

# 写

我期望在类似和格的化数据上行一个来返回精确匹配的文。然而，好的全文搜索不是完全相同的限定。相反，我可以大以包括可能的匹配，而根据相关性得分将更好的匹配推到果集的前部。

事实上，只能完全匹配的全文搜索可能会困的用。道不希望在搜索 quick brown fox 匹配一个包含 fast brown foxes 的文，搜索 Johnny Walker 同匹配 Johnnie Walker，搜索 Arnold Schwarzenegger 同匹配 Arnold Schwarzenegger？

如果存在完全符合用的文，他出在果集的前部，而弱的匹配可以被包含在列表的后面。如果没有精确匹配的文，至少我可以示有可能匹配用要求的文，它甚至可能是用最初想要的！

我已在 [token-normalization] 看自由音匹配，[stemming] 中的干，[synonyms] 中的同，但所有些方法假定写正，或者个写只有唯一的方法。

Fuzzy matching 允匹配写的，而音元器可以在索引用来行近似音匹配。

## 模糊性

模糊匹配待“模糊”相似的一个似乎是同一个。首先，我需要我的模糊性行定。

在1965年，Vladimir Levenshtein 出了 Levenshtein distance，用来度量从一个到一个需要多少次字符。他提出了三型的字符：

- 一个字符替一个字符：\_f\_ox → \_b\_ox
- 入一个新的字符：sic → sic\_k\_
- 除一个字符：b\_l\_ack → back

Frederick Damerau 后来在些操作基上做了一个展：

- 相个字符的位：\_st\_ar → \_ts\_ar

个例子，将 beiber 成 beaver 需要下面几个：

1. 把 b 替成 v：bie\_b\_er → bie\_v\_er
2. 把 i 替成 a：b\_i\_ever → b\_a\_ever
3. 把 e 和 a 行位：b\_ae\_ver → b\_ea\_ver

三个表示 Damerau-Levenshtein edit distance 距 3。

然，从 <code>beaver</code> 成 <code>bieber</code> 是一个很的程；他相距甚而不能一个的写。Damerau 80% 的写距 1。句，80% 的写可以原始字符串用 <em> 次 </em> 行修正。

Elasticsearch 指定了 fuzziness 参数支持最大距的配置，2。

当然，每次字符串的匹配取决于字符串的长度。每次只能生成 `mad`，所以一个只有 3 个字符长度的字符串允许的次数虽然太多了。`fuzziness` 参数可以被置为 `AUTO`，将导致以下的最大距离：

- 字符串只有 1 到 2 个字符 是 `0`
- 字符串有 3、4 或者 5 个字符 是 `1`
- 字符串大于 5 个字符 是 `2`

当然，可能会距离 `2` 然是太多了，返回的结果似乎并不相关。把最大 `fuzziness` 置为 `1`，可以得到更好的结果和更好的性能。

## 模糊

[{ref}/query-dsl-fuzzy-query.html](#) [`fuzzy`] 是 `term` 的模糊等。也很少直接使用它，但是理解它是如何工作的，可以帮助在更高的 `match` 中使用模糊性。

了解它是如何工作的，我首先索引一些文档：

```
POST /my_index/my_type/_bulk
{ "index": { "_id": 1 } }
{ "text": "Surprise me!" }
{ "index": { "_id": 2 } }
{ "text": "That was surprising." }
{ "index": { "_id": 3 } }
{ "text": "I wasn't surprised." }
```

在我可以 `surprise` 行一个 `fuzzy`：

```
GET /my_index/my_type/_search
{
  "query": {
    "fuzzy": {
      "text": "surprise"
    }
  }
}
```

`fuzzy` 是一个简单的，所以它不做任何分析。它通过某个以及指定的 `fuzziness` 到字典中所有的。 `fuzziness` 置为 `AUTO`。

在我的例子中，`surprise` 比 `surprise` 和 `surprised` 都在距离 `2` 以内，所以文档 `1` 和 `3` 匹配。通过以下，我可以少匹配度到匹配 `surprise`：

```
GET /my_index/my_type/_search
{
  "query": {
    "fuzzy": {
      "text": {
        "value": "surprise",
        "fuzziness": 1
      }
    }
  }
}
```

## 提高性能

`fuzzy` 的工作原理是 定原始 及 造一个 `<em>自 机</em>`。像表示所有原始字符串指定 距 的字符串的一个大 表。

然后模糊 使用 个自 机依次高效遍 典中的所有 以 定是否匹配。一旦收集了 典中存在的所有匹配 , 就可以 算匹配文 列表。

当然, 根据存 在索引中的数据 型, 一个 距 2 的模糊 能 匹配一个非常大数量的 同行效率会非常糟 。下面 个参数可以用来限制 性能的影响 :

### prefix\_length

不能被 '模糊化' 的初始字符数。大部分的 写 生在 的 尾, 而不是 的 始。例如通 将 'prefix\_length 置 3, 可能 著降低匹配的 数量。

### max\_expansions

如果一个模糊 展了三个或四个模糊 , 些新的模糊 也 是有意 的。如 果它 生 1000 个模糊 , 那 就基本没有意 了。 置 max\_expansions 用来限制将 生的模糊 的数量。模糊 将收集匹配 直到 到 max\_expansions 的限制。

## 模糊匹配

match 支持 箱即用的模糊匹配 :

```
GET /my_index/my_type/_search
{
  "query": {
    "match": {
      "text": {
        "query": "SURPRIZE ME!",
        "fuzziness": "AUTO",
        "operator": "and"
      }
    }
  }
}
```

字符串首先进行分析，会生成 `[surprise, me]`，并且会根据指定的 `fuzziness` 行模糊化。

同样，`multi_match` 也支持 `fuzziness`，但只有当行类型是 `best_fields` 或者 `most_fields`：

```
GET /my_index/my_type/_search
{
  "query": {
    "multi_match": {
      "fields": [ "text", "title" ],
      "query": "SURPRIZE ME!",
      "fuzziness": "AUTO"
    }
  }
}
```

`match` 和 `multi_match` 都支持 `prefix_length` 和 `max_expansions` 参数。

#### TIP

模糊性 (Fuzziness) 只能在 `match` and `multi_match` 中使用。不能使用在短匹配、常用或 `cross_fields` 匹配。

## 模糊性 分

用 喜 模糊。他会魔法般的到正写合。很憾，效果平平。

假我有1000个文包含 `Schwarzenegger`，只是一个文的出写 `Schwarzeneger`。根据 [term frequency/inverse document frequency](#) 理，个写文比写正的相度更高，因写出在更少的文中！

句，如果我待模糊匹配似其他匹配方法，我将偏的写超了正的写，会用狂。

#### TIP

模糊匹配不用于参与分—只能在有写大匹配的。

情况下，`match` 定所有的模糊匹配的恒定分 1。可以在足在果列表的末尾添加潜在的匹配，并且没有干非模糊的相性分。

#### TIP

在模糊最初出很少能独使用。他更好的作一个 [bigger](#) 景的部分功能特性，如 [search-as-you-type](#) [{ref}/search-suggesters-completion.html](#)[\[完成 建\]](#) 或 [did-you-mean](#) [{ref}/search-suggesters-phrase.html](#)[\[短 建\]](#)。

## 音匹配

最后，在任何其他匹配方法都无效后，我可以求助于搜索音相似的，即使他的写不同。

有一些用于将成音的算法。[Soundex](#) 算法是些算法的鼻祖，而且大多数音算法是

Soundex 的改 或者 版本, 例如 [Metaphone](#) 和 [Double Metaphone](#) ( 展了除英 以外的其他言的 音匹配), [Caverphone](#) 算法匹配了新西 的名称, [Beider-Morse](#) 算法吸收了 Soundex 算法了更好的匹配 和依地 名称, [Kölner Phonetik](#) 了更好的 理 。

得一提的是, 音算法是相当 的, 他 初衷 的 言通常是英 或 。 限制了他 的 用性。不 , 了某些明 的目 , 并与其他技 相 合, 音匹配能 作 一个有用的工具。

首先, 需要从 <https://www.elastic.co/guide/en/elasticsearch/plugins/current/analysis-phonetic.html> 取 音分析 件并在集群的 个 点安装, 然后重 个 点。

然后, 可以 建一个使用 音 元 器的自定 分析器, 并 下面的方法 :

```
PUT /my_index
{
  "settings": {
    "analysis": {
      "filter": {
        "dbl_metaphone": { ①
          "type": "phonetic",
          "encoder": "double_metaphone"
        }
      },
      "analyzer": {
        "dbl_metaphone": {
          "tokenizer": "standard",
          "filter": "dbl_metaphone" ②
        }
      }
    }
  }
}
```

① 首先, 配置一个自定 **phonetic** 元 器并使用 **double\_metaphone** 器。

② 然后在自定 分析器中使用自定 元 器。

在我 可以通 **analyze** API 来 行 :

```
GET /my_index/_analyze?analyzer=dbl_metaphone
Smith Smythe
```

个 **Smith** 和 **Smythe** 在同一位置 生 个 元 : **SM0** 和 **XMT** 。通 分析器播放 **John** , **Jon** 和 **Johnnie** 将 生 个 元 **JN** 和 **AN** , 而 **Jonathon** 生 元 **JN0N** 和 **ANTN** 。

音分析器可以像任何其他分析器一 使用。首先映射一个字段来使用它, 然后索引一些数据 :

```

PUT /my_index/_mapping/my_type
{
  "properties": {
    "name": {
      "type": "string",
      "fields": {
        "phonetic": { ①
          "type": "string",
          "analyzer": "dbl_metaphone"
        }
      }
    }
  }
}

PUT /my_index/my_type/1
{
  "name": "John Smith"
}

PUT /my_index/my_type/2
{
  "name": "Jonnie Smythe"
}

```

① `name.phonetic` 字段使用自定义 `dbl_metaphone` 分析器。

可以使用 `match` 来进行搜索：

```

GET /my_index/my_type/_search
{
  "query": {
    "match": {
      "name.phonetic": {
        "query": "Jahnnie Smeeth",
        "operator": "and"
      }
    }
  }
}

```

这个请求返回全部文档，演示了如何进行音匹配。使用音算法计算分是没有意义的。音匹配的目的不是为了提高精度，而是要提高召回率——以展示足够的文档来捕捉可能匹配的文档。

通常更有意义的使用音算法是在索引到结果后，由一台计算机进行消歧和处理，而不是由人直接使用。