

相似度阈值设置为什么值，相关性判断的结果最准确

要确定相似度阈值设置为多少时相关性判断结果最准确，需要先分析当前数据中不同相似度分数对应的实际相关性情况，再结合业务场景需求综合判断。以下是具体分析过程：

一、数据相似度分数分布与相关性标签分析

当前 JSON 数据中，标签分为 `similar`（相似）和 `no_similar`（不相似），对应的相似度分数如下：

标签	相似度分数范围	样本数	典型分数示例
similar	0.6725 - 0.9729	17 组	0.9274、 0.9424、0.9729 等
no_similar	0.6268 - 0.7579	2 组	0.6268、0.7579

关键观察：

- `similar` 组的最低分数为 0.6725（如 `similar_9` 中对比文本 1），而 `no_similar` 组的最高分数为 0.7579（如 `no_similar_2`）。
- 两组数据的分数区间存在重叠（0.6725 - 0.7579），这意味着直接通过分数区分相关性存在一定难度。
- `similar` 组中大部分分数集中在 0.8 以上（如 0.8396、0.8895 等），而 `no_similar` 组分数均低于 0.76，说明高分段（ ≥ 0.8 ）与相关性的匹配度更高。

二、阈值候选方案及优缺点分析

方案 1：阈值设为 0.75

- 依据：`no_similar` 组最高分为 0.7579，若设为 0.75，可将 `no_similar_2`（0.7579）归为不相似，但 `similar_9` 的 0.6725 会被误判为不相似。
- 优缺点：
 - 优点：能正确识别大部分 `no_similar` 样本。
 - 缺点：会遗漏部分 `similar` 样本（如 0.6725-0.75 之间的相似数据），导致召回率降低。

方案 2：阈值设为 0.8

- 依据：similar 组中 80% 以上的分数 ≥ 0.8 ，且 no_similar 组分数均 < 0.8 ($0.7579 < 0.8$)。
- 优缺点：
 - 优点：
 - 能过滤掉 no_similar 组的所有样本（避免误判为相似）。
 - 保留 similar 组中高分段的强相关数据，提升判断准确性。
 - 缺点：可能将 0.8 以下的部分相似数据（如 similar_9 的 0.6725）误判为不相似，但这类数据可能属于弱相关，业务上可接受。

方案 3：阈值设为 0.7

- 依据：覆盖 similar_9 的最低分 0.6725，但会将 no_similar_1 的 0.6268 误判为相似。
- 优缺点：
 - 优点：召回率高，几乎不遗漏 similar 样本。
 - 缺点：会引入 no_similar 样本，导致误判率上升，准确性下降。

三、业务场景下的最优阈值选择

4. 若业务需求侧重“精准识别强相关性”：

- 推荐阈值：0.8
- 理由：
 - 从数据看， ≥ 0.8 的分数与 similar 标签的匹配度极高（如 similar_1 - similar_5 分数均 ≥ 0.89 ），能确保大部分强相关数据被正确识别。
 - 过滤掉 0.8 以下的弱相关数据，减少误判（如 similar_9 的 0.6725 可能属于内容相关性较低的场景，业务上可视为不相关）。

5. 若业务需求侧重“不遗漏任何可能的相关数据”：

- 推荐阈值：0.7
- 理由：
 - 覆盖 similar 组的最低分 0.6725，确保所有标注为相似的数据被召回。
 - 但需接受 no_similar_1 (0.6268) 可能被误判为相似的风险，需后续人工审核校准。

四、综合结论：阈值设为 0.8 时准确性最高

支持依据：

1. 数据匹配度：`similar` 组中 ≥ 0.8 的样本占比超过 80%，且 `no_similar` 组无样本超过 0.8，阈值 0.8 能有效区分两组数据。

2. 业务逻辑合理性：

- 高分段 (≥ 0.8) 通常意味着文本内容、促销信息、品牌关键词等高度重合（如 `similar_14` 中多个对比文本与原始数据几乎一致，分数达 0.9298），相关性更强。
- 0.8 以下的分数可能对应弱相关场景（如仅部分关键词重叠），业务上可根据需求进一步筛选。

3. 误判风险最低：相比 0.7 和 0.75，阈值 0.8 能在召回率和准确率之间取得平衡，避免将 `no_similar` 样本误判为相似，同时保留大部分真实相似数据。

五、建议：结合业务场景动态调整

- 若需要严格控制误判率（如用于精准营销推送），优先使用 0.8 作为阈值。
- 若需要扩大数据覆盖范围（如用于舆情监控或数据收集），可将阈值设为 0.7，并搭配人工审核机制过滤误判样本。