**(1)** | Objective of the project |

We want to Recommended
the product that user
has is currently looking

↳ Why we car about
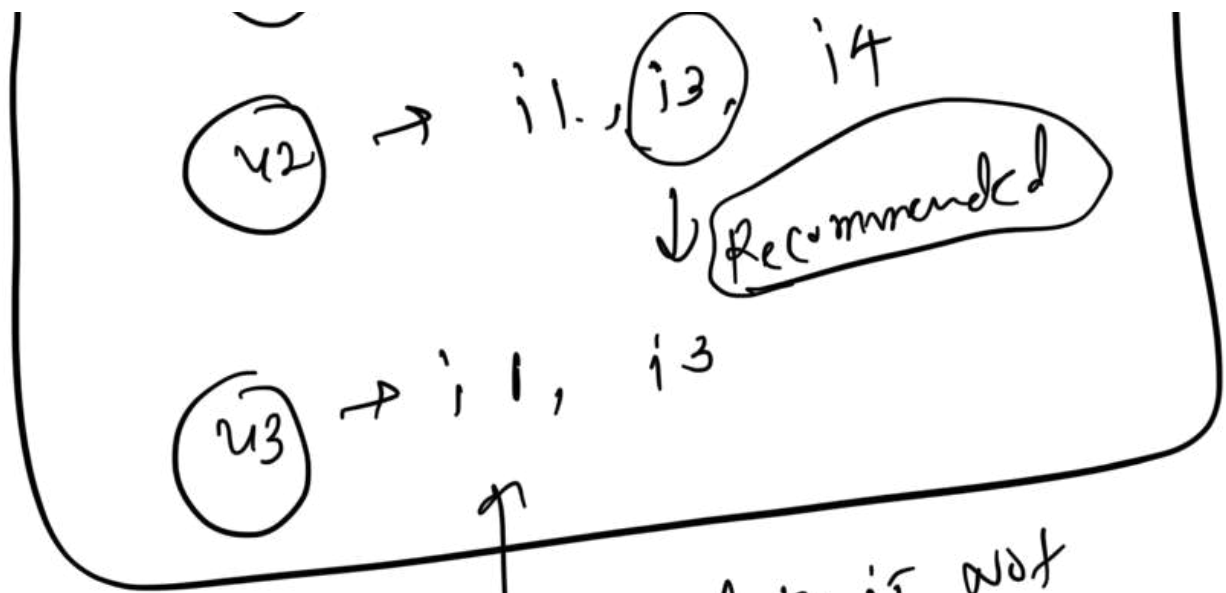  search about.

→ estimated 35% increase
  in Revenue.
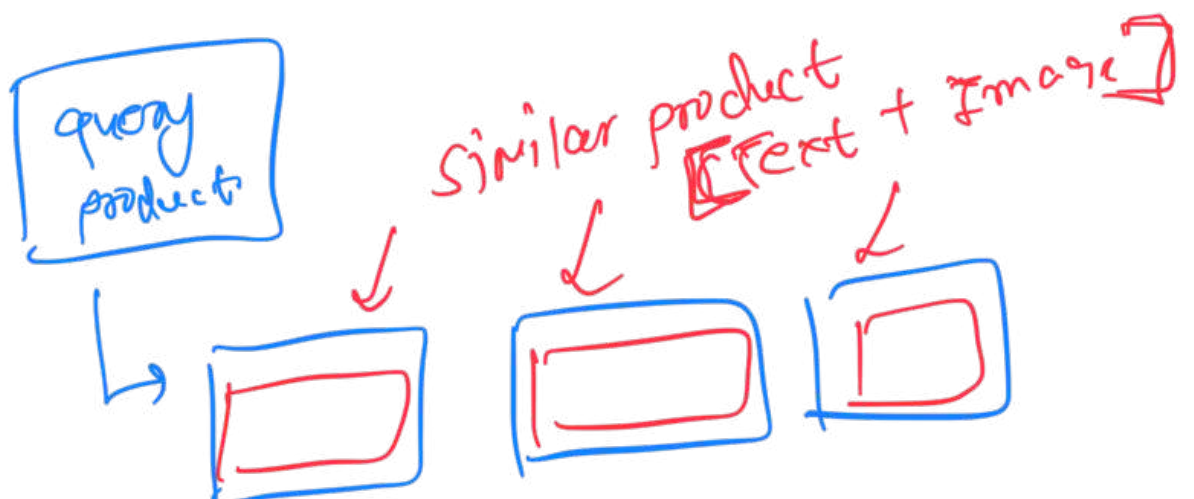
(1) content based
    Recommendation.
    | Text, Image description |

(11) | Collborative Filtering |

    (u1) → i1, i2, i3

$u2 \rightarrow i1, (i3), i4$

Recommended

$\downarrow$

$u3 \rightarrow i1, i3$

$\uparrow$

this data is not
easily available

For this project we will use
Content Sorted Recommendation.

query
product

Similar product
[Text + Image]

# Plan of Attack

1. data acquisition
2. data cleaning
3. Text-processing
4. Linear Algebra.
5. Text based product Recommendation.
   * BoW
   * W2V
   * TF-IDF
6. Image based product Recommendation
7. A/B Testing
8. Future Work

# 1 dataset Information

1 Asin
2 brand
3 Color
5 Image-url
S title
C price

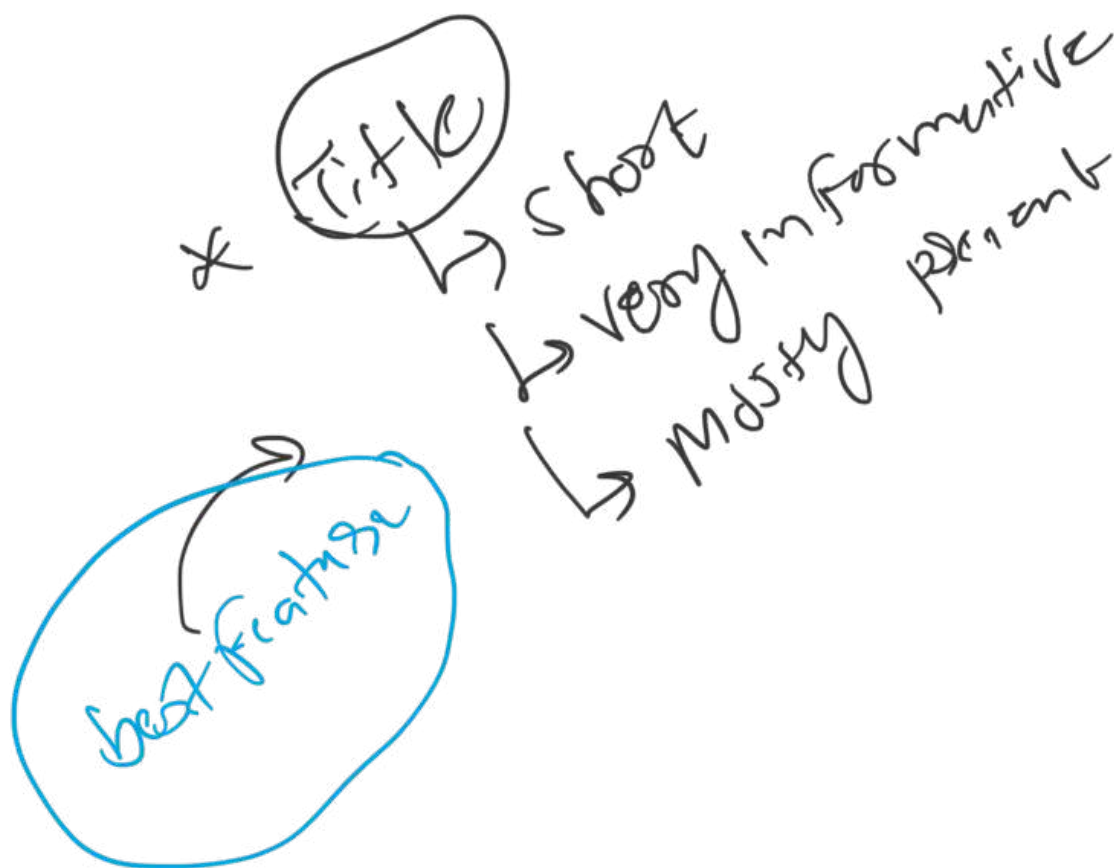## 2) data cleaning

\* Very important
overlooked

* often

* 
* basic EDA

* Top products
* Top 10 Most Frequent products

↑ you will get class imblance information

* Top 20 brands
Top 20 Frequent brands

bad or crazy

* Real world data is very

* (Title)
  ↳ Short
  ↳ Very informative
  ↳ Mostly relevant

best feature

* Remove duplicates.

* Find duplicates title

pandas dataframe. duplicated (File)

ASIN=#1

title: Red top

ASIN= #2

title: Red top

only color difference

ASIN = #3

title - Blue top

* Remove products with short title

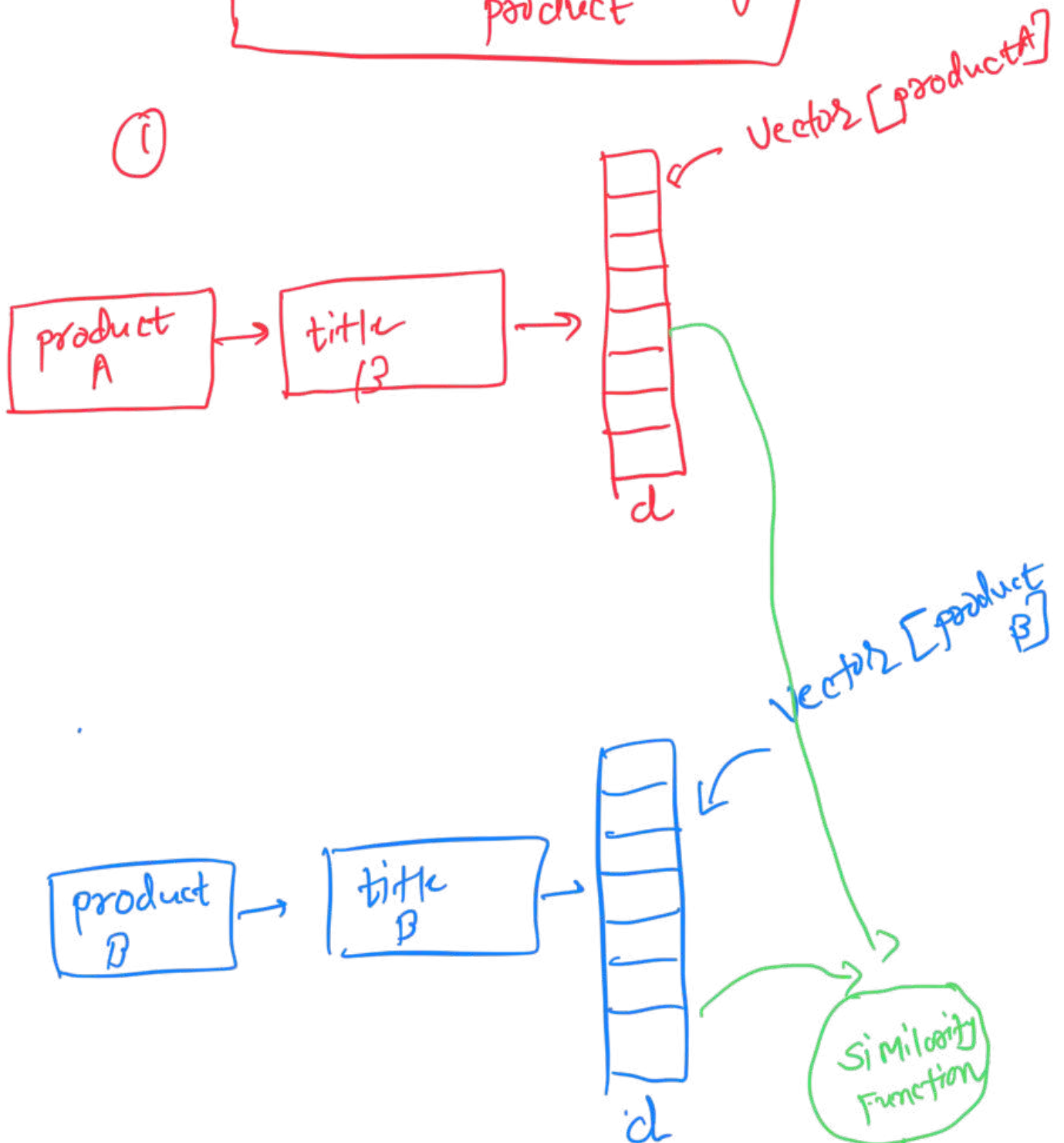Text pre processing: Tokenization and stop word Removal

Text pre processing

(i) STOP word Removal

(ii) Stemming [experiment and check if it works]

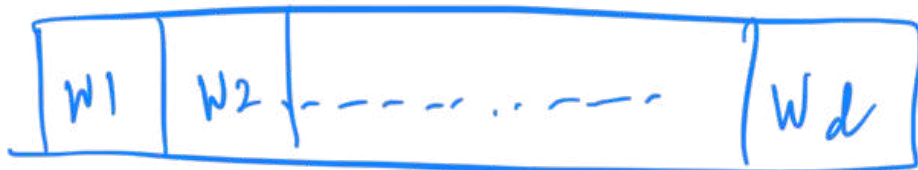First art approach

(i) Text based similarity product

(i)

product A → title B →

Vector [productA]

d

product B → title B →

Vector [product B]

d

Similarity Function

Similarity score

① Simple BOW Vectors.

T1, T2, T3 ..... Tm (Title)

① build vocabulary set [S]

S = all the words in all the title.

$S \in \mathbb{R}^d$

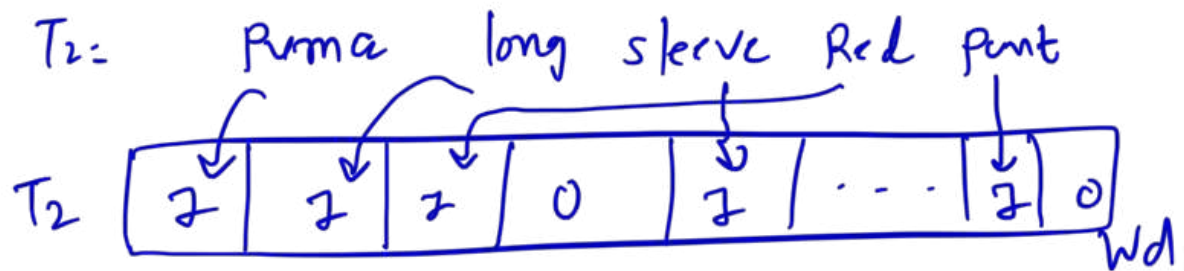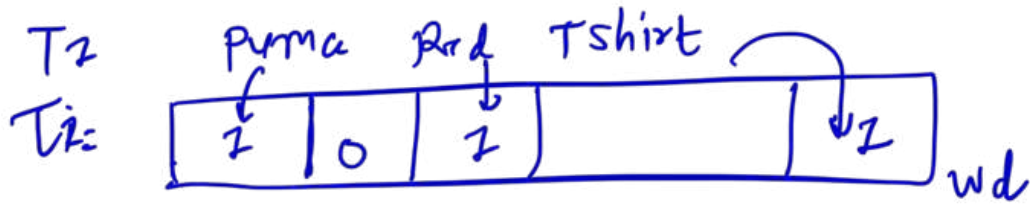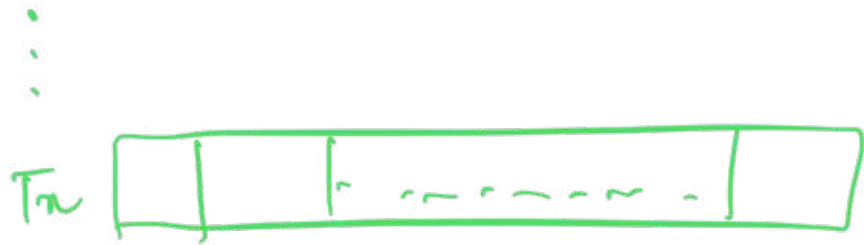| W1 | W2 | - - - - - - - . - - - - | Wd |

d is very large

② Generate vectors using BOW

Ti | W1 | V2 | W3 ........ | Wd |

T2 | W1 | V1 | W4 . .... | Wd |

:

$T_n$

$T_2$     Puma    Red    Tshirt

$T_1:$

| 1 | 0 | 1 | | | 1 |
|---|---|---|---|---|---|

wd

$T_2:$     Puma    long    sleeve    Red    Pant

$T_2$

| 1 | 1 | 1 | 0 | 1 | $\cdots$ | 1 | 0 |
|---|---|---|---|---|---|---|---|

wd

## Problem

①   sequence information

②    counter

product $\rightarrow$ Title $\rightarrow$ 

**Vectorizer**

BoW
TF IDF
TF
embedding
Matrix

$\rightarrow$ $R^d$ vector

①    BoW $\rightarrow$ vectors

↳ Very sparse Vector



This very big matrix store in sparsed Form

① get vector → BIW

② use Eucledeon distance For pair wise distance

if most of the words commonly occuring then they similar

$$\text{dist} [t_i, t_j] = \sum_{k=1}^{n} [T_{ik} - T_{jk}]^2$$

deduplication
can be improved

2. TF IDF based
approach

Term Frequency    TF

Inverse Document    IDF
Frequency

$$TF [w_j, t_i] = \frac{\# \text{ times } w_j \text{ occures in } T_i}{\# \text{ words in } T_i}$$

$$IDF(Wi, D) = log\left[\frac{\text{\# Document in corpus}}{\text{\# document in } D \text{ having word } Wj}\right]$$

entire corpus

How rare is a word is in corpus

TF IDF Vector

$T_i \longrightarrow W_1, W_2, V_3 \ldots W_d$

$W_1 \longrightarrow TF(W_1, T_i) * IDF(W_1, D)$

BOW Vector

| W1 | W2 | | | | Wd |
|----|----|---|---|---|----|
| 3 | 2 | . | . | . | .- |

TF IDF Vector

| W1 | W2 | | | Wd |
|-----|------|--|--|------|
| 0.2 | 0.02 | | | 0.05 |

TF IDF (Wi) ↑ → TF ↑

$*$ TF IDF results were better than IDF.

Problem with TF IDF  $*$

Observation ① words don't Repeat Mostly in our dataset

$$TF(W_i, T_j) = \frac{\# \ W_i \ in \ T_j}{\# \ words \ in}$$

$$\frac{\#\text{ words in}}{T_j}$$

$T_1 \rightarrow \quad W1, \ W2, \ W3.. \ W4, \ W5, \ W6... \quad W10$

$T_2 \rightarrow \quad W1, \ W2, \ W3, \ W4$

※ if all word Frequency is 1 then

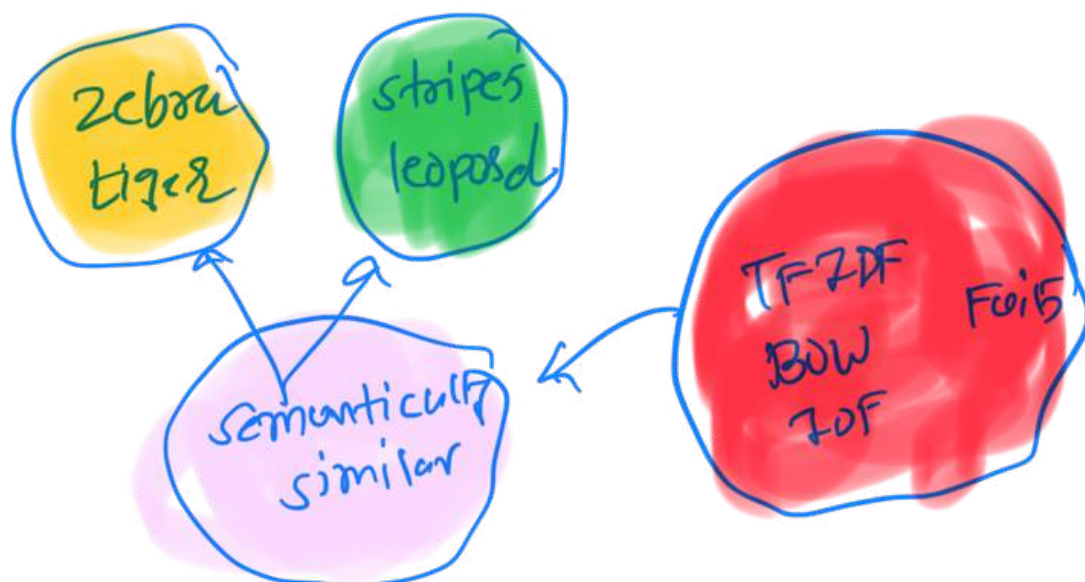For $T_1$ $\qquad TF(W_i, T_1) = \dfrac{1}{10} = 0.1$

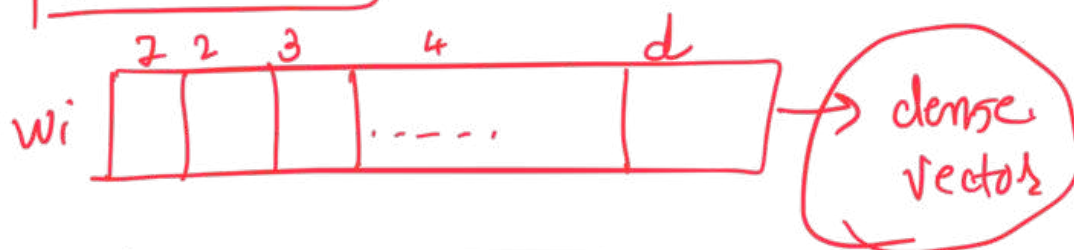For $T_2$ $\qquad TF(W_j, T_2) = \dfrac{1}{4} = 0.25$

this algorithm will favour th short title

✗ We can use directly IDF

This experiment we should Try, it may work as it may Not.

semantic similarity
and W2V

zebra
tiger

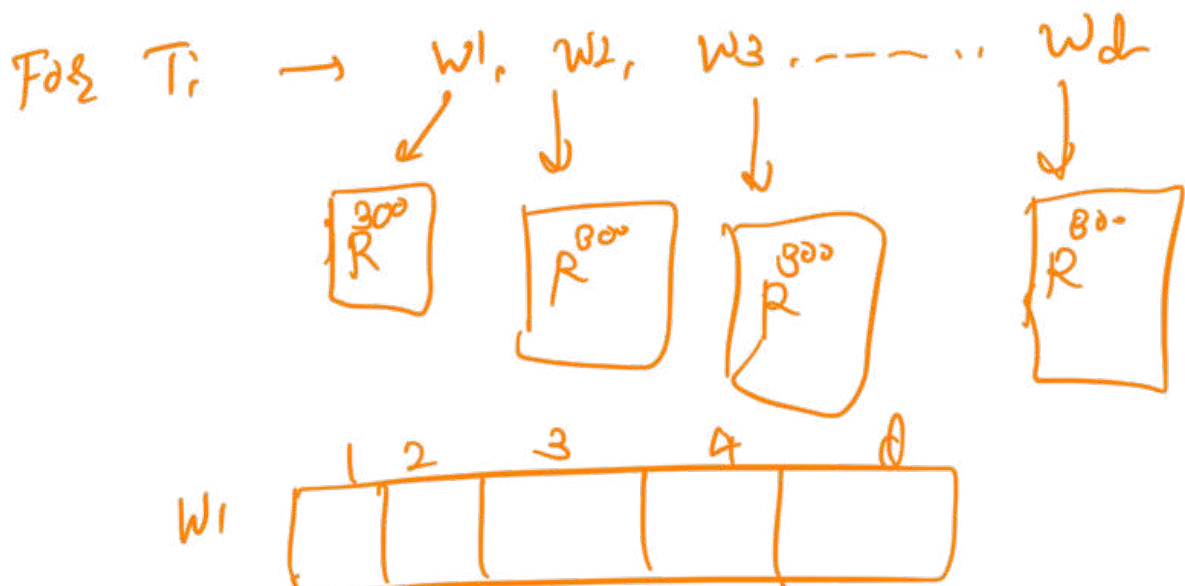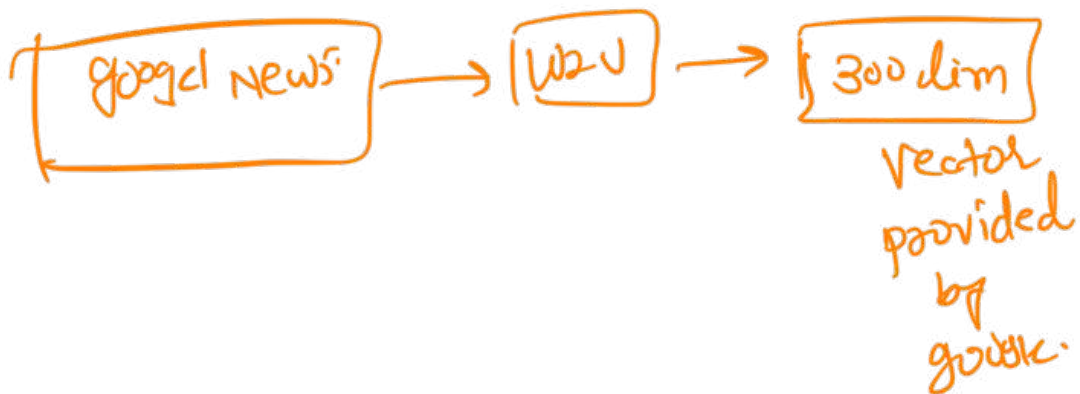stripes
leopard

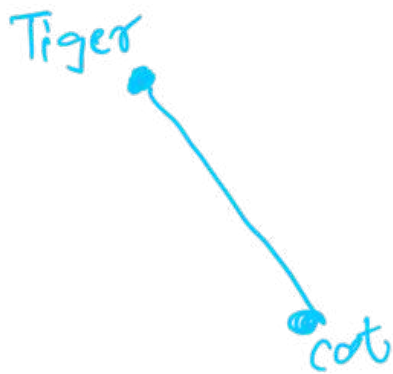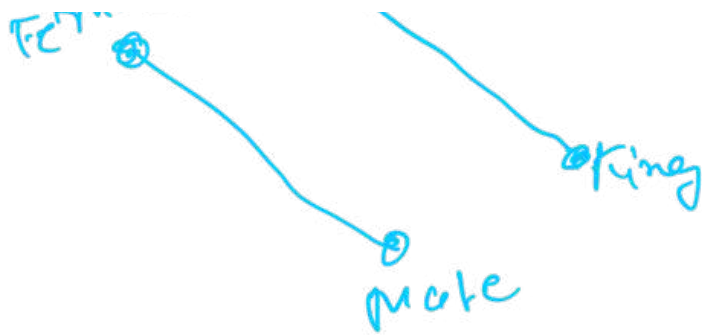TF-IDF
BOW
7OF

Fois

Semantically
similar

# Word2vec

$$w_i \quad \boxed{7 \; 2 \; 3 \quad 4 \quad \cdots \quad d}$$

$w_i$ → dense vector

D → Very large corpus required

geometric intuition of W2V

queen

male

T.C... King

Male

Tiger

cat

google News $\rightarrow$ W2V $\rightarrow$ 300 dim

vector
provided
by
google.

For $T_i \rightarrow W_1, W_2, W_3 \cdots \cdots W_d$

$R^{300}$   $R^{300}$   $R^{300}$   $R^{300}$

|  | 1 | 2 | 3 | 4 | d |
|---|---|---|---|---|---|
| W1 |  |  |  |  |  |

$w_2$

| 1 | 2 | 3 | n | d |
|---|---|---|---|---|

. . .

$W_n$

| 1 | 2 | 3 | 4 | | l |
|---|---|---|---|---|---|

W2V(Ti)   $\dfrac{1}{K}$

| 1 | 2 | 3 | 4 | | d |
|---|---|---|---|---|---|

Avg W2V

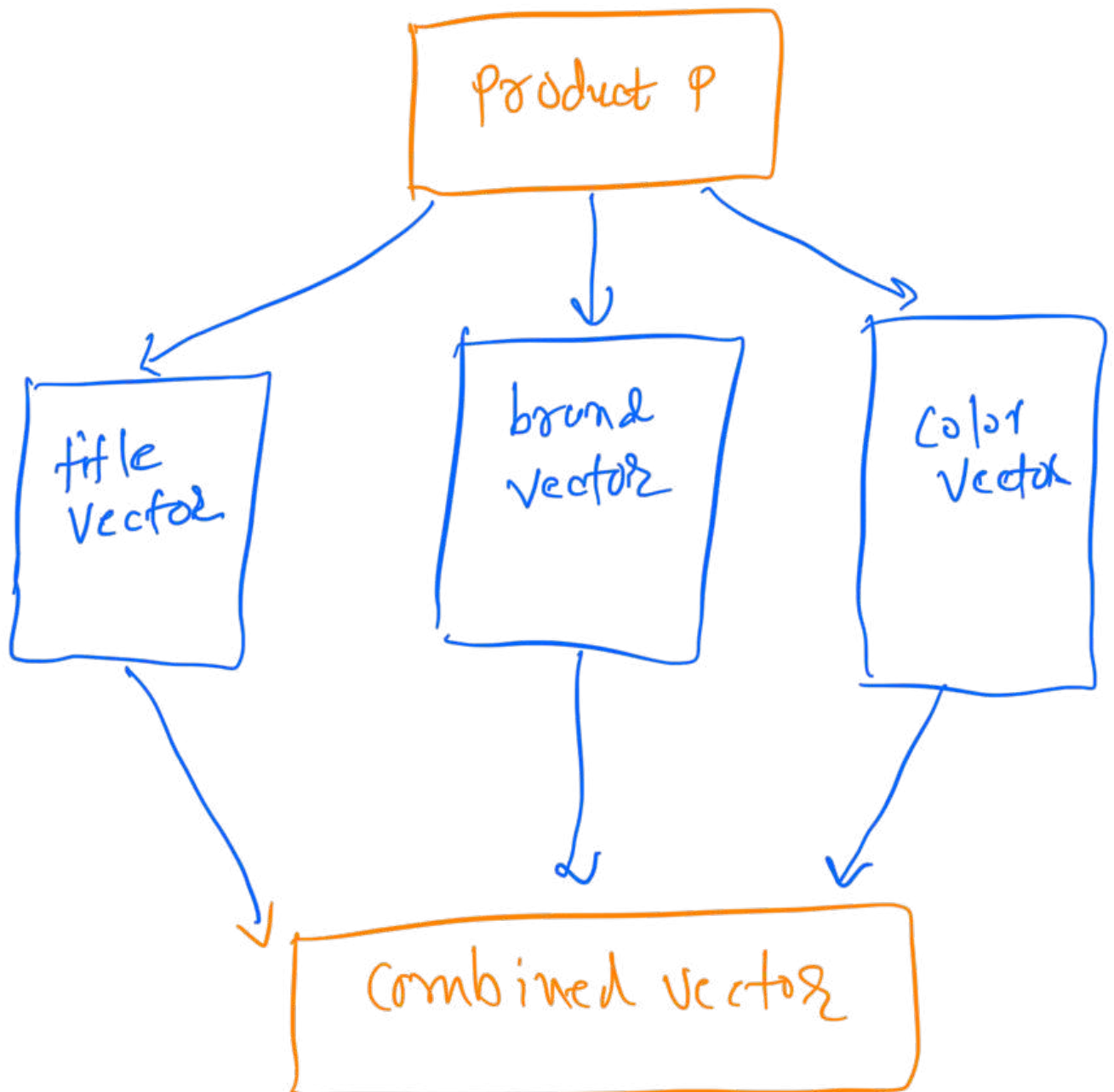↳ We get better result with
W2V → this take cares
about ↑ similarity
semantic

TF IDF weighted
W2 V

↳ word importance is useful

↳

semantic similarity.
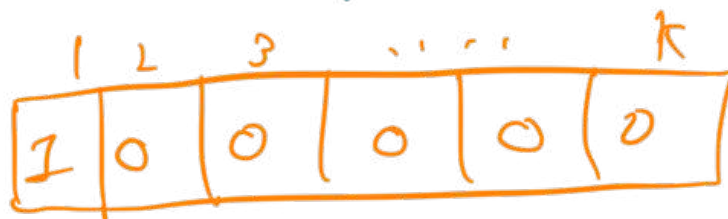
using brand and
color and other features

Product P

title
vector

brand
vector

color
vector

combined vector

brand vector $\longrightarrow$ one hot encoded

b1.  b2.  b3.....  bm

| 1 | 2 | 3 | . . . | . | m |
|---|---|---|---|---|---|
| 7 | 0 | 0 | 0 | 0 | 0 |

b1

one hot encoded vector

| 1 | 2 | 3 | . . . . . | k |
|---|---|---|---|---|
| I | 0 | 0 | 0 | 0 | 0 |

C1, C2, C3 . . . . . Ck  color

K $\longleftarrow$ 300 $\longrightarrow$ | $\longleftarrow$ m $\longrightarrow$ | $\longleftarrow$ n $\longrightarrow$ |

| Title | brand | color |
|---|---|---|

use Euclidean distance

(1)  Prefer showing customer to

specific   boxand   or color

use   Weighted Euclidean distance

$W_{title} = 1$ ,   $W_{brand} = 3$ , $W_{color} = 7$

$distance = W_{title} * title\ feature$

$+ W_{brand} * brand\ feature$

$+ W_{color} * color\ feature.$

building Real world
solution.

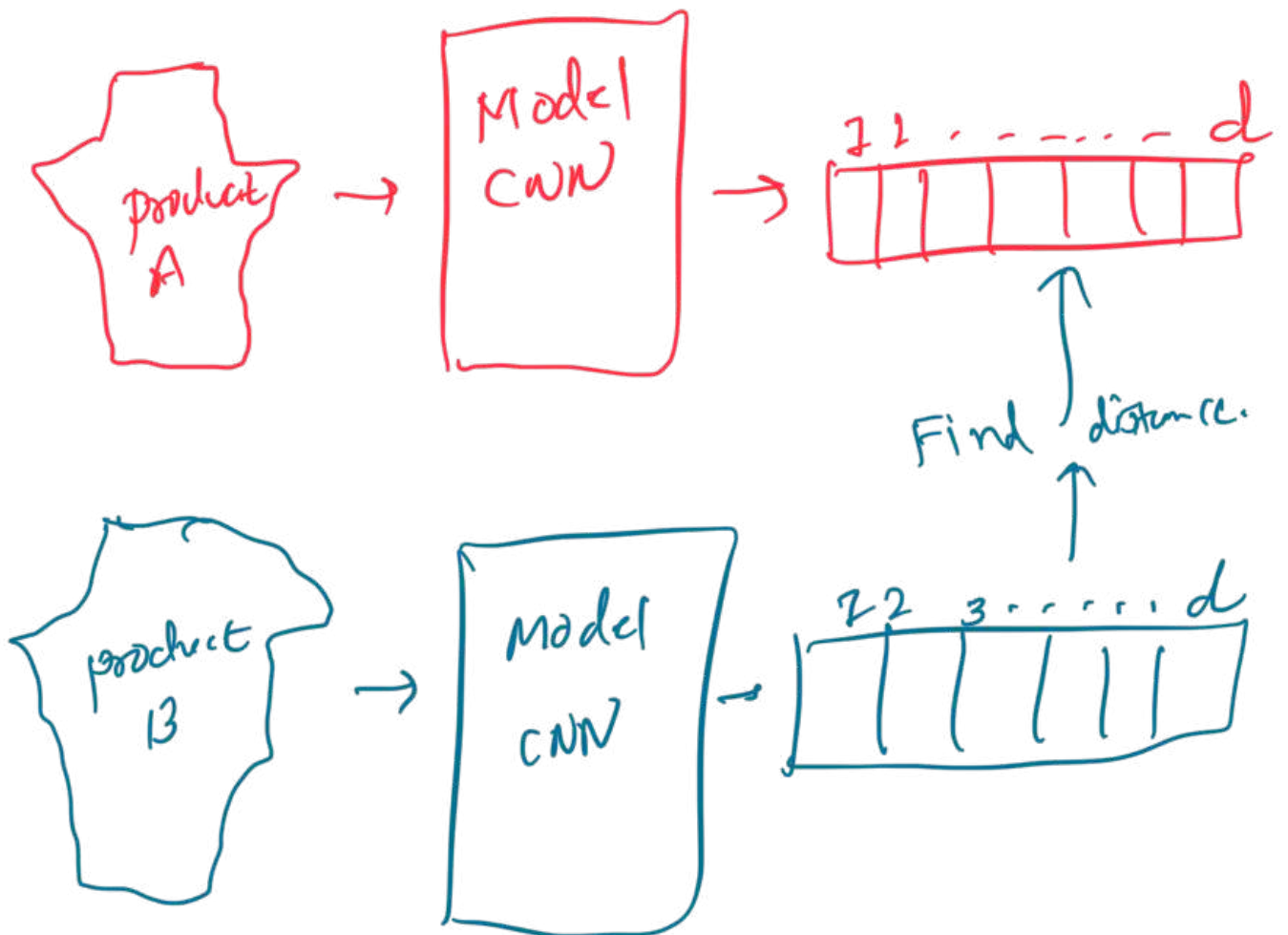| Product query | BOW | 7 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| | TFIDF | 7 | 2 | 5 | ··· 7 |
| | wav | 3 | 7 | 5. |

Apply business Rule

business is Most important

Deep Learning based
visual similarity.

Algorithm
Remains
Same

Color

product → Product Image → shirt / pant / tees → Simi larity vector

Embedding

Product A → Model CNN → z1 ...... d

Product B → Model CNN → z2 3 ...... d

Find distance

## A|B Testing

$P_q$ → Solution 1 → $P_1$ $B$ $P_5$ $P_7$

$P_q$ → Solution 2 → $P_3$ $P_5$ $P_7$ $P_r$.

How should we evaluate solution 1 and solution 2 quantatively?

↳ A|B Testing.

$u$

$\subset$ users

Solution 7

Measure sales.

Sales $[u_1]$ > Sales $[u_2]$

solution 7 is good