

核57, 2020.12.12.

回顾:

线性回归 $\xrightarrow{\text{Sigmoid } (\frac{1}{1+e^{-x}})}$ LR 逻辑回归; 二分类.

多分类, 二分类
回归问题

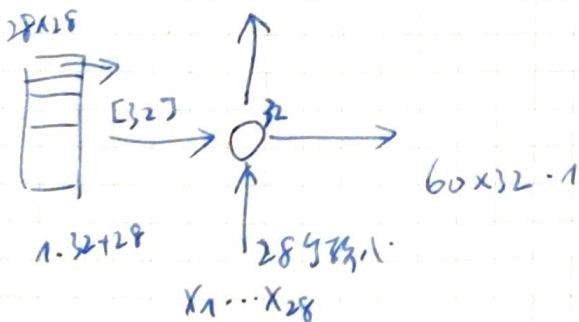
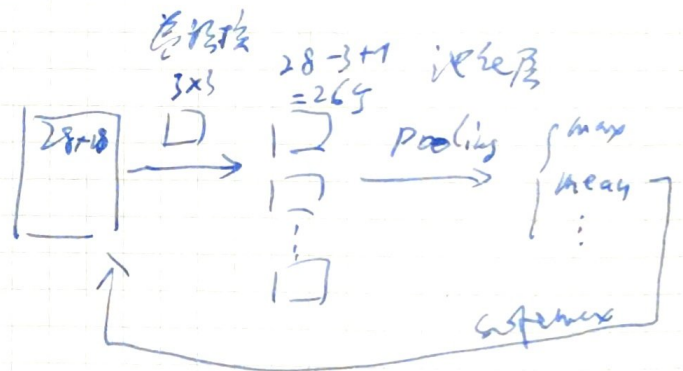
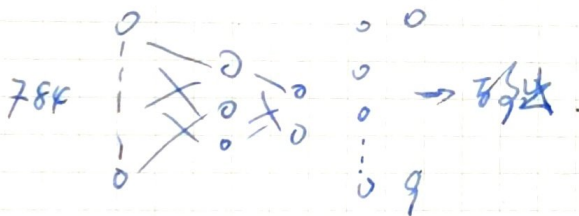
ID3/2045 树 \rightarrow CART $\xrightarrow{\text{基尼系数}} \text{GINI}$, 找最优分裂点

$\left\{ \begin{array}{l} \text{bagging} \\ \text{boosting} \end{array} \right\} \Rightarrow \text{RT} \rightarrow \text{GBDT} \Rightarrow \text{XGB} \rightarrow \text{Lgb}$

时序 \rightarrow 模型
 $\left\{ \begin{array}{l} \text{Arims} \\ \text{LSTM (RNN)} \end{array} \right.$
 BP $\xrightarrow{\text{反向传播}}$ LSTM (RNN) \rightarrow Embedding \uparrow
 梯度下降法, 迭代问题

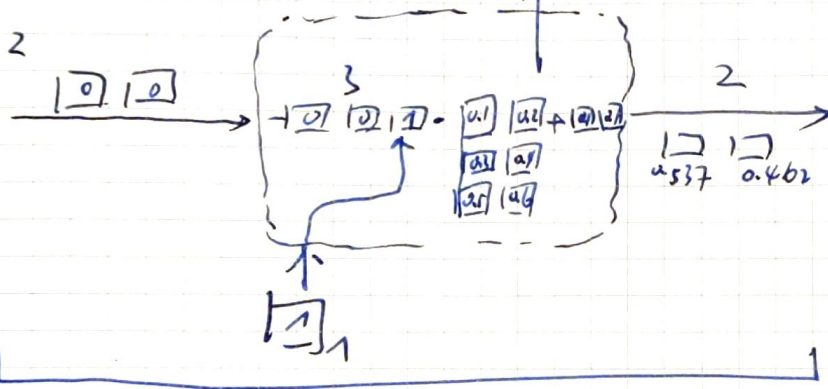
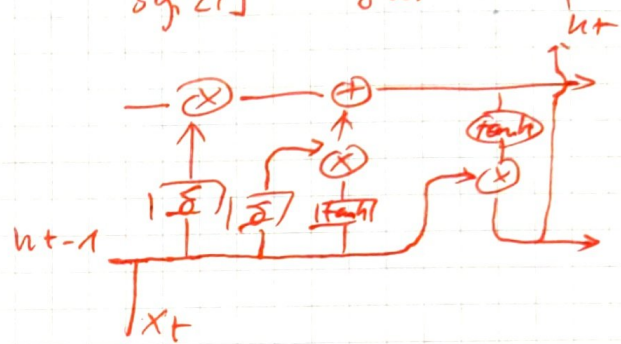
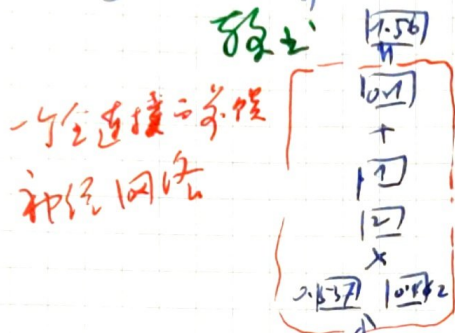
推荐系统: User \times Item $\xrightarrow{\text{召回}}$ GBDT + LR

• 图像, MNIST, $28 \times 28 = 784$



LSTM: Long short Term 长短期记忆.

- 应该将哪些信息状态 遗忘门
- 应该加入哪些新的状态 输入门
- 根据当前状态和输入, 输出是什么? 输出门



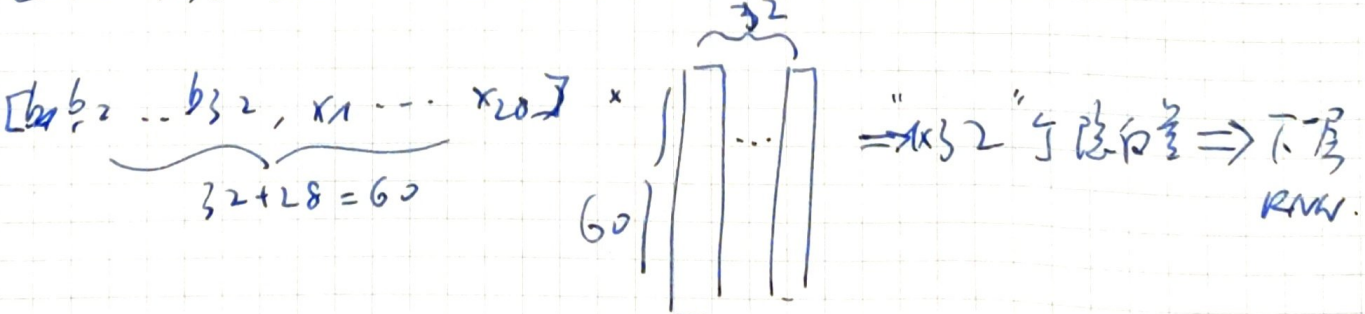
RNN 接收信息: $t_0 \rightarrow t_1 \rightarrow t_2 \dots t_n$

I am a student.

t_0 t_1 t_2

~~conv~~ MNIST: 手写识别: 28×28 .

LSTM 层 32 个节点. \Rightarrow 32 个隐变量.



Count Vectorizer: 文本特征提取方法.

- 只考虑每个词汇在语料库文本中出现的频率.

每篇文章 \rightarrow [— — —] - n words.

语料库 \rightarrow [— — —]

TF-IDF 原理: 一个词在当前文章中出现的频率, 在所有文档中出现的倒数.

表示词代表这篇文章.

词频: $TF_x = \frac{\text{词} \times \text{出现的次数}}{\text{该文档所有词的总数}} \cdot \text{即} \frac{\text{词数}}{\text{文章总词数}}$

问题: ① 在文本中出现的词数越多.

词频: 词频越大, 表示主题越强 \uparrow
词频越小, 表示主题越弱 \downarrow

② 通用词语词频大, 但不表示主题, 反而, 某些词词频的代表文章

所有的统计文章中, 一个词只在很少几篇文章中出现, 则其主题相关.

Term Frequency - Inverse Document Frequency

语料库文章总数

$$IDF = \log \left(\frac{\text{语料库文章总数}}{\text{包含词w的文章总数} + 1} \right)$$

逆文件频率.

← IDF

$$TF-IDF = TF \cdot IDF$$

在文章中重要性

权重, 则词频高, 权重小.

逆文件频率.

截断, 截取...

Tokenizer & pad-sequence 原理

词频, 统计词频.

填充后, 词频为在序列中的位置 (keras 只接受数字输入)

from sklearn.feature_extraction.text import CountVectorizer,

TFidfVectorizer.

① 先 Count 统计 CountVectorizer(). fit_transform(texts)

TFidfTransformer

② vector.get_feature_names() / vector.vocabulary_

TF-IDF

③ Count to array

[— — —] ← 文章1
[— — —] ← 文章2
[— — —] ← 文章n.

```
transformer = TfidfTransformer()
tfidf_matrix = transformer.fit_transform(count)
```

Tokenizer: from keras.preprocessing.text import Tokenizer

```
tok = Tokenizer().fit_transform(text)
```

tok.word-index / tok.word-counts

tok.texts-to-matrix(text) → $\begin{bmatrix} 0 & 1 & 1 & 0 & 1 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$ numpy 矩阵

tok.texts-to-sequences(text) → $\begin{bmatrix} [3, 4, 1, 2, 2], [\dots], [\dots] \end{bmatrix}$ 句子的词索引 (list)

Pad-sequences: from keras.preprocessing.sequence import pad_sequences

```
pad_sequences(text, maxlen=100)
```

系列!! 流程: 先用 IOF / Tokenizer / pad-sequences. Lgb ← numpy sparse (translating)

定义信息: 让每个句子的长度一致, 然后再用给定的词频子模型进行 Count-Vector.

通过问题及相互包含的关联特性 ⇒ 是否定义.

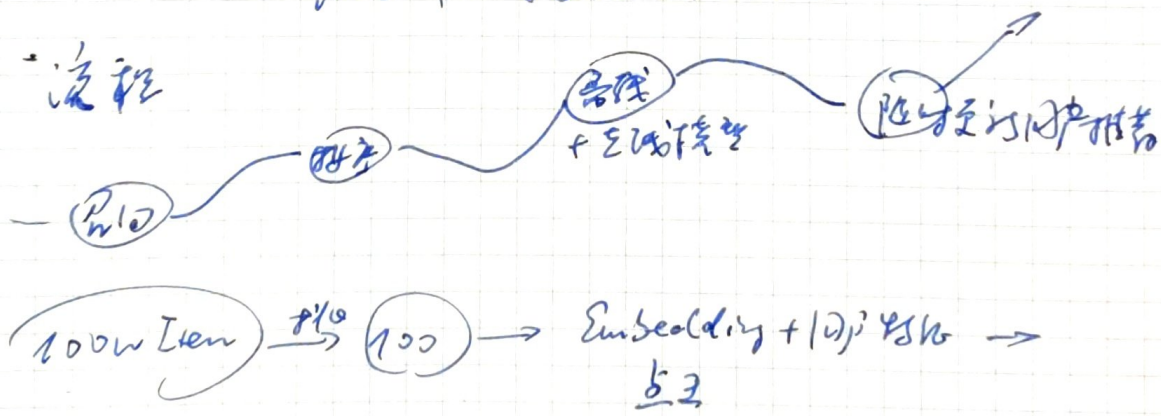
推荐系统, 文本. 词频包含 ⇒ 是否感兴趣.

推荐系统

- 用户: 解决“信息过载”, 高效获取感兴趣的信息.
- 企业: 吸引, 留存用户, 增加粘性, 提高用户转化率.
 - 视频网站: 观看时长 (CVR)
 - 电商: 购买转化率
 - 新闻: 点击率 CTR
- base: $\frac{User}{人} / \frac{Item}{物} / \frac{Context}{场景}$
 - 时间, 地点, 网络环境
 - 社交网络/人际关系/关系网络

- 热门推荐
- 个性化推荐
- 广告曝光推荐
- 搜索推荐

- 因子: 相关性 / 新颖性 / 多样性 / 多样性
- 目的:
 - 用户满意度, 以及企业业务增长.
 - 损失函数: $loss / 召回率 / 准确率$
 - top N: 召回和排序
 - 覆盖率: 挖掘“长尾”, 多样性



- 特点:
 - 数据量之大
 - 用户个性化
 - 实时更新
 - 冷启动问题

显式反馈: 明确告知

隐式反馈: 不告知

召回和排序

召回和排序

召回和排序

← UserCF / ItemCF 协同过滤