



Databáze ◊ poznámky k přednášce

12. Pokročilé podmínky

verze z 16. prosince 2024

V příkladech budeme uvažovat indexovou proměnnou `movie` danou: „Film d má název D .“ s hodnotou I_1 :

```
{<1, "The Fellowship of the Ring">,<br> <2, "The Two Towers">,<br> <3, "The Return of the King">}
```

1 Konjunkce termů

Vezměme podmínky $\theta_1, \dots, \theta_n$. Pak

$$\theta_1 \wedge \dots \wedge \theta_n$$

je podmínka, kterou dokument splňuje, jestliže splňuje všechny podmínky $\theta_1, \dots, \theta_n$. Například dokument "The Fellowship of the Ring" splňuje podmínku `ring` \wedge `fellowship`. Dále $\sigma_{\text{of} \wedge \text{ring}}(I_1)$ je rovno:

```
{<1, "The Fellowship of the Ring">}
```

a dáno: „Film d má název D , který obsahuje slovo `of` a slovo `ring`.“

Vezměme dokument Q , pak řetězec:

```
{<br>  "match": {<br>    "text": {<br>      "query": Q,<br>      "operator": "AND"<br>    }<br>  }<br>}
```

zapisuje podmínku $q_1 \wedge \dots \wedge q_n$, kde $T(Q) = (q_1, \dots, q_n)$.

Například:

```
GET /movie/_search<br>{
```

```

"query": {
  "match": {
    "text": {
      "query": "OF King",
      "operator": "AND"
    }
  }
}
}

{
  "hits": {
    "hits": [
      {
        "_id": "3",
        "_source": {
          "text": "The Return of the King"
        }
      }
    ]
  }
}
}

```

Kvalita splnění podmínky $\theta_1 \wedge \dots \wedge \theta_n$ je dána:

$$\text{score}(I, D, \theta_1 \wedge \dots \wedge \theta_n) = \sum_{i=1}^n \text{score}(I, D, \theta_i)$$

Například dokument v předchozím příkladu splnil podmínku s kvalitou přibližně $1.3 = 0.4 + 0.9$.

2 Obecné složené podmínky

Vezměme podmínky $\theta_1, \dots, \theta_n$. Pak podmínku $\theta_1 \wedge \dots \wedge \theta_n$ zapíšeme řetězcem:

```

{
  "bool": {
    "must": [
       $\theta_1$ ,
       $\vdots$ ,
       $\theta_n$ 
    ]
  }
}

```

Například podmínku $(\text{fellowship} \vee \text{two}) \wedge (\text{ring} \vee \text{towers})$ zapíšeme řetězcem:

```
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "text": "fellowship two"
          }
        },
        {
          "match": {
            "text": "ring towers"
          }
        }
      ]
    }
  }
}
```

Podmínku $\theta_1 \vee \dots \vee \theta_n$ zapíšeme řetězcem:

```
{
  "bool": {
    "should": [
       $\theta_1$ ,
      :
       $\theta_n$ 
    ],
    "minimum_should_match": 1
  }
}
```

Například podmínku $(\text{ring} \vee \text{fellowship}) \wedge (\text{towers} \vee \text{two})$ zapíšeme řetězcem:

```
{
  "bool": {
    "should": [
      {
        "match": {
          "text": {
            "query": "ring fellowship",
            "operator": "AND"
          }
        }
      }
    ]
  }
}
```

```

    }
  },
  {
    "match": {
      "text": {
        "query": "towers two",
        "operator": "AND"
      }
    }
  }
],
"minimum_should_match" : 1
}
}

```

Pro podmínku θ je $\neg\theta$ podmínka, kterou dokument splňuje, jestliže nespĺňuje podmínku θ . Například dokument "The Return of the King" splňuje podmínku \neg ring.

Podmínku $\neg\theta$ zapíšeme řetězcem:

```

{
  "bool": {
    "must_not":  $\theta$ 
  }
}

```

Například podmínku \neg ring zapíšeme:

```

{
  "bool": {
    "must_not": {
      "match": {
        "text": "ring"
      }
    }
  }
}
}

```

Kvalita splnění podmínky $\neg\theta$ je $score(I, D, \neg\theta) = 0$. Například máme:

$$\begin{aligned}
 &score(I_1, \text{"The Return of the King"}, of \wedge \neg ring) \\
 &= score(I_1, \text{"The Return of the King"}, of) \\
 &\approx 0.44
 \end{aligned}$$

Řetězec:

```
{
  "bool": {
    "must_not": [
       $\theta_1$ ,
      :
       $\theta_n$ ,
    ]
  }
}
```

je zkratkou za $\neg(\theta_1 \vee \dots \vee \theta_n)$.

Například řetězec:

```
{
  "bool": {
    "must_not": [
      {
        "match": {
          "text": "ring"
        }
      },
      {
        "match": {
          "text": "king"
        }
      }
    ]
  }
}
```

zapisuje podmínku $\neg(\text{ring} \vee \text{king})$.

Každý dokument splňuje podmínku 1, kterou zapíšeme řetězcem:

```
{
  "match_all": {}
}
```

Kvalita splnění podmínky 1 je $score(I, D, 1) = 1$.

3 Zobecněná disjunkce

Vezměme podmínky $\theta_1, \dots, \theta_m$ a číslo $0 \leq n \leq m$, pak

$$\bigvee_n(\theta_1, \dots, \theta_m) = \alpha_1 \vee \dots \vee \alpha_k,$$

kde α_i jsou všechny možné konjunkce tvaru $\theta_{j_1} \wedge \dots \wedge \theta_{j_n}$ takové, že $j_1 < \dots < j_n$. Platí, že $k = \binom{m}{n}$. Například:

$$\alpha_1 = \bigvee_2(\text{ring, king, return}) = (\text{ring} \wedge \text{king}) \vee (\text{ring} \wedge \text{return}) \vee (\text{king} \wedge \text{return})$$

Speciálně dostáváme, že:

$$\begin{aligned} \bigvee_0(\theta_1, \dots, \theta_m) &= 1, \\ \bigvee_1(\theta_1, \dots, \theta_m) &= \theta_1 \vee \dots \vee \theta_m, \\ \bigvee_m(\theta_1, \dots, \theta_m) &= \theta_1 \wedge \dots \wedge \theta_m. \end{aligned}$$

Pro $n \neq 0$ řetězec:

```
{
  "bool": {
    "should": [  $\theta_1, \dots, \theta_m$  ],
    "minimum_should_match":  $n$ 
  }
}
```

zapisuje podmínku $\bigvee_n(\theta_1, \dots, \theta_m)$. Například výše uvedenou podmínku α_1 zapíšeme:

```
{
  "bool": {
    "should": [
      {
        "match": {
          "text": "ring"
        }
      },
      {
        "match": {
          "text": "king"
        }
      },
      {
        "match": {
          "text": "return"
        }
      }
    ],
    "minimum_should_match" : 2
  }
}
```

Kvalita splnění podmínky $\bigvee_n(\theta_1, \dots, \theta_m)$ dokumentem D z indexu I je dána vztahem:

$$\text{score}(I, D, \bigvee_n(\theta_1, \dots, \theta_m)) = \sum_{i=1}^m \text{score}(I, D, \theta_i)$$

4 Obecná podmínka bool

Vezměme podmínky $\alpha_1, \dots, \alpha_m, \beta_1, \dots, \beta_n, \gamma_1, \dots, \gamma_k$ a číslo $0 \leq l \leq n$, takové, že $l = 0$ implikuje $m > 0$, pak řetězec:

```
{
  "bool": {
    "must": [  $\alpha_1, \dots, \alpha_m$  ],
    "should": [  $\beta_1, \dots, \beta_n$  ],
    "minimum_should_match":  $l$ 
    "must_not": [  $\gamma_1, \dots, \gamma_k$  ],
  }
}
```

je zkratkou za

$$\alpha_1 \wedge \dots \wedge \alpha_m \wedge \left(\bigvee_l(\beta_1, \dots, \beta_n) \right) \wedge \neg(\gamma_1 \vee \dots \vee \gamma_k).$$

Položka:

```
"must":  $\alpha$ 
```

je zkratkou za:

```
"must": [  $\alpha$  ]
```

Podobně pro položky "should" a "must_not". Dále chybějící položka "must" je zkratkou za:

```
"must": []
```

chybějící položka "should" je zkratkou za:

```
"should": []
```

pokud $k = 0$, jinak:

```
"should": 1
```

Chybějící položka "minimum_should_match" je zkratkou za:

```
"minimum_should_match": p
```

kde $p = 0$ pokud $m > 0$, jinak $p = 1$. Chybějící položka "must_not" je zkratkou za:

```
"must_not": []
```

Například:

```
{
  "bool": {
    "must": {
      "match": {
        "text": "of"
      }
    },
    "must_not": {
      "match": {
        "text": "ring"
      }
    }
  }
}
```

je zkratkou za:

```
{
  "bool": {
    "must": [
      {
        "match": {
          "text": "of"
        }
      }
    ],
    "should": [],
    "minimum_should_match": 0,
    "must_not": [
      {
        "match": {
```



```

    "text": "ring"
  }
}
]
}
}

```

a tedy zapisuje: $of \wedge \neg ring$.

5 Přibližné vyhledávání

Uvažujme nyní jen řetězce, které jsou termy nebo prázdný řetězec označený ϵ . Přidání znaku x k řetězci w označíme xw . Například pro $x = r$ a $w = ing$ je $xw = ring$ nebo $g\epsilon = g$.

Délku řetězce w označíme $|w|$. Například $|Ring| = 4$ a $|\epsilon| = 0$.

Vezměme řetězce w_1, w_2 . Definujeme (**editační** nebo také **Levenštejnovu vzdálenost** $d(w_1, w_2)$ řetězců w_1, w_2 následovně.

1. Pokud $w_1 = \epsilon$, pak $d(w_1, w_2) = |w_2|$,
2. pokud $w_2 = \epsilon$, pak $d(w_1, w_2) = |w_1|$,
3. pokud $w_1 \neq \epsilon$ a $w_2 \neq \epsilon$ a $xw'_1 = w_1$ a $xw'_2 = w_2$ pro nějaký znak x , pak $d(w_1, w_2) = d(w'_1, w'_2)$,
4. pokud $w_1 \neq \epsilon$ a $w_2 \neq \epsilon$ a $xw'_1 = w_1$ a $yw'_2 = w_2$ pro nějaké znaky x, y takové, že $x \neq y$, pak $d(w_1, w_2) = 1 + \min(d(w'_1, w'_2), d(w'_1, w_2), d(w_1, w'_2))$.

Například: $d(ing, ing) = d(ng, ng) = d(g, g) = d(\epsilon, \epsilon) = 0$ a proto $d(ring, king) = 1 + \min(d(ing, ing), d(ing, king), d(ring, ing)) = 1$.

Pro číslo $k \geq 0$ zavedeme ekvivalenci \approx_k na termech takovou, že položíme $w_1 \approx_k w_2$, jestliže $d(w_1, w_2) \leq k$. Například $ring \approx_1 king$, ale $ring \not\approx_1 kin$. Ekvivalence \approx_0 je relace rovnosti na termech. Tedy $w_1 \approx_0 w_2$, právě když $w_1 = w_2$.

Vezměme term q , pak

$$q \approx_k$$

je podmínka, kterou dokument D splňuje, jestliže pro $T(D) = t_1, \dots, t_n$ existuje $1 \leq i \leq n$ tak, že $q \approx_k t_i$. Například dokument "The Fellowship of the Ring" splňuje podmínku $king \approx_1$, protože obsahuje term $ring$ a $king \approx_1 ring$.

Kvalita splnění podmínky $q \approx_k$ dokumentem D z indexu I je dána vztahem:

$$score(I, D, q \approx_k) = \sum_{i=1}^n boost(t_i, q) \cdot IDF(I, t_i) \cdot tf(I, t_i, D)$$

kde t_1, \dots, t_n jsou všechny termy z $Set(T(D))$ ekvivalentní s q vzhledem k \approx_k . Číslo $boost(t_i, q)$ udává podobnost termů t_i a q a klesá od 2.2 k nule.

Například:

$score(I_1, \text{"The Fellowship of the Ring"}, \text{king} \approx_1) \approx 1.65 \cdot 0.98 \cdot 0.43$.

Pro $1 \leq k \leq 2$ řetězec:

```
{
  "match": {
    "text": {
      "query":  $D$ ,
      "fuzziness":  $k$ ,
      "fuzzy_transpositions": false
    }
  }
}
```

zapisuje podmínku $q_1 \approx_k \vee \dots \vee q_n \approx_k$, kde $T(D) = (q_1, \dots, q_n)$ a řetězec:

```
{
  "match": {
    "text": {
      "query":  $D$ ,
      "operator": "AND",
      "fuzziness":  $k$ ,
      "fuzzy_transpositions": false
    }
  }
}
```

podmínku $q_1 \approx_k \wedge \dots \wedge q_n \approx_k$. Například:

```
GET /movie/_search
{
  "query": {
    "match": {
      "text": {
        "query": "rurturn ring",
        "fuzziness": 1,
        "fuzzy_transpositions": false
      }
    }
  }
}
```

```

{
  "hits": {
    "hits": [
      {
        "_id": "3",
        "_score": 1.4610269,
        "_source": {
          "text": "The Return of the King"
        }
      },
      {
        "_id": "1",
        "_score": 0.9227538,
        "_source": {
          "text": "The Fellowship of the Ring"
        }
      }
    ]
  }
}

```

Otázky a úkoly na cvičení

1. Uvažujme indexovou proměnnou `book` s charakteristickou vlastností „Kniha d se jmenuje D .“ a hodnotou:

- 1 The Life And Opinions Of Tristram Shandy
- 2 Emma
- 3 Nightmare Abbey
- 4 One Day in the Life of Ivan Denisovich
- 5 Life After Life

Napište výrazy pro Elasticsearch, které mají následující charakteristické vlastnosti.

- (a) „Název D knihy d obsahuje slova `life` a `of`.“
- (b) „Název D knihy d obsahuje jméno `Ivan Denisovich` nebo `Tristram Shandy`.“
- (c) „Název D knihy d obsahuje slovo `life` a předložku `of` nebo `in`.“
- (d) „Název D knihy d neobsahuje slovo `life`.“
- (e) „Název D knihy d neobsahuje jména: `Emma`, `Shandy` a `Ivan`.“
- (f) „Název D knihy d obsahuje aspoň dvě ze slov: `in`, `of` a `life`.“

Spočítejte hodnoty výrazů a ověřte je s těmi, které spočítá systém Elasticsearch.

2. Určete hodnotu a charakteristickou vlastnost následujícího výrazu.

```
{
  "bool": {
    "must": {
      "match": {
        "text": "life of"
      }
    },
    "should": {
      "match": {
        "text": "day"
      }
    }
  }
}
```

Hodnotu porovnejte s hodnotou získanou ze systému Elasticsearch.

3. Spočítejte editační vzdálenost následujících dvojic termů:

- (a) day a da
- (b) day a dya
- (c) day a dy
- (d) day a of

4. Určete, které dokumenty v hodnotě indexové proměnné book splňují podmínku $one \approx_2$. Odpověď ověřte systémem Elasticsearch.
5. Určete hodnotu následujícího výrazu.

```
GET /movie/_search
{
  "query": {
    "match": {
      "text": {
        "query": "lfe ifan",
        "operator": "AND",
        "fuzziness": 1,
        "fuzzy_transpositions": false
      }
    }
  }
}
```