

# Hitungan Manual: BoW, TF-IDF, dan Cosine Similarity

## 1. Dataset Mainan

Tiga dokumen pendek (lowercase, tanpa penghapusan stopword agar mudah dihitung):

- $d_1$ : “barang cepat sampai”
  - $d_2$ : “pengiriman cepat dan rapi”
  - $d_3$ : “sangat lambat pengiriman”
- Kosakata global (urut abjad) kita tetapkan sebagai

$$\mathcal{V} = [\text{barang}, \text{cepat}, \text{sampai}, \text{pengiriman}, \text{dan}, \text{rapi}, \text{sangat}, \text{lambat}],$$

dengan ukuran korpus  $N = 3$  dokumen.

## 2. Bag-of-Words (BoW)

### 2.1 Definisi

Representasi BoW untuk sebuah dokumen  $d$  adalah vektor frekuensi:

$$\text{BoW}(d) = [f(w_1, d), f(w_2, d), \dots, f(w_V, d)],$$

dengan  $f(w_i, d)$  menyatakan jumlah kemunculan kata  $w_i$  dalam  $d$ , dan  $V = |\mathcal{V}|$ .

### 2.2 Perhitungan BoW

Panjang dokumen:

$$|d_1| = 3, \quad |d_2| = 4, \quad |d_3| = 3.$$

Susunan vektor mengikuti urutan  $\mathcal{V}$ .

$$\text{BoW}(d_1) = [1, 1, 1, 0, 0, 0, 0, 0]$$

$$\text{BoW}(d_2) = [0, 1, 0, 1, 1, 1, 0, 0]$$

$$\text{BoW}(d_3) = [0, 0, 0, 1, 0, 0, 1, 1]$$

## 3. TF, DF, dan IDF

### 3.1 Term Frequency (TF)

Gunakan TF ternormalisasi (proporsional terhadap panjang dokumen):

$$\text{tf}(t, d) = \frac{f_{t,d}}{|d|}.$$

Contoh:

- Di  $d_1$ : setiap kata yang muncul memiliki  $\text{tf} = \frac{1}{3}$ .
- Di  $d_2$ : setiap kata yang muncul memiliki  $\text{tf} = \frac{1}{4}$ .
- Di  $d_3$ : setiap kata yang muncul memiliki  $\text{tf} = \frac{1}{3}$ .

### 3.2 Document Frequency (DF)

$df(t)$  adalah jumlah dokumen yang *mengandung* kata  $t$ :

$$\begin{aligned} df(\text{barang}) &= 1 (d_1), \quad df(\text{cepat}) = 2 (d_1, d_2), \quad df(\text{sampai}) = 1 (d_1), \\ df(\text{pengiriman}) &= 2 (d_2, d_3), \quad df(\text{dan}) = 1 (d_2), \quad df(\text{rapi}) = 1 (d_2), \\ df(\text{sangat}) &= 1 (d_3), \quad df(\text{lambat}) = 1 (d_3). \end{aligned}$$

### 3.3 Inverse Document Frequency (IDF)

Pakai varian IDF yang stabil terhadap pembagi nol:

$$idf(t) = \ln\left(\frac{N}{1 + df(t)}\right), \quad N = 3.$$

Maka:

$$df(t) = 1 \Rightarrow idf(t) = \ln\left(\frac{3}{2}\right) \approx 0.4055, \quad df(t) = 2 \Rightarrow idf(t) = \ln\left(\frac{3}{3}\right) = 0.$$

Konsekuensinya:

$$\begin{aligned} idf(\text{barang}) &= idf(\text{sampai}) = idf(\text{dan}) = idf(\text{rapi}) = idf(\text{sangat}) = idf(\text{lambat}) \approx 0.4055, \\ idf(\text{cepat}) &= idf(\text{pengiriman}) = 0. \end{aligned}$$

## 4. Vektor TF-IDF

Definisi dasar:

$$tfidf(t, d) = tf(t, d) \cdot idf(t).$$

Kita tampilkan hanya komponen yang tidak nol (pembulatan 4 desimal).

### 4.1 Dokumen $d_1$ (“barang cepat sampai”)

$$\begin{aligned} tfidf(\text{barang}, d_1) &= \frac{1}{3} \cdot 0.4055 \approx 0.1352 \\ tfidf(\text{cepat}, d_1) &= \frac{1}{3} \cdot 0 = 0 \\ tfidf(\text{sampai}, d_1) &= \frac{1}{3} \cdot 0.4055 \approx 0.1352 \end{aligned}$$

Vektor TF-IDF (urutan  $\mathcal{V}$ ):

$$[0.1352, 0, 0.1352, 0, 0, 0, 0, 0].$$

### 4.2 Dokumen $d_2$ (“pengiriman cepat dan rapi”)

$$\begin{aligned} tfidf(\text{pengiriman}, d_2) &= \frac{1}{4} \cdot 0 = 0, \quad tfidf(\text{cepat}, d_2) = \frac{1}{4} \cdot 0 = 0, \\ tfidf(\text{dan}, d_2) &= \frac{1}{4} \cdot 0.4055 \approx 0.1014, \quad tfidf(\text{rapi}, d_2) = \frac{1}{4} \cdot 0.4055 \approx 0.1014. \end{aligned}$$

Vektor TF-IDF:

$$[0, 0, 0, 0, 0.1014, 0.1014, 0, 0].$$

### 4.3 Dokumen $d_3$ (“sangat lambat pengiriman”)

$$\begin{aligned} tfidf(\text{sangat}, d_3) &= \frac{1}{3} \cdot 0.4055 \approx 0.1352, \quad tfidf(\text{lambat}, d_3) = \frac{1}{3} \cdot 0.4055 \approx 0.1352, \\ tfidf(\text{pengiriman}, d_3) &= \frac{1}{3} \cdot 0 = 0. \end{aligned}$$

Vektor TF-IDF:

$$[0, 0, 0, 0, 0, 0, 0.1352, 0.1352].$$

## 5. Cosine Similarity (Contoh)

Misalkan ada *query*  $q$ : “barang cepat”. Dengan kosakata yang sama,

$$|q| = 2, \quad \text{tf(barang, } q) = \frac{1}{2}, \quad \text{tf(cepat, } q) = \frac{1}{2}.$$

Dengan  $\text{idf}(\text{barang}) \approx 0.4055$  dan  $\text{idf}(\text{cepat}) = 0$ :

$$\text{tfidf}(\text{barang, } q) = \frac{1}{2} \cdot 0.4055 \approx 0.2027, \quad \text{tfidf}(\text{cepat, } q) = \frac{1}{2} \cdot 0 = 0.$$

Vektor TF-IDF  $q$  (urutan  $\mathcal{V}$ ):

$$\mathbf{q} = [0.2027, 0, 0, 0, 0, 0, 0, 0].$$

Cosine similarity antara  $\mathbf{q}$  dan  $\mathbf{d}$  didefinisikan:

$$\cos(\theta) = \frac{\mathbf{q} \cdot \mathbf{d}}{\|\mathbf{q}\| \|\mathbf{d}\|}.$$

**Dengan  $d_1$ .**

$$\mathbf{d}_1 = [0.1352, 0, 0.1352, 0, 0, 0, 0, 0].$$

$$\mathbf{q} \cdot \mathbf{d}_1 = 0.2027 \times 0.1352 + 0 \times 0.1352 = 0.0274.$$

$$\|\mathbf{q}\| = \sqrt{0.2027^2} \approx 0.2027, \quad \|\mathbf{d}_1\| = \sqrt{0.1352^2 + 0.1352^2} \approx \sqrt{2 \times 0.0183} \approx 0.1913.$$

$$\cos(\mathbf{q}, \mathbf{d}_1) \approx \frac{0.0274}{0.2027 \times 0.1913} \approx \frac{0.0274}{0.0388} \approx 0.706.$$

**Dengan  $d_2$  dan  $d_3$ .** Karena  $\mathbf{q}$  hanya memiliki komponen “barang” dan  $d_2, d_3$  tidak memiliki “barang”,

$$\mathbf{q} \cdot \mathbf{d}_2 = 0, \quad \mathbf{q} \cdot \mathbf{d}_3 = 0 \Rightarrow \cos(\mathbf{q}, \mathbf{d}_2) = 0, \quad \cos(\mathbf{q}, \mathbf{d}_3) = 0.$$

**Kesimpulan:**  $q$  paling mirip ke  $d_1$  (nilai cosine terbesar).

## 6. Catatan Singkat tentang N-gram

Jika menggunakan *bigram* untuk menangkap negasi:

dokumen: “tidak bagus”  $\Rightarrow$  fitur bigram: “tidak\_bagus”.

BoW bigram:

$$\text{BoW}_{\text{bigram}}(d) = [f(\text{“tidak\_bagus”}, d), \dots].$$

Skor TF-IDF bigram dihitung sama: gunakan frekuensi bigram sebagai  $f_{t,d}$  lalu kalikan dengan  $\text{idf}(t)$ .

## 7. Ringkasan

- BoW menghitung *frekuensi kata* per dokumen.
- TF-IDF menurunkan bobot kata yang *terlalu umum* dan meninggikan kata *informatif*.
- Cosine similarity mengukur kedekatan vektor TF-IDF antar teks.
- Vektor TF-IDF siap dipakai sebagai fitur untuk model klasik (LogReg/SVM/NB).