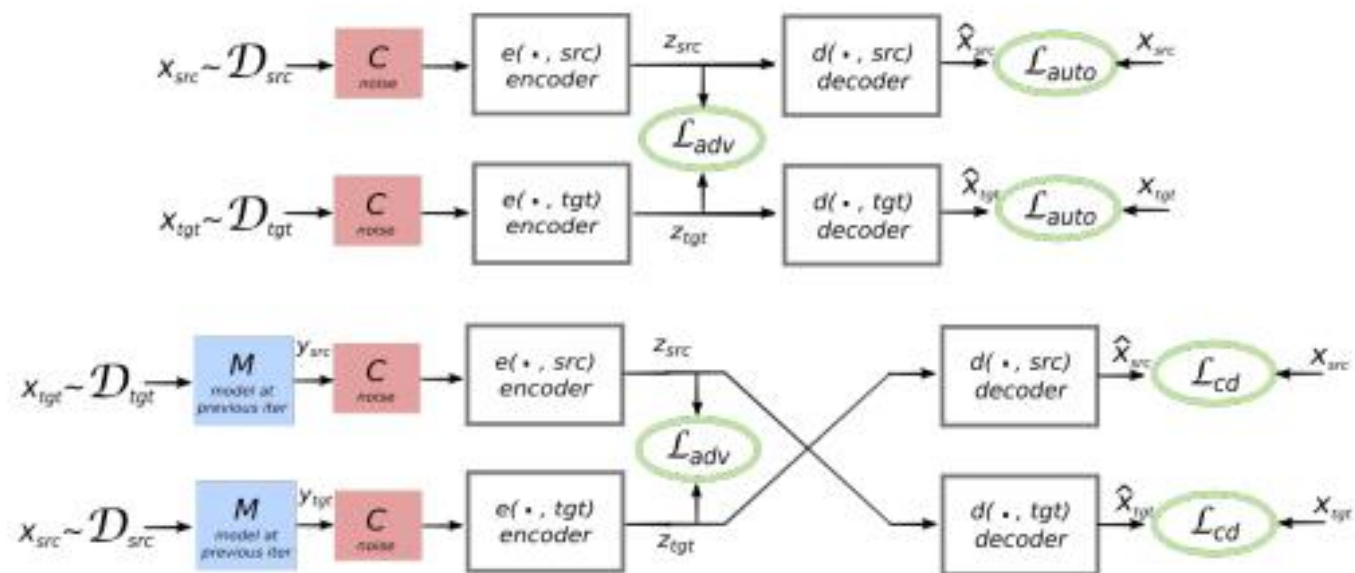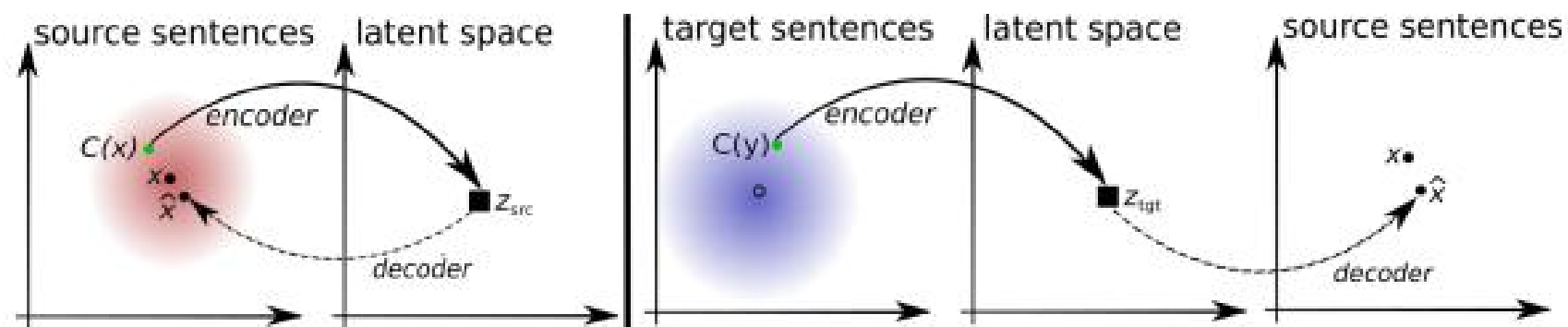
$$\lambda_{adv}L_{adv}(\theta_{enc}, Z_{emb} | \theta_D)$$

where

$$L_{auto}(\theta_{enc}, \theta_{dec}, Z_{emb}, l) = E_{x \sim D_l, \tilde{x} \sim d(e(C(x),l),l)}[\Delta(\tilde{x}, x)]$$

$$L_{cd}(\theta_{enc}, \theta_{dec}, Z_{emb}, l_1, l_2) = E_{x \sim D_{l_1}, \tilde{x} \sim d(e(C(M(x)),l_2),l_1)}[\Delta(\tilde{x}, x)]$$

$$L_{adv}(\theta_{enc}, Z_{emb} | \theta_D) = -E_{(x_i, l_i)}[\log p_D(l_j | e(x_i, l_j))]$$

$$L_D(\theta_D | \theta, Z) = -E_{(x_i, l_i)}[\log p_D(l_i | e(x_i, l_i))]$$

# 数据增强

---

**Algorithm 1** Unsupervised Training for Machine Translation

1: **procedure** TRAINING($\mathcal{D}_{src}, \mathcal{D}_{tgt}, T$)
2:     Infer bilingual dictionary using monolingual data (Conneau et al., 2017)
3:     $M^{(1)} \leftarrow$ unsupervised word-by-word translation model using the inferred dictionary
4:     **for** $t = 1, T$ **do**
5:         using $M^{(t)}$, translate each monolingual dataset
6:         // discriminator training & model training as in eq. 4
7:         $\theta_{\mathrm{discr}} \leftarrow \arg\min \mathcal{L}_D, \quad \theta_{\mathrm{enc}}, \theta_{\mathrm{dec}}, \mathcal{Z} \leftarrow \arg\min \mathcal{L}$
8:         $M^{(t+1)} \leftarrow e^{(t)} \circ d^{(t)}$ // update MT model
9:     **end for**
10:     return $M^{(T+1)}$
11: **end procedure**

---

# 结果

- 帮 我 预 订 个 <入住城市> 旁 边 的 <酒店品牌> ， <入住日期> 入 住 END | 帮 我 预 订 个 <入住城市><地址> 周 边 的 酒店， <入住日期> 入 住
- 我 打 算 查 下 <出发时间> 从 <出发城市> 出 发 的 车 票 END | 我 打 算 查 下 从 <出发城市> 出 发 的 <列车类型><座位类型> ， <出发日期> 的
- 搜 下 <日期><城市> 限 行 状 况 如 何 END | 搜 下 <城市><日期> 的 限 行 情 况
- 给 我 设 定 个 <日期> 的 起 床 闹 钟 哦 END | 给 我 设 定 一 个 起 床 闹 钟， <日期> 的 哦
- 我 想 要 把 <全部范围> 的 <设备名> 亮 度 小 一 点 END | 我 想 要 把 <房间><全部范围> 的 <设备名> 亮 度 小 一 点
- <出发城市> 到 <到达城市> 车 票 还 有 吗， <出发日期> 的 END | 我 有 没 有 <出发城市> 到 <到达城市> 票
- 帮 我 看 一 下 概 念 是 什 么 END | 帮 我 看 一 下 <定义关键词> 的 概 念 是 什 么
- 把 <全部范围><房间><设备名> 速 度 调 高 END | 把 <房间> 的 <设备名> 风 速 调 高
- 我 要 让 <房间><设备名><全部范围> 风 速 不 够 小 小 点 END | 我 要 让 <房间> 的 <设备名> 风 速 小 点
- 我 希 望 设 定 个 <日期><时间区间> 的 早 起 闹 钟 END | 我 希 望 设 定 一 个 <日期><时间区间><时间> 的 早 起 闹 钟
- 我 要 放 一 下 这 一 个 <故事内容名> 的 故 事 END | 我 要 放 一 个 这 一 段 <故事内容名> 的 故 事
- 将 <成语> 的 释 义， 近 义 词， 反 义 词 都 讲 一 下 END | 请 帮 我 查 一 下 <成语> 的 释 义 和 意 思 相 似 的 成 语
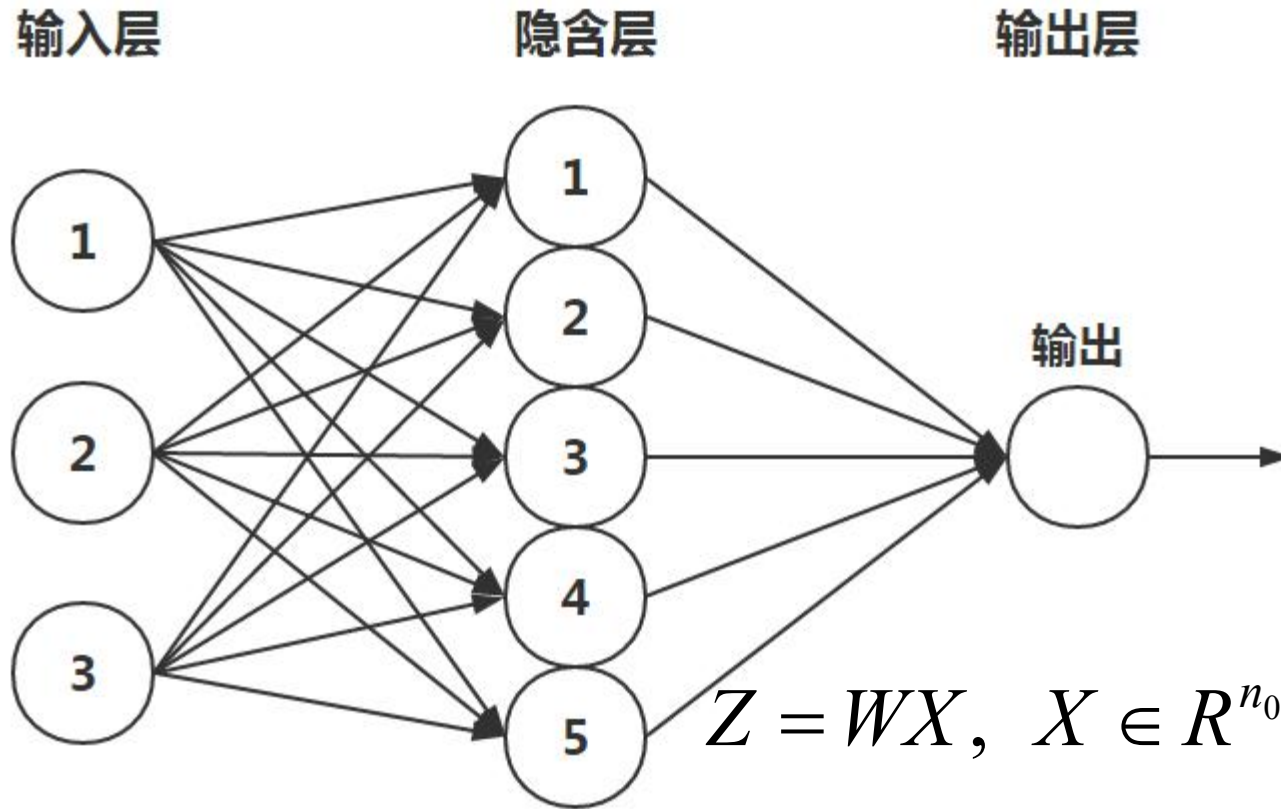- 请 为 我 设 置 <时间> 的 闹 表 END | 请 为 我 设 置 个 闹 铃

# 大纲

- 机器学习系统概览
- 我的NLP研究工作回顾
  - 端到端对话系统
  - 序列标注
  - 动态激活函数选择
  - 数据增强
- **机器学习与RMT结合探索**

# 机器学习与RMT结合探索

Borrow the idea from the statistical physics: approximate the constituents of a large complex system with random variables

- 随机特征(Random Feature)

- Hessian/Jacobian matrix of the loss function

# Basics of Neural Network



输入层　　　　　隐含层　　　　　输出层

输出

$$Z = WX, \ X \in R^{n_0 \times m} \ W \in R^{n_1 \times n_0}$$

$$Y = f(Z), \ \text{non}linear \ activation \ function \ f : R \to R$$

# Central difficulties of DNN

- Non-Convex

- High-Dimensional

# RMT for the analysis of deep learning

*Notation* :

$X \in R^{n_0 \times m}$ $W \in R^{n_1 \times n_0}$ let non*linear activation funtion* $f : R \to R$ *with zero mean and finite moments*

$W$ *and* $X$ *are Gaussian ramdom matrices* with *i.i.d elements* $X_{i\mu} \sim N(0, \sigma_x^2), W_{ij} \sim N(0, \sigma_w^2 / n_0)$

D*efine* $\phi \equiv \dfrac{n_0}{m}, \psi \equiv \dfrac{n_0}{n_1}$ *to be fixed cons* tan t*s*

*Cons*tants $\eta$ *and* $\varsigma$ *defined as* :

$$\eta = \int \frac{e^{-z^2/2}}{\sqrt{2\pi}} f(\sigma_w \sigma_x z)^2 dz \qquad \varsigma = [\sigma_w \sigma_x \int \frac{e^{-z^2/2}}{\sqrt{2\pi}} f^{'}(\sigma_w \sigma_x z) dz]^2$$

# RMT for the analysis of deep learning

Ramdom *Feature Map* :

$Z = WX$

$Y = f(Z)$

$M = \dfrac{1}{m} YY^T \in R^{n_1 \times n_1}$

*Empirical Spectral Density* :

$\rho_M(t) = \dfrac{1}{n_1} \sum\limits_{j=1}^{n_1} \delta(t - \lambda_j(M)), \delta \text{ is the Dirac delta function, } \lambda_j(M) \text{ denote the jth eigenvalue of } M$

To analyze of the eigenvalues(eigenvalues distribution) of the Gram matrix M as it propagates through a neural network.

# RMT for the analysis of deep learning

*For $z \in C \setminus \sup p(\rho_M)$ the Stieltjes transform G of $\rho_M$ :*

$$G(z) = \int \frac{\rho_M(t)}{z-t} dt = -\frac{1}{n_1} E[tr(M - zI_{n_1})^{-1}]$$

$$\rho_M(\lambda) = -\frac{1}{\pi} \lim_{\varepsilon \to 0^+} \operatorname{Im} G(\lambda + i\varepsilon)$$

But it's hard to solve the problem!

# RMT for the analysis of deep learning

$$G(z) = \int \frac{\rho_M(t)}{z-t} dt = -\frac{1}{n_1} E[tr(M - zI_{n_1})^{-1}]$$

$$\rho_M(\lambda) = -\frac{1}{\pi} \lim_{\varepsilon \to 0^+} \operatorname{Im} G(\lambda + i\varepsilon)$$

*Moment Method* :

$$G(z) = \sum_{k=0}^{\infty} \frac{m_k}{z^{k+1}} , \; m_k \; is \; the \; kth \; moment \; of \; the \; \text{distribution} \; \rho_M$$

$$\mathrm{m}_k = \int \rho_M(t) t^k dt = \frac{1}{n_1} E[trM^k]$$

$$\frac{1}{n_1} E[trM^k] = \frac{1}{n_1} E[\sum_{i_1,\ldots,i_k \in [n_1]} M_{i_1 i_2} M_{i_2 i_3} \ldots M_{i_{k-1} i_k} M_{i_k i_1}]$$

**untraceable when n is infinite**

# RMT for the analysis of deep learning

*the Stieltjes transform of the spectral density of M satisfies,*

$$G(z) = \frac{\psi}{z} P\left(\frac{1}{z\psi}\right) + \frac{1-\psi}{z},$$

*where,*

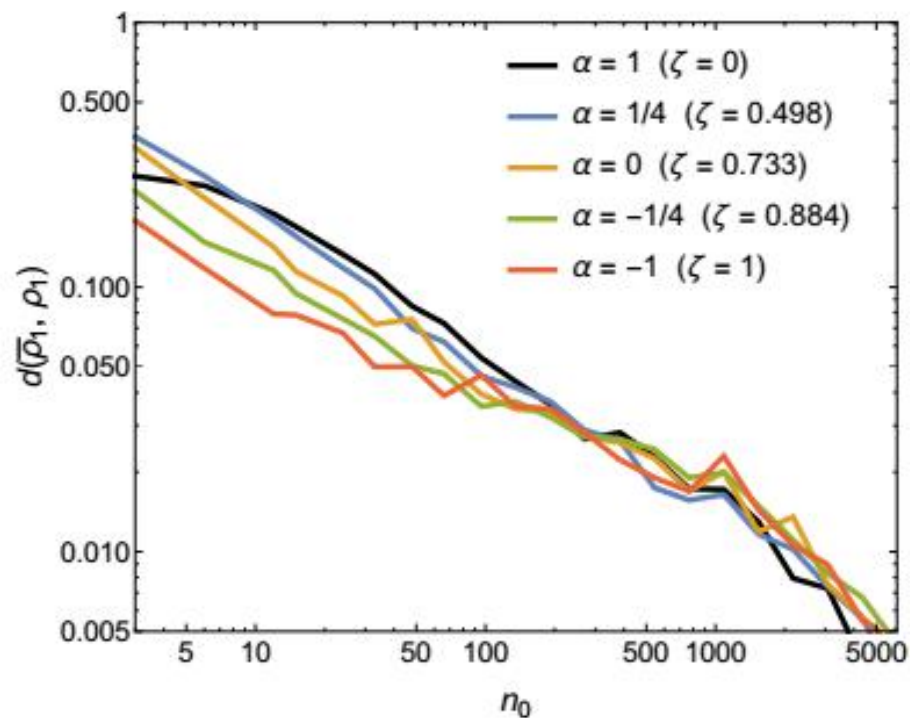$$P = 1 + (\eta - \zeta)tP_\phi P_\psi + \frac{P_\phi P_\psi t\zeta}{1 - P_\phi P_\psi t\zeta},$$

*and*
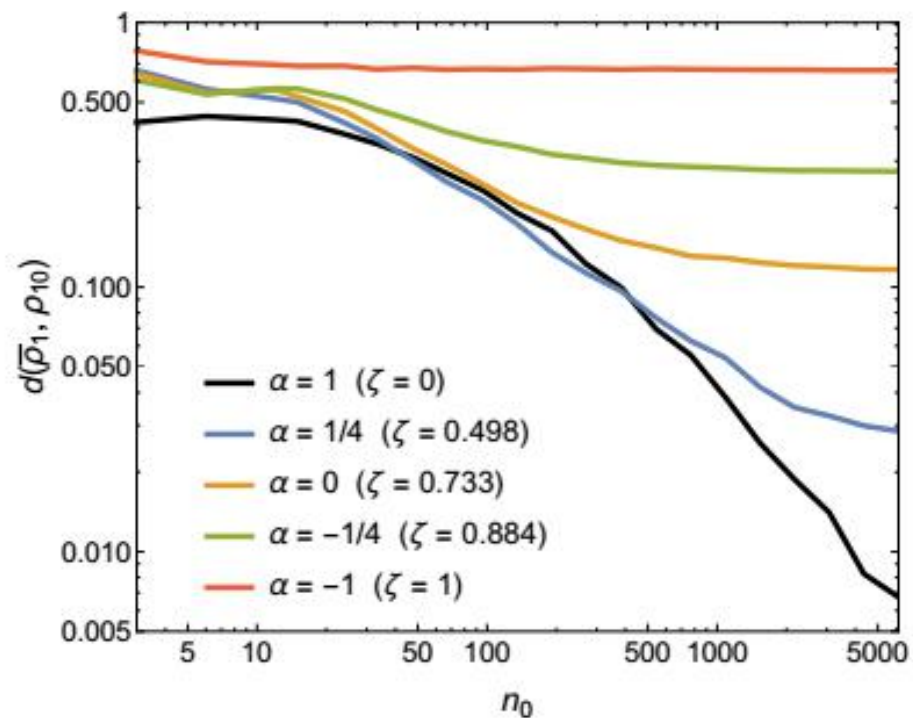
$$P_\phi = 1 + (P-1)\phi, \quad P_\psi = 1 + (P-1)\psi.$$

Another interesting limit is when $\zeta = 0$, which significantly simplifies the expression in eqn. (12). Without loss of generality, we can take $\eta = 1$ (the general case can be recovered by rescaling $z$). The resulting equation is,

$$z\,G^2 + \left(\left(1 - \frac{\psi}{\phi}\right)z - 1\right)G + \frac{\psi}{\phi} = 0, \tag{14}$$

which is precisely the equation satisfied by the Stieltjes transform of the Marchenko-Pastur distribution with shape parameter $\phi/\psi$. Notice that when $\psi = 1$, the latter is the limiting spectral distribution of $XX^T$, which implies that $YY^T$ and $XX^T$ have the same limiting spectral distribution. Therefore we have identified a novel type of isospectral nonlinear transformation. We investigate this observation in Section 4.1.

Figure 1: Distance between the (a) first-layer and (b) tenth-layer empirical eigenvalue distributions of the data covariance matrices and our theoretical prediction for the first-layer limiting distribution $\bar{\rho}_1$, as a function of network width $n_0$. Plots are for shape parameters $\phi = 1$ and $\psi = 3/2$. The different curves correspond to different piecewise linear activation functions parameterize by $\alpha$: $\alpha = -1$ is linear, $\alpha = 0$ is (shifted) relu, and $\alpha = 1$ is (shifted) absolute value. In (a), for all $\alpha$, we see good convergence of the empirical distribution $\rho_1$ to our asymptotic prediction $\bar{\rho}_1$. In (b), in accordance with our conjecture, we find good agreement between $\bar{\rho}_1$ and the tenth-layer empirical distribution $\zeta = 0$, but not for other values of $\zeta$. This provides evidence that when $\zeta = 0$ the eigenvalue distribution is preserved by the nonlinear transformations.

# RMT for the analysis of deep learning

Highly skewed distributions indicate strong anisotropy in the embedded feature space, which is a form of poor conditioning that is likely to derail or impede learning.

# Limitations

Strong Requirements:

- Activation function f : R→R with zero mean and finite moments
- Gaussian assumption
- $\varsigma = [\sigma_w \sigma_x \int \frac{e^{-z^2/2}}{\sqrt{2\pi}} f^{'}(\sigma_w \sigma_x z) dz]^2 = 0$

# Hessian matrix of the loss function

*Decompose Hessian matrix* at critical points int *o two pieces* : $H = H_0 + H_1$

*whe* re

$$[H_0]_{\alpha\beta} \equiv \frac{1}{m} \sum_{i,\mu=1}^{n,m} \frac{\partial \hat{y}_{i,\mu}}{\partial \theta_\alpha} \frac{\partial \hat{y}_{i,\mu}}{\partial \theta_\beta} \equiv \frac{1}{m} [JJ^T]_{\alpha\beta} \text{ , } J \text{ is Jacobian matrix}$$

$$[H_1]_{\alpha\beta} \equiv \frac{1}{m} \sum_{i,\mu=1}^{n,m} e_{i,\mu} \left( \frac{\partial^2 \hat{y}_{i,\mu}}{\partial \theta_\alpha \partial \theta_\beta} \right)$$

$$e_{i,\mu} = \hat{y}_{i,\mu} - y_{i,\mu}$$

# Two assumptions

## Primary assumptions:

1. The matrices $H_0$ and $H_1$ are *freely independent*, a property we discuss in sec. 3.

2. The residuals are i.i.d. normal random variables with tunable variance governed by $\epsilon$, $e_{i\mu} \sim \mathcal{N}(0, 2\epsilon)$. This assumption allows the gradient to vanish in the large $m$ limit, specifying our analysis to critical points.

3. The data features are i.i.d. normal random variables.

4. The weights are i.i.d. normal random variables.

## Secondary assumption:

The elements of $J$ and $H_1$ are i.i.d normal random variables.

# RMT for the analysis of deep learning

*Take $\sigma_{H_0} = 1$ and $\sigma_{H1} = \sqrt{2\varepsilon}$, where $2\varepsilon$ is the variance of $e_{i,\mu}$, $e_{i,\mu} \sim N(0, 2\varepsilon)$*

*then have $\rho_{H_0}(\lambda) = \rho_{MP}(\lambda; 1, \phi)$, $\rho_{H_1}(\lambda) = \rho_{SC}(\lambda; \sqrt{2\varepsilon}, \phi)$*

*where $\phi = 2n/m$*

$$\rho_{\mathrm{MP}}(\lambda; \sigma, \phi) = \begin{cases} \rho(\lambda) & \text{if } \phi < 1 \\ (1 - \phi^{-1})\delta(\lambda) + \rho(\lambda) & \text{otherwise} \end{cases} \qquad \rho_{\mathrm{SC}}(\lambda; \sigma, \phi) = \begin{cases} \frac{1}{2\pi\sigma^2}\sqrt{4\sigma^2 - \lambda^2} & \text{if } |\lambda| \leq 2\sigma \\ 0 & \text{otherwise} \end{cases}$$

where $\phi = n/p$ and,

$$\rho(\lambda) = \frac{1}{2\pi\lambda\sigma\phi}\sqrt{(\lambda - \lambda_-)(\lambda_+ - \lambda)}$$

$$\lambda_\pm = \sigma(1 \pm \sqrt{\phi})^2.$$

# RMT for the analysis of deep learning

For :

$$G(z) = \int_R \frac{\rho(t)}{z - t} dt \quad (\textit{Stieltjes transform})$$

$$R(G(z)) + \frac{1}{G(z)} = z \quad (\textit{R transform})$$

$$R_{H_0 + H_1} = R_{H_0} + R_{H_1} \quad (H_0 \textit{ and } H_1 \textit{ are freely independent})$$
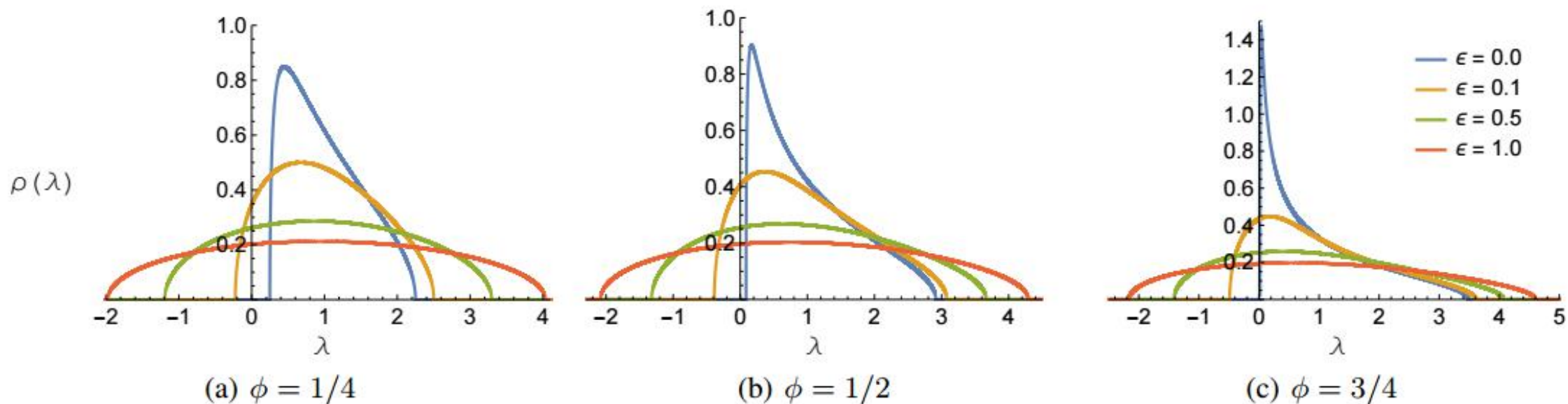
*We can have :*

$$R_{H_0}(z) = \frac{1}{1 - z\phi} \, , \, R_{H_1}(z) = 2\varepsilon z$$

$$R_H(z) = R_{H_0} + R_{H_1} = \frac{1}{1 - z\phi} + 2\varepsilon z$$

$$2\varepsilon\phi G_H^3 - (2\varepsilon + z\phi)G_H^2 + (z + \phi - 1)G_H - 1 = 0 \, , \, G_H \sim 1/z \textit{ as } z \to \infty \, (\textit{Tao}, 2012)$$

# RMT for the analysis of deep learning



Figure 1. Spectral distributions of the Wishart + Wigner approximation of the Hessian for three different ratios of parameters to data points, $\phi$. As the energy $\epsilon$ of the critical point increases, the spectrum becomes more semicircular and negative eigenvalues emerge.

# RMT for the analysis of deep learning

What does that mean?

# RMT for the analysis of deep learning

$H = U \sum U^T = \sum \mathbf{u}_i u_i^T \sigma_i$ , *where* $(\sigma_i, u_i)$ *are eigenvalue and eigenvector*

$l(W) \approx l(W^*) + (W - W^*)^T \nabla l(W^*) + \dfrac{1}{2}(W - W^*)^T H(W^*)(W - W^*), W^*$ *is a critical po*int

$\approx l(W^*) + \dfrac{1}{2}(W - W^*)^T H(W^*)(W - W^*)$

$\approx l(W^*) + \dfrac{1}{2}(W - W^*)^T U_H \sum_H U_H^T (W - W^*)$

$\approx l(W^*) + \sum \sigma_i \dfrac{1}{2}(W - W^*)^T u_i u_i^T (W - W^*)$

$\Rightarrow$

$l(W) - l(W^*) \approx \sum \sigma_i \dfrac{1}{2}(W - W^*)^T u_i u_i^T (W - W^*)$

# RMT for the analysis of deep learning

So we can conclude:

- As the energy ε of the critical point increase, the critical point's probability of being saddle point will increase!
- Increase the ratio Φ(#parameters/#samples), it could be more probably to escape from the saddle points.

# References

- Sutskever I, Vinyals O, Le Q V. Sequence to sequence learning with neural networks[J]. 2014, 4:3104-3112.

- Sordoni A, Bengio Y, Vahabi H, et al. A Hierarchical Recurrent Encoder-Decoder for Generative Context-Aware Query Suggestion[C]// ACM International on Conference on Information and Knowledge Management. ACM, 2015:553-562.

- Pennington J, Worah P. Nonlinear random matrix theory for deep learning[C]//Advances in Neural Information Processing Systems. 2017: 2637-2646.

- Pennington J, Bahri Y. Geometry of neural network loss surfaces via random matrix theory[C]//International Conference on Machine Learning. 2017: 2798-2806.

- Tao, Terence. Topics in random matrix theory, volume 132. American Mathematical Society Providence, RI, 2012.

- https://zhenyu-liao.github.io/pdf/tutorial_Eusipco_handout.pdf

- https://zhuanlan.zhihu.com/p/38777140

谢　谢！