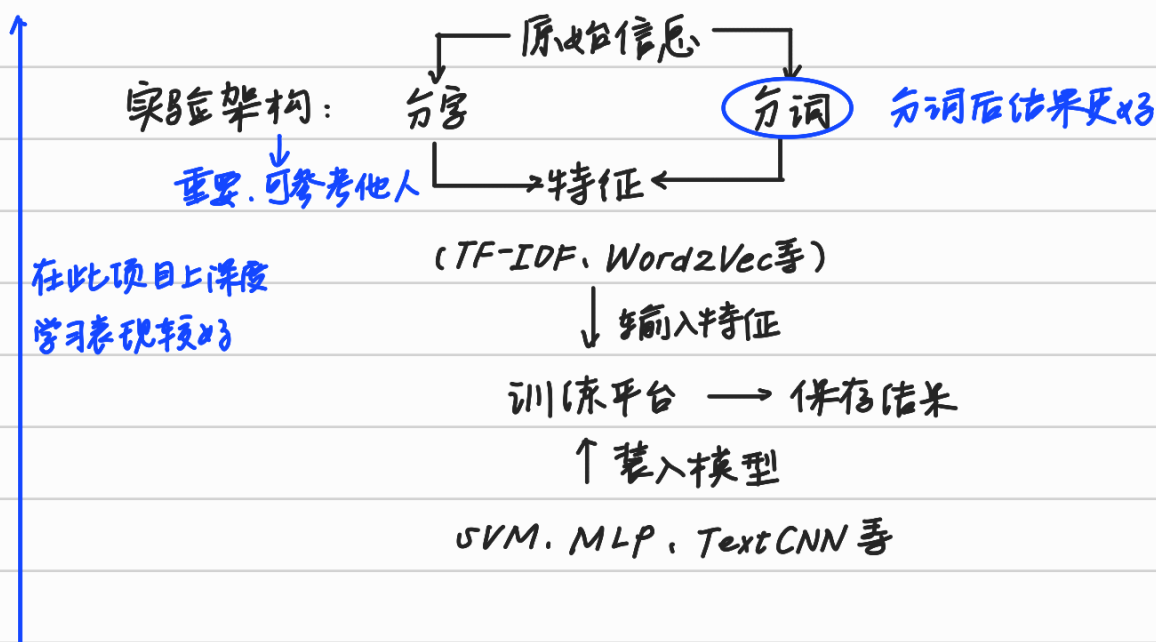


## 1. 文本数字化



## 2. 机器学习分类器



## 3. 深度学习实现文本分类

→ 在小数据集上使用预训练效果不错

(1) FastText: 模型简单, 在简单文本分类任务中效果不错, 直接 Softmax

(2) Text CNN: 用卷积提取局部特征, 用池化对特征进行选择

→ 局部信息, 并行度好

(3) Text RNN: 将上一个词的隐藏层也作为输入传给下一层, “整个文本形成链”

→ 全局信息, 并行度差

(4) Attention 机制

→ 既有全局信息, 并行度也好。

## 4. 模型融合: 关键在多个“好而不同”的模型

## 4) 模型的多样性

结论 ① 方字和方词的差异性最大 ② 传统ML和DL差异性较大

### ③ 数据增强也有差异性

#### 4.1.3 模型多样性度量

模型多样性度量是指度量模型融合中单模型的多样性，用来刻画单模型的多样化程度，它在模型融合和集成学习中是一个比较重要的问题。比较典型的做法是，考虑单模型的两两相似性或不相似性。总体上，模型多样性度量主要分为两大类：成对的多样性度量和非成对的多样性度量。

首先，假设二分类任务有 $m$ 个样本，分类器 $h_j$ 和分类器 $h_i$ 对样本的预测结果组合情况如表4-1所示。

表4-1 分类器预测结果组合情况

	$h_i = 1$	$h_i = 0$
$h_j = 1$	$a$	$c$
$h_j = 0$	$b$	$d$

成对多样性度量主要有如下参数。

不一致度量：

$$\text{dis}_{i,j} = \frac{b+c}{m} \quad (4.12)$$

相关系数：

$$\rho_{ij} = \frac{ad-bc}{\sqrt{(a+b)(a+c)(c+d)(b+d)}} \quad (4.13)$$

5. BERT: 效果很好, 只是训练耗费较大