

机器学习（本科生公选课）GEC6531

第12节 机器学习调试 Diagnosing a Learning Algorithm

计算机科学与技术学院

张瑞 教授

邮箱: ruizhang6@hust.edu.cn

签到 & 思考

■ 微助教签到（学校要求）

1. 加入课堂：微信扫码或者通过微助教公众号



二维码有效期至: 2024-11-16

课堂名称: GEC6531 机器学习 (公选课)

课堂编号: OA628

1、扫码关注公众号: 微助教服务号。

2、点击系统通知: “[点击此处加入【GEC6531 机器学习 \(公选课\)】课堂](#)”, 填写学生资料加入课堂。

*如未成功收到系统通知, 请点击公众号下方“学生” - “全部(A)” - “加入课堂” --- “输入课堂编号”手动加入课堂

2. 微信扫码签到

回顾线性回归的损失函数、神经网络

今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

调试机器学习算法

- 假设你实现了一个线性回归算法

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

- 但是发现做预测时误差大 --- 测试误差大，该怎么办

- 常见的办法如下

- 收集更多的训练数据
- 尝试用更少的特征
- 尝试用更多的特征
- 尝试增加多项式特征 ($x_1^2, x_2^2, x_1x_2, \text{etc.}$)
- 减小 λ 值
- 增大 λ 值

机器学习算法的诊断

- 机器学习算法的诊断： 你可以运行的一项测试，从而了解你的算法哪些方面有效或者无效，并获取如何最好地提高其性能的指导。
- 诊断可能需要一些时间来实现，但可能非常值得。

判断算法的好坏：将数据分为训练、验证、测试

■ 诊断的第一步我们需要判断算法的好坏

- 通常将所有数据 D 分成三个子集: D_{TR} 为训练数据, D_{VA} 为验证数据, D_{TE} 为测试数据
- 然后用交叉验证法等来获得训练误差和测试误差 (也叫泛化误差)



今天的目录

■ 调试机器学习算法

- 常见办法
- 机器学习算法的诊断

■ 估计量的偏置 (Bias) 和方差 (Variance)

- 偏置和方差与欠拟合/过拟合的关系
- 偏置和方差与正则化的关系

■ 学习曲线

■ 调试机器学习算法总结

- 神经网络和过拟合

■ 期末考核：结课报告

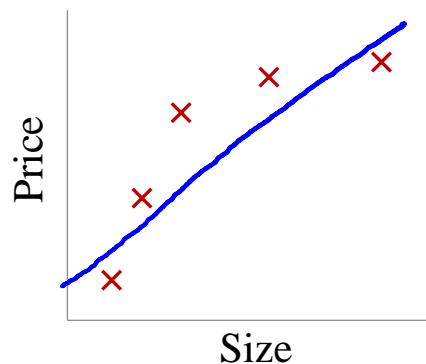
■ 机器学习/人工智能科研

■ 下一步？

估计量的偏置 (Bias) 和方差 (Variance)

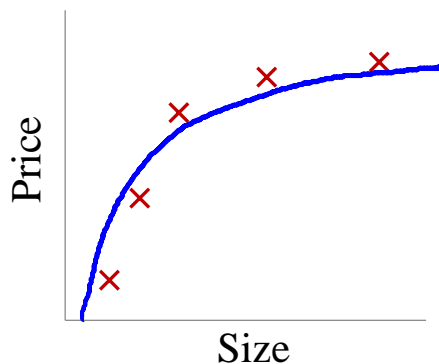
- **估计量的偏置 bias:** $B_{\theta}(\hat{\theta}) = E_{\theta}[\hat{\theta}(X_1, \dots, X_n)] - \theta$
 - 数据真的来自于 p_{θ}
 - 从 p_{θ} 采样出多个数据集, 从每个数据集可以按照模型估算出一个 $\hat{\theta}$, 所以 $\hat{\theta}$ 也服从某个随机分布、是一个随机变量
 - 在多个数据集 S_i 上估算出 $\hat{\theta}_i$, 这些 $\hat{\theta}_i$ 的均值, 也就是 $\hat{\theta}$ 的期望 $E_{\theta}[\hat{\theta}(X_1, \dots, X_n)]$
 - 估计量的偏置定义为 $E_{\theta}[\hat{\theta}(X_1, \dots, X_n)]$ 与 θ 的差, 一般希望偏置越小越好, 当这个差为0时, 我们说这个估计量是无偏的, 简称为**无偏估计**
- **估计量的方差 variance:** $\text{Var}_{\theta}(\hat{\theta}) = E_{\theta}[(\hat{\theta} - E_{\theta}[\hat{\theta}])^2]$
 - 多个 $\hat{\theta}_i$ 的方差的期望, 简单理解为 $\hat{\theta}_i$ 的均方差
 - 一般希望方差越小越好, 一般不为0, 有一个大于0的下界

偏置和方差与欠拟合/过拟合的关系



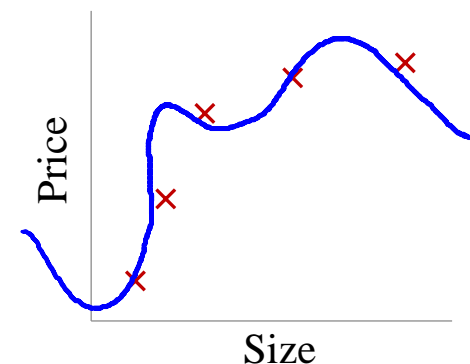
Size
 $\theta_0 + \theta_1 x$

High bias
(underfit)
 $d=1$



Size
 $\theta_0 + \theta_1 x + \theta_2 x^2$

“Just right”
 $d=2$

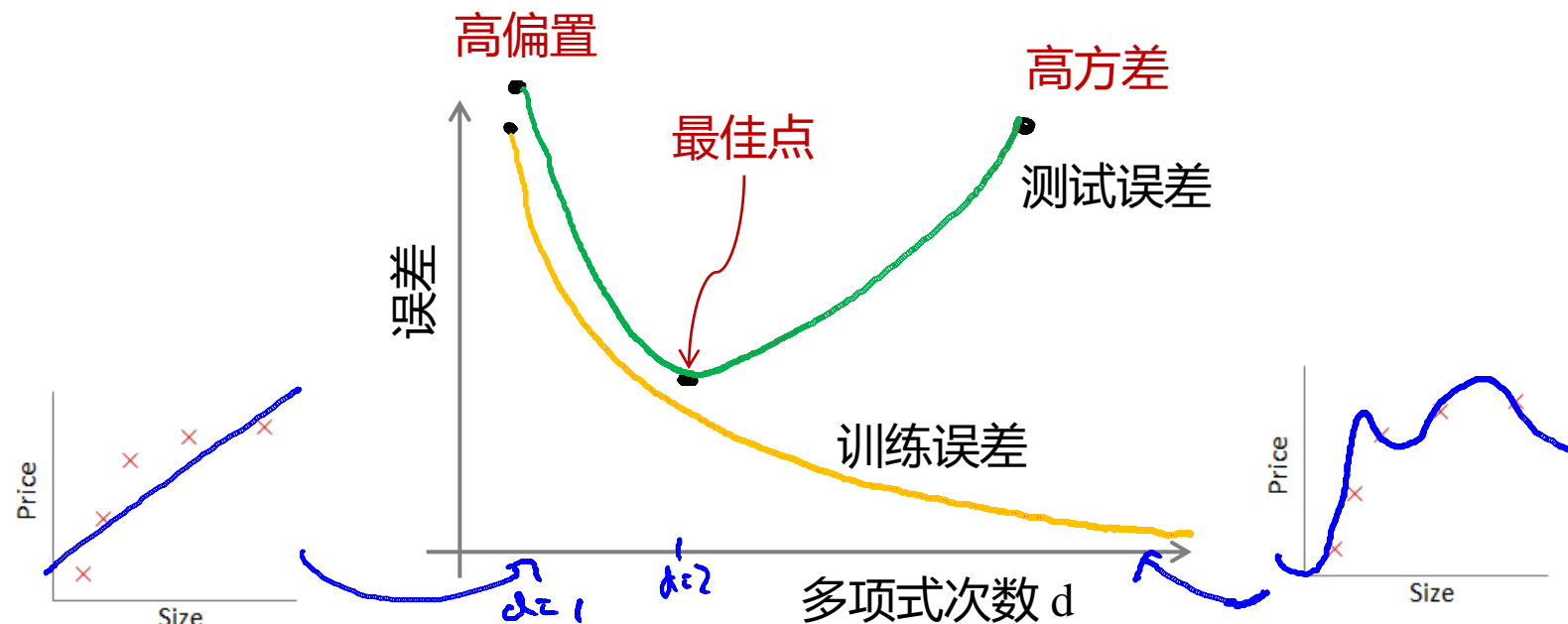


Size
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

High variance
(overfit)
 $d=4$

- 简单的模型：欠拟合，高偏置
- 复杂的模型：过拟合，高方差

误差 vs 回归函数的多项式次数

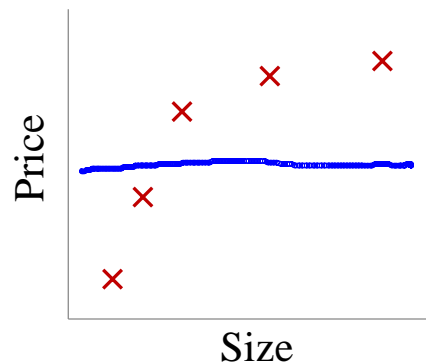


模型测试误差大，是偏置问题还是方差问题？

- 简单的模型：欠拟合，高偏置
- 复杂的模型：过拟合，高方差

如果偏置过高，说明模型欠拟合，需要提高模型复杂度
如果方差过高，说明模型过拟合，需要降低模型复杂度

偏置和方差与正则化的关系

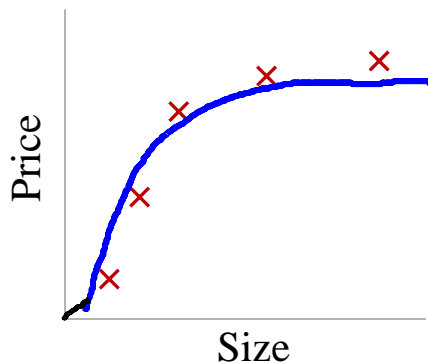


Large λ

High bias (underfit)

$\lambda = 10000$. $\theta_1 \approx 0, \theta_2 \approx 0, \dots$

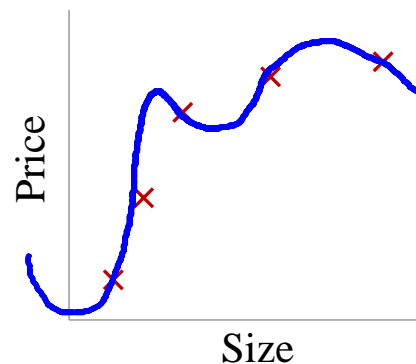
$h_{\theta}(x) \approx \theta_0$



Size

Intermediate λ

“Just right”



Size

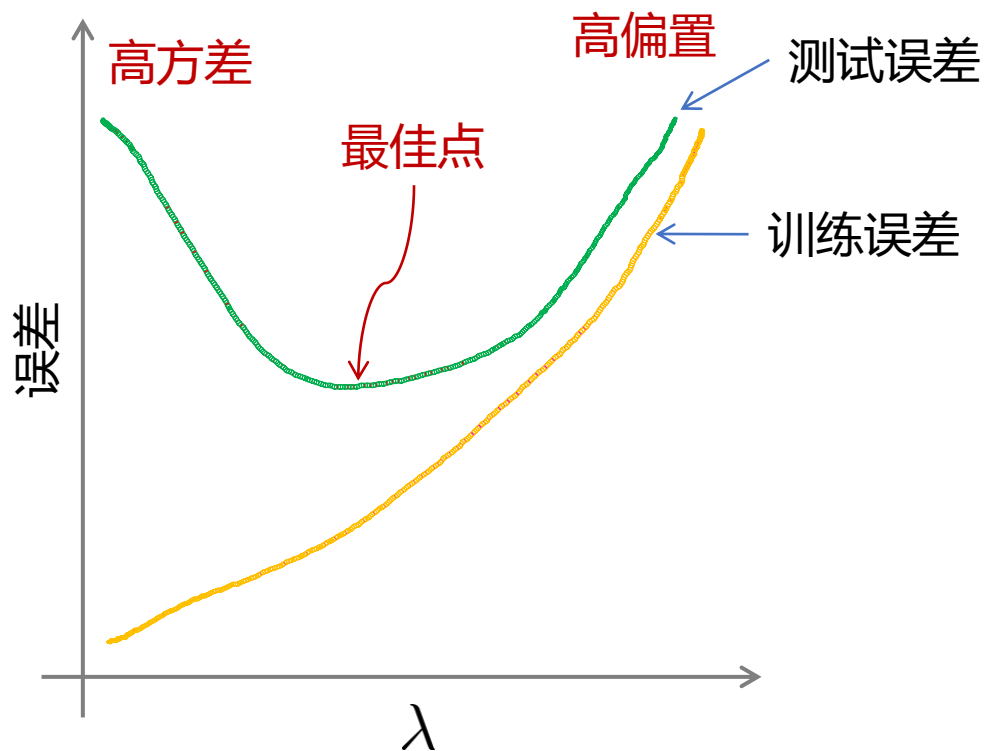
Small λ

High variance (overfit)

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

- 小的 λ : 过拟合, 高方差
- 大的 λ : 欠拟合, 高偏置

误差 vs 正则化参数 λ



- 小的 λ : 过拟合, 高方差
- 大的 λ : 欠拟合, 高偏置

模型测试误差大, 是偏置问题还是方差问题?

如果偏置过高, 说明模型欠拟合, 需要减小 λ

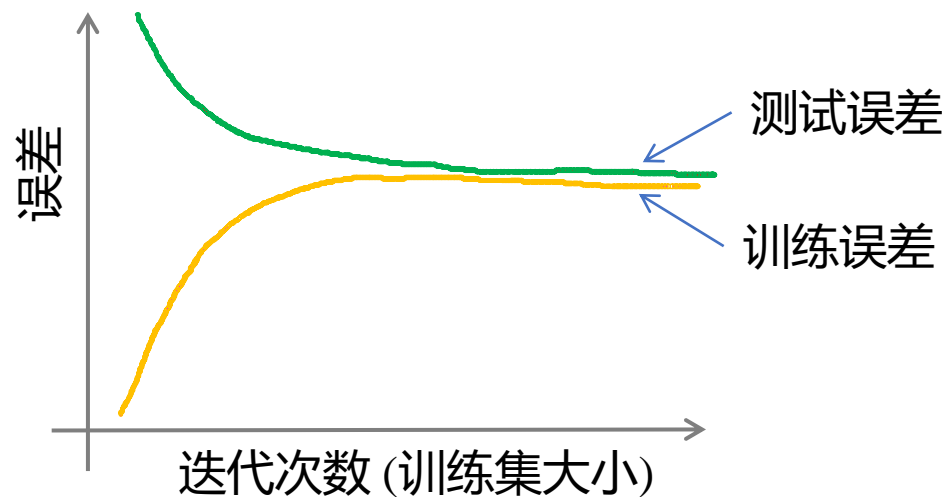
如果方差过高, 说明模型过拟合, 需要增大 λ

今天的目录

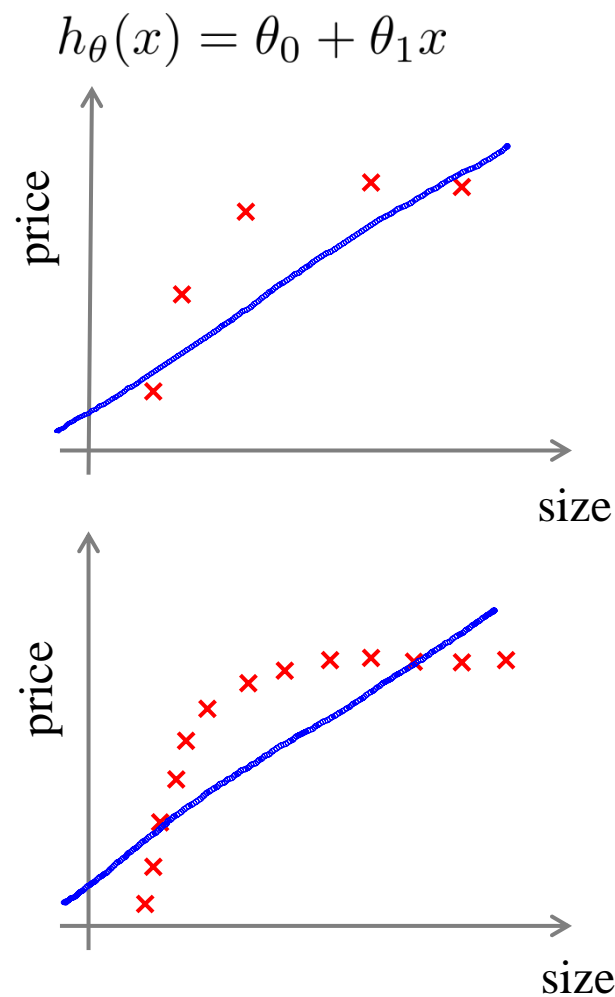
- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

学习曲线：高偏置

高偏置

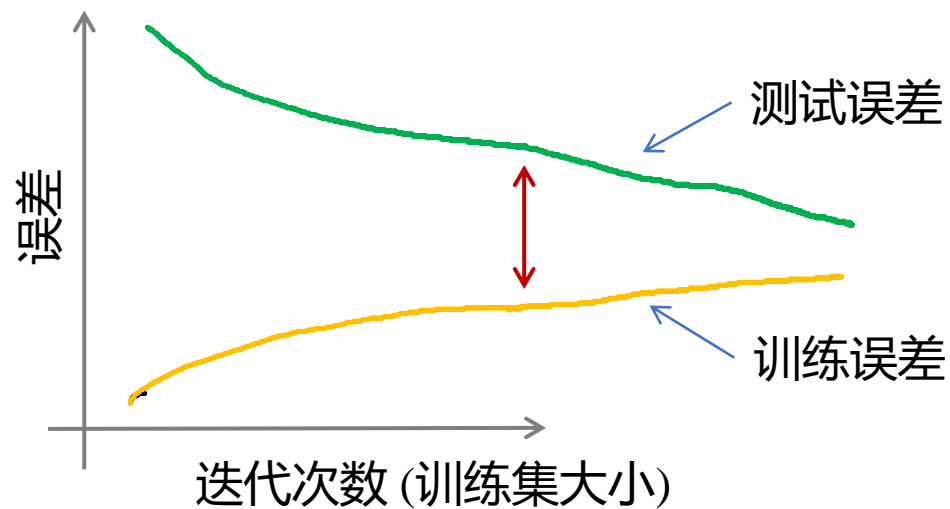


如果是高偏置问题，那么获取更多数据不会带来多少提高



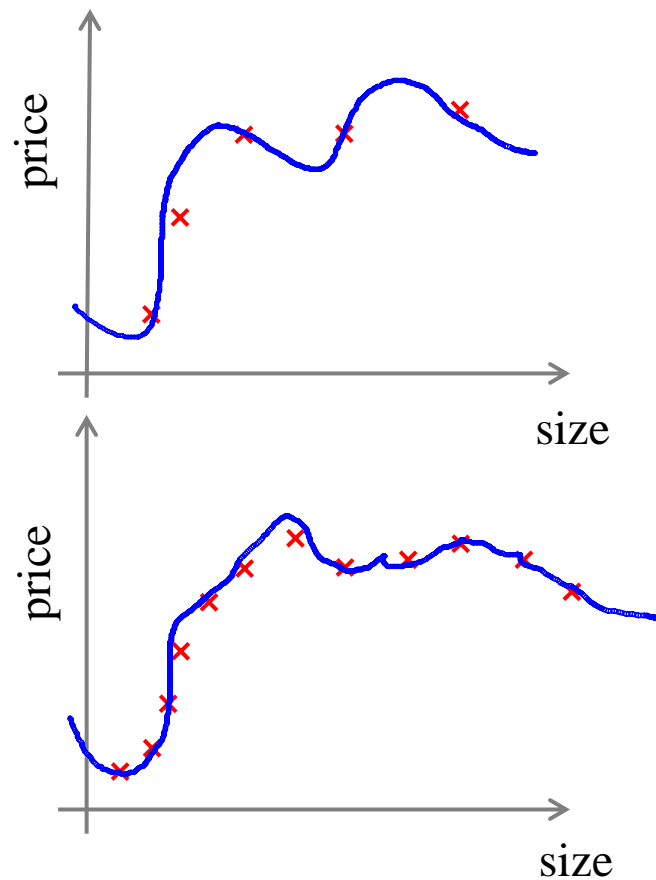
学习曲线：高方差

High variance



如果是高方差问题，那么获取更多数据可能带来提高

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_{100} x^{100}$$



今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

总结：调试机器学习算法

- 假设你实现了一个线性回顾算法

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

- 但是发现做预测时误差大 --- 测试误差大，该怎么办

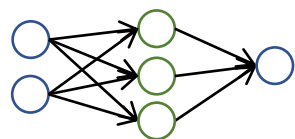
- 常见的办法如下

- 收集更多的训练数据 --- 解决高方差
- 尝试用更少的特征 --- 解决高方差
- 尝试用更多的特征 --- 解决高偏置
- 尝试增加多项式特征 ($x_1^2, x_2^2, x_1x_2, \text{etc.}$) --- 解决高偏置
- 减小 λ 值 --- 解决高偏置
- 增大 λ 值 --- 解决高方差

神经网络和过拟合

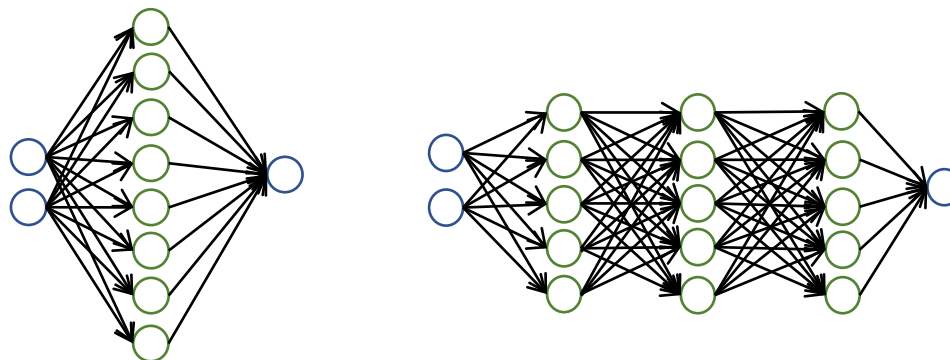
小的神经网络：

- 参数少，容易欠拟合
- 计算更便宜



大的神经网络：

- 参数多，容易过拟合
- 计算更昂贵
- 用**正则化**，**dropout** 等来解决过拟合



今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

期末考核：结课报告（60分）

- 《结课报告》撰写说明，见华中科技大学课程平台，资料
- 网页版：

<https://hustiibd.feishu.cn/wiki/MRs4wEt9kiVRzfksOdPc2CJZnhc/>

今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步？**

机器学习/人工智能科研

- **深度学习 (Deep Learning, DL)**
- **计算机视觉 (Computer Vision, CV)**
- **自然语言处理 (Natural Language Processing, NLP)**
- **图学习 (Graph Learning)**
- **推荐模型 (Recommender Systems)**
- **强化学习 (Reinforcement Learning, RL)**
- **具身智能 (Embodied AI)**
- **AI for Science**

多模态生成 (Multimodal Generation)



midjourney

<https://www.midjourney.com/>



Prompt: Several giant woolly mammoths approach treading through a snowy meadow, their long woolly fur lightly blows in the wind as they walk, snow covered trees and dramatic snow capped... +

Sora

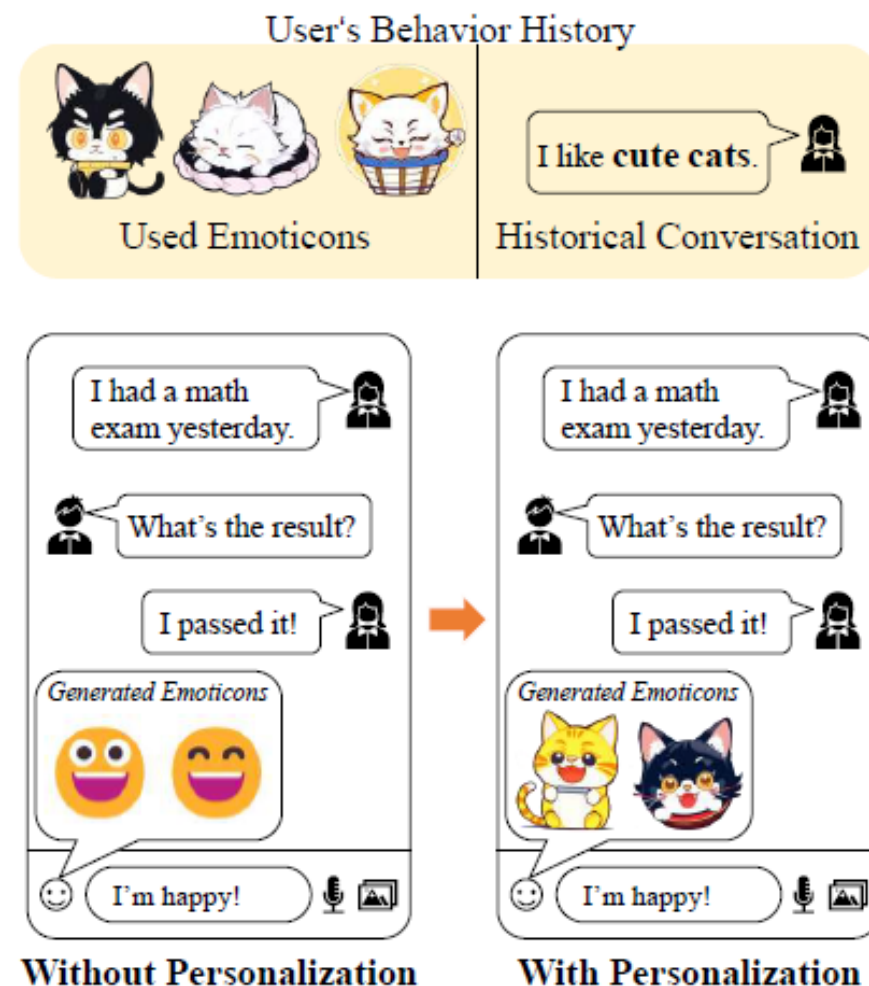
<https://openai.com/index/sora/>

是否能让生成的内容个性化？

个性化多模态生成 (Personalized Multimodal Generation)

■ 个性化多模态生成技术 PMG

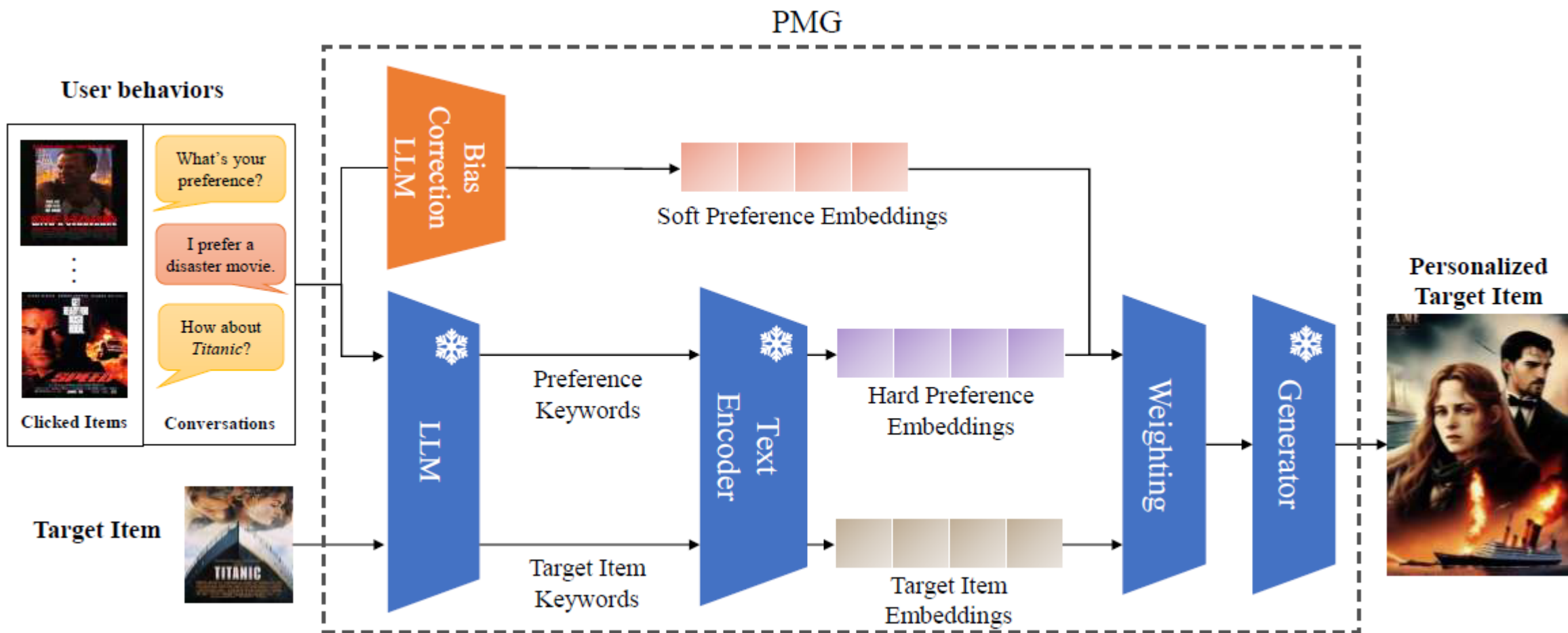
- 将用户行为（对话、点击等）转换为自然语言
- 提取用户偏好描述，包括硬偏好和软偏好嵌入
- 用偏好作为条件控制大模型进行多模态生成
- 在个性化度量方面提高了8%



PMG : Personalized Multimodal Generation with Large Language Models, The Web Conference 2024

https://mp.weixin.qq.com/s/Ysqa_XSXL7nb11q-ZOF6jA 量子位报道

个性化多模态生成 (Personalized Multimodal Generation)



个性化多模态生成 (Personalized Multimodal Generation)

User Behaviors Candidate Target Item	Without Personalization	With Personalization			
	N/A	Business Style	Girl's Style	Boy's Style	Cartoon Style
shoes					
shirt					

$$d_p = \frac{e_M \cdot e_p}{\|e_M\|_2 \|e_p\|_2},$$

$$d_t = \frac{e_M \cdot e_t}{\|e_M\|_2 \|e_t\|_2}.$$

Finally, our objective is to optimize the weighted sum of d_p and d_t .

$$z = \alpha \cdot \log d_p + (1 - \alpha) \cdot \log d_t.$$



(a) $w_p : w_t = 0 : 4$



(b) $w_p : w_t = 1 : 3$



(c) $w_p : w_t = 2 : 2$



(d) $w_p : w_t = 3 : 1$



(e) $w_p : w_t = 4 : 0$

Figure 7: Generated poster of movie *Titanic* with different weights of conditions. w_p is the weight of preference conditions, which prefer disaster movie. w_t is the weight of target item conditions, which consider it as a romantic movie. When $w_p : w_t = 1 : 3$ it achieves the highest z score and the generated poster is a combination of romance and disaster.

知识库/知识图谱 (Knowledge Base/Graph)

□ 一些信息的集合，通常包含大量的非结构化数据，以图结构的形式存储。

- 与传统关系数据库对比
- 广泛应用于搜索引擎（如谷歌的知识图谱）、问答系统、以及信息提取等一般数据挖掘任务中
- 开源的知识库: YAGO, freebase, Wikidata



Bill Gates
American business magnate

 gatesnotes.com

William Henry Gates III is an American business magnate, software developer, investor, and philanthropist. He is best known as the co-founder of Microsoft Corporation. [Wikipedia](#)

Born: 28 October 1955 (age 64 years), [Seattle, Washington, United States](#)

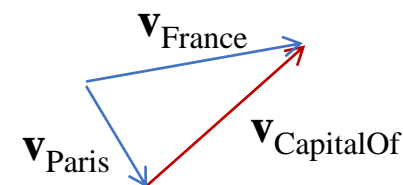
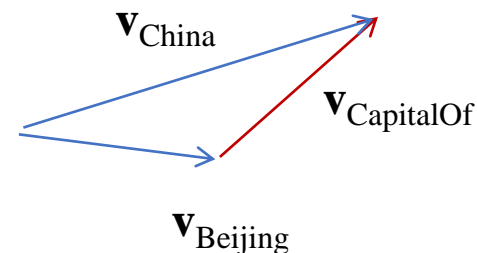
Net worth: 107.4 billion USD (2019)

Spouse: [Melinda Gates](#) (m. 1994)

Children: [Phoebe Adele Gates](#), [Rory John Gates](#), [Jennifer Katharine Gates](#)

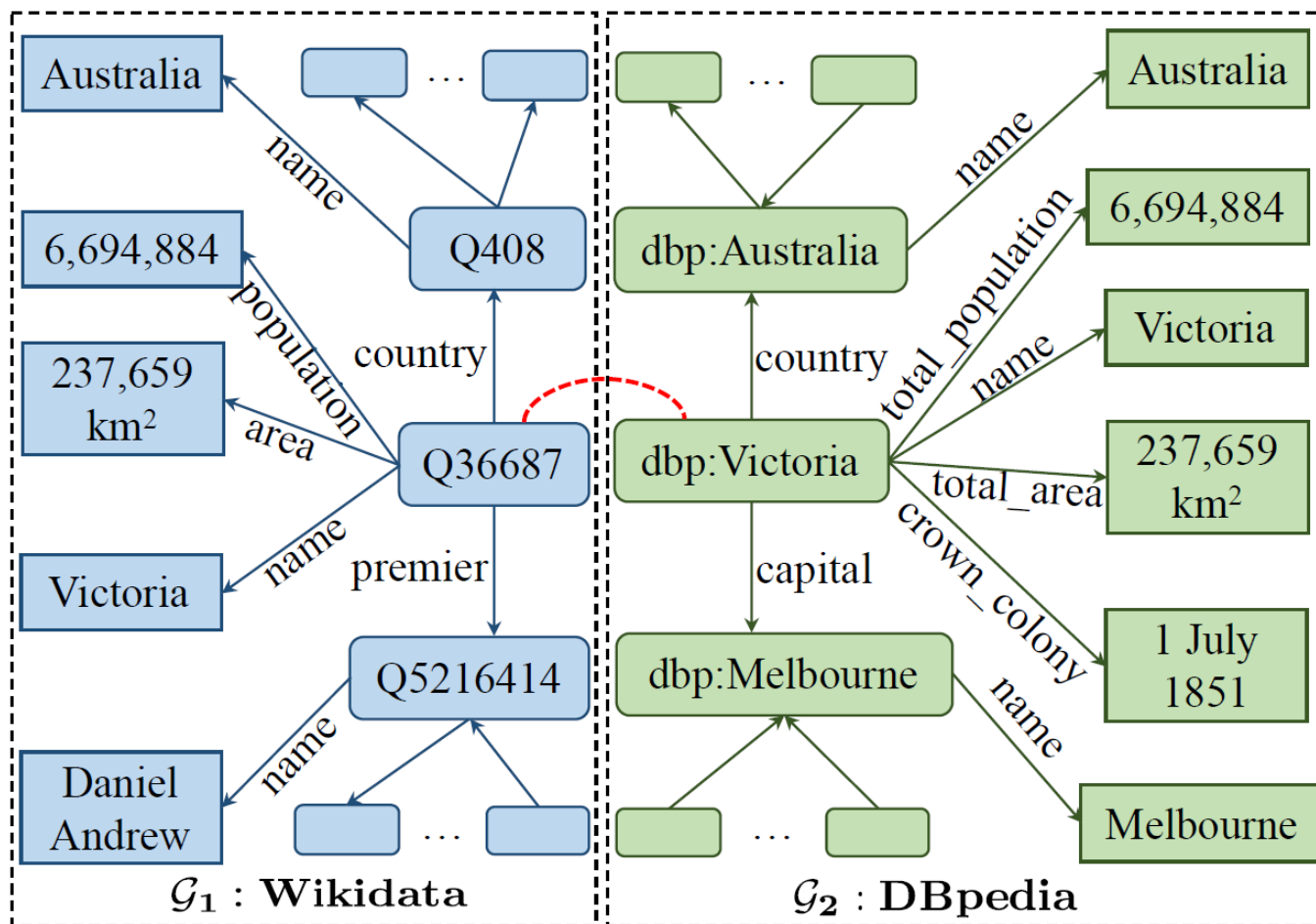
知识图谱的表征 (Representation)

- 知识图谱由三元组组成
 - ◆ $\langle \text{头实体}, \text{关系}, \text{尾实体} \rangle$
 - ◆ $\langle \text{Beijing}, \text{CapitalOf}, \text{China} \rangle$
- 每个实体和关系都用一个向量 \mathbf{v} 来表征
- 三元组 $\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$ 的实体满足 $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$
也就是 $\mathbf{r} \approx \mathbf{t} - \mathbf{h}$
- 所以, $\mathbf{v}_{\text{China}} - \mathbf{v}_{\text{Beijing}} \approx \mathbf{v}_{\text{France}} - \mathbf{v}_{\text{Paris}} \approx \mathbf{v}_{\text{Germany}} - \mathbf{v}_{\text{Berlin}} \approx \dots$
- 随机初始化每个实体和关系, 然后利用上述等式来反复训练



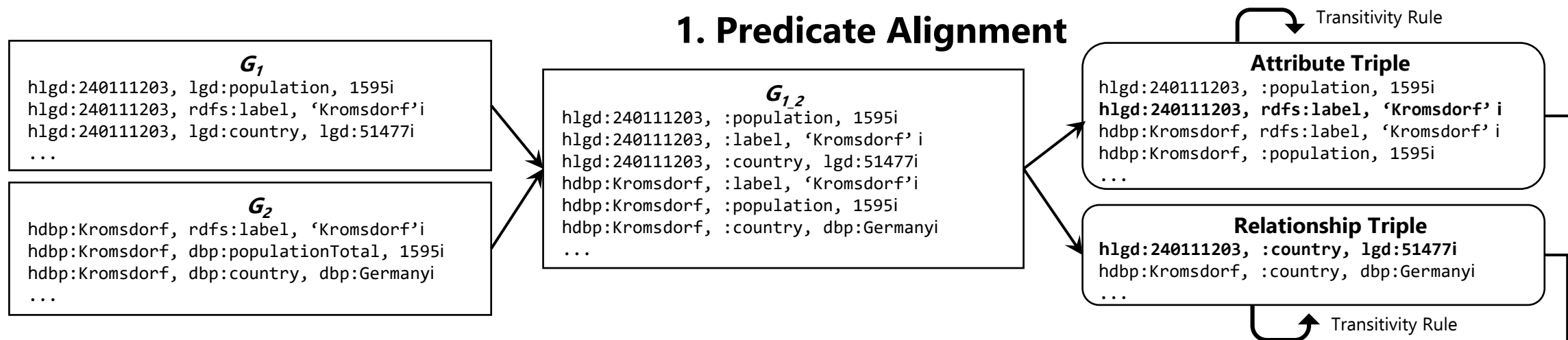
知识图谱的对齐 (Knowledge Graph Alignment)

知识图谱的对齐：找出在不同知识图谱里代表现实中相同实体的对应项

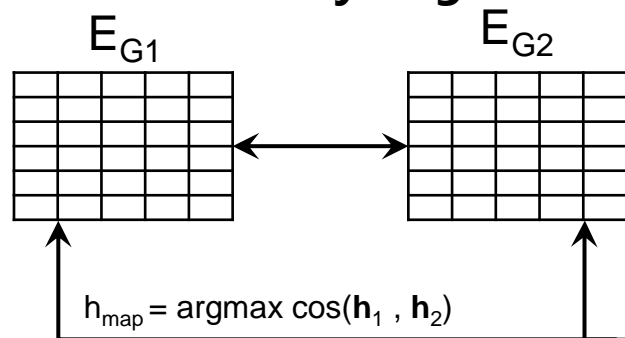


知识图谱的自动对齐算法 (Automatic Knowledge Graph Alignment)

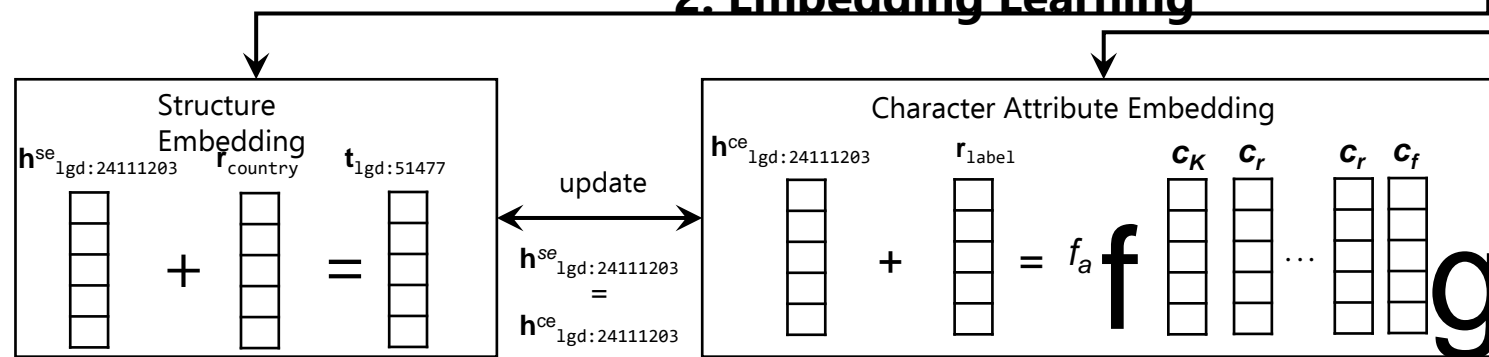
1. Predicate Alignment



3. Entity Alignment



2. Embedding Learning



准确率提高50%

AutoAlign: Fully Automatic and Effective Knowledge Graph Alignment enabled by Large Language Models, IEEE Transactions on Knowledge and Data Engineering (TKDE), 2024.

今天的目录

- **调试机器学习算法**
 - 常见办法
 - 机器学习算法的诊断
- **估计量的偏置 (Bias) 和方差 (Variance)**
 - 偏置和方差与欠拟合/过拟合的关系
 - 偏置和方差与正则化的关系
- **学习曲线**
- **调试机器学习算法总结**
 - 神经网络和过拟合
- **期末考核：结课报告**
- **机器学习/人工智能科研**
- **下一步?**

想进一步学习研究AI的同学

■ 低年级同学：

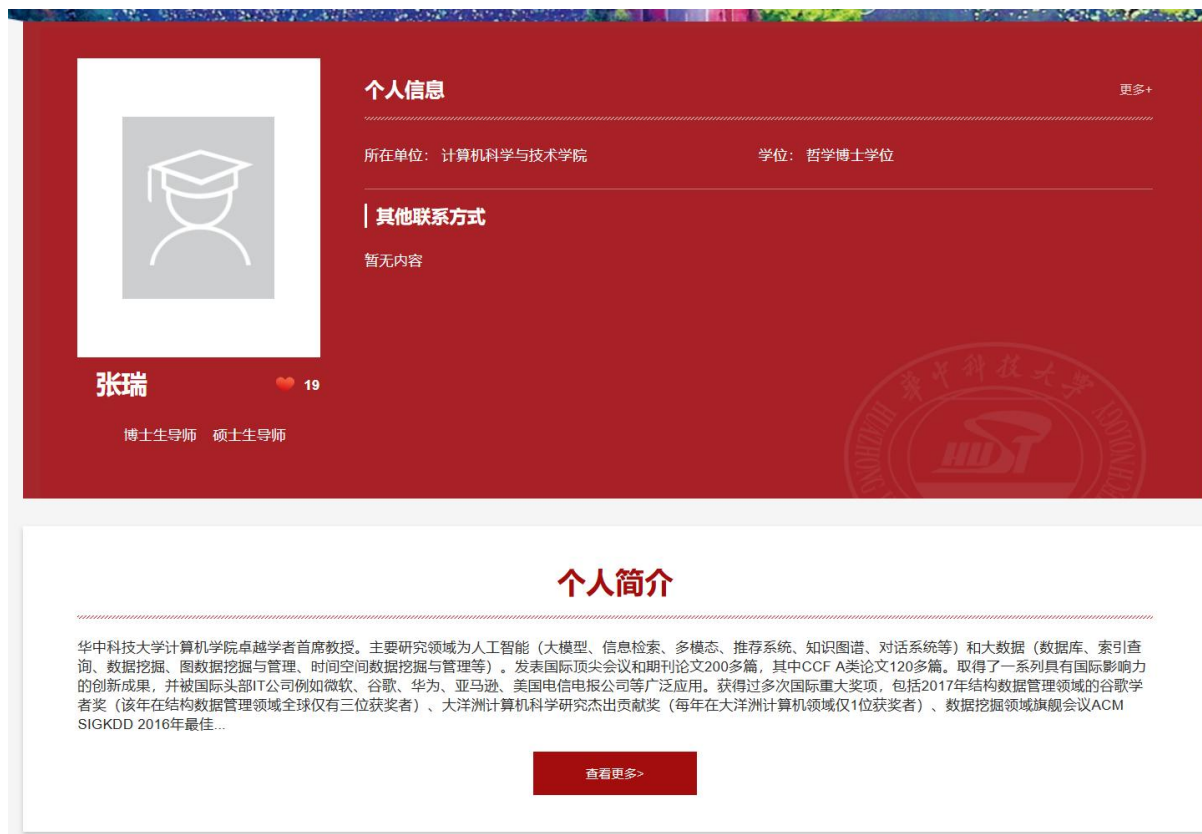
- 进一步学习机器学习/AI相关课程

■ 较高年级同学

- 参与实验室的一些研究
- 毕业设计项目
- 继续读硕士/博士
- 欢迎与我联系 <ruizhang6@hust.edu.cn>, 附上简单的简历&成绩单

如果你们从这门课中收益，想了解更多

■ http://faculty.hust.edu.cn/zr/zh_CN/index.htm



■ 教学调查问卷

祝大家考试好运，获得好成绩！