

→ So both alg-w2v and tfidf-w2v } weighting schemes
Sentences \rightarrow Vec

21.11 Bag of words (code sample)

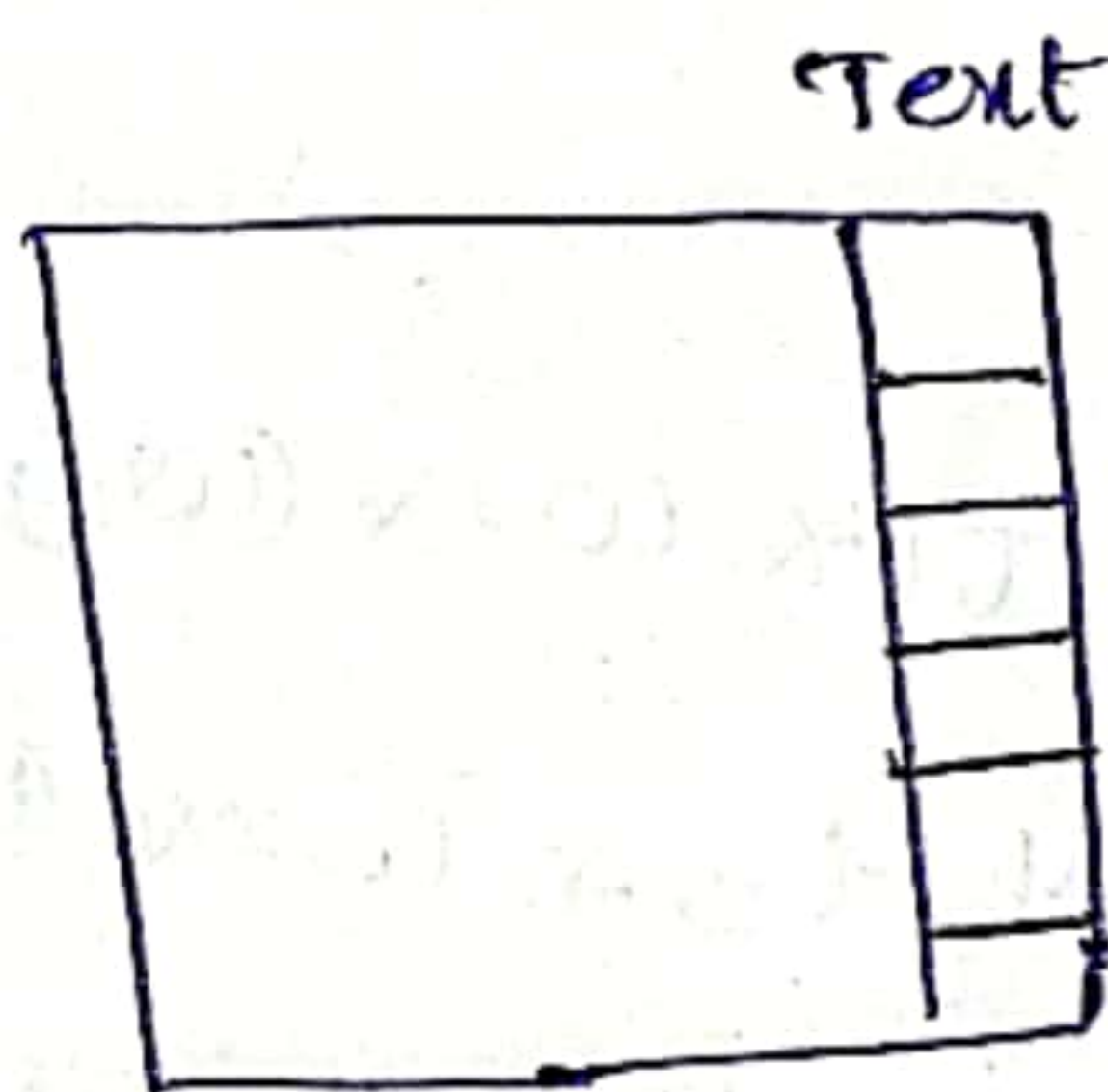
→ for Bag of words we use Scikit learn

→ we use function called `CountVectorizer()`

$X = \text{CountVectorizer}()$

$Y = X.\text{fit_transform}(\text{final}['\text{Text}']).\text{values}$

`final['Text']` basically means final data frame
and with a text



→ getting Text column and convert them into
values

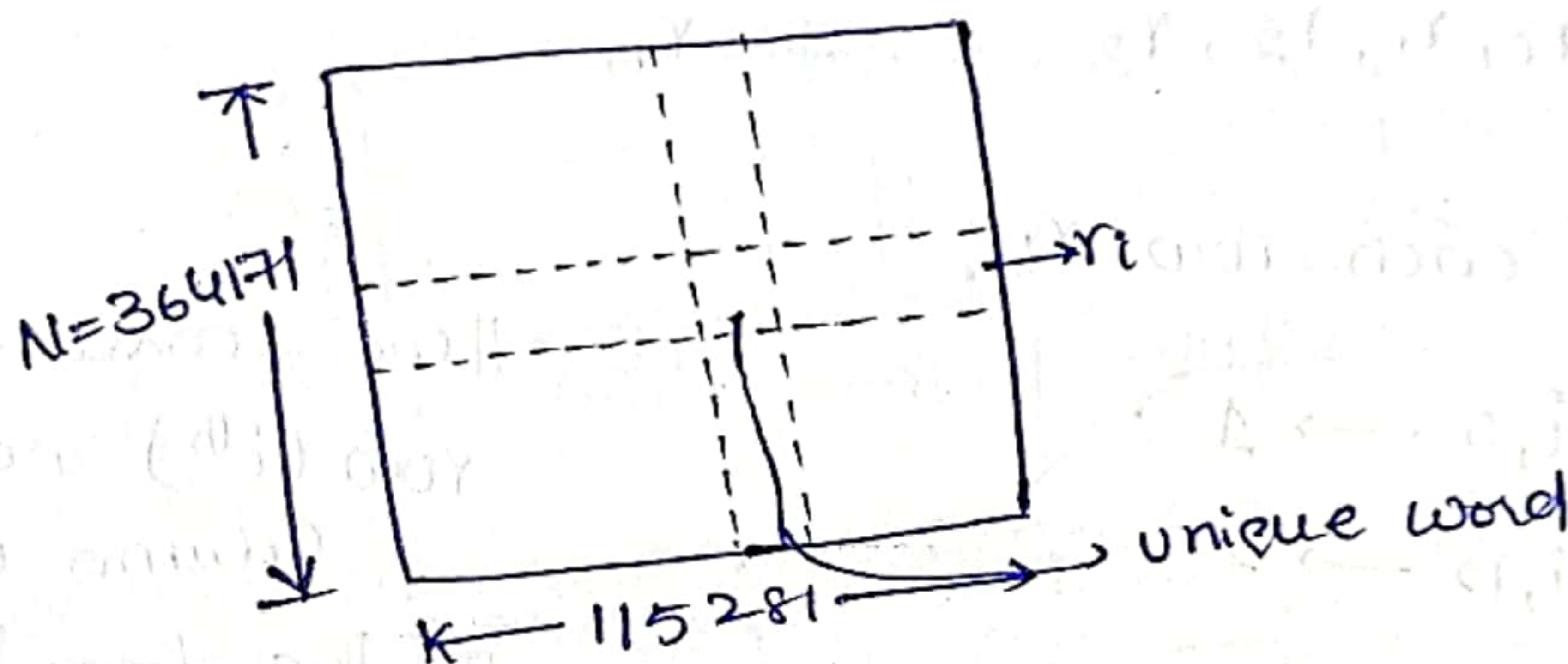
→ So `final['Text'].values` is converting final
Text \rightarrow values (vectors)

→ `type(Y)` gives `Scipy.sparse.csr.csr_matrix`
as this is sparse matrix

→ `y.get_shape()`

`get_shape()` fn gives the shape of the sparse matrix.

here of $(364171, 115281)$
 $N = \text{no of reviews or documents}$
 $\rightarrow \text{no of unique words.}$

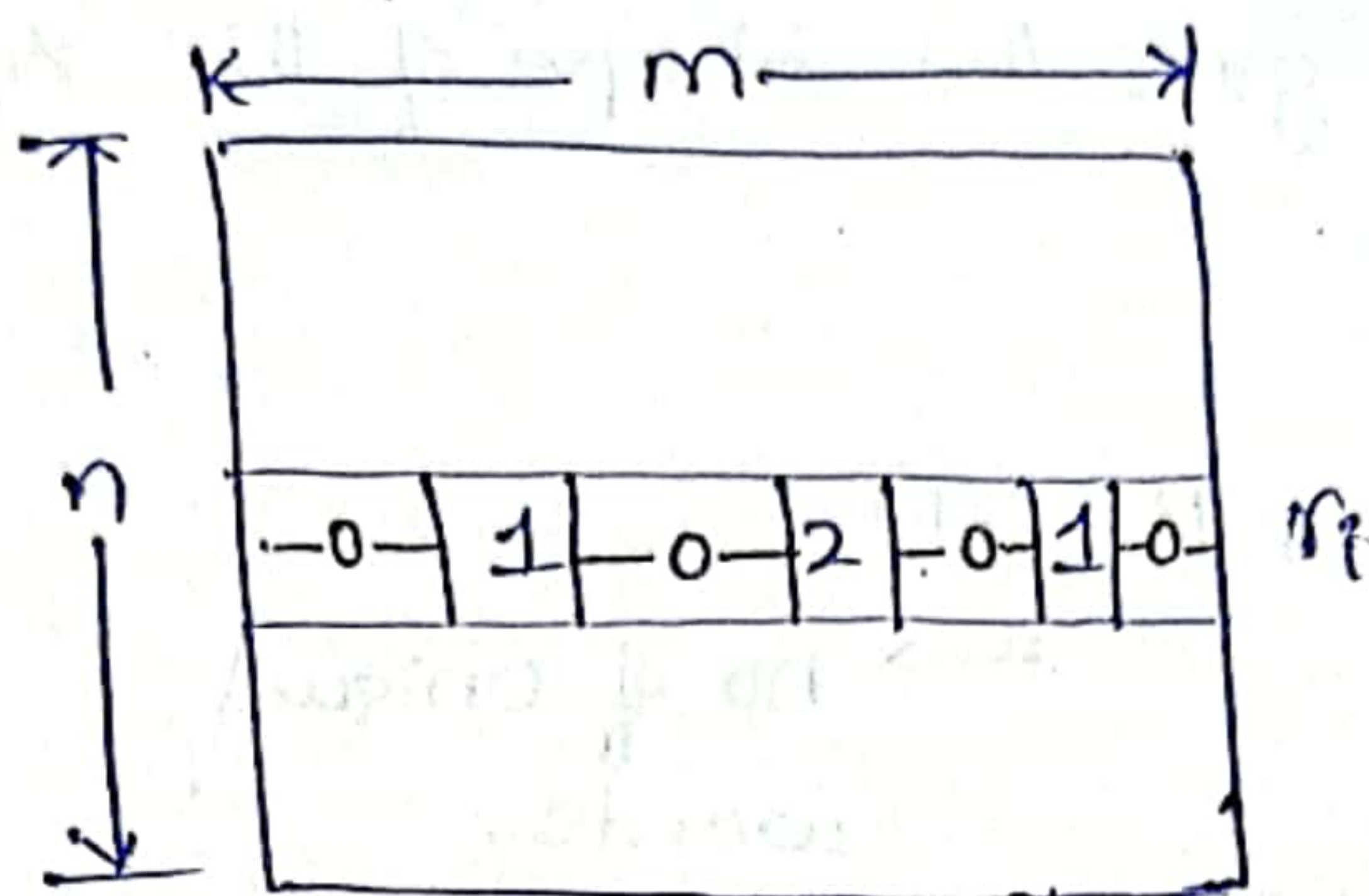


→ In General if we have a matrix of m rows and n columns the time/space complexity would be

$$O(m \times n)$$

→ But in Bow we get the matrix which is a sparse matrix. The space complexity of a sparse matrix can be reduced compared to a normal matrix or even low.

→ let us consider a sparse matrix, with n rows and m columns.



let $m=100$

$r_0, r_1, r_2, r_3, \dots, r_n$

for each row r_i :

$(i, 5 \rightarrow 1)$

$(i, 12 \rightarrow 2)$

$(i, 36 \rightarrow 1)$

we consider the row (i^{th}) & element column with value in the form key value pairs

$\underbrace{\text{row, column}}_{\text{Key}}, \rightarrow \text{Value}$

like we try to store the non zero values in the python dictionaries with the help of row and column

→ So the space complexity will be reduced to $\text{no. of rows} \rightarrow \text{total no. of non zero cells}$
3

why 3 because for each data point we are storing row no, column no & data point

r_i :- $O(m)$ space $m=100$.

as we are storing

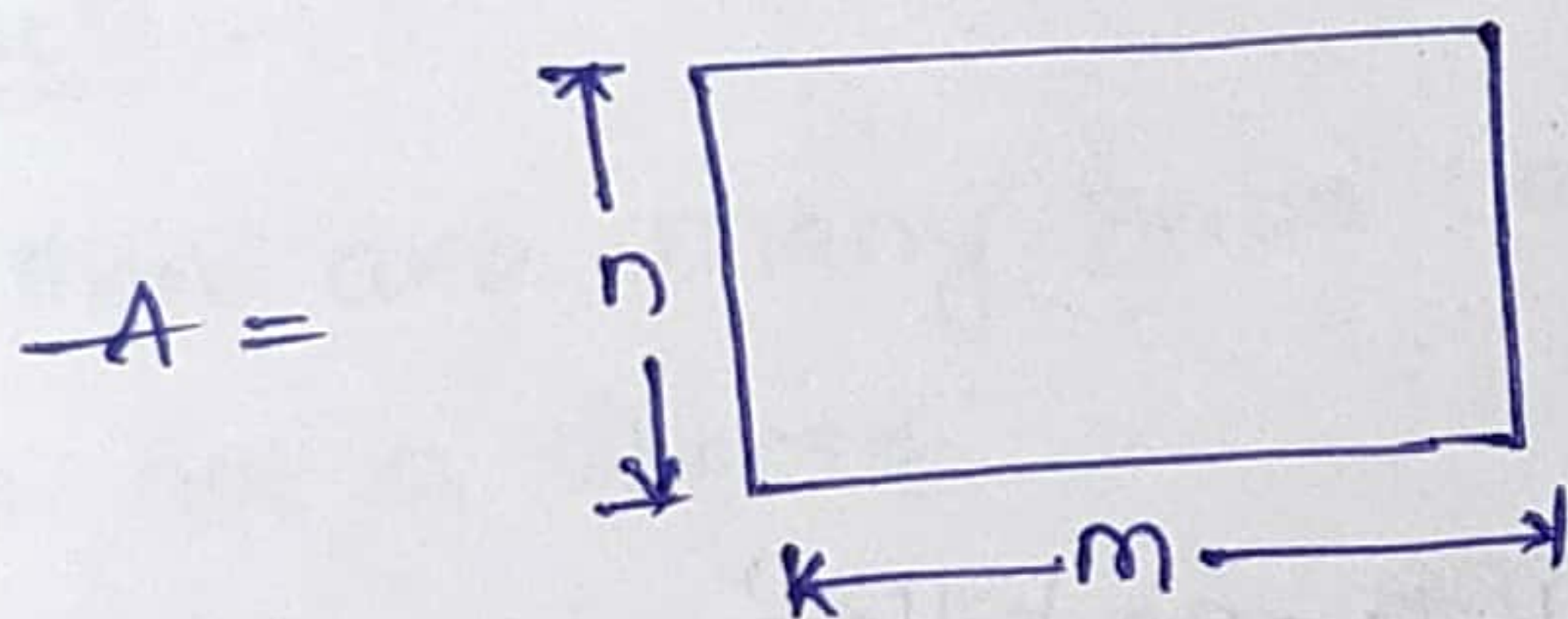
row no, col no, Value \rightarrow 3 * No of non zero cells

= 9 Values

Instead of 100 values we are storing only 9 values

\rightarrow Sparsity of a Matrix :-

for a matrix A with m rows and n columns



Total no of cells = ~~nm~~ $n \times m$

k cells have non-zero values rest of them have '0's

$$\text{Sparsity of } A = \frac{k}{n \times m}$$

where $k =$ no of cells with non-zero values

$n \times m = \text{rows} \times \text{col.}$

\rightarrow more sparse a matrix is the more efficient is

sparse matrix representation technique.

20.12 Text preprocessing:-

Code:-

→ get set of Stop words

⇒ $x = \text{Set}(\text{Stopwords}, \text{words}('english'))$

⇒ To Initialize Snowball Stemmer

$y = \text{nltk.Stem.SnowballStemmer}('english')$

regular expressions are ~~the~~ to recognize patterns in the sentences or words

Ex:- identifying 'that'

'you can call that kite as killer'

Ex- abaa identifying

'abaaaaaabaaaaabaaaa'

20.13 Bi-grams and n-grams

We begin analysis by getting the frequency distribution of the words

~~for~~ $x = \text{nltk.FreqDist}(\text{all_positive_words})$
 $y = \text{nltk.FreqDist}(\text{all_negative_words})$

$\text{print}('Most common^{+ve} words.', x.\text{most_common}(20))$
" " " negative " " $y.\text{most_common}(20)$

Observation

- In this if we see if we perform like this that the most common positive and negative words overlap for eg. 'like' could be used as 'not like' etc, so it is good idea to consider pairs of consequent words (bi-grams) or q sequence of n .

Note

- If there are many zeros values in your matrix it is not a sparse.
- And if there are many non-zero values in your matrix it does not refer as dense matrix.
- Sparse matrix or dense matrix is a way of storing your matrix
- when to store a matrix as a sparse matrix & when to store a matrix as a dense matrix.
- Given a dense matrix

$$\begin{bmatrix} 7 & 0 & 0 \\ 4 & 5 & 0 \\ 0 & 0 & 6 \end{bmatrix}_{3 \times 3}$$

whereas in dense matrix only 9 values got stored.

Sparse representation
(row, column, value)

(0, 0, 7)

(1, 0, 4)

(1, 1, 5)

(2, 2, 6)

} 4 values

• If your matrix is having many non zero values then there will not be much differences if you store it in either a sparse matrix form or in dense form.
→ Use sparse matrix format mainly when your matrix has several zero values.

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$