

## Lec 1

NLP - Natural Language Processing

Natural lang - related by humans

Software controls Hardware

Natural lang - defines how humans perceive things

the ball smashed through the table,

because it was glass

↓  
prior knowledge  
say table is glass

steel?  
10<sup>th</sup> hard is it ball / table?

CV uses pattern recognition

But NLP is diff.

Images are same, lang diff (Korean, Chinese, English)

## Text Representation & Classification

### 1. NLP Fundamentals

→ Preprocessing techniques

### 2. N-gram models

→ Smoothing techniques

(HW)

### 3. ML models / techniques (to classify texts)

→ Naive Bayes

→ Sentiment Classification

→ Logistic Regression

### 4. Vector Semantics, Embeddings, Neural language models Numerically represent lang.

## Natural Language Understanding

(Parts of Speech)

### 1. POS Tagging,

Sequence Processing;  
RNNs

→ time series / sequence  
data

→ markov model

### 2. Encoders - Decoders ,

Attention ,

Contextual Embeddings

### 3. Constituency Parsing

TOP

Advanced way of  
representing lang.

### 4. Dependency Parsing ,

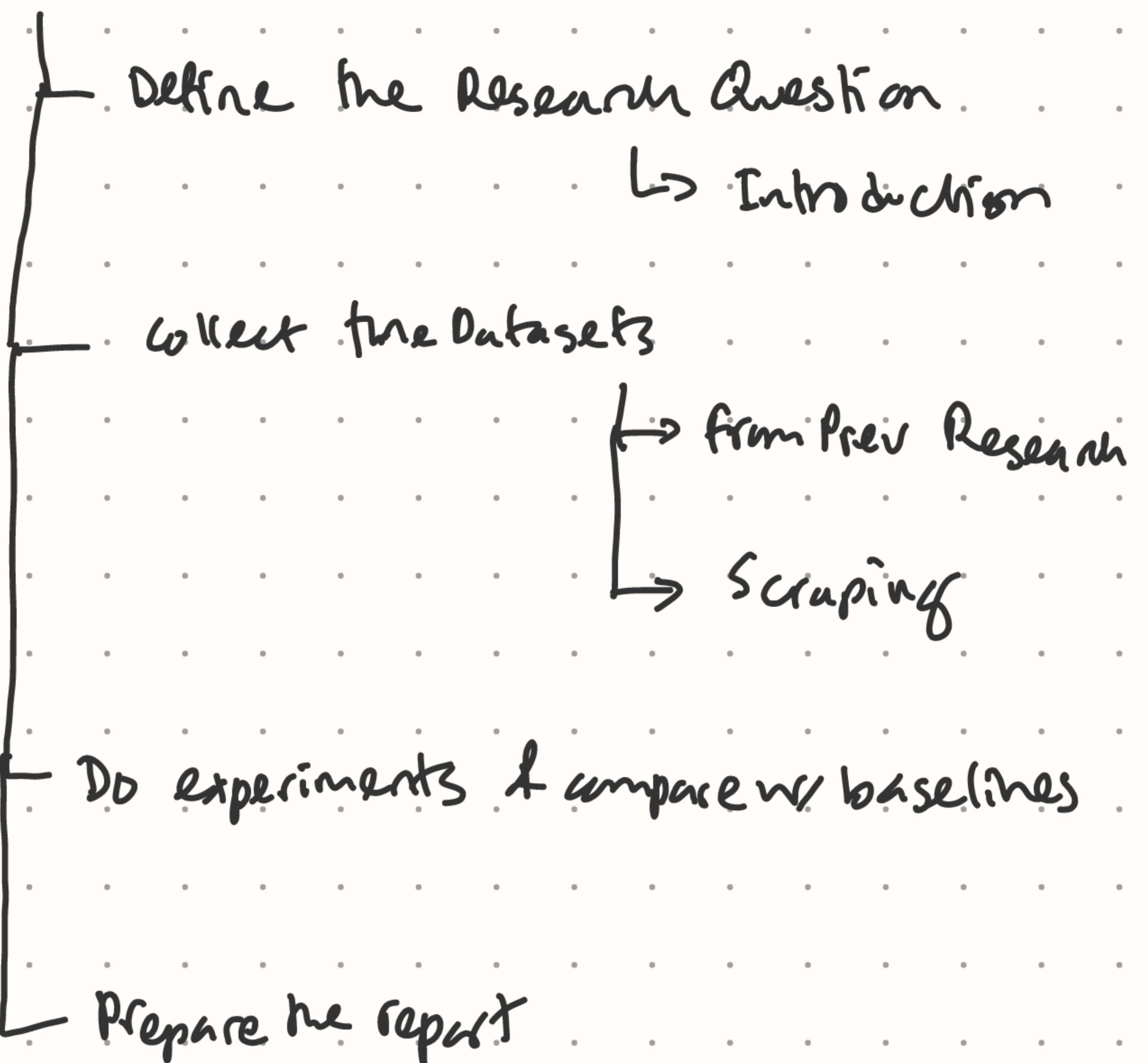
Logical representation

## Discourse & Dialogue :

### 1. Information Extraction, word sense, wordnet

### 2. Question-Answering, Dialogue systems, Chatbots

## Project



Zipt's law : frequency  $\propto \frac{1}{\text{rank}}$   
↓ word → most used ranking (the in 1)

Properties of language

Regular expression

Text Normalization

Edit Distance

3000 words will cover  
95% of everyday  
writing - Robert Charles

Synsets

Stop words : removing them, meaning remains same

Why regular expression ?

(

String search in Google doc - Primitive string result X

b              n : String search ✓ → Using regular exp.  
↓              ↓  
start          end

re python library : regular exp.

>> import re

re.findall(regex, text) → find bubblegum

re.sub(regex, repl, text) → Replace (search substituting)

$s = "Jack is a boy"$

Index of a

Ja/\_a ]

whitespace

re.search("( ) | \s)a", s)

↳ Ja/\_a

re.search('a', s).span() → Return index of 'a'  
(only 1 instance)

.findall + search → Return all indexes of 'a'  
(all instances)

[ { "indices": m.span() } for m in (re.finditer('a', s)) ]

\d → any digit

\d+

metacharacters → For this reason regular exp is used  
instead of string search

$s = \text{"Price? US or 55.33?"}$

To specify any digit using regular exp  $\rightarrow \backslash d$

grouping  $( )$

$\backslash d + (\backslash . \backslash d+)?$

digit      More      special      other      decimal might/  
occurrence      char      occurrences      might not  
                    (.)      of digit      be there  
                    (at)

regex101.com

$\rightarrow$  Find all numbers from string  $s$

rc. findall("(\d+(\.\d+)?", s)  $\rightarrow$  US 55.33

$s = \text{"Price? The price of } \underline{1} \text{ turkey is } \underline{30} \text{ \$."}$

Many product prices  $\rightarrow \$$

$(\backslash d + (\backslash . \backslash d+)? \backslash \$)$

Positive

Lookahead

$(\$ 30)$

Negative lookahead

does not end w/ \$  
(unit)

Positive  
lookbehind

$(\$ 30)$

Negative  
lookbehind

Negative lookahead  $\rightarrow$  Do not want str that  
ends w/ sth (\$)

## Regular exp. (cont.)

1 turkey costs 30 \$ → positive lookahead

\d+ (?= \\$)



positive lookahead

(30)

Is there a digit  
followed by a \$

\d+ (? != \\$) → negative lookahead



(1)

? <= \\$) \d+ → positive lookbehind  
pos/none

(? < ! \\$) \d+ → negative lookbehind

## Text Normalization

(

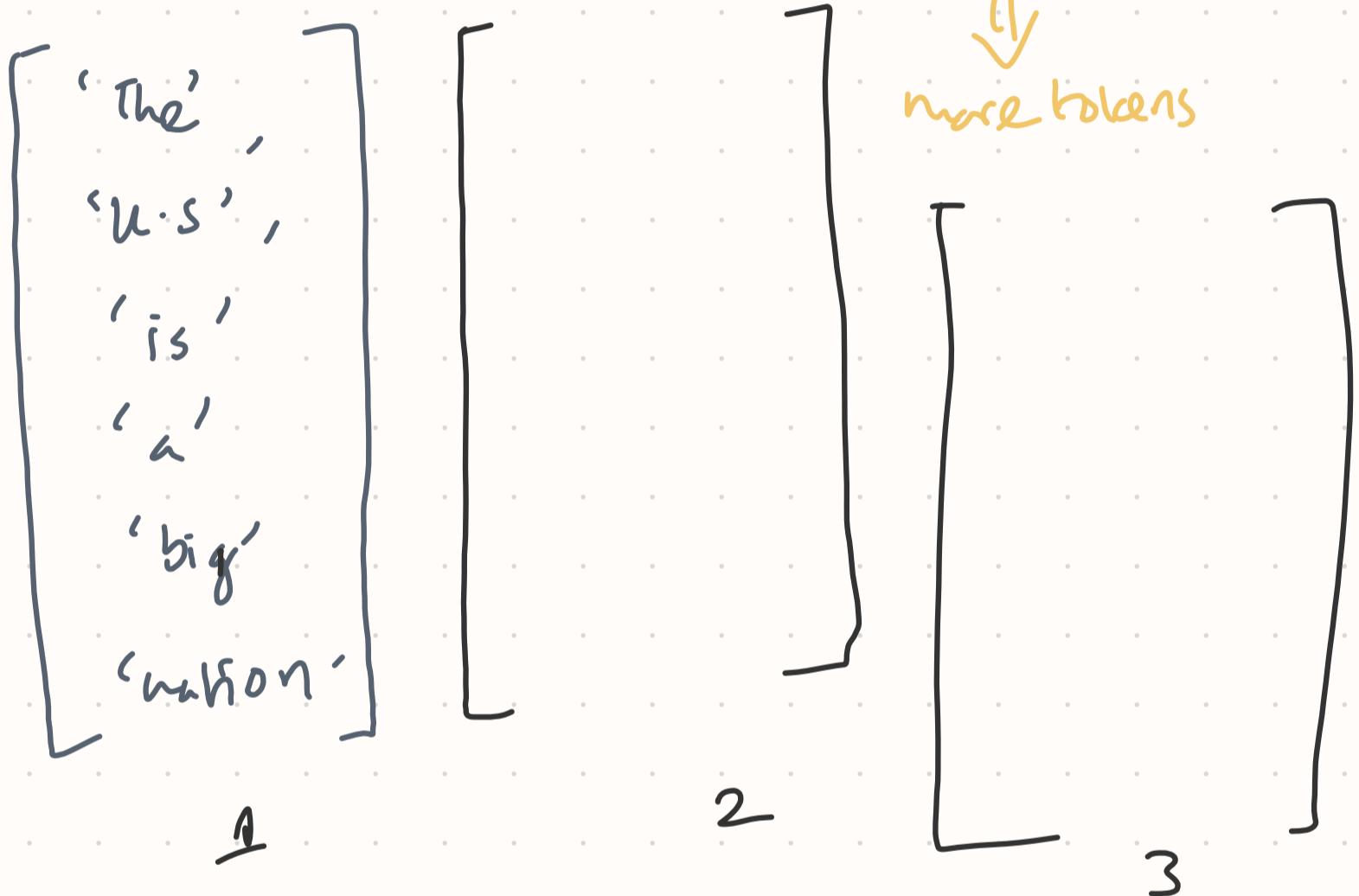
A set of collective tasks

Primary pre-processing tasks before feeding  
to model

1. Tokenization
2. Normalizing word formats
3. Segmenting sentences

- 1 "The U.S. is a big nation."
- 2 The U.S.A poster-print costs 30\$.
- 3 The driving speed unit here is m.p.h, not km."

Embedded vectors



Tokenization : segmenting a running/text into words  
(tokens)

Represent tokens(words) as numbers

3 diff size vectors

## Rule based tokenization — prev.

Each sentence end w/ . → for 1st sentence you get  
3 diff sentences

abbreviations ]  
periods ] .

white space split →

prob ↴

soln ↴ Penn Treebank standard → not all . are sentence  
(PTS) breakers

post-start-point — 1 single token

'30', '£'

PTS + regex

Based on tokenization, it depends how effective  
embedding will be.

Before seq/transformer based models — RNN earlier

bert-based model - embedded vector

Power of transformer / pretraining

Knows

which words are similar

King → queen → Cluster 1

student → teacher → Cluster 2

country → patriot

} cluster of tokens

corpus

'nation'

: similar vector value like

country

(embedding vector)

)

RNN

doesn't have this  
info  
only context

bert tokenizes / bert embedding

better representation of words ✓ → do training parameters

## automatic token parsing techniques : (Tokenization)

1. Byte Pair Encoding (BPE)

2. unigram language modeling (ULM)

3. wordpiece →  $b_{\text{er}}/X$

[sentence piece : BPE + ULM]

Corpus → { low, new, never, wider, lowest }

2 parts :

{ token learner  
token segmenter }

PTs fails to understand what to do w/ unknown words.

BPE : Corpus - 18 tokens

Divide all words based on whitespace and word frequency

Count

Corpus

all unique chars

↖ vocabulary

5

low -

↘ whitespace

- , b, e, i, l,

n, o, s, t,

k, w, e, l,

er-, ne,

new, lo,

low, never,

low -

2

lowest -

6

never -

3

wider -

2

new -

morphemes

Task: create a vocab

List most frequent pair of symbols together

$$ef = 6 + 3 = 9$$

n e w  
w i d

n e w er -

w i d er -

morphemes; smallest unit of

BPE → lower  $\xrightarrow{\text{low}}$  er } 2 tokens

(learns on unsupervised  
New word:  
lower rate

low, er -

present in teacher vocab  
low er - diff token  
(merged)

Subwording: Dividing the tokens into smaller parts

Subword = morpheme

ne w er -  
ne w -

BPE : K parameter  $\Rightarrow$  how many cycles you want to run

BPE (corpus, k)  
a

Normalizing :

am is are — same word containing same meaning

case folding :

(

lower case  
everything

Woodhouse  
woodhouse

} Same value for  
embedding vector

→ Info retrieval ✓

Text/Sentiment classification

& machine translation

(  
not doing lower casing

US → country

us → we

Lemmaization : Determine which words are the same root

addition  
additional } same root

I have some additional information.

I can do some addition on information.

same



replace each word w/ roots → Lemmatization ✓  
 using pre-existing dictionary  
 singer, sang, sung (better)

Do Lemmatization before feeding corpus to language model.

Stemming : Remove affixes from words

X additional

Changed actual meaning

(dang → doe) → error prone (stemming)

(dang → do) → not error prone (lemmatization)

happiness → happy  
 ^ ^  
 h h

Lemmatization over stemming ✓

# Term Paper

## Performance Review Research

Language model optimized with least amount of data

(child language acquisition)

Bert, Roberta - in traditional sense, they are not LLMs

(at least abillim parameter) ↩

ChatGPT - General task solvers

Benchmark of ChatGPT in summarization datasets

General task solvers are better than specific task solvers

Performance of LLM in low research language

only good at previously seen data

task contamination

SOTA

## Class Notes

Edit distance : way to quantify string similarity

spelling mistake: graffle

giraffe	graf	grail	spell checking
---------	------	-------	----------------

Harvard President Claudine | go reference

↓  
if two strings  
contain same  
entity

Define cost — LEVENSHTEIN

The apple was red

1. The apple was ~~big~~ <sup>2\$</sup> <sup>2\$</sup> ~~red~~ <sup>2\$</sup> → <sup>cost</sup> 4\$ i insertion - 2\$ d deletion - 2\$

2. The ~~fish~~ <sup>2\$</sup> was ~~blue~~ <sup>2\$</sup> <sup>red</sup> → 8\$ s substitution  
apple <sup>2\$</sup>

1st was more similar to the string.

↳ minimum Edit distance (check test similarity)

dss-i s      cost=8

INTENTION → EXECUTION

↓↓↓↓↓  
X E Y E C U T I O N

min Edit Distance?

Can start from  
anywhere

$$\begin{aligned}
 i &= 1 \\
 d &= 1 \\
 s &= d + i
 \end{aligned}$$

use DP (Dynamic Programming)

- search problem

Find most cost efficient approach

## Edit Distance Algo

- ✓ 1. Levenshtein
- ✓ 2. Hamming Distance
- 3. Jaro-winkler

## Token-based

- ✓ 1. Cosine-Similarity (Similarity in text/may)
- ✓ 2. Jaccard Index
- 3. Tversky Index

## Sequence-based (programming)

- 1. LCS (longest common subsequence)
- 2. Ratcliff-Obershelp Similarity

## Cosine-Similarity

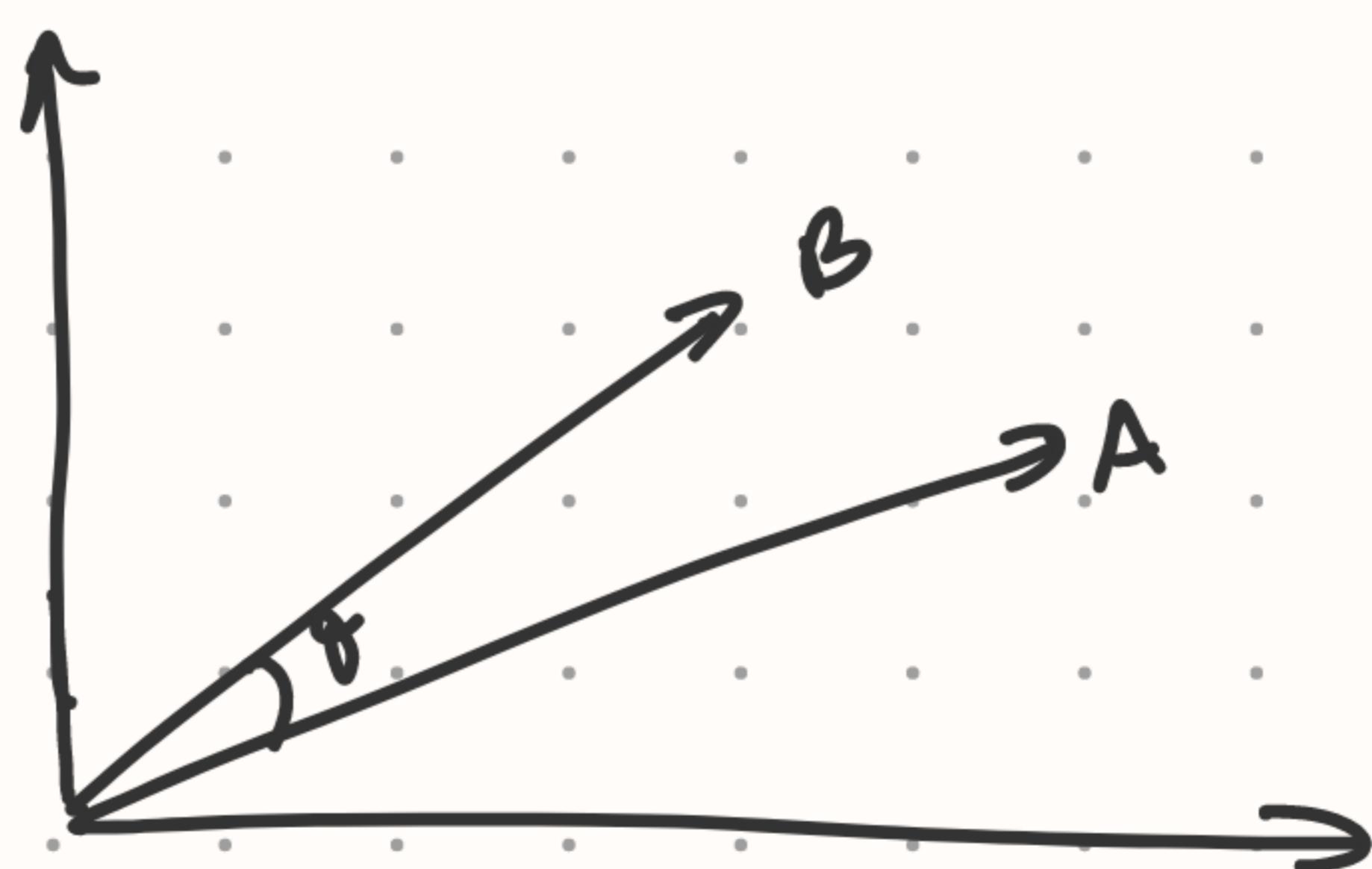
converts the whole corpus into same dim. vectors and projects it.

similar sentence  $\rightarrow$  similar vectors

more similar  $\rightarrow$  less distance

2 same sentence  $\rightarrow$  project on each other

$$\theta = 0$$



Cosine Similarity ( $\text{doc1}, \text{doc2}$ )

Document 1: The cat in the hat.

Document 2: A black cat in a black hat,

DOC → vectors : Find all unique tokens from two docs.

	The	cat	in	hat	a	black	
DOC1	2	1	1	1	0	0	Frequency
DOC2	0	1	1	1	2	2	

$$A = \begin{bmatrix} 2 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 2 \\ 2 \end{bmatrix}$$

coseine similarity (doc1, doc2)

$$= \frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{2 \cdot 0 + 1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1 + 0 \cdot 2 + 0 \cdot 2}{\sqrt{7} \cdot \sqrt{11}}$$

$$= 0.342$$

similar → close to 1

diff. → close to -1

⊥ vectors → close to 0  
(no similarity / dissimilarity)

— completely diff set of words  
(> frequency)

There should be a ( $\neq$ ) value theoretically.

### Insight

- 1) The patient died before the doc came.
- 2) The doc came after the patient died.

Cosine similarity cannot detect semantic meaning of sentence.

### (Cons)

Computing time is very long (Pros)

- sensitive to sentence length.  
long vector (not necessary) | cons

sparsity - having a lot of 0s in vectors / matrix  
cosine similarity creates sparse vectors that will dominate the entire value.

Jaccard Index is more sensitive to context in finding text similarity.

$$J.F = \frac{A \cap B}{A \cup B}$$

$$\text{Jaccard similarity } (\text{doc1}, \text{doc2}) = \frac{\text{doc1} \cap \text{doc2}}{\text{doc1} \cup \text{doc2}}$$

if 2 docs talking about same context

docs  $\rightarrow$  tokens

$$J.S. = \frac{\{\text{'The', 'cat', 'in', 'the', 'hat'}\} \cap \{\text{'A', 'black', 'cat', 'in', 'a', 'black', 'hat'}\}}{\text{U}}$$

$$= \frac{3}{8}$$

$$= \frac{3}{6} \approx 50\% \text{ similarity.}$$

Imp. of tokenization/ lemmatization.

throw, threw, thrown  $\rightarrow$  w/o lemmatization, consider  
(verb forms) the m diff.

## N-gram language models

Lang models that assign probabilities to N-grams

Please turn in your homework.

### word bigrams (2)

Please turn  
turn in  
in your  
your homework

### word trigrams (3)

Please turn in  
turn in your  
in your homework

### character bigrams

pl, le, ea ...

### character trigrams

ple, lea, eas, ase ...

### Skip grams: Nm continuous

### Skip bigrams

please in  
turn your  
in homework

} skip 1 element

### Language models : Probabilistic model

Assign probabilities to lang values  
(sequence of value)

chart least probabilities of words

Assign probabilities to the possible ans.

Highest probability ans chosen ✓

Bigram model : Trained w/ simple bigrams  
Probability of next word to occur

$P(w_n | w_{n-1}) \xrightarrow{\text{bigram}} \text{Markov model}$  (Not look too far into past)  
Ans closest to you.

$P(\text{turn} | \text{please}) = ? \xrightarrow{\text{Probability that turn will come after please}}$

$P(w_n | w_{n-1}, w_{n-2}) \xrightarrow{\text{trigram markov model}}$

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(w_{n-1})} \xrightarrow{\substack{\text{prev word occurrence} \\ w/ this word}}$$

conditional probability eqn

$$P(a|b) = P$$

$\hookrightarrow$  occurrence of prev word

$\hookrightarrow$  MLE (maximum likelihood estimation)

$\langle s \rangle$  I am. Abrar  $\langle \backslash s \rangle$

$\langle s \rangle$  Abrar I am  $\langle \backslash s \rangle$

$\langle s \rangle$  I do not like eggs.  $\langle \backslash s \rangle$

$\swarrow$   
start of string

need to define  $\langle s \rangle$  and  $\langle \backslash s \rangle$  to find  
occurrences.

$$P(I | \langle s \rangle) = \frac{C(\langle s \rangle^I)}{C(\langle s \rangle)} = \frac{2}{3}$$

$$P(\text{am} | \text{I}) = \frac{C(\text{I am})}{C(\text{I})} = \frac{2}{3}$$

$$P(<\text{s}> | \text{Abrar}) = \frac{C(\text{Abrar } <\text{s}>)}{C(\text{Abrar})} = \frac{1}{2}$$

↓

Probability that Abrar is the end of string:

Probability that string ends w/ Abrar



$$P() = \frac{C(\text{Abrar } <\text{s}>)}{C(<\text{s}>)}$$

Long models learn from data that you provide.

Longer spans → sparse data → capture more context  
 $(N \uparrow)$       (less matches)

shorter spans → many matches → less context  
 $(N \downarrow)$

[large corpus → trigram / 5-gram as features]

Feature extraction → NLP

count vectorizer → trigrams, 5-grams  
Bag of words

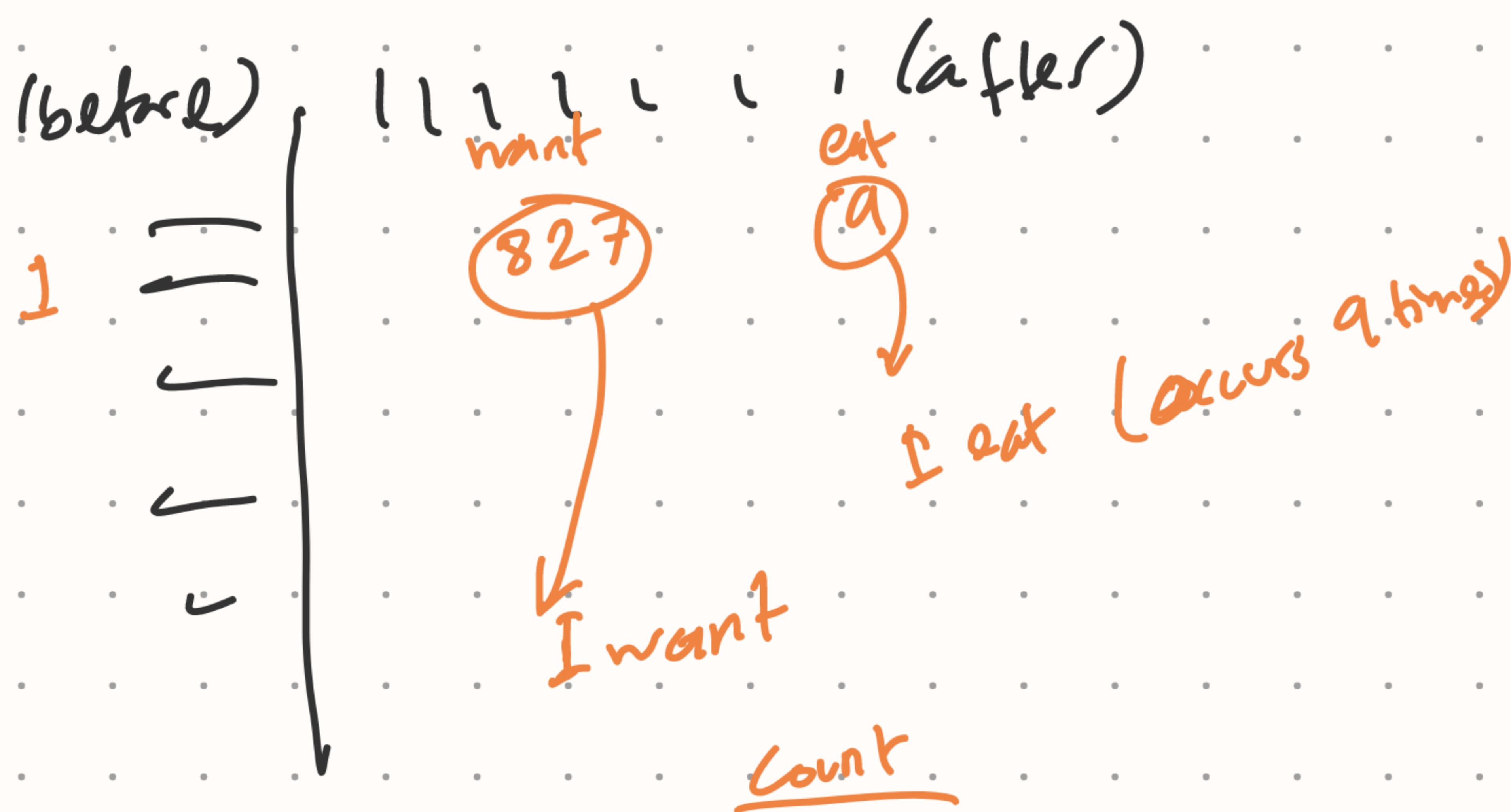
use character/next N-gram combinedly  
word

fasttext embedding → N-gram words + chars  
↓  
neural networks  
(not looking any features)

Rule based vs feature based

↓                    ↓  
fixed cases        all cases

Berkeley Restaurant Project - Dataset



Probability       $P(v_{want}|I) = \frac{c(I_{want})}{c(I)} =$  0

I want

<s> I want chinese food <\s>

$$P(i | \langle s \rangle) = 0.25$$

$$P(\text{food} | \langle 15 \rangle) = 0.68$$

bigrum  
probabilis  
given

$$= p(i \langle s \rangle) p(\text{want} | i) p(\text{Chinese want}) p(\text{food}) \cancel{p(\text{Chinese})}$$
$$p(\langle s \rangle | \text{food})$$

$$= 0.25 \times 0.33 \times 0.0065 \times 0.52 \times 0.68$$

$$= 0.0002$$

$$\log(0.25) + \log(0.33) + \log(0.0065) \\ + \log(0.52) + \log(0.68) \\ = \{ \dots \} \text{ large number}$$

$$\log(P_1 \times P_2 \times P_3) = \log P_1 + \log P_2 + \log P_3$$

$$\Rightarrow P_1 \times P_2 \times P_3 = \exp \left( - \frac{E}{kT} \right)$$

working w/ small number(s)  $\Rightarrow$  log

## Evaluating language models (LM)

- Evaluate the system which embeds the LM = **Extrinsic Evaluation**  
(expensive)
  - bug fixing expensive
- **Intrinsic Evaluation**: simulate the environment where users use LM
  - (use test set to evaluate)
  - train model using training corpus

\* one lm fits the test set better than the other one.

How to define this performance?

(Higher probability to the words

Assign probabilities better  $\rightarrow$  better LM

(Probability of occurrence of a sentence in a corpus).

Dev set / validation set - optimize hyperparameter

unseen test set - test performance

abnormally  
Assign higher probabilities to the seen data.  
but poor on unseen data.

↑  
**Overfitting** : Higher bias towards training set  
highly unlikely to get 99% accuracy (memorize rather than generalize)

**Underfitting** : performs poor on the training set as well as on seen data

Model cannot find any generalized pattern in training data.

- cannot assign probabilities to seen data

### Large Dataset

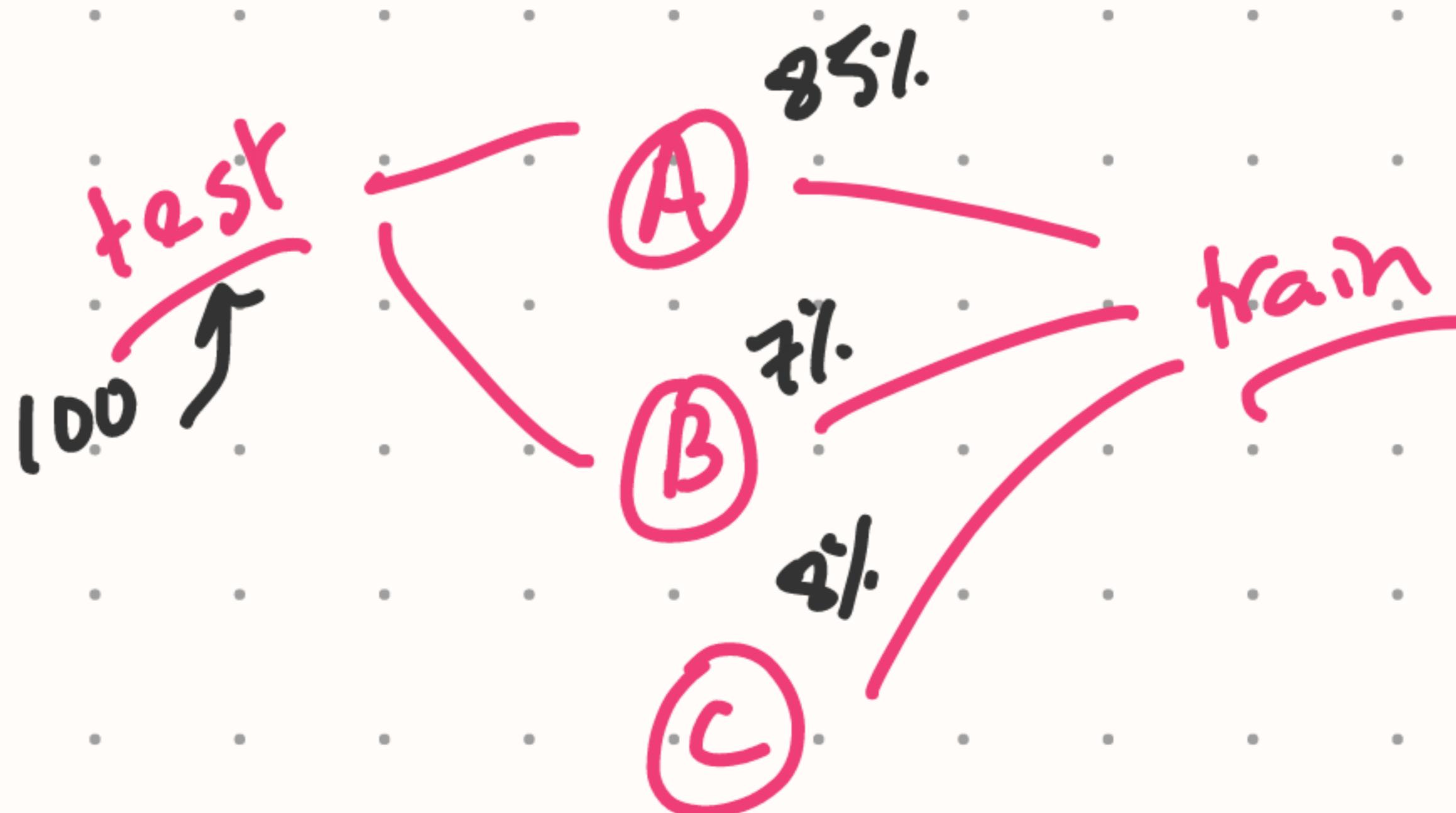
Training: dev: test

80 : 10 : 10

70 : 15 : 15

60 : 20 : 20 → why not this?

low amount data — unrepresentative of actual scenario  
1000 samples



maintain ratio:  
80 - training  
7 - dev  
5 - test

biased test set — cannot evaluate LM

train-test-split

## Evaluation matrix:

For general task:

Next word / sentence prediction — doctor away  
doctor far away

Summarization

Paraphrasing

Dif eval metrics for dif tasks

**perplexity**: To assess performance of lm in next word prediction

$$PPL(w) = \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

$$= 2^{-\frac{1}{N} \sum_{i=1}^N \log_2 P(x_i)}$$

$N$  = length of any test sample

PPL = logarithmic exponential of loss function  
(cross entropy)

"The cat is on the mat" → Test sample

$$P(\text{"the"}) = 0.1$$

$$P(\text{"cat"}) = 0.2$$

$$P(\text{"is"}) = 0.15$$

$$P(\text{"on"}) = 0.1$$

$$P(\text{"me"}) = 0.2$$

$$P(\text{"mat"}) = 0.25$$



Learned from  
training corpus

N=6 words/tokens in the  
test set

$$\text{PPL} = 2^{-\frac{1}{6} \left[ \log_2(0.1) + \log_2(0.2) + \log_2(0.15) + \log_2(0.1) + \log_2(0.2) + \log_2(0.25) \right]}$$

$$= 2^{-\frac{1}{6} \times (-15.02)}$$

$$= 2^{.25}$$

$$\approx 5.66$$

6 branches in training corpus when  
these probabilities will be found  
more branches → more  
non-deterministic  
→ more difficult

PPL → weighted avg branching factor of a LM

PPL ↓ ✓ (less branches preferred)

## \* Assignment

task + documentation

edit distance - brief explanation / defn

Notebook

Documentation - hard copy

Islamophobia bias  
Selection bias  
Ukraine vs Palestine

## \* Don't break accuracy!

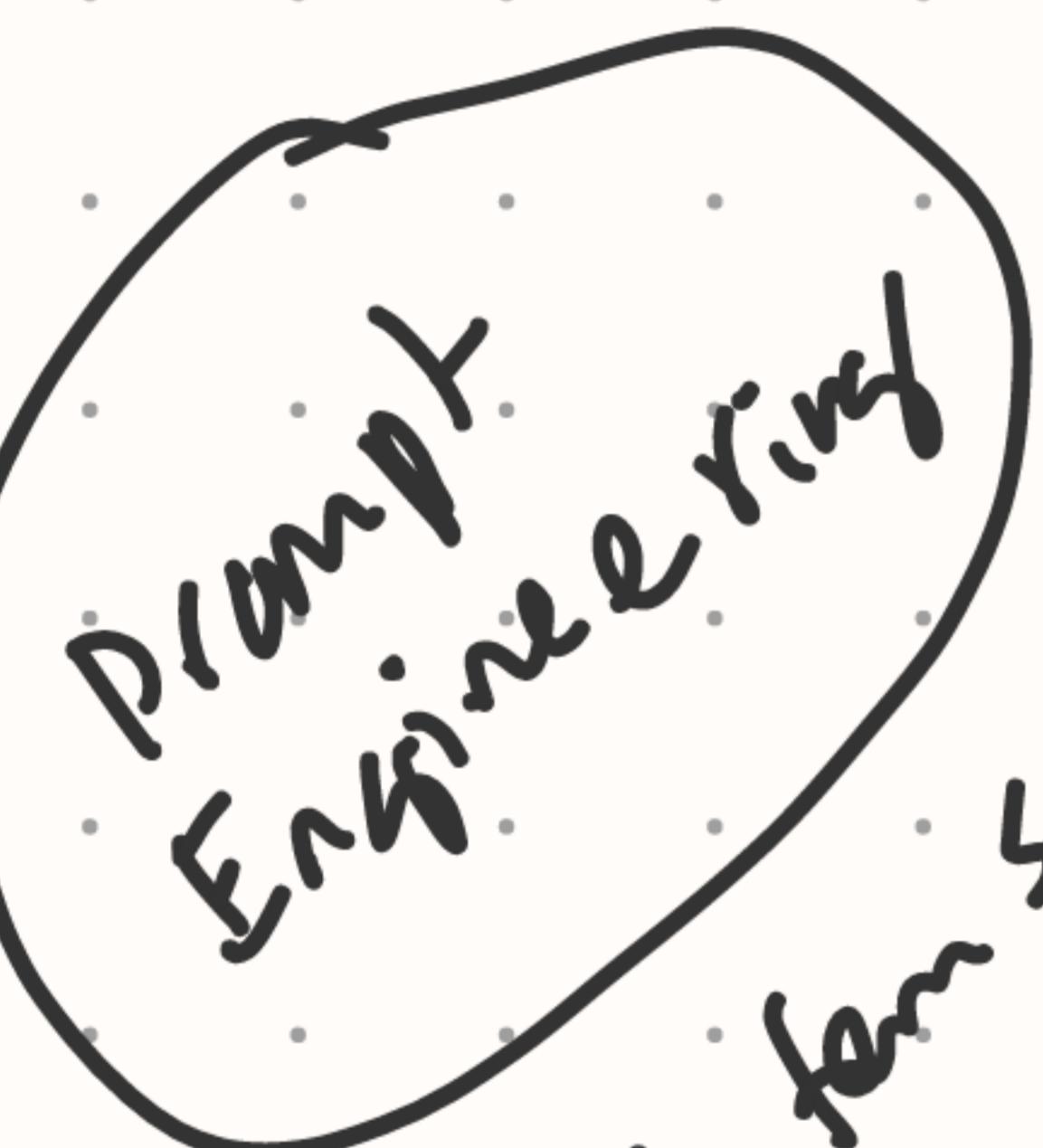
interesting prob.

transliteration  
express different intent  
alphabet  
~~brayligh~~

Sentiment analysis - classification

context / culture

+ the sentiment in Britain will differ in BO  
negative sentence ism good → ism regard  
(not same in Bangla)



→ parsing info from samples  
extracting info from samples  
how do re-pattern work in Bangla compared to English?

Linguistic pattern, not model

Design experiment

- mental health datasets - sample labels

Um

depressing

- kids

- elderly people Um  
lonely and  
DO provide support --  
mental health

Perplexity  $\rightarrow$  Probability



left  
represents form  
of variability  
Project idea by



Large N-gram to model lang., what are the problems?

- less matches  $\rightarrow$  parsing

(many rows 0 in probability distribution matrix)

Berkeley  $\rightarrow$  Probability of a word given another word

Restaurant

count

Project Corpus

Probabilities = bigram count

unigram count

$$P(i|i) = \frac{\text{Count (ii)}}{\text{Count (i)}} = \frac{5}{2533}$$

Smoothing/Discounting : Take from other probabilities  
and give to the 0 probability

collective probability = 1

Laplace Smoothing : Add 1s to all the 0 counts

$$P = \frac{\text{Count (w, v)} + 1}{\text{All counts of bigrams} + 1}$$

After Laplace smoothing :  $P(i|i) = \frac{6}{253}$

I want food lunch.

$$P = P(I | \langle s \rangle) \times P(\text{want} | I) \times P(\text{food} | \text{want}) \times \\ P(\text{lunch} | \text{food}) \times P(\langle \rangle | \text{lunch})$$

For a particular bigram whose probability becomes 0.

Smoothing  
(Add 1)

Laplace/Add-K (another variant of Laplace)

$$P = \frac{\text{count}(w, v) + k}{\text{all counts of bigrams} + K}$$

try diff K.

For which K, we get max probability

Good for text classification but not lang model

## Backoff

$P(w_n | w_{n-1} w_{n-2}) \rightarrow$  trigram probability

$P(w_n | w_{n-1}) \rightarrow$  bigram probability

If you don't find probability of a trigram, back off to bigram.

----- unigram

$P(\text{lunch} | \text{food}) = 0$  Bigram

↓ backoff

$P(\text{lunch}) =$  unigram

Interpolation : Advanced technique of backoff

$$P(w_n | w_{n-1} w_{n-2}) = \lambda_1 \cdot P(w_n) + \lambda_2 \cdot (\text{Bigram Prob.}) + \lambda_3 \cdot (\text{Trigram Prob.})$$

$$\sum_{i=1}^n \lambda_i = 1$$

Kneser-Ney smoothing: sophisticated way

I can't see without my reading

✓ glasses

✗ Francisco (error based on prob. count)



if you only depend on prob. count

Though Francisco is very frequent, only appropriate after

San:

San Francisco ✓

$P_{\text{continuation}}$  = the no. of diff contexts the word has appeared in.

$P_{\text{W continuation}} \propto \{ v : \text{unique } v \text{ before w} \}$

$P_{\text{W continuation}} \propto \{ v : C(v, w) > 0 \}$

glasses - 2

Francisco - 1

San Francisco - 1

$$\frac{P_C(\text{glasses})}{P_C(\text{Francisco})} = \frac{2}{1}$$

occurs in  
1 context & only glasses ✓

(higher prob of occurring in diff contexts)

Preprocessing  
n-gram LM

LM calculates probability distribution.

Stopwords — high probability?



Shouldn't be only  
depend on word count

Text / sentiment classification

Linear ML models — Naive Bayes, Logistic Regression

Evaluation metrics

## Sentiment Classification

(

sentiment extraction

what particular document is about  
positive/negative

spam detection : email is spam/not

previously - hand rules (if else ...)

Soln: supervised ml

↳ input data w/ output labels  
ML learn from these samples

Classifications:

→ binary (+/-, spam/non-spam)

→ multinomial (series of emotion - anger, joy, pride)  
↳ discrete

→ ordinal (child, youth, teen, elder)  
↳ series

→ multiclass (happy, angry)

↳ n. no. of classes

classify k out of n classes

from multiple classes

→ multilevel (happy, angry)

multiple level for a particular element

## Naive Bayes

$$P(x|y) = \frac{P(x) \times P(y|x)}{P(y)}$$

Naive assumption about how features interact

Bag of Words

↓  
frequency based

Does not consider order of the elements

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c|d)$$

Calculate prob of all classes  $\rightarrow$  max posterior prob. choose

$$\hat{c} = \operatorname{argmax}_{c \in C} \frac{P(c) \times P(d|c)}{P(d)}$$

win a free vacation  
document

document

$P(d) = 1$

have aduc.

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c) \times P(d|c)$$

Prior probability

likelihood estimation  
Prob. of doc give that  
it's inside that class

100 training samples - 50 spam, 50 non-spam

$$P(\text{spam}) = \frac{1}{2}$$

$$P(\text{non-spam}) = \frac{1}{2}$$

Even in real scenario, model assumes get  $\frac{1}{2}$  always.

$$P(d|c) = P(f_1 f_2 f_3 \dots f_n | c)$$

↳

All possible combinations of features  
inside the class.

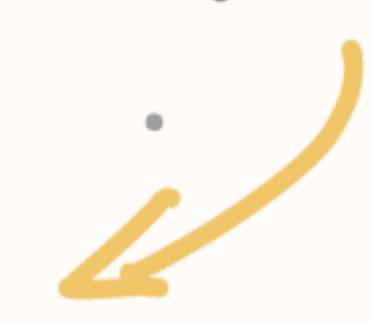
Expensive - need to simplify

Naive simplification

All features are independent

$$P(d|c) = P(f_1 f_2 f_3 \dots f_n | c) = P(f_1 | c) \cdot P(f_2 | c) \cdot \dots \cdot P(f_n | c)$$

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} \ P(c) \prod_{f \in F} P(f | c) = \underset{c \in C}{\operatorname{argmax}} \ P(c) \prod_{i \in \text{position}} P(w_i | c)$$



go to every position in doc  
check the probability  
of that doc in that class

features - words w/~~fix~~ sentiments

$$= \underset{c \in C}{\operatorname{argmax}} \ \log P(c) + \sum_{i \in \text{pos}} \log (w_i | c)$$

Naive Bayes eq<sup>n</sup> for text

generative classifier

add features (linear classifiers)

Logistic Regression — Discriminative classifier  
classification (not generation)

Task: If incoming email is spam/not  
sentiment classification:

prob. of the class that we're looking for

$$P(c) = \frac{N_c}{N_{\text{doc}}}$$

$N_{\text{doc}}$  = # training samples

$N_c$  = # samples in each class

$$P_{\text{spam}} = \frac{2}{5} \quad \left| \begin{array}{l} \text{classes = spam, non-spam} \\ \text{total # doc} \end{array} \right.$$

$$P_{\text{Nonspam}} = \frac{3}{5}$$

/ particular word

Likelihood  $P(w_i | c) = \frac{\text{Count}(w_i, c)}{\sum_{w \in V} \text{Count}(w, c)}$

(  
all words in  
that class)

vocabulary

word you're looking for is not inside spam class  
but in non-spam class

Soln Laplace smoothing

$$P(w_i | c) = \frac{\text{Count}(w_i, c) + 1}{\sum_{w \in V} (\text{Count}(w, c) + 1)}$$

$$= \frac{\text{Count}(w_i, c) + 1}{\sum_{w \in V} \text{Count}(w, c) + |V|}$$

all words in vocab

## Spam Detection using Naive Bayes Classifier

Feature Extraction using bag of words

test sample,  $s = \text{"Win a free vacation"}$

$$P(s \text{ spam}) = \frac{3}{6}$$

$$P(\text{ham}) = \frac{3}{6}$$

unknown word

(didn't occur in training sample) or, randomly assign a probability

6 doc  
3 spam  
3 ham

ignore (cannot apply Laplace smoothing)

Laplace smoothing  $\rightarrow$  when word is here but particular bigram/trigram not there  
(Add 1)

bigram/trigram not there

Prob. of occurrence of each word in each of the classes:

$$P(\text{win} | \text{spam}) = \frac{1+1}{10+16}$$

$$P(\text{free} | \text{spam}) = \frac{2+1}{10+16}$$

$$P(\text{vacation} | \text{spam}) = \frac{0+1}{10+16}$$

$$P(\text{win} | \text{Ham}) = \frac{0+1}{8+16}$$

$$P(\text{free} | \text{Ham}) = \frac{0+1}{8+16}$$

$$P(\text{vacation} | \text{Ham}) = \frac{1+1}{8+16}$$

16 unique words - vocab (unique ignore duplicate)  $|V|=16$

2 bags of words (1 for spam  
1 for ham)

Count (not unique) - count duplicate words

$$P(s \text{ in spam}) \quad \hat{c} = P(c) \cdot \prod P(w_i | c)$$

$$\hookrightarrow P(\text{spam}) \times P(S|\text{spam})$$

$$= \frac{3}{6} \times \frac{2 \times 3 \times 1}{26^3} \approx 0.000171 \rightarrow \checkmark \text{ (bigger)}$$

$P(s \text{ in ham})$

$$\hookrightarrow P(\text{Ham}) \times P(S|\text{Ham}) = \frac{3}{6} \times \frac{1 \times 1 \times 2}{24^3}$$

$$\approx 0.0000172$$

model classifies  $s$  as spam (Naïve Bayes)

I really like this book

I don't like this book.

positive class

not negative meaning

bag of words

In case of negation, using bag of words as feature extractors cannot handle these cases.

Prepending (Preprocessing)

I don't NOT-like NOT-this NOT-book.

inside negative class

until punctuation

Prepending (preprocessing) should be done on test samples too.

→ Insufficient word in training sample

Pretrained, pre annotated corpus → the, the words

LWIC (Linguistic Inquiry Word Count)

## MPQA lexicon root word

	(0.9)	(0.1)
ture	: admirable, beautiful	~ catastrophe(0.1)
-ude	: catastrophe, envious	

Catastrophe might be in the sample with low probability.  
lets, assume 0.1

spam)) compared to ham

correctly classify ham (less in number in real life)

# 100 emails      unbalanced dataset - only know spam  
(seen more spam)

99%	1%
(spam)	(ham)

Predict everything as spam  
mistake

Evaluating lang. model (Lm) :

Accuracy = 99%  $\rightarrow$  1% classifying ham as spam  
useless Lm      ↴  
risky Evaluation metric

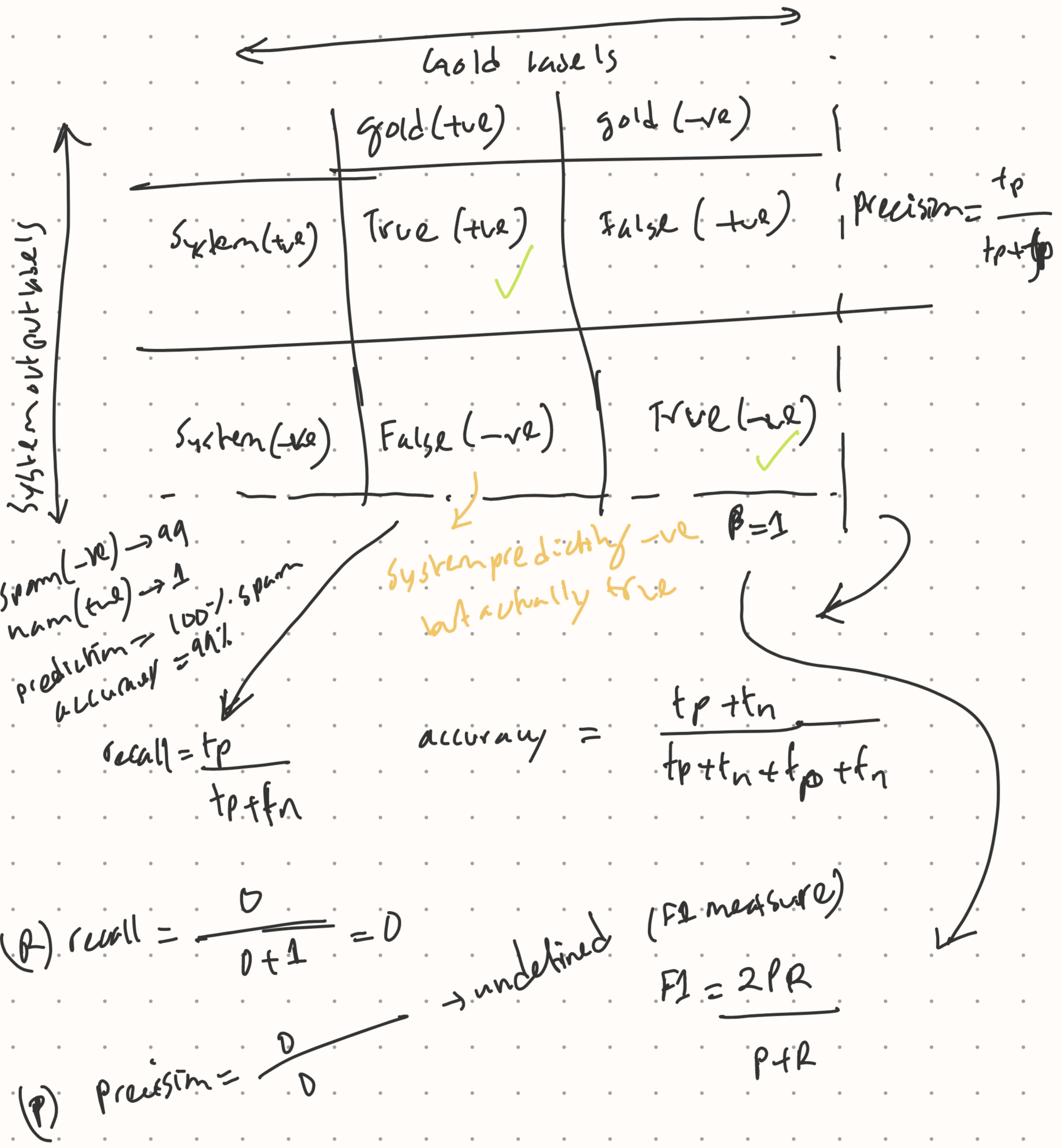
Sohn confusion matrix

(visualization of how an algo is performing  
in diff classes)

Use category for evaluation

spam  
ham

Gold labels / Human Gold labels — in the test doc,  
gold labels are used for evaluation metrics



$$F\text{-measure} = \frac{(\beta^2 + 1) \cdot \text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

$(\beta > 1)$  → more weight to Recall (R)

$(\beta < 1)$  → more weight to Precision (P)

$(\beta = 1)$  → equal weight P and R

## Precision & Recall Calculation:

Confusion matrix for 3 classes:

Urgent }  
spam  
ham }

$P = \text{Precision}$   
 $R = \text{Recall}$

	u	s	h
u	✓ 8	10	1
s	5	60	50
h	3	30	200

$$P_u = \frac{8}{8+10+1}$$

$$R_u = \frac{8}{8+5+3}$$

Performance of whole model?

\* macro-averaging

\* micro-averaging

class-1 : urgent

8	11
8	340

$$Pr_1 = \frac{8}{8+11}$$

class-2 : spam

60	55
40	212

$$Pr_2$$

class-3 : ham

200	33
51	83

$$Pr_3$$

macro-avg precision =  $\frac{Pr_1 + Pr_2 + Pr_3}{3} = 0.60$

↓  
all classes equally imp.

268	99	
99	635	

} from 3 classes  
(3 matrixes)

Micro-avg precision =  $\frac{268}{268 + 99} = 0.73$

↓  
dominating class

Determining span → micro avg

[Naive Bayes - generative classifier  
Logistic Regression - discriminative]

classify cats & dogs:  
**(Naive Bayes)** — generate as much features as possible

$$\hat{P}(c) = P(c) \cdot \underbrace{P(d|c)}_{\text{likelihood}}$$

- given cat, what are the features he eyes might have
- given dog, — — — tail might have

single dog/cat feature by itself

## (Logistic Regression)

↳

calculates  $P(d|c)$  directly

Put more weight one feature

### Model and Prediction

1. Feature Representation → find discriminative feature

2. Classification function

binary / multinomial

(based on need)

softmax

soft void

3. Objective function

loss, objective / log loss function: how your prediction is diff  
from goal

prediction and gold label

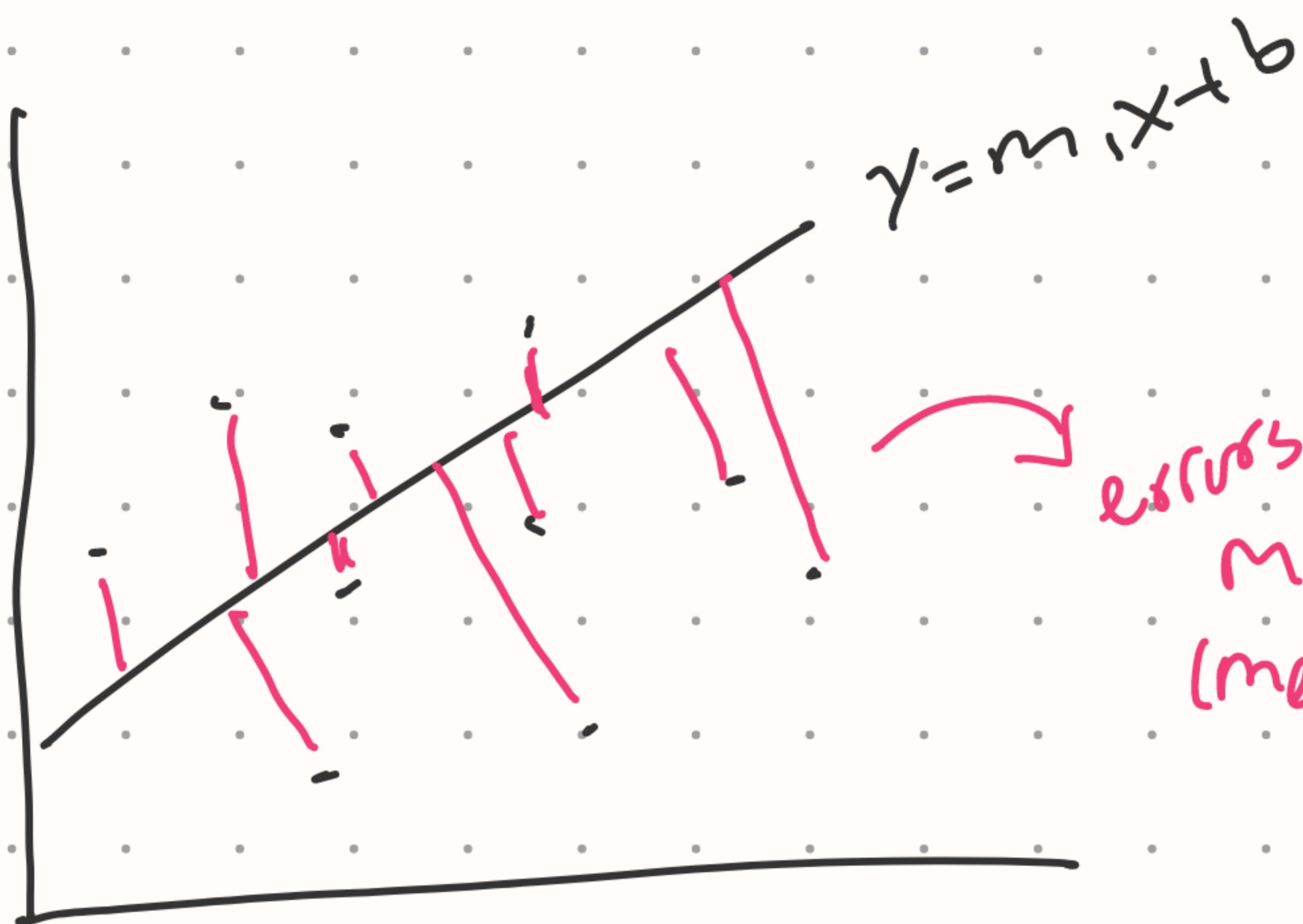
min. distance

4. Optimization function : to optimize objective func.

Gradient Descent

Linear Regr.

st. line (bad algo)



errors  
MSE

(mean square error)

→ L2 loss function

$$= \frac{1}{n} \sum_{i=1}^n (y_i - y_i^{\text{pred}})^2$$

MAE (Mean Absolute Error)

$$= \frac{\sum |y_i - y_i^{\text{pred}}|}{n}$$

L1 loss

How to optimize loss?

↳ m and b value

Gradient Descent : matrix multiplication formula  
Partial Derivative } 2 ways

Random Hill climbing Algo : random m and b  
(global)  
Calculate MSE, MAE  
(like 20 times randomly)

minimizing error → go in that way  
(decreasing loss func.)  
stop when loss func. increasing again

Outliers - skewed data → consider outliers  
Avoid outliers - L2 loss pays big price  
(high penalty)  
use L1 loss

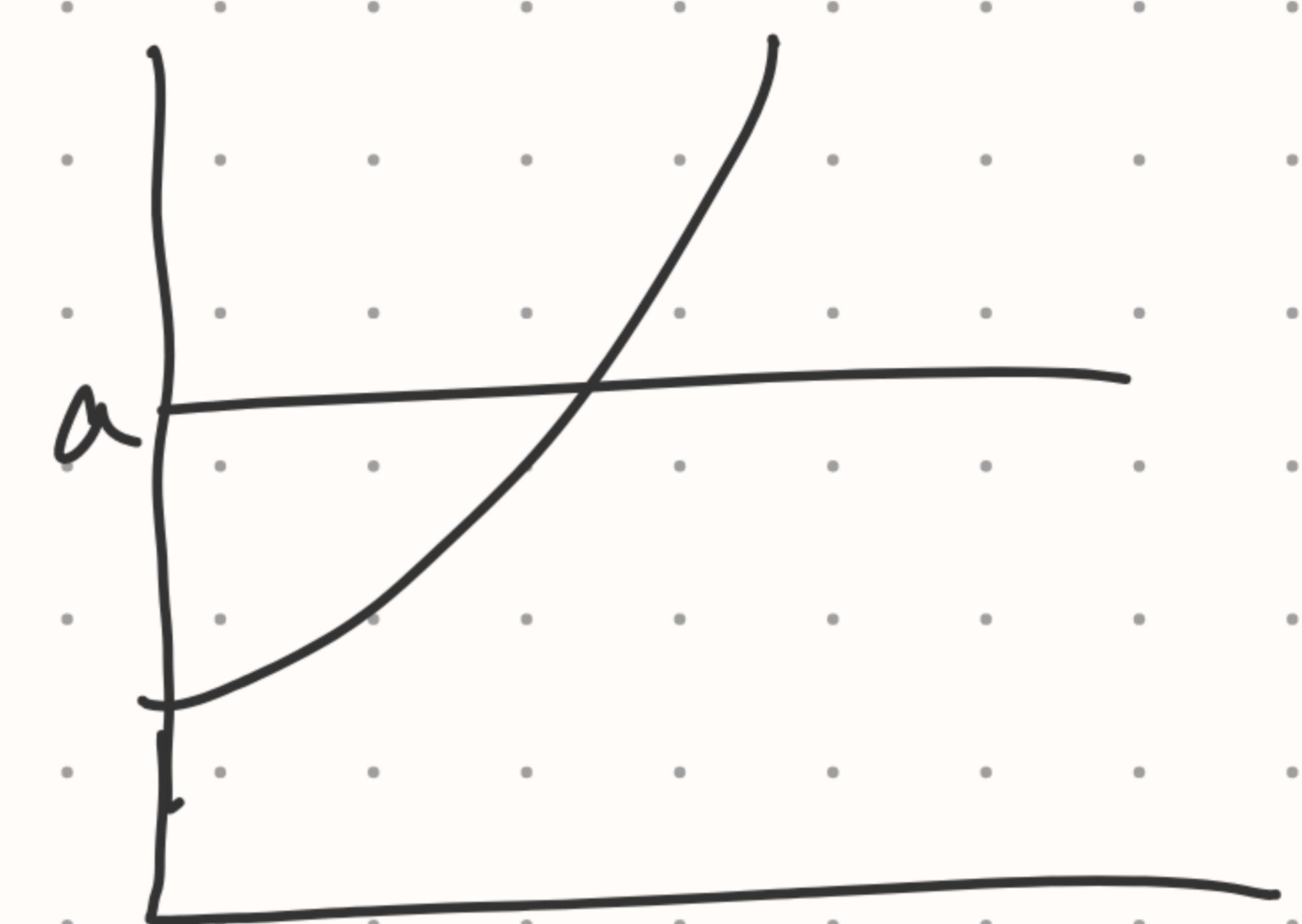
flattens out graph data points - Logistic Regs. sln  
(use such a loss function that  
flattens out w/ datapoints)

# Logistic Regression

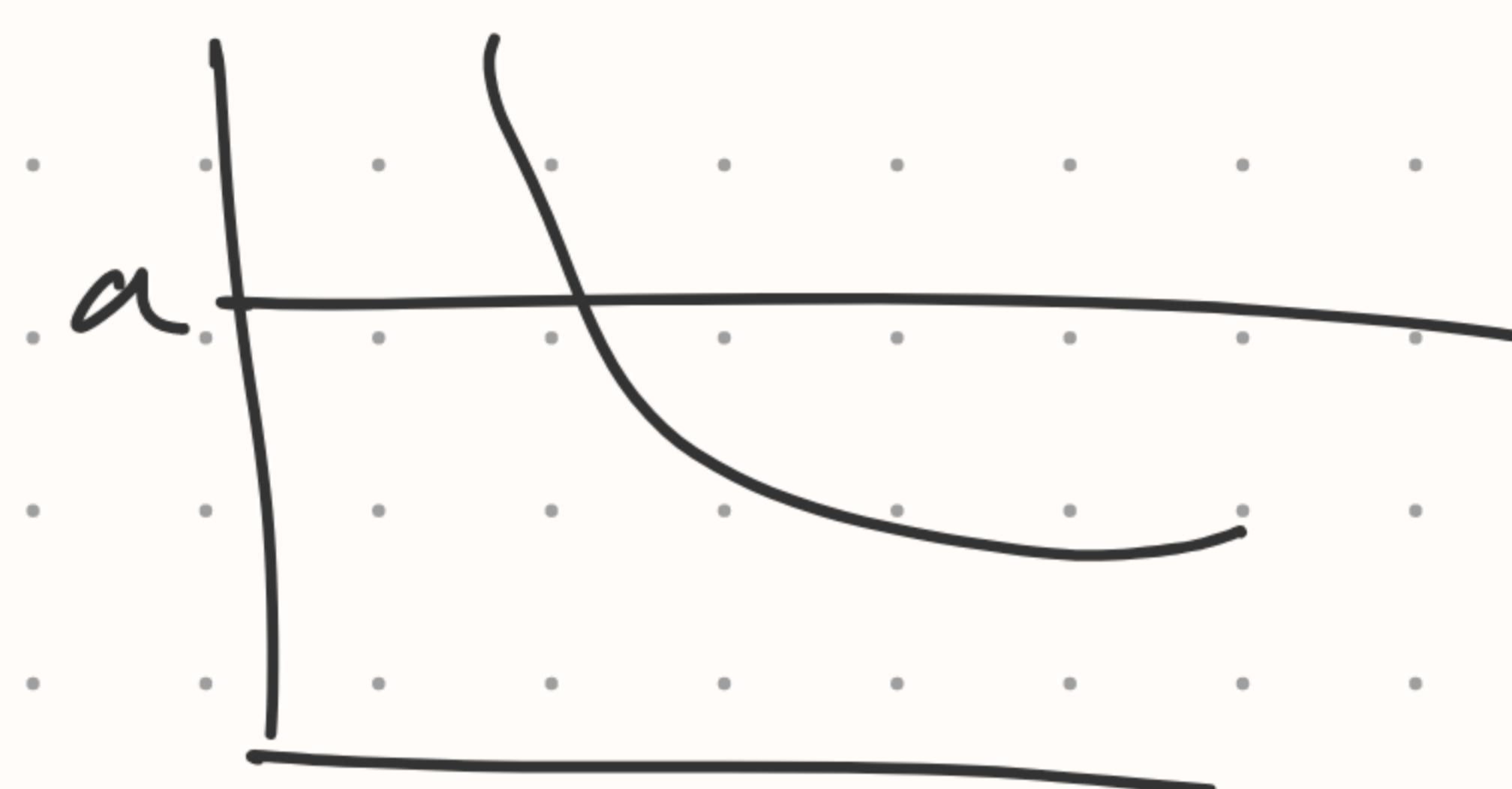
$$y = mx + b$$

$$y = a$$

$$y = e^{mx+b}$$



$$y = e^{-(mx+b)}$$



$$y = \frac{a}{1 + e^{-(mx+b)}}$$

$$y = \frac{1}{1 + e^{-(mx+b)}}$$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$m$  = weight

$b$  = bias

$$z = \sum_{i=1}^n w_i x_i + b = w x + b$$

Sigmoid function

Sigmoid func.

= fitting  
= utilities ✓  
- differentiable (optimizable)

\* 6 features

weight more for +ve words -1  
~ (pos) ~ -ve ~ -0

$$x = [3, 2, 1, 3, 0, 4.19]$$

$$w = [2.5, -5, -1.2, 0.5, 2, 0.7] \quad [b = 0.1]$$

(weight)  
high-value

Bias - prevent overfitting data

$$p(+ | \text{doc}) = \delta(w \cdot x + b)$$

$$= 3(2.5) + 2(-5) + \dots + 0.1$$

$$= \delta(0.833) = \frac{1}{1 + e^{-0.833}} \\ \approx 0.7$$

$$p(- | \text{doc}) = 1 - \delta(z) = 1 - 0.7 = 0.3$$

Autoencoders  
(Basic)  
(Advanced)

Feature template  $\rightarrow$  Representation learning  
(prev)  $\downarrow$  (now)  $\uparrow$  appropriate features that represent data

Dimensionality reduction  
(PCA)

SVD  
(Singular Value Decomposition)

$\hat{y}$  =  $\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix}$  [m × 1]       $\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_f \end{bmatrix}$  [f × 1]       $\begin{bmatrix} b^{(1)} \\ b^{(2)} \\ \vdots \\ b^{(m)} \end{bmatrix}$  [m × 1]

↓ Feature matrix      ↓ weight      ↓ Bias

$\hat{y} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_f^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_f^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} & \dots & x_f^{(m)} \end{bmatrix}$  [m × f dimension]

$$z = w \cdot x + b$$

$z = \text{weight} \times \text{feature} + \text{bias}$

- Tradeoffs between Naive Bayes & Logistic Regression

matrix multiplication → costly  
 small datasets ✓  
 - enough datasets  
 discriminate between classes ✓  
 (discriminative features)

- Binary classification → use sigmoid function

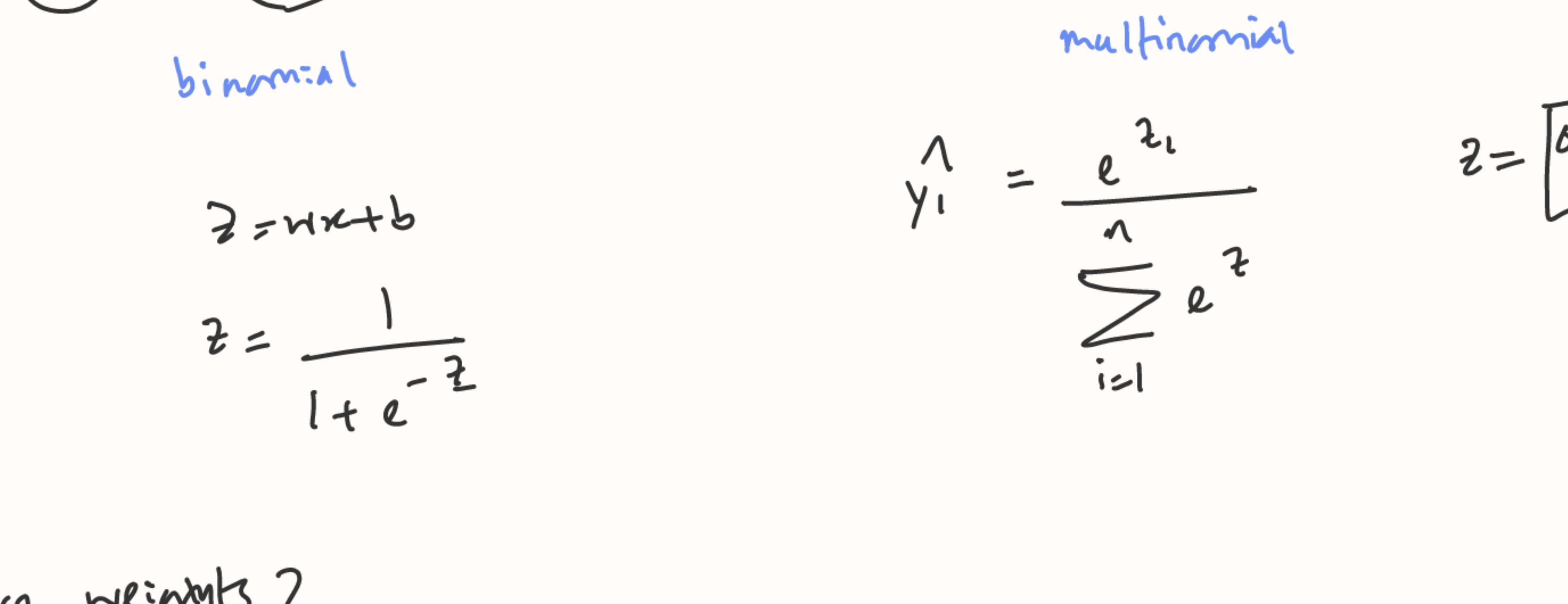
$\hat{y} \rightarrow \text{positive}$   
 $1 - \hat{y} \rightarrow \text{negative}$

- Multinomial logistic regression

Sigmoid function designed to classify between 2 classes: 0 and 1.

Change objective function 1st to classify 2 classes.

Use softmax



$$z = w \cdot x + b$$

$$\hat{y} = \frac{1}{1 + e^{-z}}$$

$$\hat{y}_i = \frac{e^{z_i}}{\sum_{i=1}^n e^{z_i}}$$

$$z = \begin{bmatrix} 0.7, 0.15, 0.02, \dots \end{bmatrix}$$

n number of classes

- How to learn weights? (w and b)

Objective function: How likely the predictions are to my gold labels (true)

$L(\hat{y}, y) = \text{How much } \hat{y} \text{ differs from true } y$

$\downarrow$   
 pred      gold labels

Cross Entropy

outcome: 0 or 1

Bernoulli distribution  
 $P(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$

only 2 possible outcomes  
 prob.      feature of doc

$(y=1)$   
 if true label true i.e. 1 class  $\rightarrow P(y|x) = \hat{y}$   
 if true label -ve i.e. 0 class  $\rightarrow P(y|x) = 1 - \hat{y}$   
 $(y=0)$

$$\log P(y|x) = y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$L_{CE}(\hat{y}, y) = - \left[ y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

$$(Cross \ Entropy \ Loss) = - \left[ y \log G(w \cdot x + b) + (1-y) \log (1 - G(w \cdot x + b)) \right]$$

Target: minimize loss function

If the classifier is perfect, you have zero loss.

... perfectly imperfect? gold label = 0

classifying it as 1.

$$\log(0) = \infty$$

$$\log(1) = 0$$

range:  $[-\infty, 0]$

(of loss function)

bad classifier close to  $-\infty$

good ... - - - 0

no. of samples

Loss function  $\rightarrow \theta$

$$\hat{\theta} = \arg \min_{\theta} = \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}, \theta), y^{(i)})$$

( $w, b$ )

For each sample, minimize loss func.

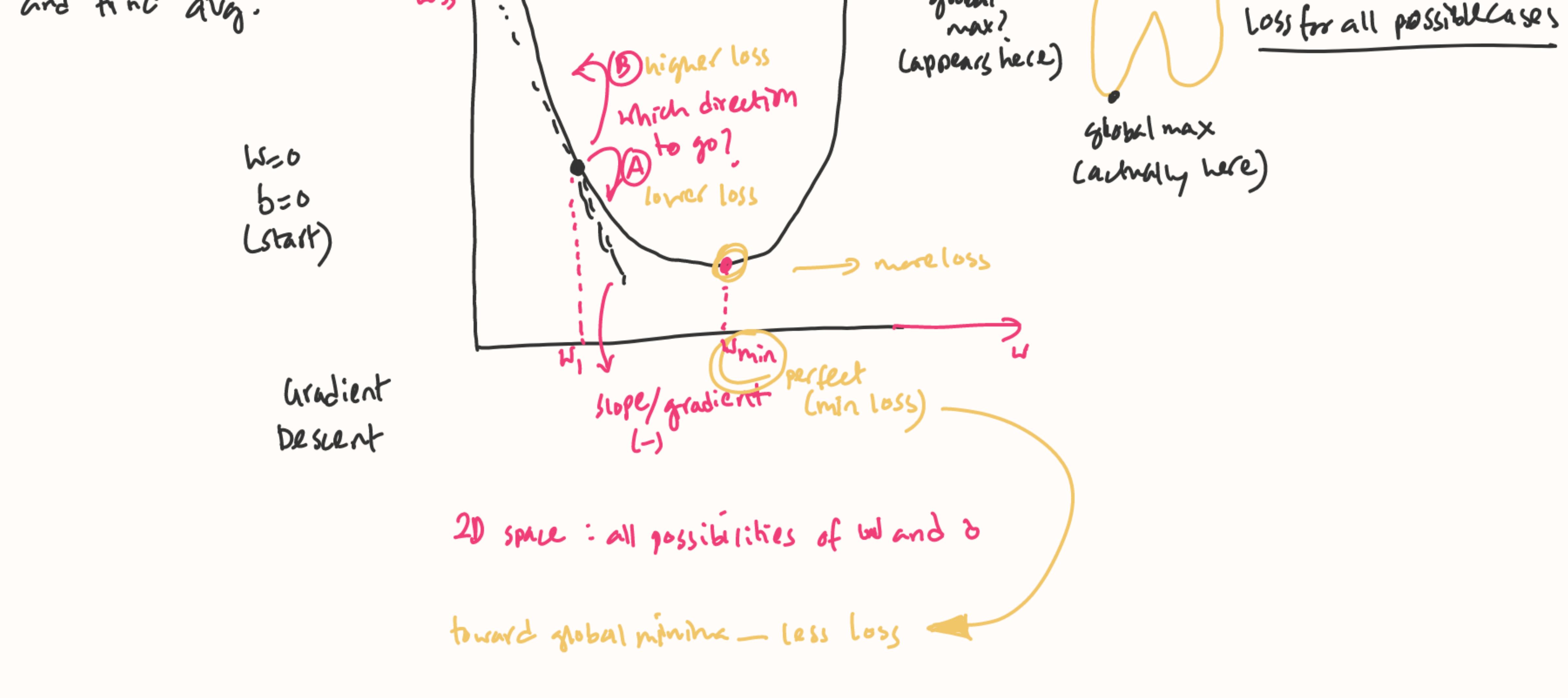
$$Loss \uparrow$$

$$\downarrow$$

$$(w, b)$$

$$\downarrow$$

$$Loss \downarrow$$



Future weight prediction:

$$\theta^{t+1} = \theta^t - \eta \frac{d}{dw} (L(f(x, \theta), y))$$

(η) Learning rate: 0.01

Gradient Descent algo decide which way to go — go down

How much to go? — Decided by η

η very high — loses global minima

η very low — take long time to find global minima

Hyperparameter testing: take large η and decrease slowly

- How to calculate gradient descent?

2 feature:  $x = [x_1, x_2]$

$$x_1 = (+ve) \text{ lexicons} = 3$$

$$x_2 = (-ve) \text{ lexicons} = 2$$

$$w_1 = w_2 = b = 0 \quad (\text{starts } w)$$

$$\eta = 0.1 \quad (\text{learning rate})$$

gradient descent of loss func.  
(inverse delta)

$$\theta^{t+1} = \theta^t - \eta \nabla L(f(x_i, \theta), y_i)$$

$$\nabla_{w,b} L = \left[ \begin{array}{c} \frac{\partial}{\partial w_1} L_{CE}(\hat{y}, y) \\ \frac{\partial}{\partial w_2} L_{CE}(\hat{y}, y) \\ \frac{\partial}{\partial b} L_{CE}(\hat{y}, y) \end{array} \right] \quad \text{partial derivatives wrt } w_1$$

loss func. for logistic regression:

$$L_{CE} = - \left[ y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right]$$

$$\frac{\partial}{\partial w_j} L_{CE} = (\hat{y} - y) x_j$$

$G(w \cdot x + b)$  sigmoid func

$$\nabla_{w,b} L = \left[ \begin{array}{c} (G(w \cdot x + b) - y) x_1 \\ (G(w \cdot x + b) - y) x_2 \\ G(w \cdot x + b) y \end{array} \right] = \left[ \begin{array}{c} -0.5 x_1 \\ -0.5 x_2 \\ -0.5 \end{array} \right] = \left[ \begin{array}{c} -1.5 \\ -1.0 \\ -0.5 \end{array} \right]$$

$$G(z) = \frac{1}{1+e^{-z}}$$

$$\theta^t = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix}$$

$$\theta^{t+1} = \theta^t - \eta \nabla L(f(x_i, \theta), y_i)$$

↓  
Stochastic grad. desc. =  $\begin{bmatrix} w_1 = 0 \\ w_2 = 0 \\ b = 0 \end{bmatrix} \rightarrow 0.1 \begin{bmatrix} -1.5 \\ -1.0 \\ -0.5 \end{bmatrix}$   
(one step at a time)

- online algo

$$= \begin{bmatrix} 0.15 \\ 0.1 \\ 0.05 \end{bmatrix} \rightarrow \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \text{ after 1 step}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}, \theta), y^{(i)})$$

convex optimization — gradient descent always points towards global min.  
optimize weights and biases

Perfectly fit weight on features → model tries

Cannot generalize unseen data → overfitting  
performs well on training data

### • How to tackle overfitting in optimization problem?

— Regularization

1. L1 norm : Absolute distance from origin

2. L2 norm : Euclidean distance from origin  
(degree=2)

For L2 norm, eqn:  $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}, \theta), y^{(i)}) - \alpha \sum_{j=1}^n \theta_j^2$

$\downarrow$   
distributes weights so that no weight becomes very large (break down weight into smaller weight)

For L1 norm, eqn:  $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(f(x^{(i)}, \theta), y^{(i)}) - \alpha \sum_{j=1}^n |\theta_j|$

$\downarrow$   
sparsity  
big weight randomly make some weights 0

solt randomly make some weights 0

L2 better than L1 ✓

all eq at a time: time and resources ↑

solt batch-size

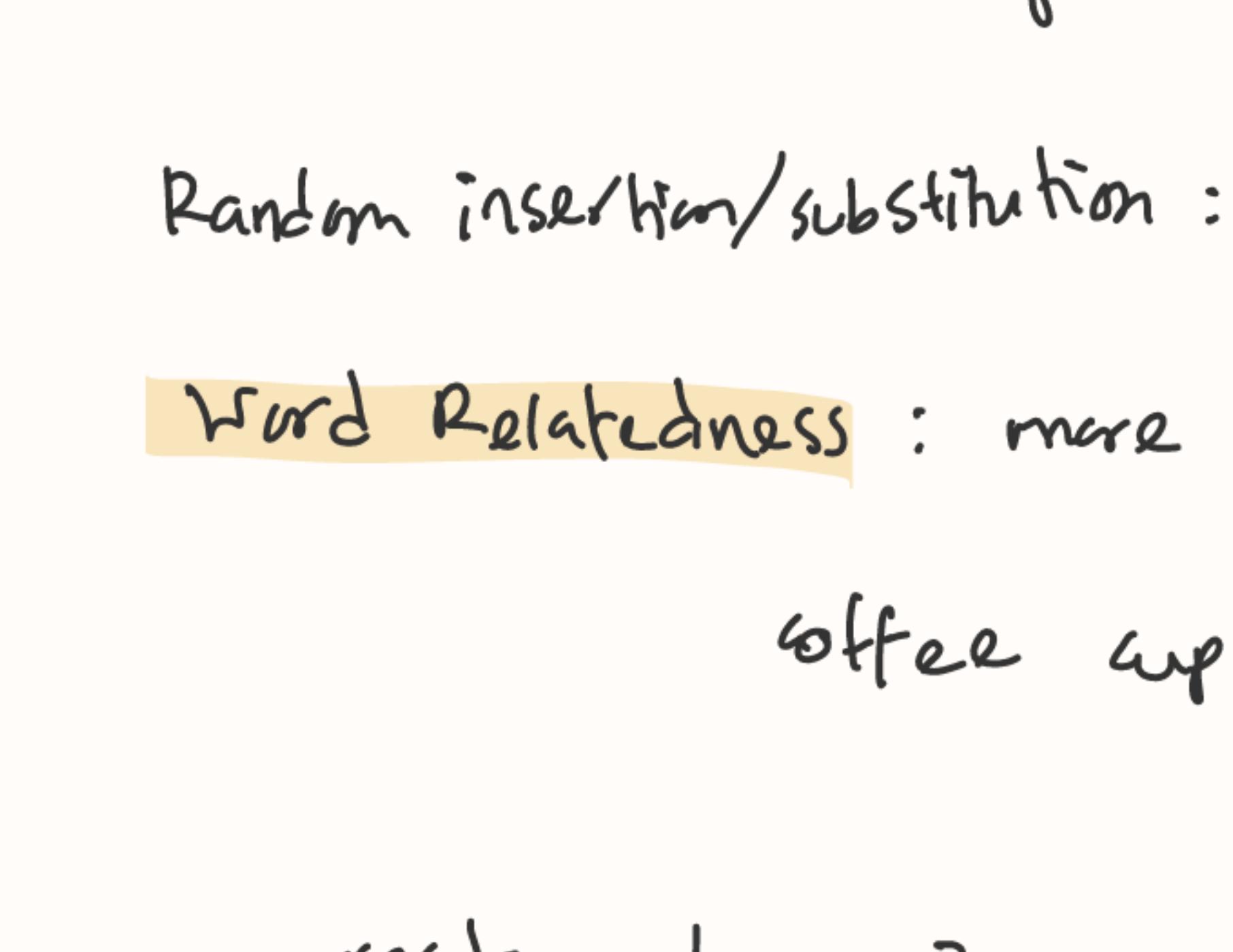
Stochastic gradient descent : one by one eq  
(online algo)

### Embeddings

Feature Engineering : hand crafted features  
not good for large doc

frequency — prev. evaluation metric to represent features

Represent words in vector space using their meanings → vector semantics  
(Representation learning)



Synonymy : couch, sofa  
water, H<sub>2</sub>O

"No two words are perfectly synonyms" — synonyms are not always substitutable

Word Similarity : principle of

no perfect synonyms but have similarity

Cat, dog — both animals

Random insertion/substitution : cat → dog

Word Relatedness : more broadened scope

coffee cup — related

(used in same context / semantic field)

restaurant : waiter, menu, recipe — related

### LDA : Latent Dirichlet Allocation

$\downarrow$   
models topics

uses a context window

makes a topic based on sentences occurring in a context window

document: Doctor nurse hospital OT

Connotations : related to sentiments

innocent, naive

$\downarrow$

positive

$\downarrow$

negative

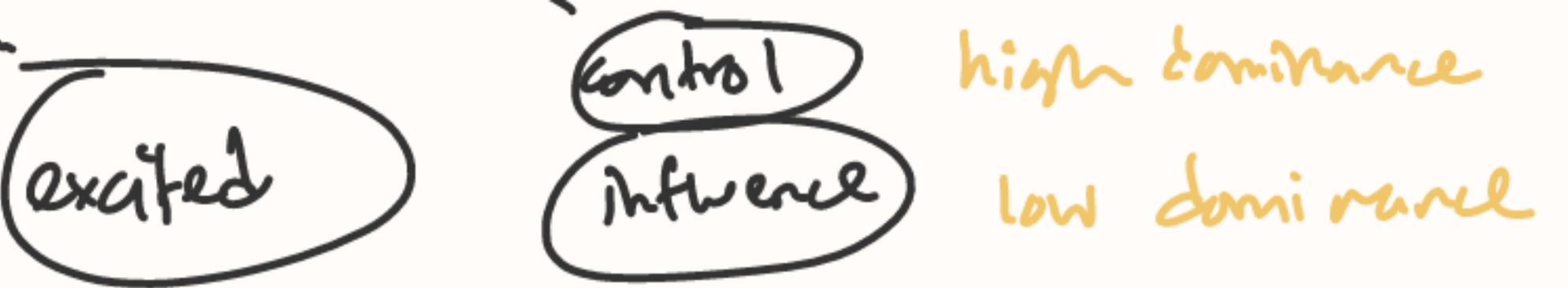
→ connotation

Every word has 3 particular characteristics:

1) Balance (pleasantness)

2) Arousal (how intensely tries to stimulate)

3) Dominance (degree of control)



Excited (—, —, —)  
 ↓      ↓      ↓  
 balance    arousal    dominance

Projecting a word in 3D space  
 (multidim)

Embeddings

Word can be represented as a point in multidim space.

Document similarity - how similar doc is to query  $\rightarrow$  search engine

Info retrieval (IR) - query  $\rightarrow$  docs provided  
 (most appropriate)

Cosine similarity  $\rightarrow$  smaller angle - more relevant to query

Col vector: docs as cols

Row vector: words as cols

<u>hypothesis</u>	<u>decision</u>
similar words $\rightarrow$ similar docs	occur in
similar docs $\rightarrow$ similar words	contains

Distributional semantics : Projecting words in vector space based on meanings

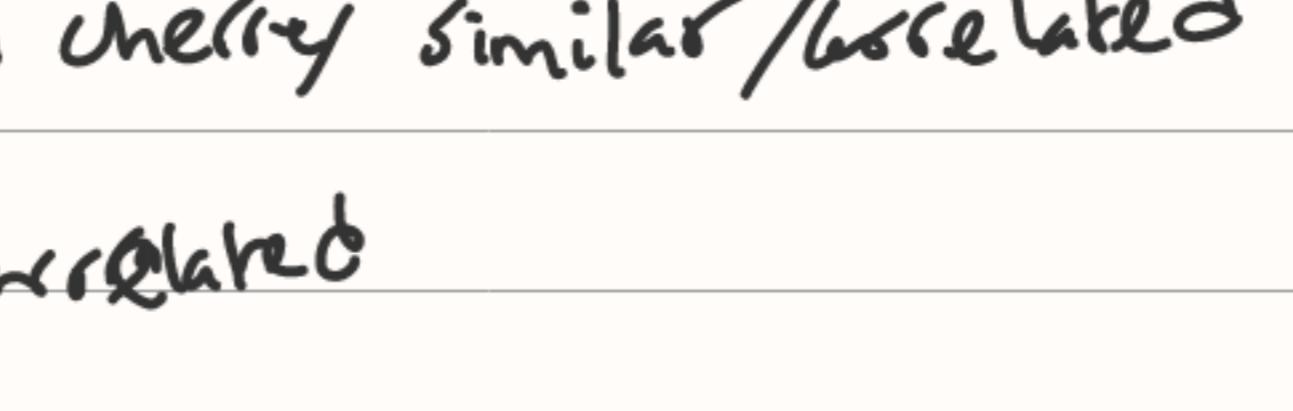
- till now, uses frequency

### Term document frequency

# rows → # words

# cols → # docs

represent words in terms of doc



Represent words in terms of words

context : 4 words before, 4 words after

for each word : context window = 4

which words occur frequently in same context window

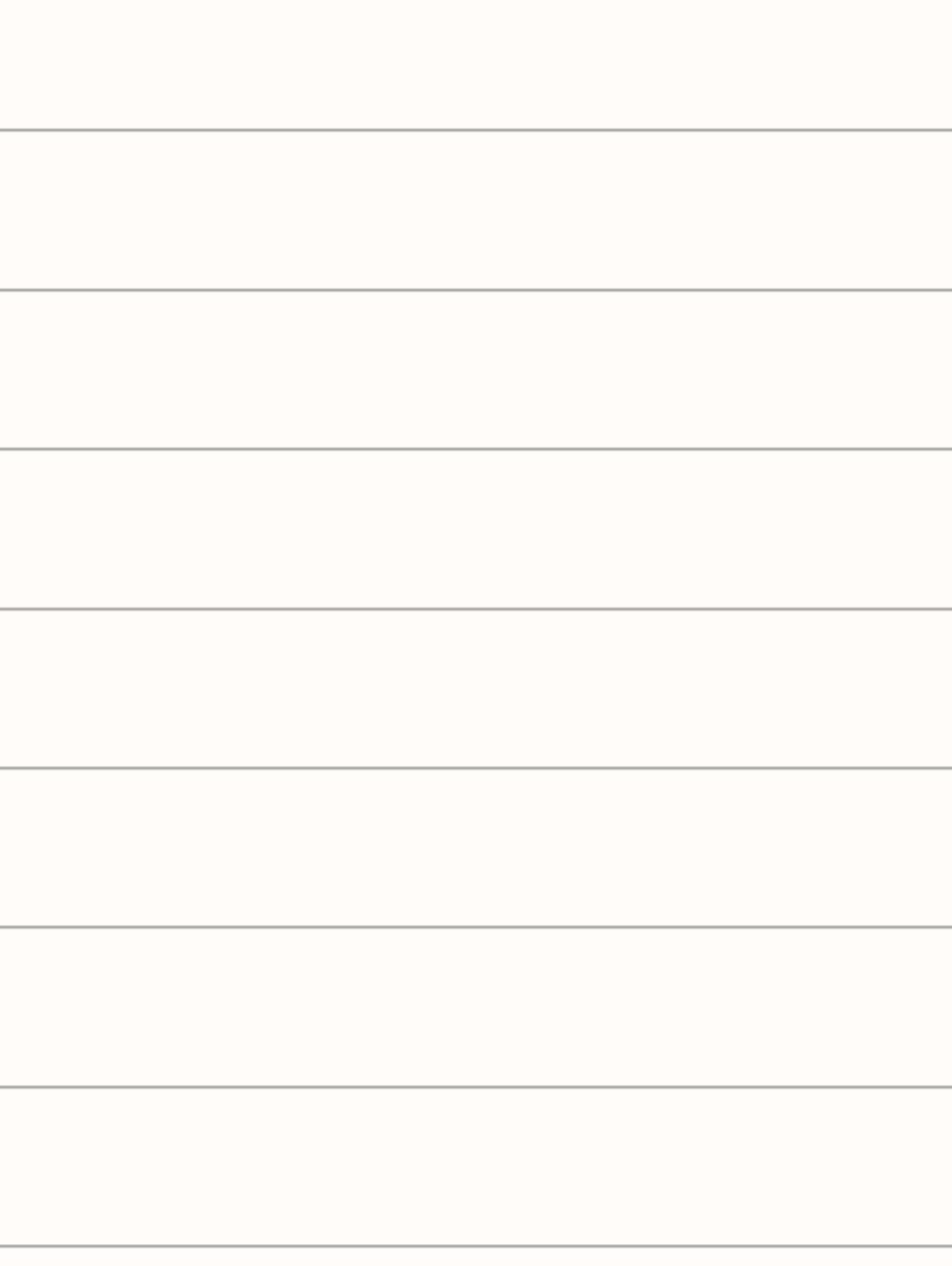
digital computer

strawberry pie

cherry pie

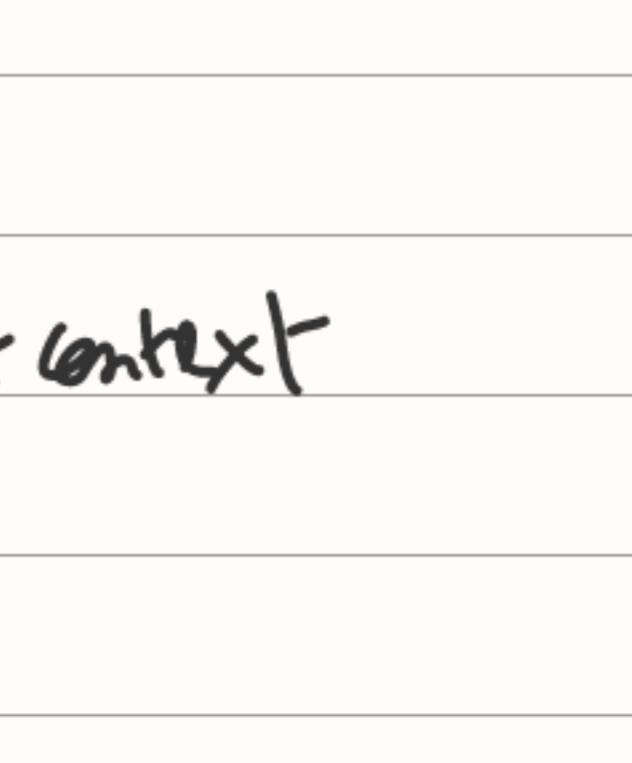
strawberry and cherry similar/correlated

digital and info correlated



Embedding : Projecting those words in terms of meanings

If two words appear in same context window, they are similar in meaning. (Almost the same vector)



### Cosine similarity

measure similarity between 2 vectors  
(in 2D)

But in multiDim?

$$Dm = \sqrt{v \times D}$$

$\downarrow$   
# vocab    # docs

If vectors are similar, dot product will be higher

$$\cos(\text{cherry, info}) = 0.007$$
$$\cos(\text{digital, info}) = 0.998$$

digital and info more similar

NOT

In 1 doc, high value → negative doc (rare type doc)

Term frequency says about context

Sparse vector (0 at most entries)

Freq. words are not good discriminators of context.

Stopwords - highest occurring terms → Doesn't say anything about context

### TF-IDF Term frequency Inverse document frequency

tf-idf → suppress the imp of frequency →  $1 \rightarrow 1$

$$10 \rightarrow 2$$

$$100 \rightarrow 3$$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

tf-idf →  $100 \rightarrow 3$

tf-idf →  $1 \rightarrow 1$

tf-idf →  $10 \rightarrow 2$

## Word2vec

Embedding vector : dense vector  
(not sparse)  
frequently occurring 0s

$$\begin{aligned} |v| \times f &\rightarrow \text{embedding for tf-idf} & \text{doc/terms} \\ |v| \times |v| &\rightarrow \text{term-term vector} & \text{the - - - (f)} \\ && \text{cherry} \\ && \text{data} \\ && \vdots \\ && \text{all unique words} \\ && (|v|) \end{aligned}$$

Prob: tf-idf and PMI  
→ frequency dependent  
→ sparse vector

Soln Dense vector : learn weights  
(not frequency dependent)

$$\begin{aligned} |v| \times d &\rightarrow \text{Dimension - random number} \\ &(300) \\ &\text{Optimal result} \end{aligned}$$

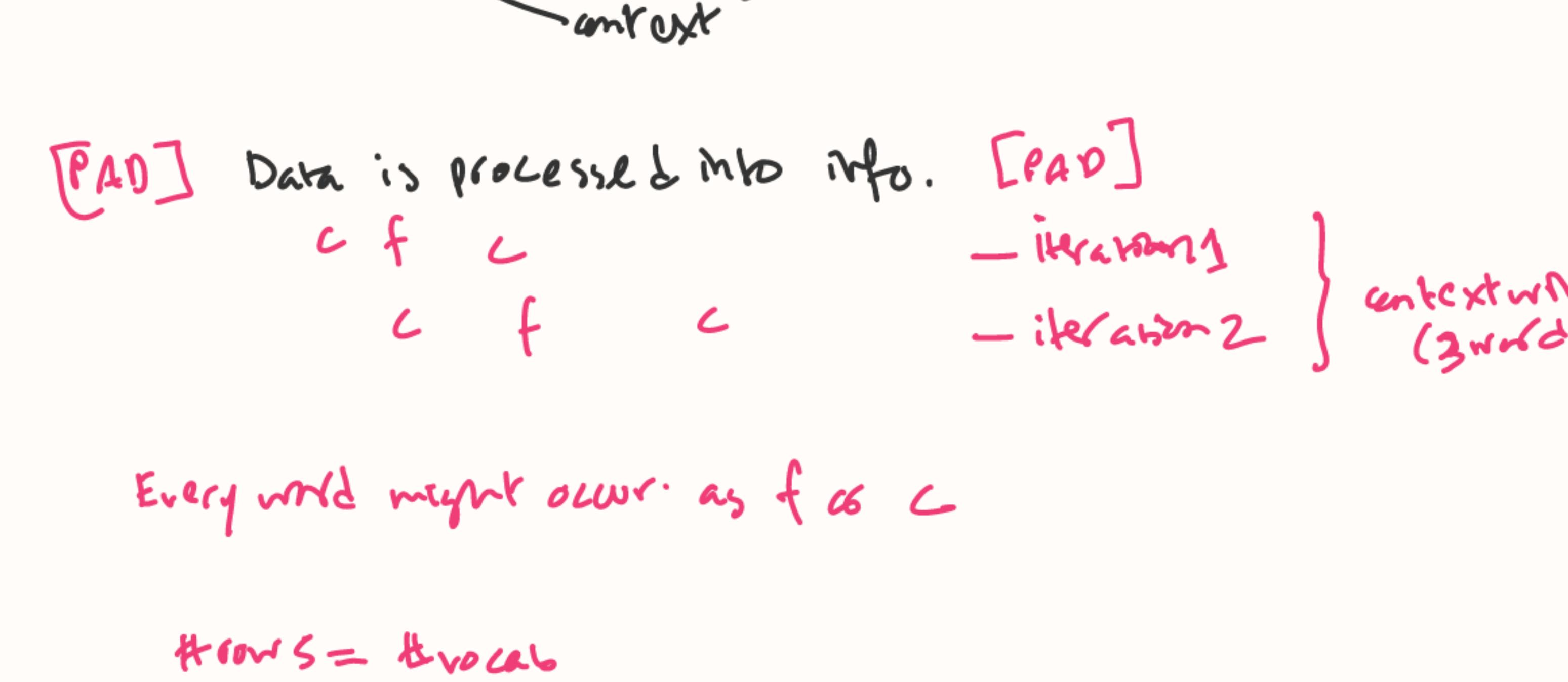
target word vs context word prediction

context words  
focus word  
Data is processed into information.

context window = 4

want to  
When you know context words, find focus words  $\Rightarrow$  CBOW  
(continuous Bag of words)

When you know focus word, - - - context words  $\Rightarrow$  Skip-Ngram

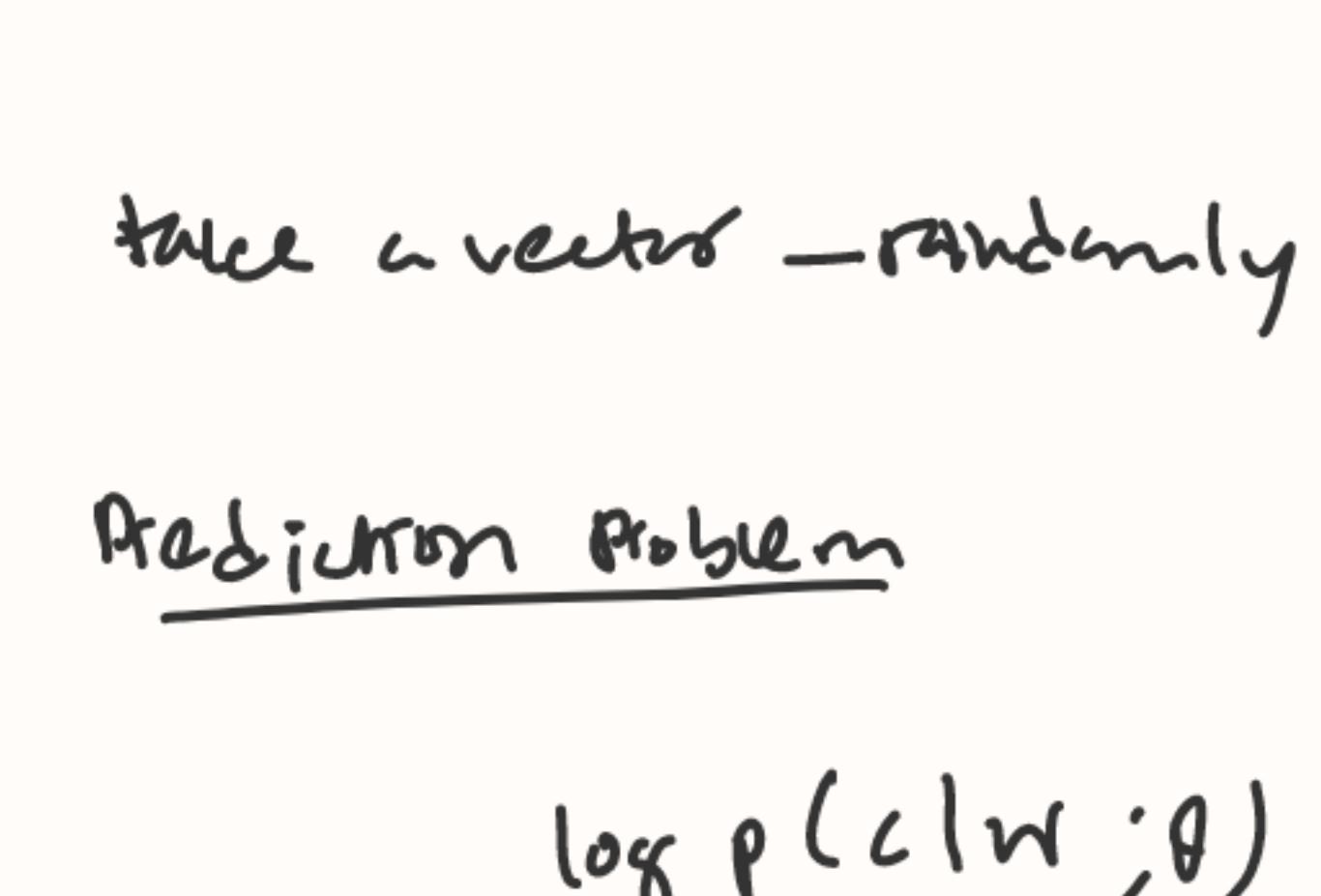


similar words into similar classes — target of embeddings

Skip-Ngrams is more suitable approach for finding embeddings of rare words.  
(diversity more)

1m words in vocab

# probabilities :  $1m \times 2$  — not feasible



[P4D] Data is processed into info. [P4D]

$C \ f \ c$   
 $c \ f \ c$

— iterations 1  
— iteration 2 } contextwindow = 1  
(3 words)

Every word might occur as f or c

# rows = # vocab

# cols = # defined by you (traditionally 300)

Data is processed into info.

is

processed

into

info

cat

dog

⋮

⋮

⋮

$V$

$C$

Brings down frequent words

Brings up infrequent words

use grad descent to optimize

$$P = \frac{c}{\sum c^d}$$

$$d = 0.75 \rightarrow \frac{3}{4}$$

{

use grad descent to optimize

Embeddings  $\Rightarrow$  weights

softmax, loss function

wrd (w)  
if actually occurring in context window - true (actual)  
 $w \cdot c_1$  - predicted

loss function:  $1 - w \cdot c_1$

$1 - w \cdot c_2$

global minima - min. loss

similar contextual words in similar windows

word2vec Jordan video

$$(-) \rightarrow \text{to minimize}$$
$$\begin{matrix} w \\ \xrightarrow{\text{SVD}} \\ V \times 300 \end{matrix} \quad C^T \quad 300 \times |V| = \begin{matrix} C \\ \downarrow \\ m \end{matrix} \quad w$$

word2vec  $(300 \times 300)$  ?

$$A = S \Lambda S^T$$

SVD  
(singular value decomposition)  
Dimensionality reduction

→ 50000 to 300 dim  
(SVD cannot run in my PC  
can run word2vec)

Evolution of words through embeddings  
(1800s to 2000s)

Aff word2vec embeddings from same doc

Inherent variability - embeddings are variable

by randomizing weights at first iteration

multiple weights - avg

bootstrapping  
(randomize)

Word2vec - static embedding

↳ Every word has same embedding regardless of context

Now - Dynamic Embedding

↳

one word has diff embedding based on context

'King' - 'man' + 'woman' = 'Queen'

man : woman :: King : Queen

Gender stereotyping: male doc — nurse

(not female doc)

Misogynist

1) Allocational  
bias

homemaker

↓  
unfairly for diff people

2) Representational bias  
(racial bias)

American names: David — Vandya

African names: —  
racial remarks

Internet data contains bias - positively biased towards white

Debiasing techniques

↳ reduce bias, not remove

Remove bias from embeddings

(not training sample)

m10

→ Probability related Ques. ✓

| Derivation X

| Application ✓