

Chapter 10: Recurrent Neural Network

- 应用
 - 机器翻译
 - 序列标注
 - 图像描述
 - 推荐系统
 - 智能聊天机器人
 - 自动作词作曲等

RNN VS CNN

- CNN
 - 传统文本处理任务的方法中一般讲TF-IDF向量作为特征输入，丢失了输入的文本序列中单词间的顺序关系
 - 接受定长的向量作为输入，通过滑动窗口加活化，可以捕捉到局部信息，但是单词间的长距离依赖关系很难学习到
- RNN
 - 能够很好地处理文本数据变长且有顺序的输入序列，将输入中的有用信息编码到神经网络状态变量中去，从而拥有了一定的记忆能力
 - 在h, T后面直接接一个Softmax层，输出文本所属类别的预测概率y，r可以选择Tanh函数或者Relu函数，通过最小化损失误差(即输出的y与真实类别之间的距离)，就可以实现文本分类
- 文本分类问题
 - $$net_t = Ux_t + Wh_{t-1}$$
$$h_t = f(net_t)$$
$$y = g(Vh_T)$$
$$h_t, h_{t-1} \text{ 的计算公式}$$
 - Relu为激活函数，U为输入层到隐含层的权重矩阵，W为隐含层从上一时刻到下一时刻状态转移的权重矩阵

激活函数

- Relu
 - 将net_t的表达式展开，会发现式中包含t个W连乘，如果W不是单位矩阵，最终网络在T时刻的输出结果将会趋于0或者无穷，引发严重的数值问题
 - 对于梯度来说，只要W不是单位矩阵，梯度还会出现消失或者爆炸的现象
 - 但是在CNN中没有该现象，因为CNN中每一层的权重矩阵W都是不同的，并且在初始化时它们是独立分布的，因此可以相互抵消，在多层之后不会出现严重的数值问题
 - 因此只有当W的取值在单位矩阵附近时才能取得比较好的结果，需要将W初始化为单位矩阵
- Sigmoid与Tanh
 - 都是饱和的，在输入达到一定值得情况下，输出就不会发生明显变化，可以达到门控的效果
 - 使用非饱和的Relu函数则难以达到门控的效果

训练

- RNN的训练可以采用BPTT(Back Propagation Through Time)，基于时间的反向传播算法实现
- 梯度爆炸
 - 预测误差沿着神经网络的一层反向传播，当雅克比矩阵的最大特征值大于1时，随着离输出越来越远，每层的梯度大小会呈指数增长，导致梯度爆炸
- 梯度消失
 - 当雅克比矩阵的最大特征值小于1时，随着离输出越来越远，每层的梯度大小会呈指数减小，产生梯度消失
- 需要对本模型本身进行改进
 - 上短时记忆模型LSTM
 - 加入门控机制，很大程度上弥补了梯度消失所带来的损失

LSTM: Long Short Term Memory

- 输入门*i*
 - 控制当前计算的新状态以多大程度更新到记忆单元中
- 遗忘门*f*
 - 控制前一步记忆单元中的信息有多大程度被遗忘掉
- 输出门*o*
 - 控制当前的输出有多大程度上取决于当前记忆单元
- 内部记忆单元*c*
 - 仍然基于*x*, *h*和*h*_(t-1)来计算*h*_t
- Tanh函数
 - $$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$
$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$
$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$
$$c'_t = \tanh(W_c x_t + U_c h_{t-1})$$
$$c_t = f_t \odot c_{t-1} + i_t \odot c'_t$$
$$h_t = o_t \odot \tanh(c_t)$$
- Sigmoid函数

Seq2Seq: Sequence to Sequence

- 核心思想
 - 选取一种度量标准后，每次都在当前状态下选择最佳的一个结果，直到结束
 - 计算代价地，获得的是一个局部最优解，对于实际问题，往往不能取得最好的结果
- 贪心法
 - 启发式方法，会保存b(beam size)个当前的较佳选择，然后解码时每一步根据保存的选择进行下一步扩展和排序，接着选择前b个进行保存，循环迭代，直到结束时选择最佳的一个作为解码结果
- 束搜索
 - b取1，则退化为贪心法，随着b的增大，其搜索的空间增大，最终效果会有所提升，但计算量也增大，实际中选择b=8-12
- 解码方法
 - 堆叠RNN
 - 与解码器之间建立残差连接
 - 增加Dropout机制
 - 采用记忆网络，从外界获取知识
- 问题
 - 实际使用中Seq2Seq模型会随着输入序列的增长，模型的性能发生了显著下降，因为解码时输入序列的全部信息压缩到了一个向量表示中，序列越长，句子前面的词的信息丢失就越严重
 - 模型的输出序列中，常常会丢失部分输入序列的信息，因为在解码时，当前词语对应的源语言词的上下文信息和位置信息在解码过程中丢失了

改进

注意力机制Attention Mechanism

- 仍然使用普通的RNN对输入序列进行编码，得到隐状态，但是在解码时，每一个输出词都依赖于前一个隐状态以及输入序列每一个对应的隐状态
- $$s_i = f(s_{i-1}, y_{i-1}, c_i)$$
$$p(y_i | y_1, y_2, \dots, y_{i-1}) = g(y_{i-1}, s_i, c_i)$$
$$c_i = \sum_{j=1}^T \alpha_{ij} h_j$$

其中语境向量*c_i*是输入序列全部隐状态*h₁, h₂, ..., h_T*的一个加权和

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$$
$$e_{ij} = a(s_{i-1}, h_j)$$
- 注意力权重参数*a_{ij}*不是固定权重，而是另一个神经网络计算得到
- 使用双向RNN能够保证*x_t*时刻的前后隐状态信息都被考虑到

2018.12.23 Sunday. By KuKuXia@github.com

ML/DL/RL深度学习交流群: 622545311, 加群时请备注: 姓名-方向-公司or高校