

### 3. 负样本筛选

当前要做的人工筛选分为两个部分：**判断语法点正确与否、检查负样本构建是否合理。**

[illegible]

2: 该语法点判断正确。

(5) op: 样本相应的错误位置、类型和修正值;

op数组含义: [(错误起始index, 错误终止index, 待修改字符), 标签, 错误类型, (修改起始index, 修改终止index, 修改值)]

错误类型说明: 'W': Word Ordering Errors, 错序;

'M': Missing Words, 遗漏;

'S': Word Selection, 误用。

下图给出相应例子:

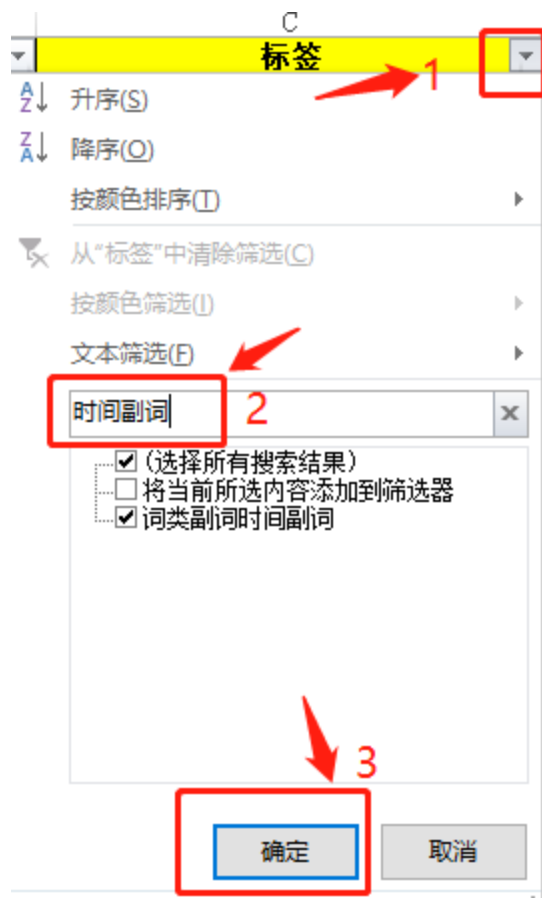
source: 种植物正这个季节长得很快, 经过短短一个星期, 它就长满了这面墙, 叶子很厚, 绿绿的。
target: 这种植物在这个季节长得很快, 经过短短一个星期, 它就长满了这面墙, 叶子很厚, 绿绿的。
[(0, 0, '种'), 词类代词指示代词, 'M', (null, null, "这")]
[(3, 3, '正'), 词类副词时间副词, 'S', (null, null, "在")]
source: 刚才听广播说明天能会下大雨, 足球比赛怕恐要推迟了。
target: 刚才听广播说明天可能会下大雨, 足球比赛恐怕要推迟了。
[(8, 8, '能'), 词类动词能愿动词, 'M', (null, null, '可')]
[(18, 18, '怕'), 词类副词情态副词, 'W', (19, 19, '恐')]
source: 刚才恐怕听广播说明天可能会下大雨, 足球比赛要推迟了。
target: 刚才听广播说明天可能会下大雨, 足球比赛恐怕要推迟了。
[(2, 3, '恐怕'), 词类副词情态副词, 'W', (21, 21, 'None')]

## 2. 语法点检查

根据op和label定位target中的语法点, 并判断该语法点是否正确, 并在“label得分”一栏给定相应分数。

如, 下表所示的样例中, 将“在”字对应的语法点错误匹配成“词类副词时间副词”, 与实际不符, 则其对应“label得分”应为0。





### 3. 负样本筛选

观察source，筛掉其中构造不合理的、生活中不可能犯的、或是正确没有语病的样本，并在“source得分”一栏给定相应分数。

如下表中的样例，给出了对部分负样本的打分（打分的主观性比较强，但尽可能确保0分和3分的“稳定”）

source	target	source得分	op
这样做，只能暂时解决问题，所要想完全解决这个难题，再需要找更好的办法。	这样做，只能暂时解决问题，所以要想完全解决这个难题，还需要找更好的办法。	3	[[ (14, 14, '要'), '句子的类型复句因果复句', 'M', ('null', 'null', '以') ], [ (25, 25, '再'), '词类副词关联副词', 'S', ('null', 'null', '还') ]]
这样做，只能暂时解决问题，以所要想完全解决这个难题，还需要找更好的办法。	这样做，只能暂时解决问题，所以要想完全解决这个难题，还需要找更好的办法。	3	[[ (13, 13, '以'), '句子的类型复句因果复句', 'W', (14, 14, '所') ], [ (26, 26, '以'), '词类副词关联副词', 'S', ('null', 'null', '还') ]]

大这个难题，也而去找更好的办法。	这个难题，也而去找更好的办法。		也)，词类副词关联副词，S，('null', '还'))]
还这样做，只能暂时解决问题，所以要想完全解决这个难题，需要找更好的办法。	这样做，只能暂时解决问题，所以要想完全解决这个难题，还需要找更好的办法。	1	[[ (15, 15, '要'), '句子的类型复句因果复句', 'M', ('null', 'null', '以')], [(0, 0, '还'), '词类副词关联副词', 'W', (26, 26, 'None')]]
这样做，只能暂时解决问题，也要想完全解决这个难题，就需要找更好的办法。	这样做，只能暂时解决问题，所以要想完全解决这个难题，还需要找更好的办法。	1	[[ (13, 13, '也'), '句子的类型复句因果复句', 'S', ('null', 'null', '所以')], [(25, 25, '就'), '词类副词关联副词', 'S', ('null', 'null', '还')]]
他年轻很，可是遇到问题很冷静，和相同年龄得人更成熟。	他很年轻，可是遇到问题很冷静，比相同年龄的人更成熟。	3	[[ (1, 1, 'None'), '词类副词程度副词', 'W', (3, 3, '很')], [(15, 15, '和'), '词类介词引出对象', 'S', ('null', 'null', '比')], [(20, 20, '得'), '词类助词结构助词', 'S', ('null', 'null', '的')]]
很他年轻更，可是遇到的问题很冷静，比相同年龄人成熟。	他很年轻，可是遇到问题很冷静，比相同年龄的人更成熟。	0	[[ (0, 0, '很'), '词类副词程度副词', 'W', (2, 2, 'None')], [(10, 10, '的'), '词类助词结构助词', 'W', (22, 22, 'None')], [(4, 4, '更'), '词类副词程度副词', 'W', (23, 23, 'None')]]
他年轻，可是遇到问题很冷静，相同年龄地人更成熟。	他很年轻，可是遇到问题很冷静，比相同年龄的人更成熟。	3	[[ (1, 1, '年'), '词类副词程度副词', 'M', ('null', 'null', '很')], [(14, 14, '相'), '词类介词引出对象', 'M', ('null', 'null', '比')], [(18, 18, '地'), '词类助词结构助词', 'S', ('null', 'null', '的')]]
他更真年轻，可是遇到问题很冷静，比相同年龄所人成熟。	他很年轻，可是遇到问题很冷静，比相同年龄的人更成熟。	0	[[ (1, 1, '更'), '词类副词程度副词', 'S', ('null', 'null', '很')], [(21, 21, '所'), '词类助词结构助词', 'S', ('null', 'null', '的')], [(1, 1, '更'), '词类副词程度副词', 'W', (23, 23, 'None')]]

同时，在人工进行评估时，我们**应尽量从标准的、语法规范的角度去评判**（如“相同年龄人”这种表达在生活中确实会出现，似乎也不构成一个错误的点。但从更规范的语法角度去理解，“年龄”和“人”之间缺少助词，不能直接相连，应改成“相同年龄的人”。）

PS：在进行数据筛选的同时，顺便检查一下其op是否正确（有很少一部分样本会有点bug=.=），有错误的地方人工修正一下~。

比如下表样例1，对比source和op后发现[(0, 1, '多久'), '词类代词疑问代词', 'S', ('null', 'null', '谁')]这个错误修正有误（究其原因，是因为在构造负样本时先将“谁”替换成“多久”，而后续在对“多”进行操作时错误定位到了“多久”而非“最多”上），对于这种构造出错的样本，直接删去这一行source、target、op所有数据。

又如样例2所示，它对“又脏又乱”中的2个“又”交换了位置，虽然[(4, 4, '又'), '固定格式', 'W', (6, 6, '又')]这个错误修正是对的，但它没有实质性表现和意义，故将其修改成样例3所示（删去对应多余操作）。

id	source	target	op
1	久在规定的时间内接到的球最多，谁就赢了比赛，所以这种游戏十分简单。	谁在规定的时间内接到的球最多，谁就赢了比赛，所以这种游戏十分简单。	[[ (0, 1, '多久'), '词类代词疑问代词', 'S', ('null', 'null', '谁') ], [ (0, 0, '久'), '词类代词疑问代词', 'M', ('null', 'null', '多') ]]
2	那个房间又脏又乱，星期六我去打扫了、整理一下。	那个房间又脏又乱，星期六我去打扫、整理了一下。	[[ (4, 4, '又'), '固定格式', 'W', (6, 6, '又') ], [ (16, 16, '了'), '句子的类型复句紧缩复句', 'W', (20, 20, 'None') ]]
3	那个房间又脏又乱，星期六我去打扫了、整理一下。	那个房间又脏又乱，星期六我去打扫、整理了一下。	[[ (16, 16, '了'), '句子的类型复句紧缩复句', 'W', (20, 20, 'None') ]]