

RS&NLP选修 Lesson-02

前情提要

- 根据业务场景不同，推荐系统可以演化得差异很大。因此，要注意学习算法背后的思想，做到融会贯通
- 实例：视频相关推荐。拆解和抽象为机器学习问题，进而用SGNS的方式去解决
- NLP和RS是算法的两个重要应用方向，其背后的思想可以找到很多相通和借鉴之处
- 理论提升：和矩阵分解的等价性



回顾：视频相关推荐

- 分为3个步骤：
 - Step1 构建序列：预处理用户的行为日志，使之适配算法
 - Step2 求解相似性：利用skip-gram with negative sampling计算item之间的相似性
 - Step3 Faiss建立索引：利用索引库为item建立索引，供线上推荐使用

Step1 构建序列

- Input: 曝光/点击日志
- Output: 点击序列
- Code: 实操

Step2 求解相似性

- Input: 点击序列
- Output: item \rightarrow vector
- Code: 实操

Step3-1 建立索引

- Input: item -> vector
- Output: 索引文件
- Code: 实操

Step3-2 用索引查找

- Input: 索引文件, key
- Output: key的最近邻
- Code: 实操

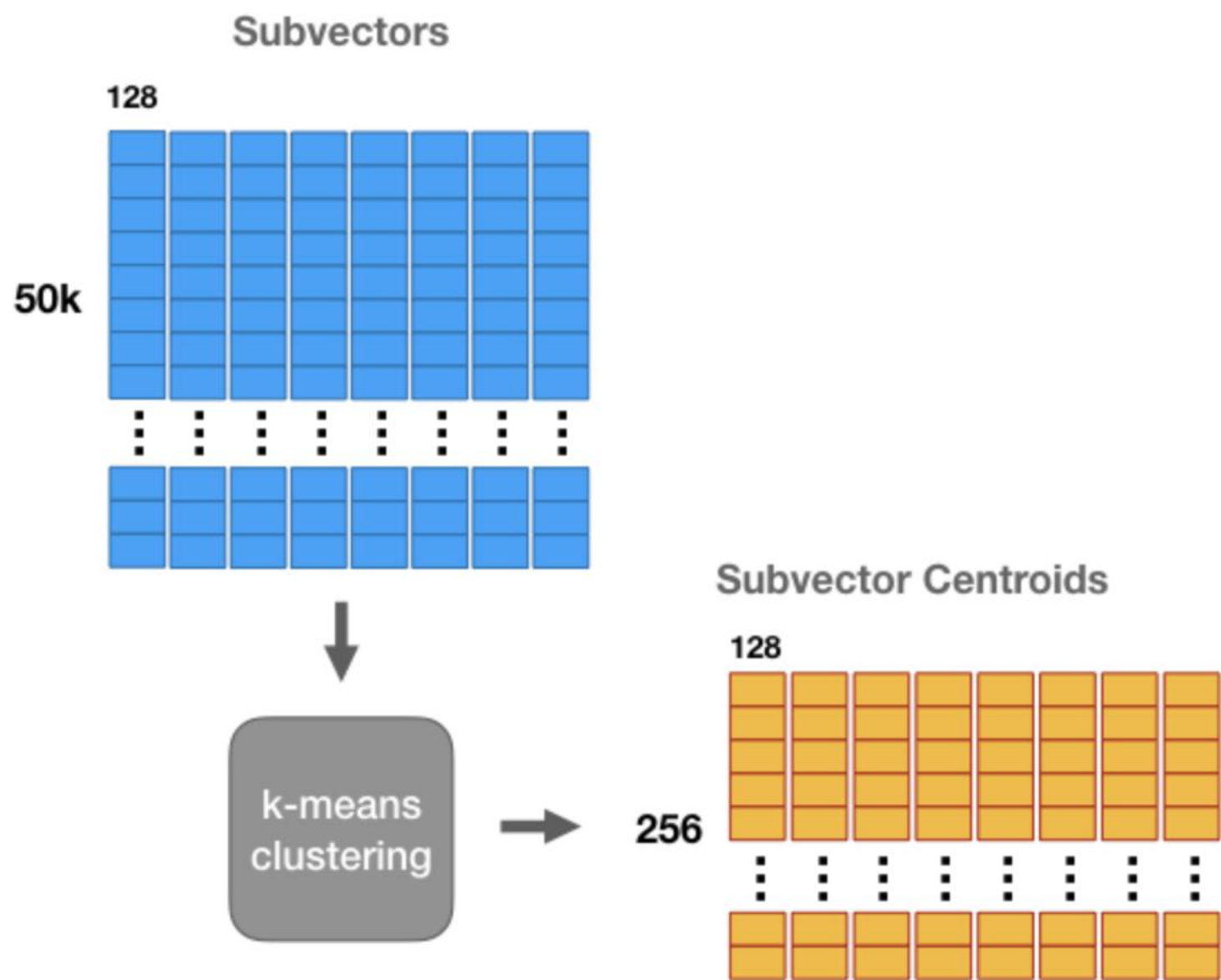
Index factory

Method	Class name	index_factory	Main parameters	Bytes/vector	Exhaustive	Comments
Exact Search for L2	IndexFlatL2	Flat	d	4*d	yes	brute-force
Exact Search for Inner Product	IndexFlatIP	Flat	d	4*d	yes	also for cosine (normalize vectors beforehand)
Hierarchical Navigable Small World graph exploration	IndexHNSWFlat	'HNSWx,Flat'	d, M	4*d + 8 * M	no	
Inverted file with exact post-verification	IndexIVFFlat	IVFx,Flat	quantizer, d, nlists, metric	4*d	no	Take another index to assign vectors to inverted lists
Locality-Sensitive Hashing (binary flat index)	IndexLSH	-	d, nbits	nbits/8	yes	optimized by using random rotation instead of random projections
Scalar quantizer (SQ) in flat mode	IndexScalarQuantizer	SQ8	d	d	yes	4 bit per component is also implemented, but the impact on accuracy may be unacceptable
Product quantizer (PQ) in flat mode	IndexPQ	PQx	d, M, nbits	M (if nbits=8)	yes	
IVF and scalar quantizer	IndexIVFScalarQuantizer	IVFx,SQ4 "IVFx,SQ8"	quantizer, d, nlists, qtype	SQfp16: 2 * d, SQ8: d or SQ4: d/2	no	there are 2 encodings: 4 bit per dimension and 8 bit per dimension
IVFADC (coarse quantizer+PQ on residuals)	IndexIVFPQ	IVFx,PQy	quantizer, d, nlists, M, nbits	M+4 or M+8	no	the memory cost depends on the data type used to represent ids (int or long), currently supports only nbits <= 8
IVFADC+R (same as IVFADC with re-ranking based on codes)	IndexIVFPQR	IVFx,PQy+z	quantizer, d, nlists, M, nbits, M_refine, nbits_refine	M+M_refine+4 or M+M_refine+8	no	

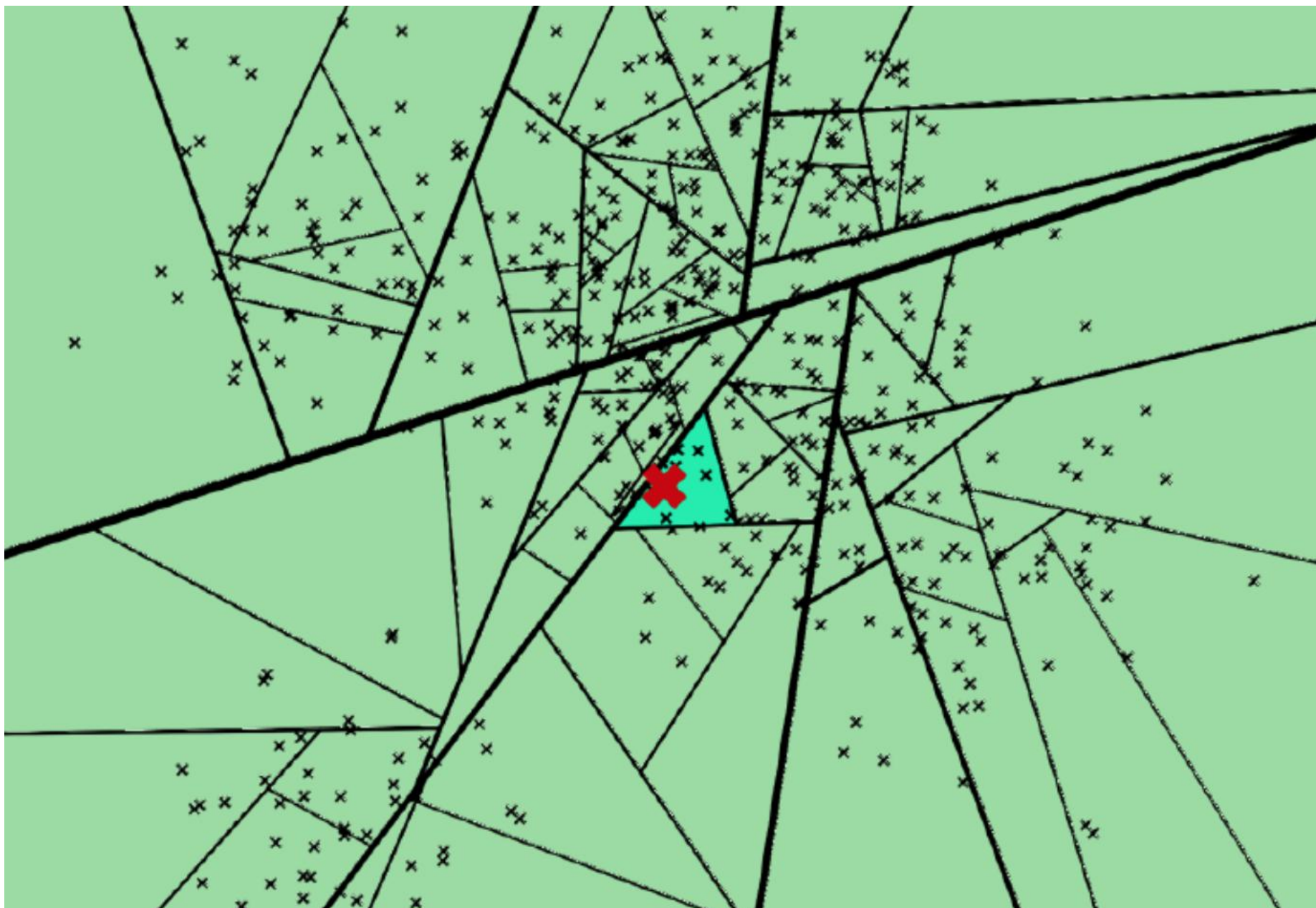
根据需求去选择
具体算法



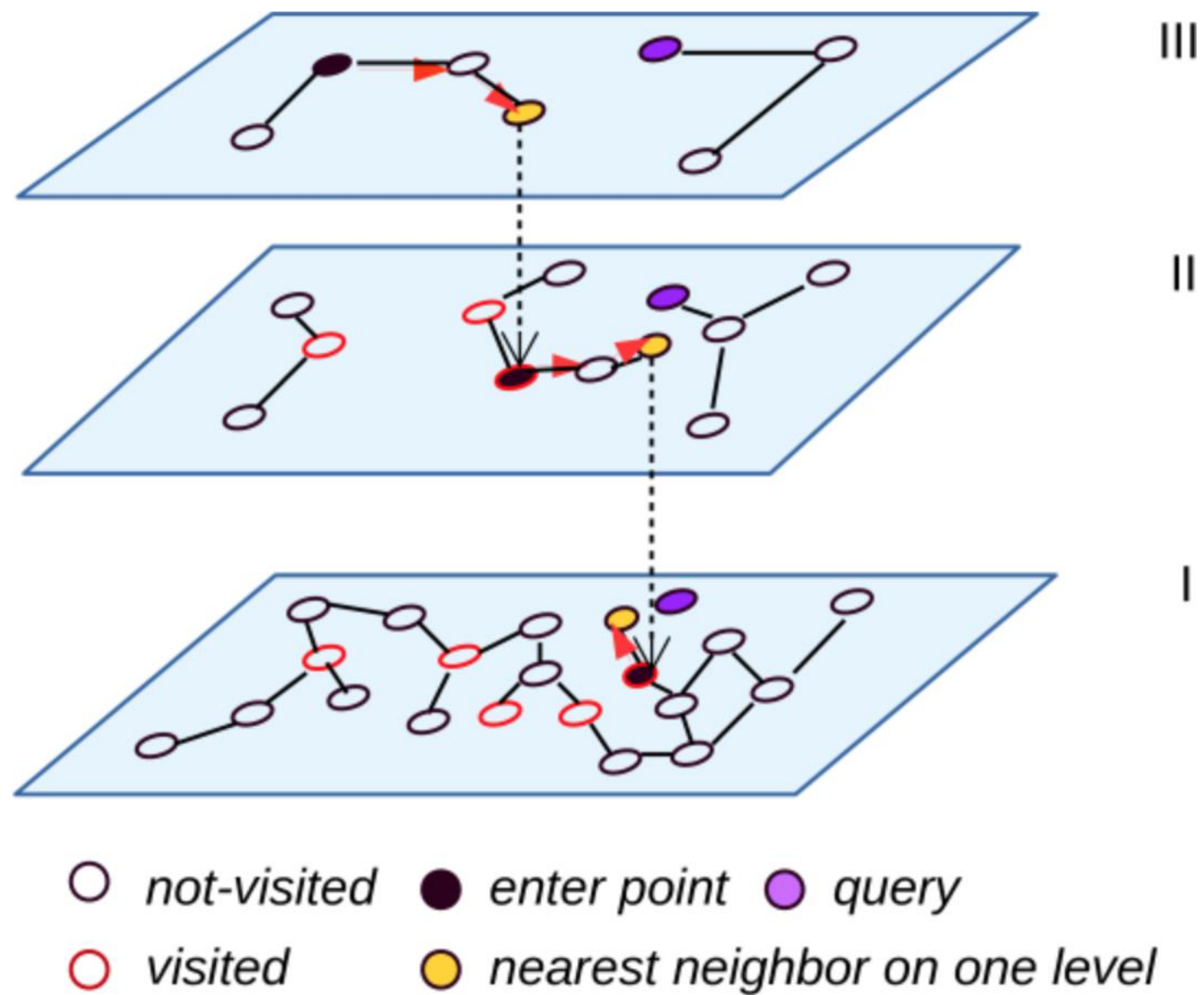
PQ算法



Annoy算法



HNSW算法



选择索引算法时的一些指导原则

- 有多少向量需要建索引
- 索引database更新的频率需要如何
- 召回率是否满足需求
- 索引算法所支持的距离度量是否适配

“有没有更好的方式？”



基于SGNS做召回的局限性

- 无监督方式，不是直接优化目标
- 没法结合side information

深度语义匹配模型DSSM

- 主要贡献
 - 利用点击数据（有监督）
 - 利用深度学习（学习能力强）
 - Word hashing（解决大词表问题）



纸上得来终觉浅， 绝知此事要躬行

TensorFlow还用介绍吗

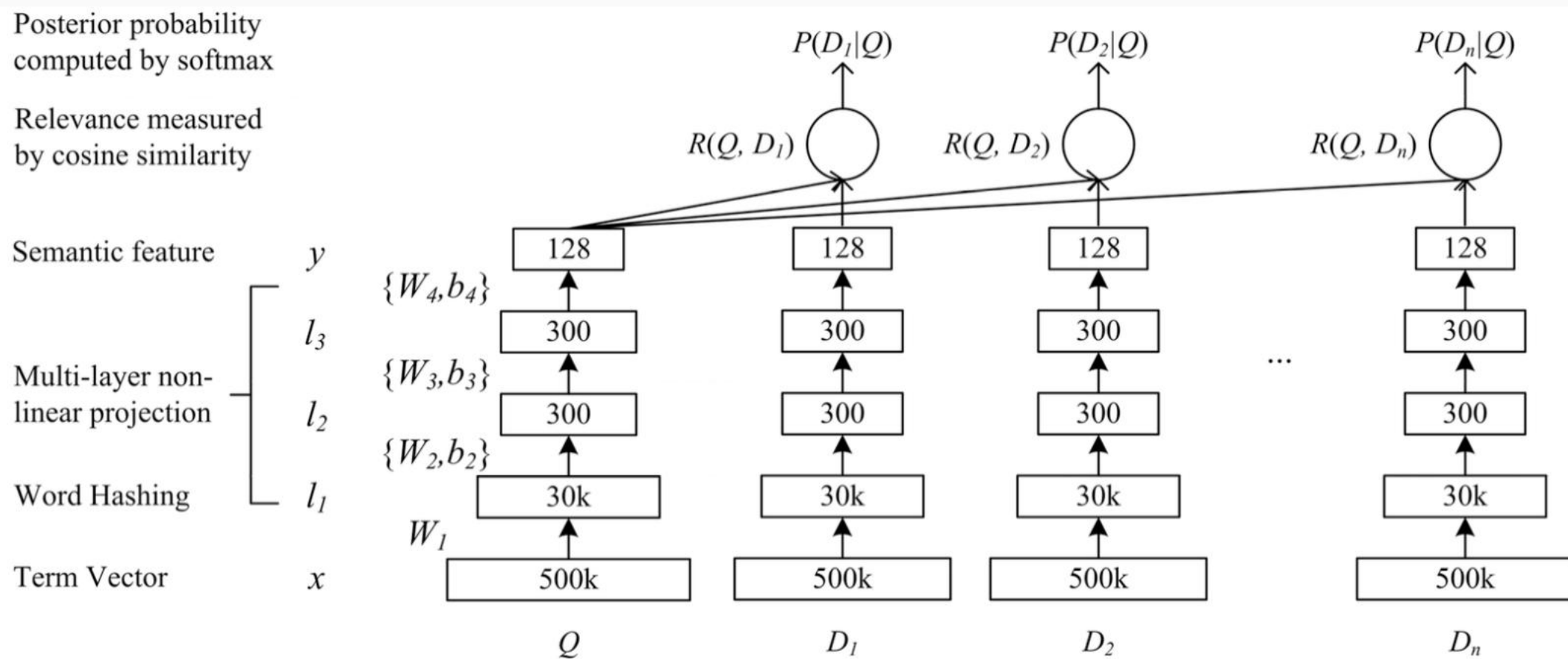
TensorFlow 是一个端到端开源机器学习平台



用TF实现DSSM都需要做什么

- 设计和定义模型结构
- 构造训练数据
- 使用GPU进行训练
- 导出训练好的模型
- 导出embedding向量

设计和定义网络结构



构造训练数据

- 数据下载
- 格式介绍
- 处理成什么样

一些必要的辅助函数（代码）





如何使用GPU（代码）

导出模型和向量（代码）



“关于模型的讨论。”

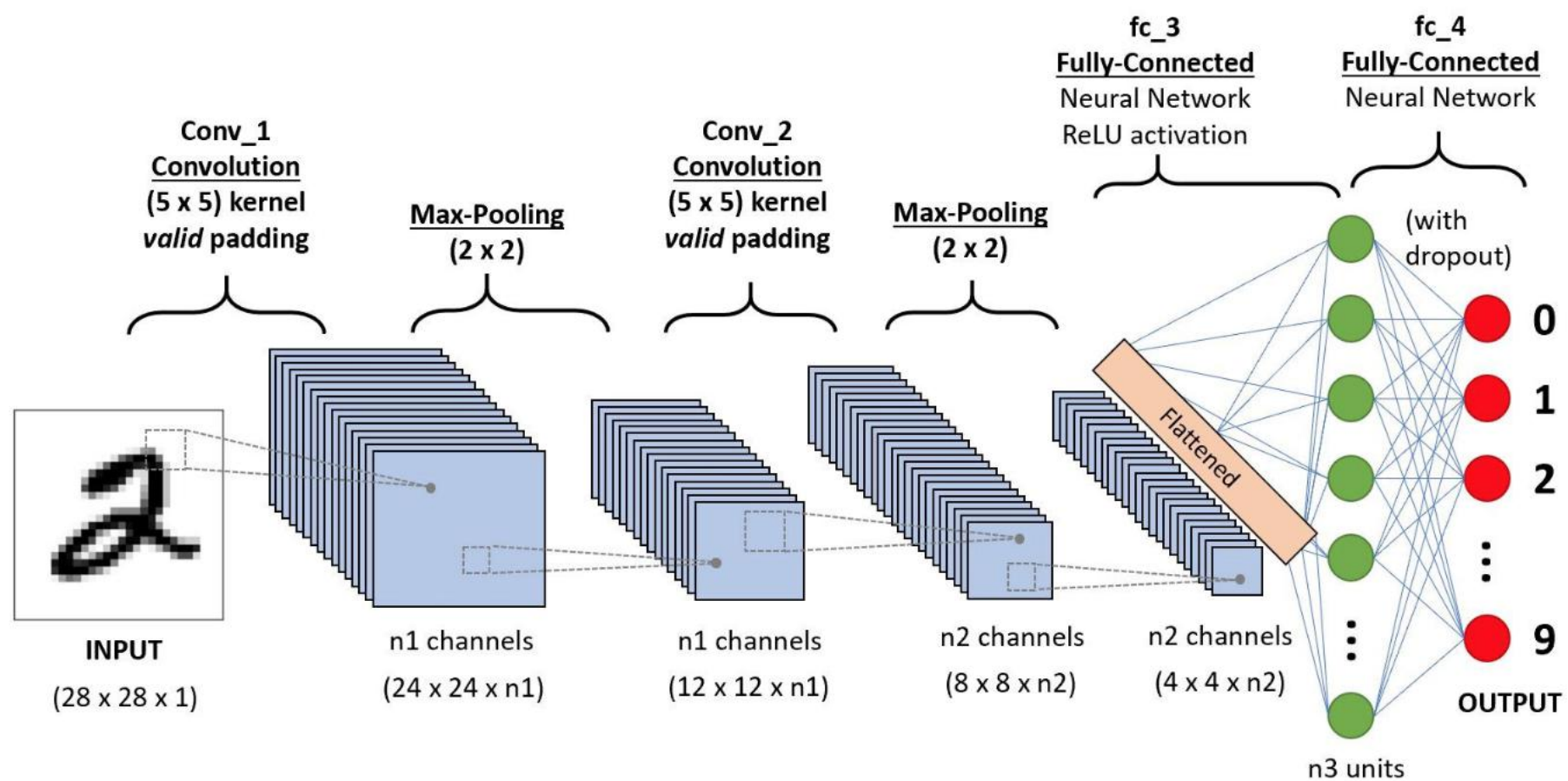
基于DNN的模型有什么缺点

- 同等看待用户的行为，而没有考虑到上下文
- 无法充分反映用户的兴趣

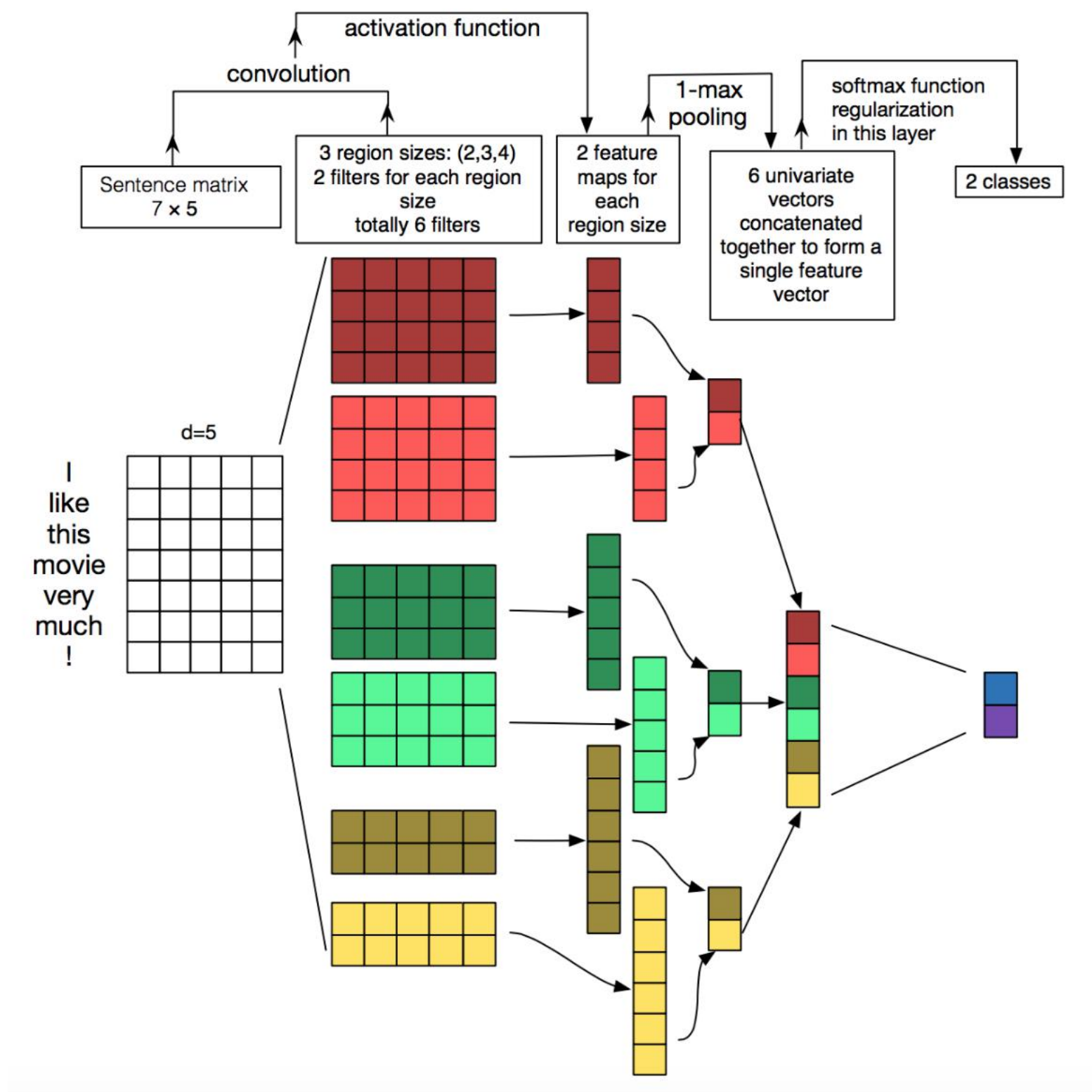
什么叫卷积

$$(f * g)(n) = \sum_{\tau=-\infty}^{\infty} f(\tau)g(n - \tau)$$

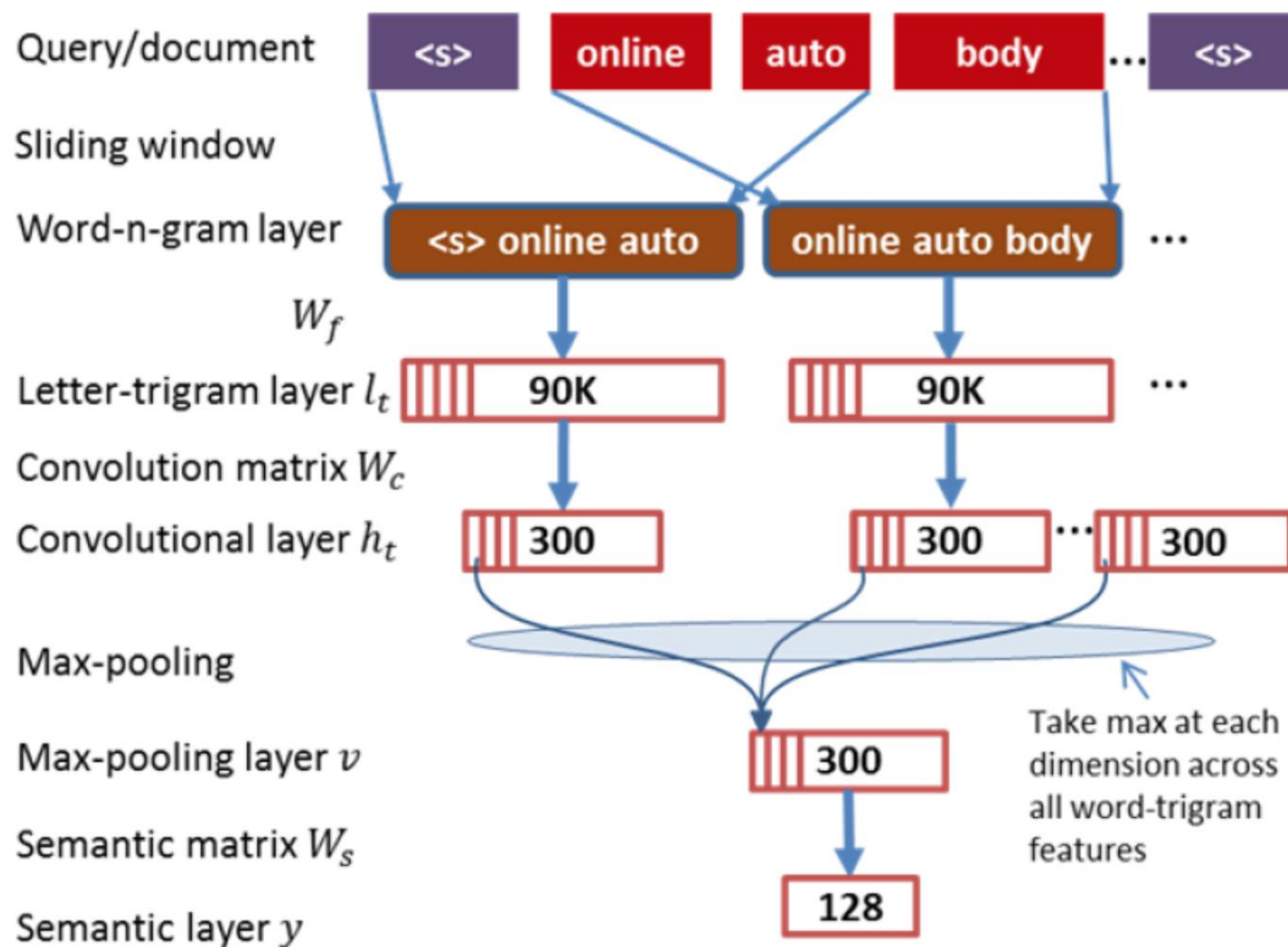
CNN用于图像



CNN用于NLP



CNN改进后的结构



“在**TF**里修改模型结构。”

下一步做什么

- 模型serving
- 确保线上线下数据分布的一致性
- badcase分析
- 调整方案

思考

- 为何深度学习有效
- 如何对行为的时序建模
- 不用深度学习怎么做

总结&回顾

- 承接上回：高维空间向量搜索算法（PQ/Annoy/HNSW）
- SGNS框架的局限性
- 用TF实现一个DSSM模型
- 如何利用CNN结构进行改进
- 理论提升