

- Topics:
- Naive Bayes (contd)
 - 'practical' problems
 - New words
 - Laplace smoothing + Bias-variance
 - Underflow & computational challenges
 - Multinomial NB vs Bernoulli NB
 - Feature-importance & interpretability

- Outliers & treating them
- Gaussian NB for numerical features
- Code - sample } Spam vs Not
 - for SVM & NB
 - ↓
 - M-NB
- ↑
- Lr-SVM
- RBF-SVM

PoEV:

$$P(y=1 \mid \text{text}) = p(y=1) \prod_{i=1}^d p(w_i \mid y=1) / \cancel{K}$$

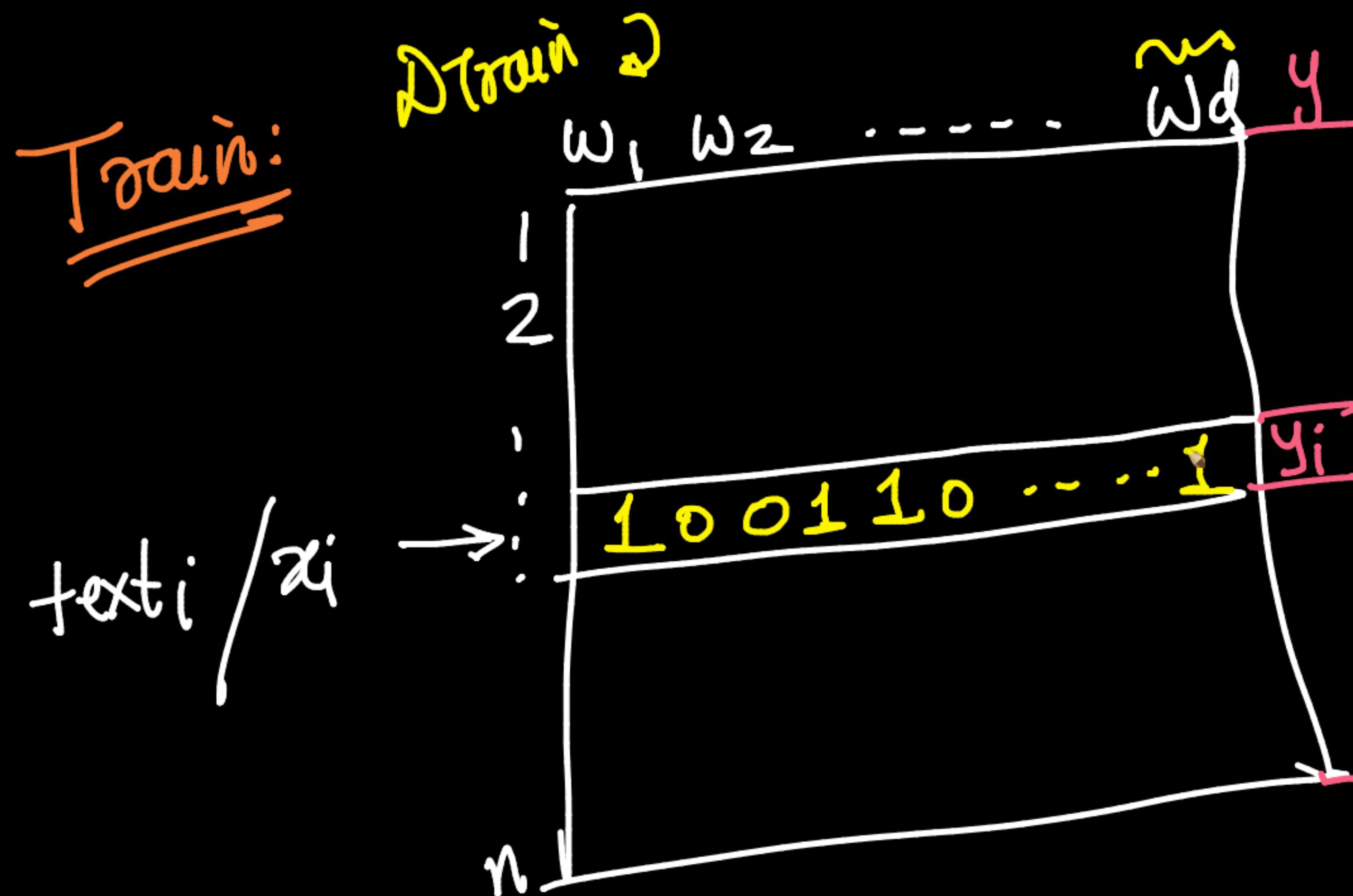
Naive Bayes

$$P(y=0 \mid \text{text}) = p(y=0) \prod_{i=1}^d p(w_i \mid y=0) / \cancel{K}$$

$$P(w_1, w_2, w_3, \dots, w_d | y=1)$$

$$= P(w_1 | y=1) \cdot P(w_2 | y=1) \cdot \dots \cdot P(w_d | y=1)$$

↳ independence of each feature / word
conditioned on the class



Class-prior

$P(y=1)$ ✓

$= \frac{\# \text{ of } y_i = 1}{n}$

$P(y=0) = \frac{\# \text{ with } y_i = 0}{n}$ ✓

②

likelihoods:

$$\hat{h}_j : l \rightarrow d$$

$$P(w_j | y_i = 1) = \frac{P(w_j \cap y_i = 1)}{P(y_i = 1)}$$

$$P(w_j | y_i = 0) = \frac{P(w_j \cap y_i = 0)}{P(y_i = 0)}$$


Run Test-time

$$x_{qj} = \underbrace{w_1 w_3 w_6 w_{12} w_{15}}_{K-\text{words}}$$

$y_{qj} = 1$

$$P(y=1 | \underbrace{w_1 w_3 w_6 w_{12} w_{15}}_{\text{words}}) \propto P(\underbrace{y=1}_{\text{test}}) \times \prod_{j=1}^K P(w_j | y_i=1)$$

$2K+2$

$$P(y=0 | \underbrace{w_1 w_3 w_6 w_{12} w_{15}}_{\text{words}}) \propto P(\underbrace{y=0}_{\text{test}}) \prod_{j=1}^K P(w_j | y_i=0)$$

(Q)

@ Run-time

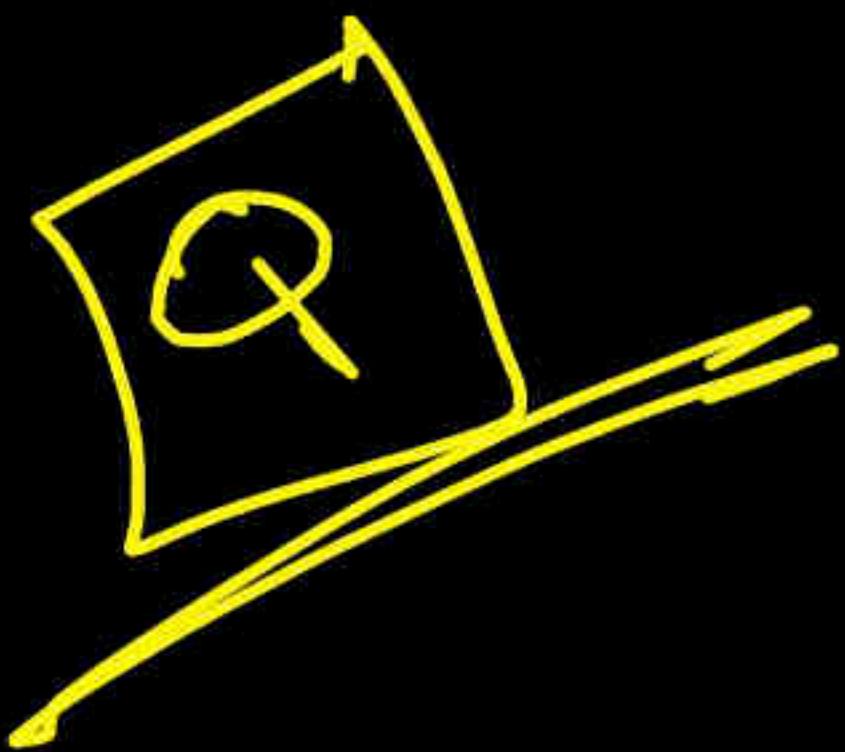
Space-complex of NB = $\mathcal{O}(d)$

likelihoods: 2^d

class priors: 2^d

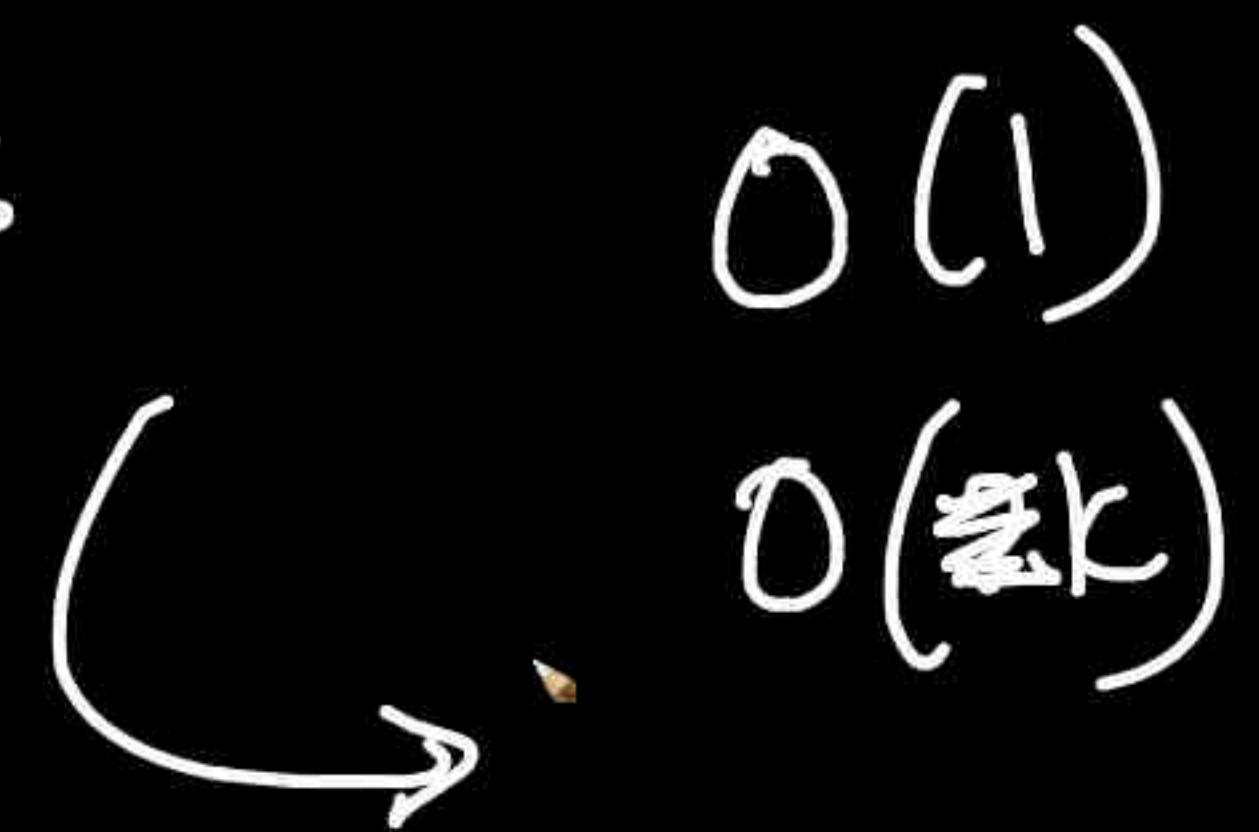
2^{d+2} probabilities

total # words in the total dataset



Run time- complx:

$x_{qj} \rightarrow k\text{-ways}$



Q

$$P(y=1 \mid w_1, w_2, w_3, \dots, w_k) \propto p(y=1) \prod_{j=1}^k p(w_j \mid y=1)$$

~~$P(w_1, w_2, w_3, \dots, w_k)$~~

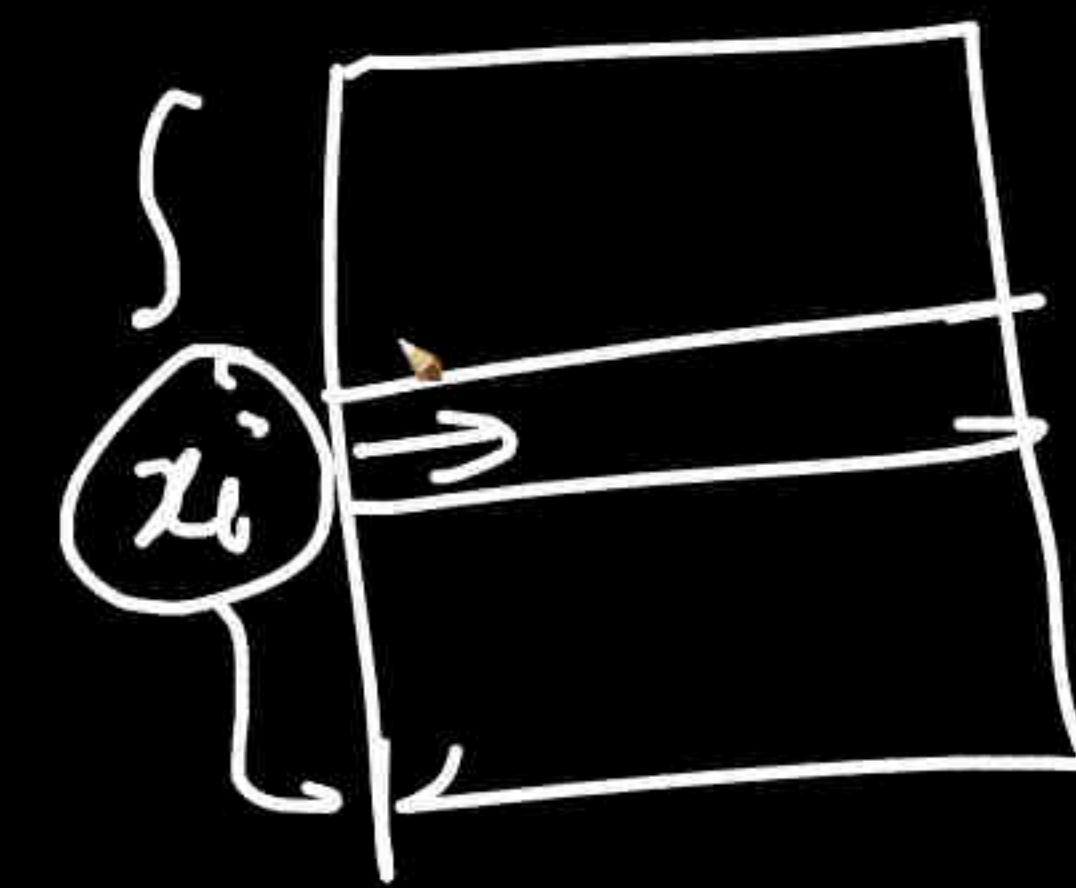
$$P(y=0 \mid w_1, w_2, \dots, w_k) \propto p(y=0) \prod_{j=1}^k p(w_j \mid y=0)$$

~~$p(w_1, w_2, \dots, w_k)$~~

popIn

den: $p(w_1, w_2, \dots, w_k)$

x_i



$$\frac{1}{n}$$

sample \rightarrow
 D_{train}

practical - issues

Run-time
 $x_{qj} =$

$\underline{\omega_1} \underline{-} \underline{\omega_2} \underline{-} \underline{\omega_6} \underline{-} \underline{\omega_7} \underline{\omega'}$

$\omega' \notin \{\omega_1, \omega_2, \dots, \omega_7\}$

↓
Set of words in
Dtrain

$$P(y=1 \mid \omega_1, \omega_2, \omega_6, \omega_7, \omega) = P(y=1) \times P(\omega_1 \mid y=1) \times \dots \times P(\omega' \mid y=1)$$



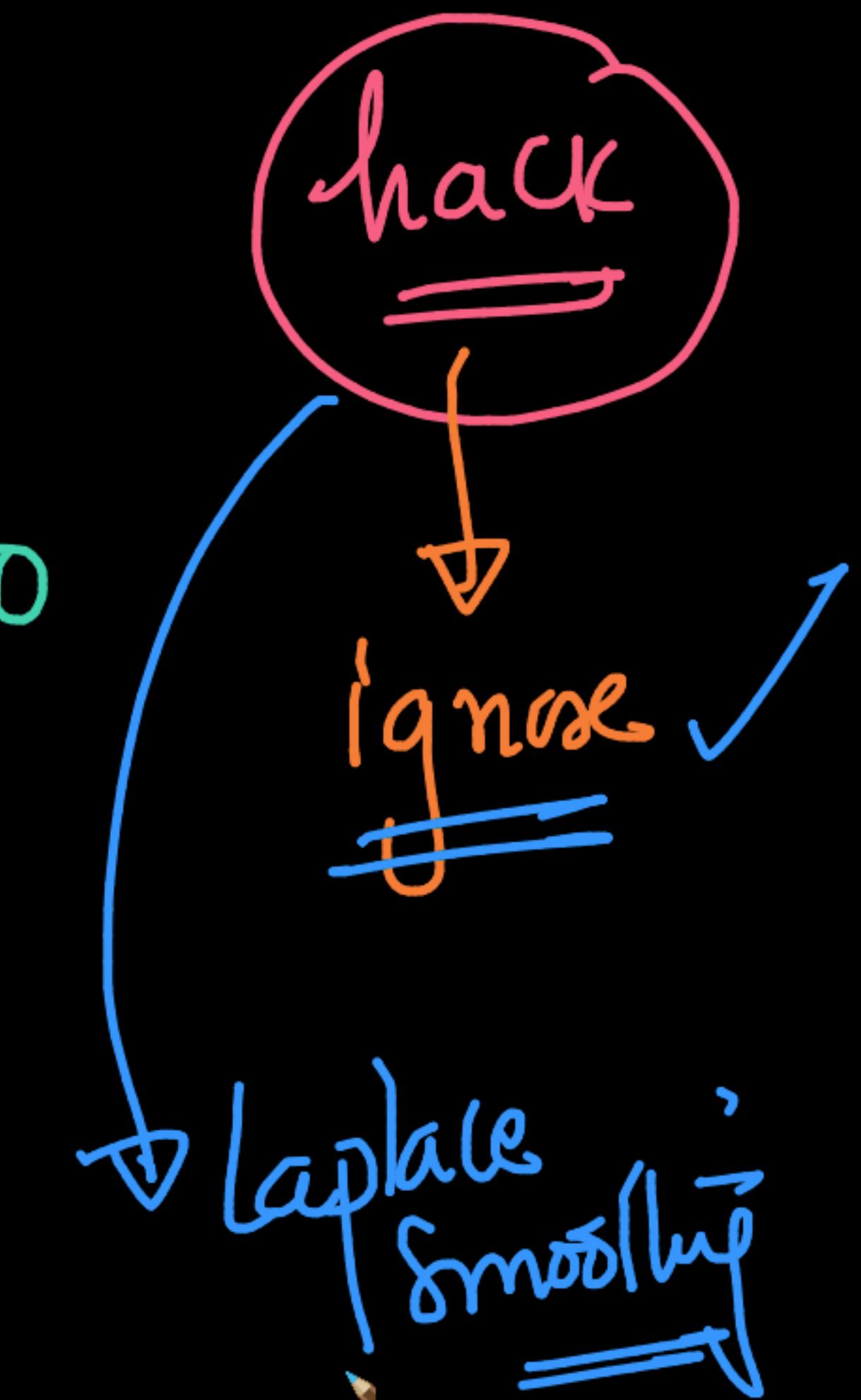
$\omega' \notin \{\omega_1, \dots, \omega_d\}$ pain

$$P(y=0 \mid \omega_1, \omega_2, \omega_6, \omega_7, \omega) = 0$$

Likelihoods

$$P(w_j | y=1) = \frac{n_{j1}}{n_1} \rightarrow D_{train}$$

$$P(w'_j | y=1) = \frac{n'_{j1}}{n_1} = \frac{0}{n_1} = 0$$



$$P(w_j | y=1) = \frac{n_{ij} + \alpha}{n_i + \alpha k}$$

Laplace smoothing
(Additive)

distinct values

than w_j can take

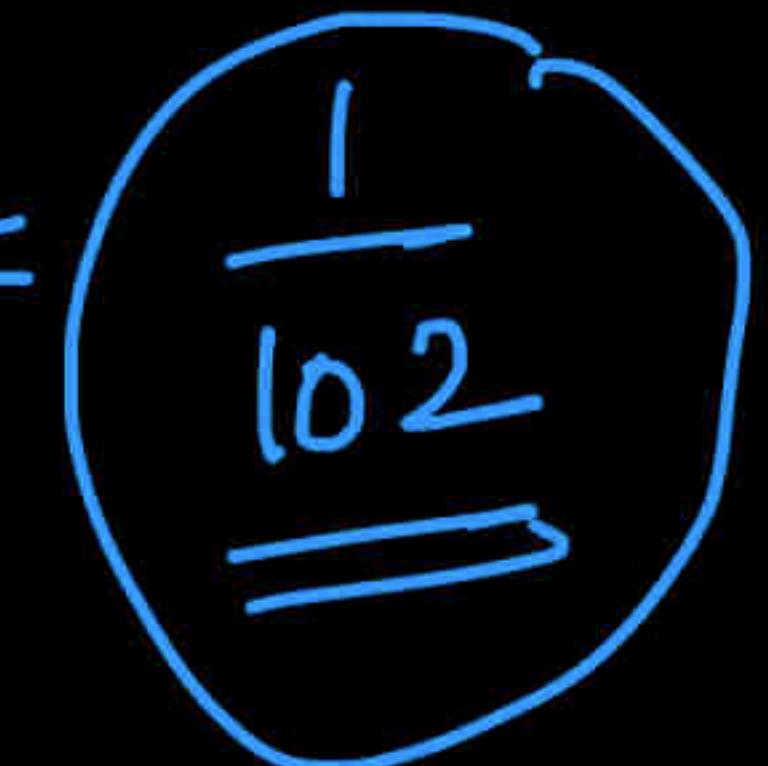
~~Simplifying:~~

$$\alpha = 1$$

$$k = 2$$

Bernoulli
2

$$P(\underline{\omega}' | y=1) = \frac{n'_1}{\underline{n}_1} = \frac{0+1}{(00+2)} = \frac{1}{102}$$



will not break my
math

NB + Laplace-smoothing

for all likelihoods:
 (train-time)
 & test-time

$$P(\omega_j | y=1) = \frac{n_{ij} + \alpha}{n_i + 2\alpha}$$

(let)

$$P(\omega_j | y=1) = \frac{1}{2}$$

(∴ we don't have any data)

$n_1 \approx 0$
 $d = 1$
 $n_1 = 100$

$$P(\omega^1 | y=1) = \frac{n'_1 + d}{n_1 + 2d} = \frac{1 + 100}{100 + 200} = \boxed{\frac{1}{3}}$$

$d = 10,000$
 $n_1 = 100$

$$P(\omega^1 | y=1) = \frac{0 + 10000}{100 + 20000} \approx \boxed{\frac{1}{2}}$$

$\hookrightarrow d$

$d_1 \gg n_1$

Training time



$$n_{j1} = 10; n_j = 100 \text{ (let)}$$

$$\text{No Laplace Smoothing} \Rightarrow -\frac{n_{j1}}{n_1} = 0.1$$

Case 1:

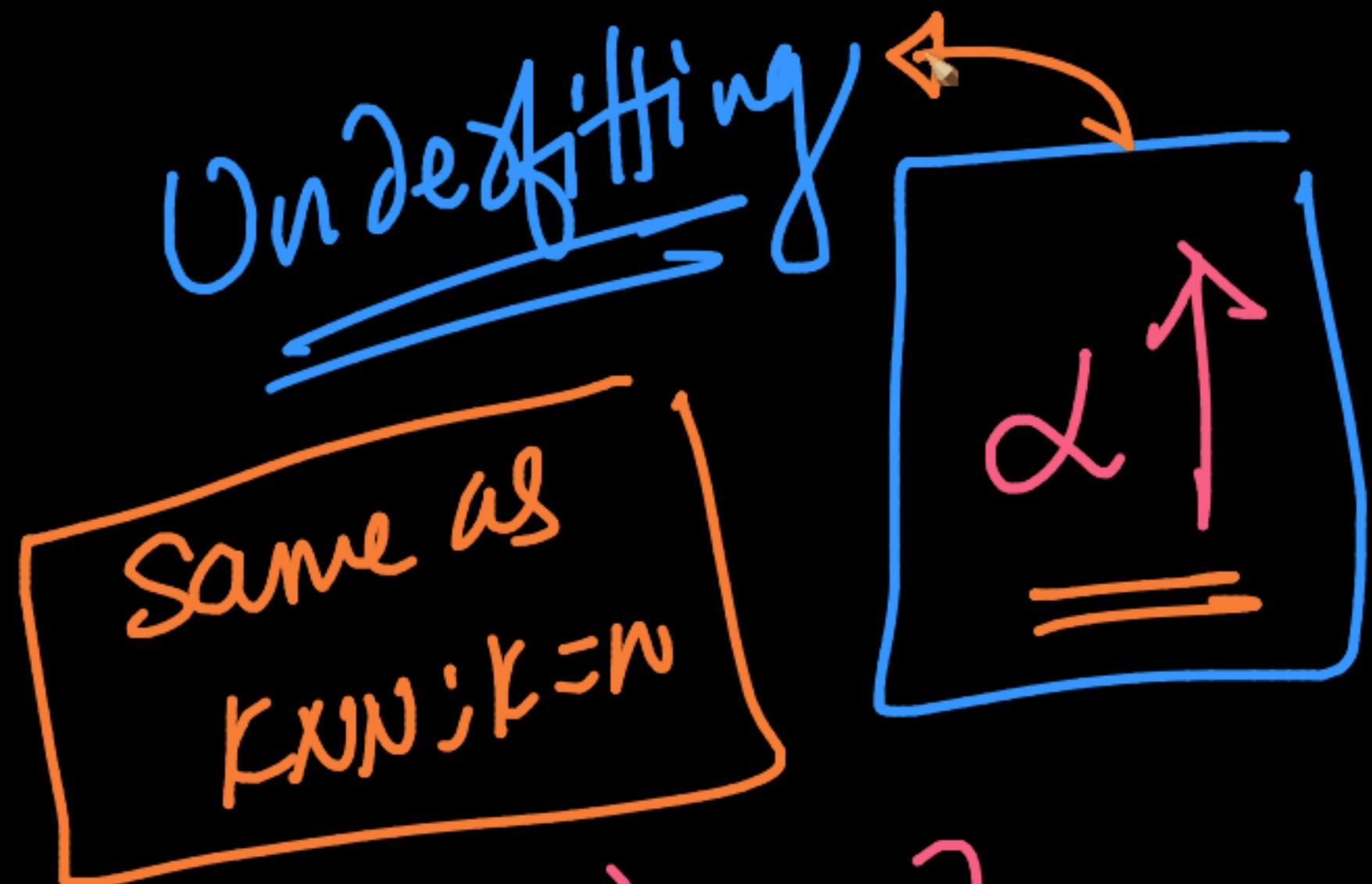
$$P(w_j | y=1) = \frac{n_{j1} + \alpha}{n_1 + 2\alpha} = \frac{(0+1)}{10+2} \approx 0.1$$

$$\alpha = 1$$

Case 2:

$$\alpha = 10,000$$

$$P(w_j | y=1) = \frac{10 + 10000}{100 + 20000} \approx \frac{1}{2}$$



$$P(Y=1) = 0.7$$

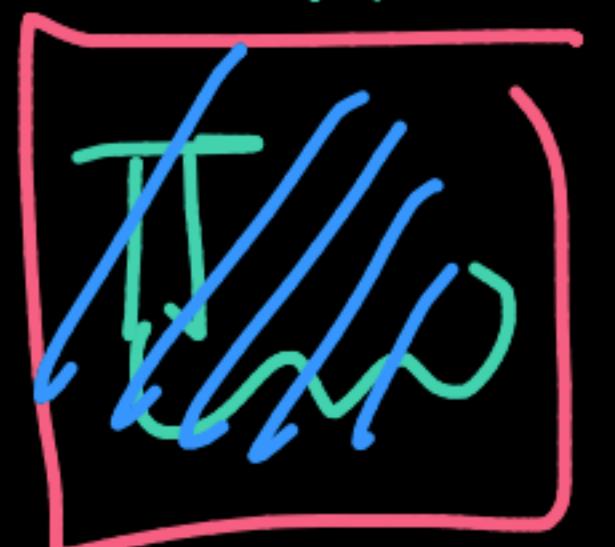
$$P(Y=0) = 0.3$$

\downarrow
for all x_q : $y_q = 1$

all my likelihoods will be
close to y_2
($\because K=2$)

$$P(Y=1 | \underbrace{w_1 \dots w_d}_{x_q}) = P(Y=1)$$

$$P(Y=0 | \underbrace{w_1 \dots w_d}_{x_q}) = P(Y=0)$$



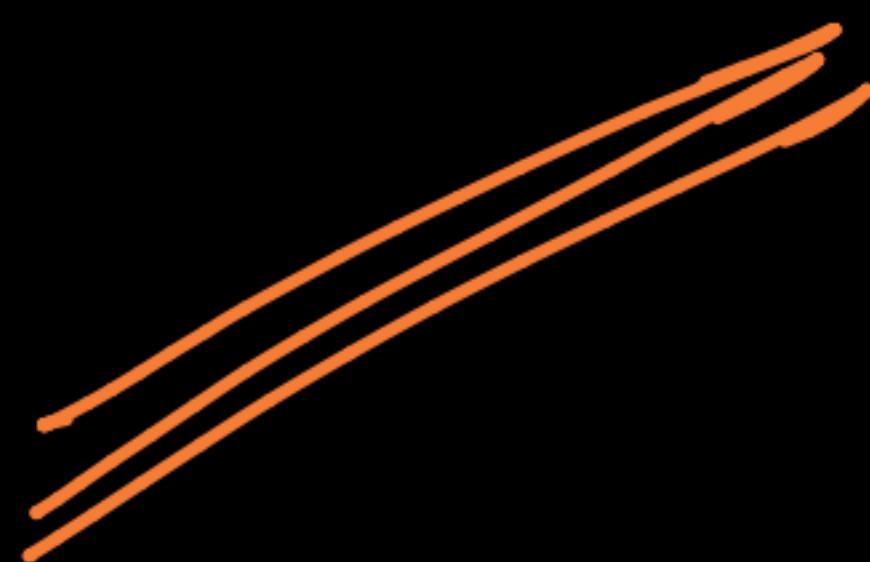
$\alpha \downarrow$

overfit

likelihoods \rightarrow

$$\frac{n_{j1}}{n_1} \propto \frac{n_{j0}}{n_0}$$

$\underbrace{\quad}_{\gamma}$



$$P(\omega | y=j) = 1$$



X



no hyper-parameter to play with



Naïve Bayes + Laplace Smoothing

HP: $\left\{ \begin{array}{ll} \alpha \uparrow & \text{underfit} \\ \alpha \downarrow & \text{overfit} \end{array} \right.$

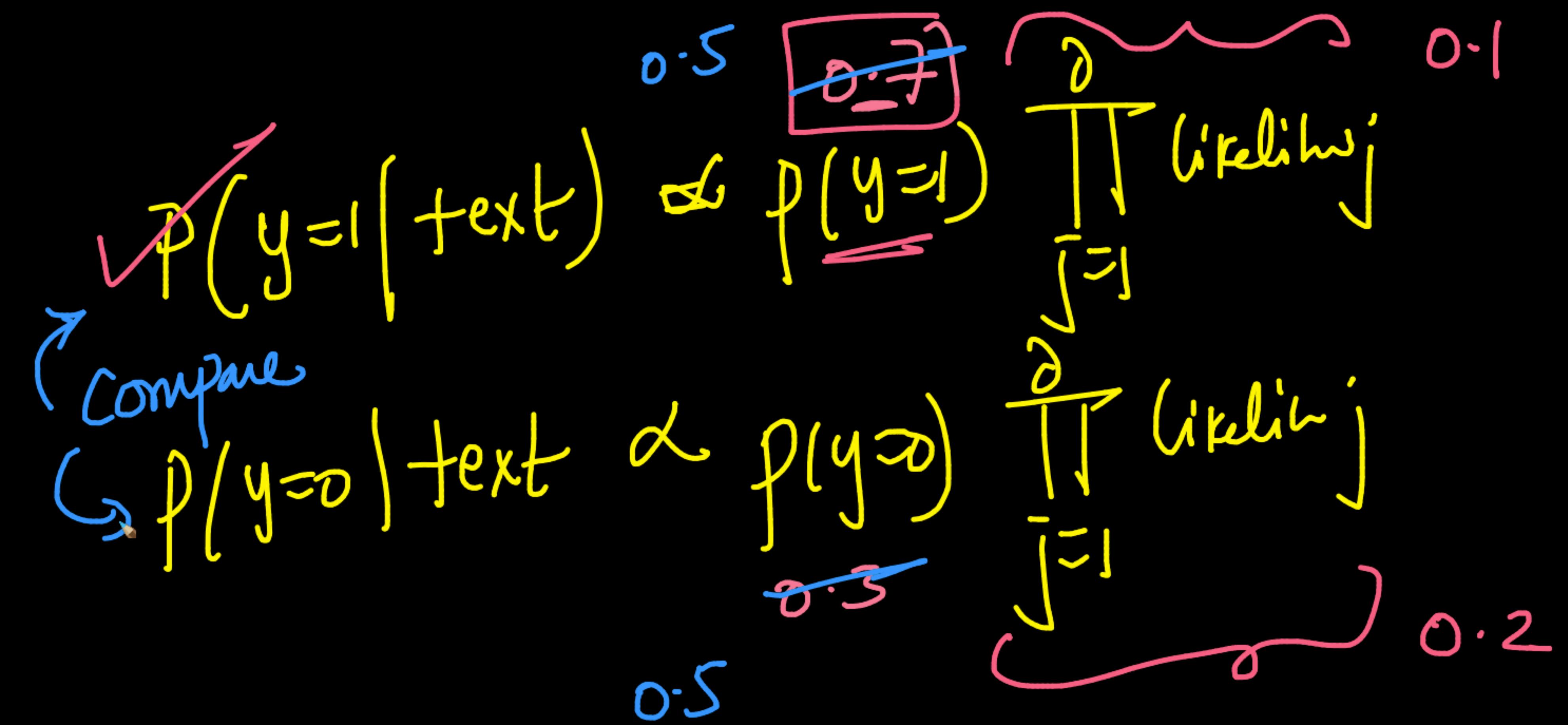
(Q)

$$P(y=1) = \underbrace{0.7}_{\sim} \quad \left. \right\} \Rightarrow \text{imbalance}$$

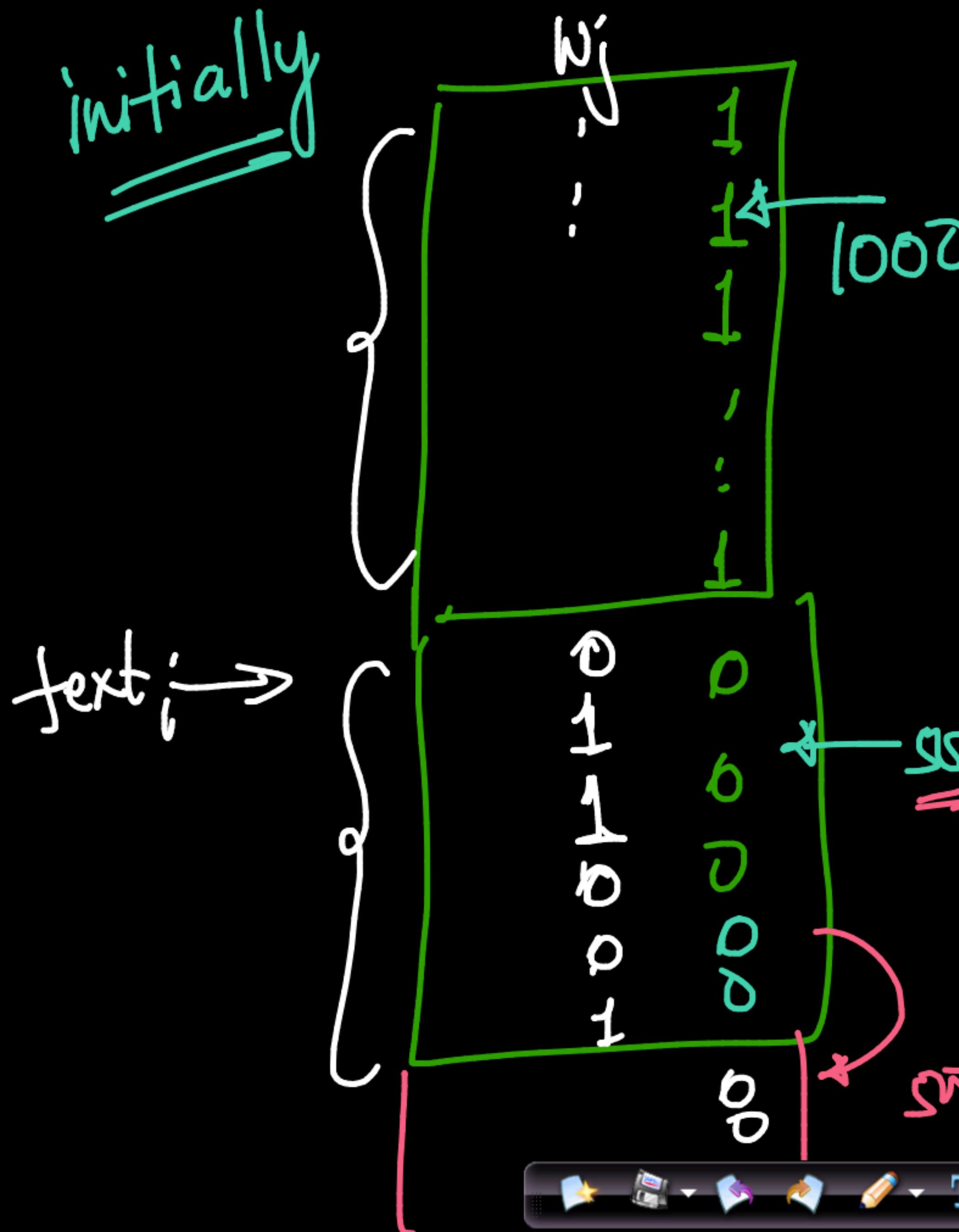
$$P(y=0) = \underbrace{0.3}_{=} \quad \begin{array}{l} \text{Yes, it is good to rebalance} \\ \hline \end{array}$$



rebalance → Upsampling



initially

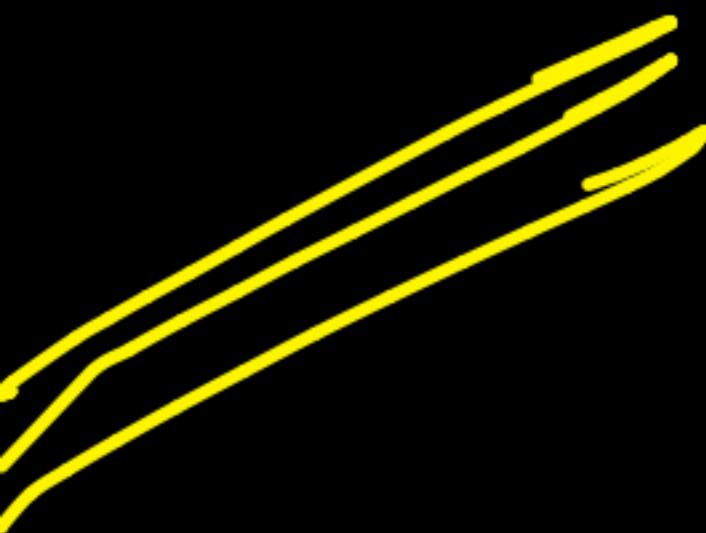


before rebalancing:

$$P(\omega_j | y^D) = \frac{n_j^D}{n^D} = \frac{50+2}{50+12}$$

rebalancing by weight:

$$P(w_j | y=2) = \frac{n_{j0}}{n_0} = \frac{600}{1000} = 0.6$$



$$p(y=1 | w_1, w_2, \dots, w_d) \propto p(y=1) \prod_{j=1}^d p(w_j | y=1)$$

Computational
problem } \rightarrow

Code

$$P(y=1 | w_1, w_2, \dots, w_d) \propto P(y=1) \prod_{j=1}^d P(w_j | y=1)$$

Computational
problem

$(d+1)$ [0 - 1]

d ↑

$$0.1 \times 0.01 \times 0.2 \times 0.3 \times \dots$$

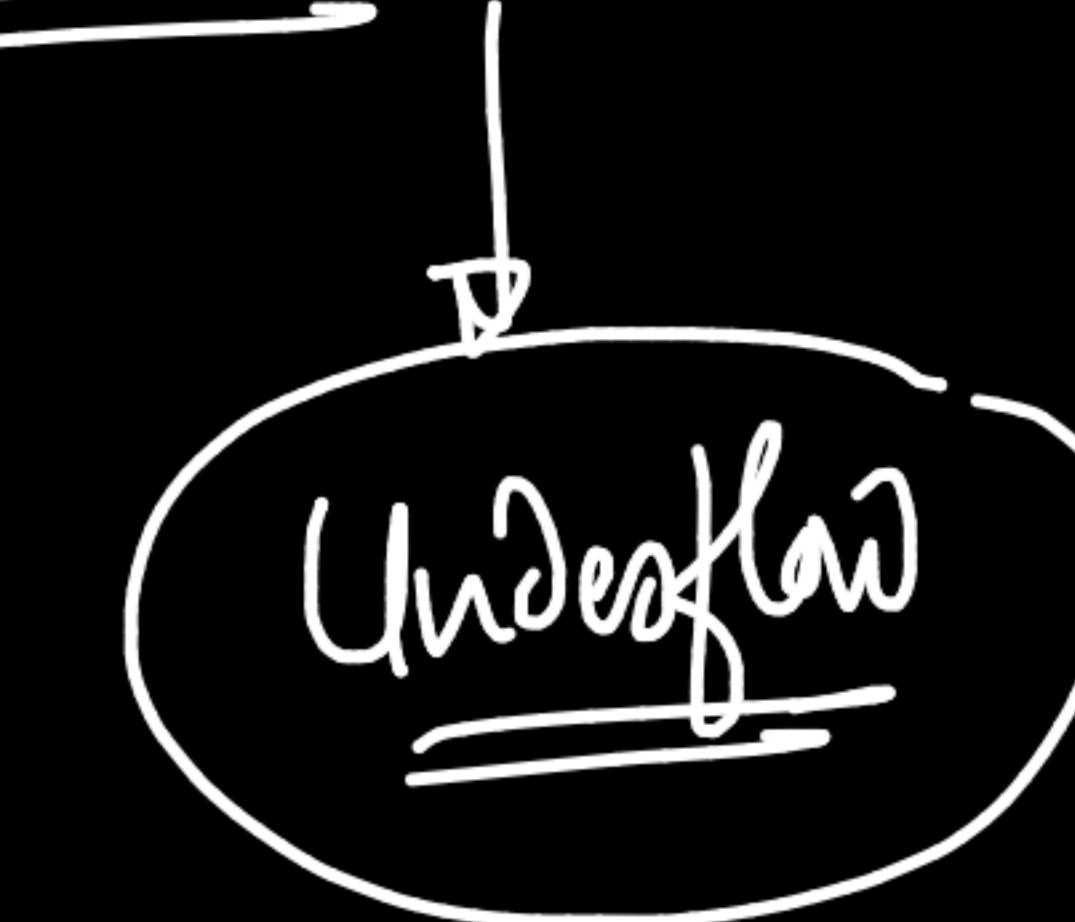
0.000123

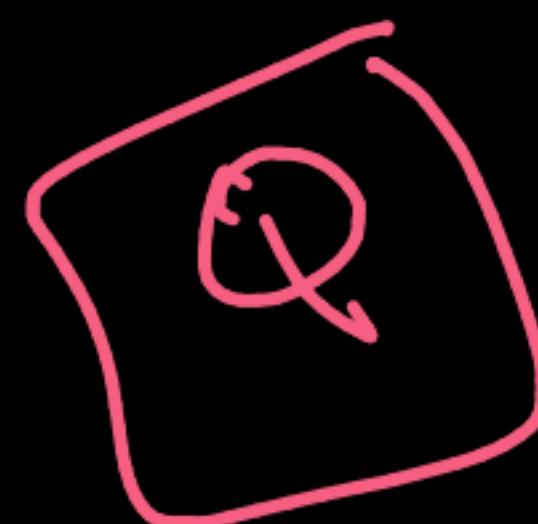
smaller

→ IEEE
FLOAT64

So small

that we
can't store it





① reduce
 $d \rightarrow \times$

How to solve
underflow problem

- ② Scaling:-
- ③ Log ✓

$$\log \left(P(y=1 \mid w_1, w_2, \dots, w_d) \right) = \underbrace{P(y=1)}_{\text{bias}} \cdot \prod_{j=1}^d P(w_j \mid y=1)$$

$$= \text{lp}(y=1) + \sum_{j=1}^d \text{lp}(w_j \mid y=1)$$

$$\begin{aligned} \log(ab) &= \log a + \log b \\ &\quad \text{---} \end{aligned}$$

$$\rightarrow \log(P(y=1 | w_1, \dots, w_d))$$

Compare

$$\hookrightarrow \log(P(y=0 | w_1, \dots, w_d))$$

- $P(y=1 | w_1, \dots, w_d)$
- ↑ Bayes + Naivety
- $P(w_1, w_2, \dots | y=1) = \prod P(w_i | y=1)$
- Recap Bernoulli NB
- $w_j \rightarrow 0$
- $w_j \rightarrow 1$
- Likelihoods & Class prior (P_{train})
- $P(y=1 | x_w) : P(y=0 | x_w) \rightarrow \underline{\underline{\text{NB}}}$

$\rightarrow P(w_j | y=1) \quad w_j \notin D_{\text{train}}$

$$\hookrightarrow \frac{n_j + d}{n_j + kd}$$

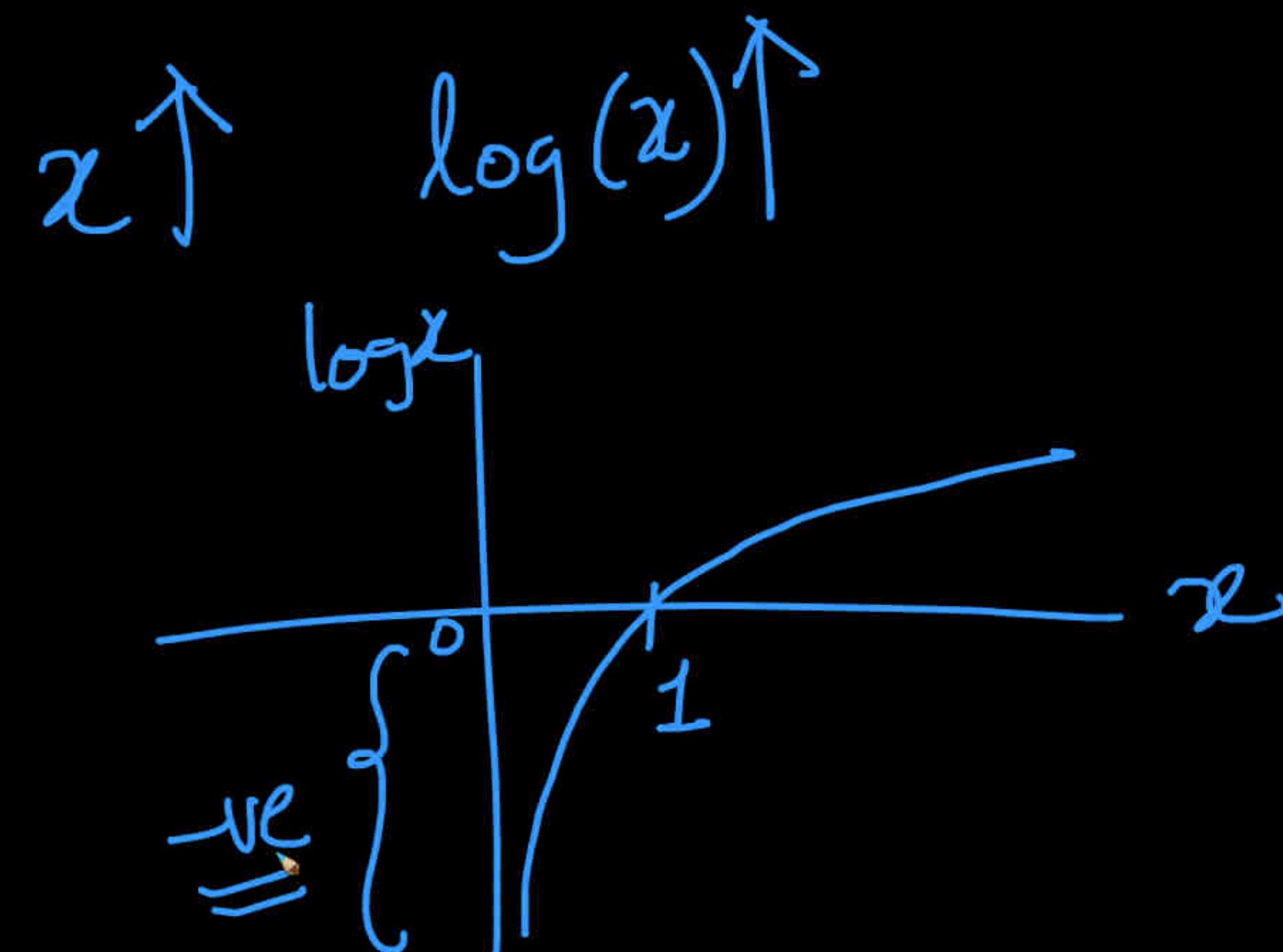
t_2 as w_j takes 2 values

$\rightarrow d \uparrow$ underfit
 $d \downarrow$ overfit

→ $\log \prod$ → replacing products with sum
(Underflow) X

ve values $\log(\prod \dots)$ ve

0 — 1



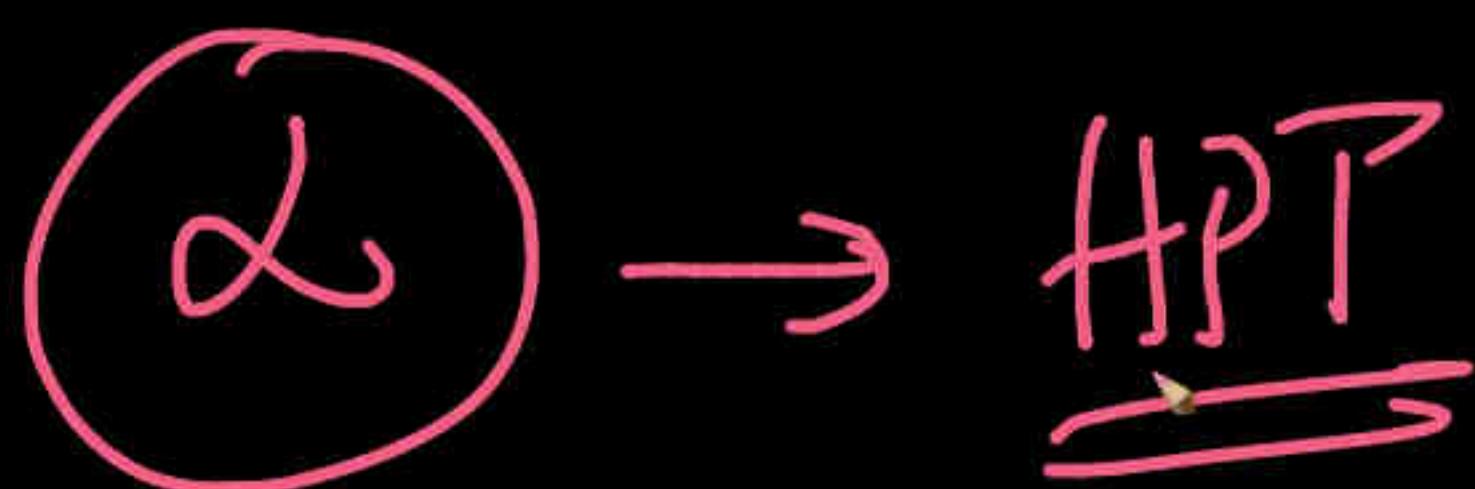
~~Ways to win~~

V.COMMON

mistake

Naive \rightarrow wi & wj are indep of
one another

$$\begin{aligned} p(w_i \underline{w_j} | y=1) \\ = p(w_i | y=1) p(w_j | y=1) \end{aligned}$$





$$P(\omega_j | y=j) = \frac{\eta_{jj} + 1}{\eta_1 + 2x_1}$$

Diagram illustrating the formula:

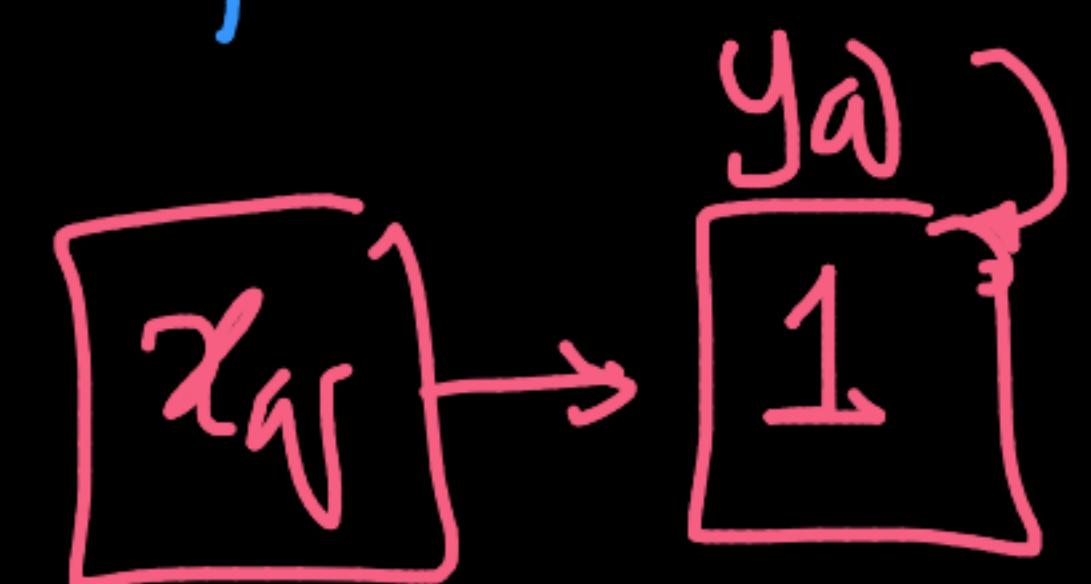
- A pink circle labeled $\alpha : HP$ is connected by arrows to:
 - A green square labeled $\eta_{jj} + 1$.
 - A green square labeled $\eta_1 + 2x_1$.
- A pink circle labeled $y=1$ is connected by arrows to:
 - A green square labeled $\eta_{jj} + 1$.
 - A green square labeled $\eta_1 + 2x_1$.
- A pink circle labeled $y=j$ is connected by arrows to:
 - A green square labeled $\eta_{jj} + 1$.
 - A green square labeled $\eta_1 + 2x_1$.

✓ Class-priors → $P(y=1)$
 $P(y=0)$

✓ likelihoods → $P(w_j | y=1)$
 $P(w_j | y=0)$

F.I & Interpretability

feature(mpg)
interpretability



top likelihoods

$$x_{ij} : w_1, w_3, \overrightarrow{w_6}, w_7 \quad \hat{y}_{ij} = 1$$

$p(w_6 | y=1) = \boxed{0.8}$

$\rightarrow 0.1$
 $\rightarrow 0.6$

Outliers:

@ Test time:-

$$x_{qj} = w_1 w_3 w_5 \underbrace{w_j}_{\text{outlier}} \rightarrow \text{mitigated by Laplace, Smoothing}$$

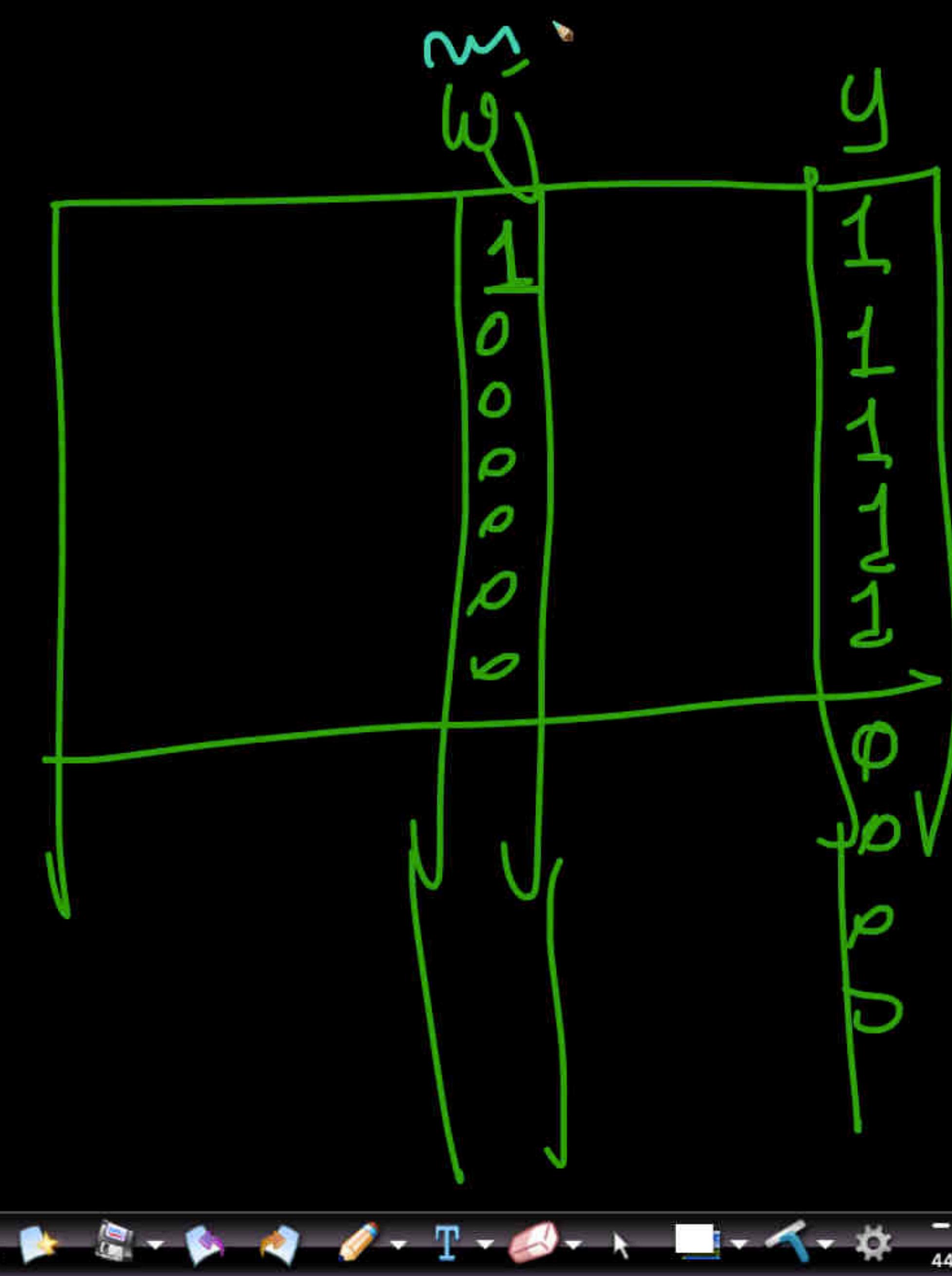
@ Train time

$$P(w_j | \underline{y} = 1) = \frac{n_{j1}}{n_1} = \frac{1 + \alpha}{1000 + 2\alpha}$$

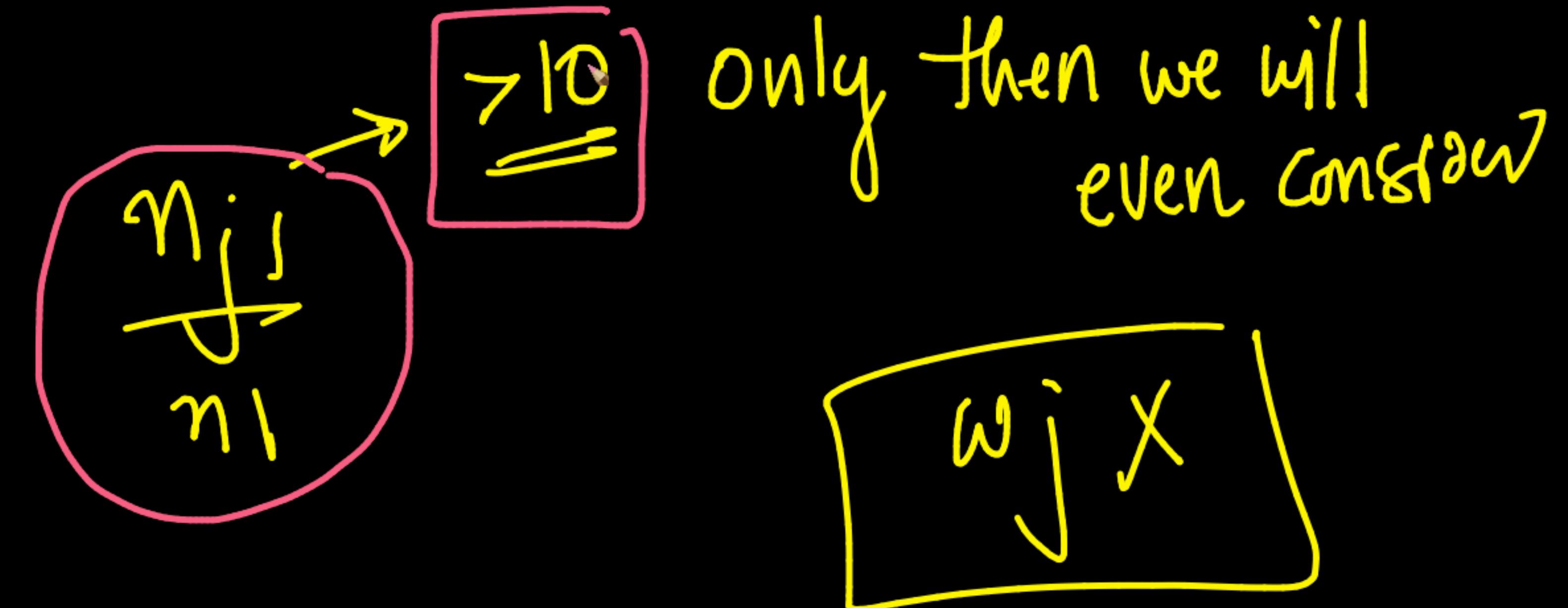
outlier? → false warn / noise

↑
100

{ Laplace smoothing
larger α



hack:

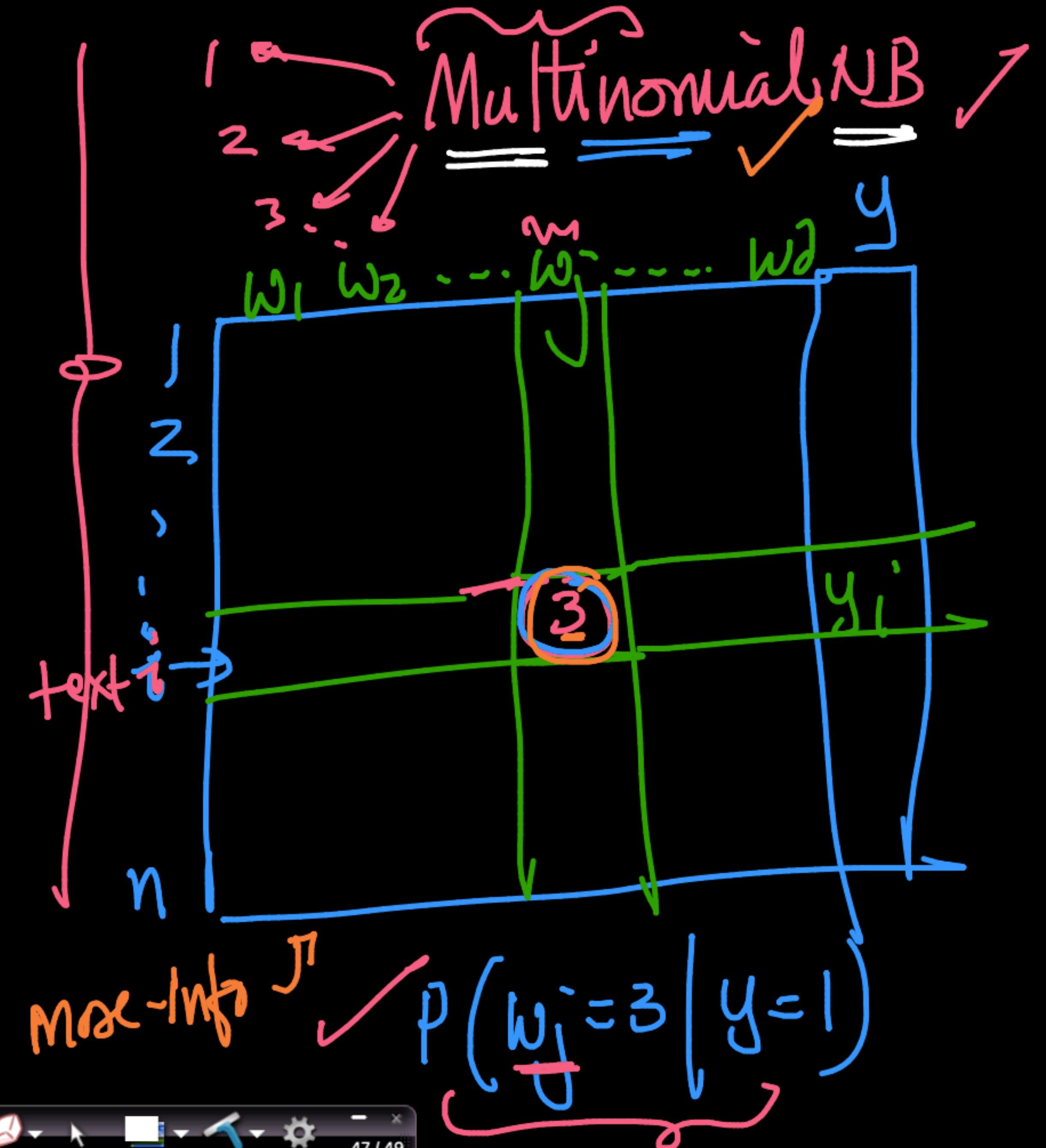
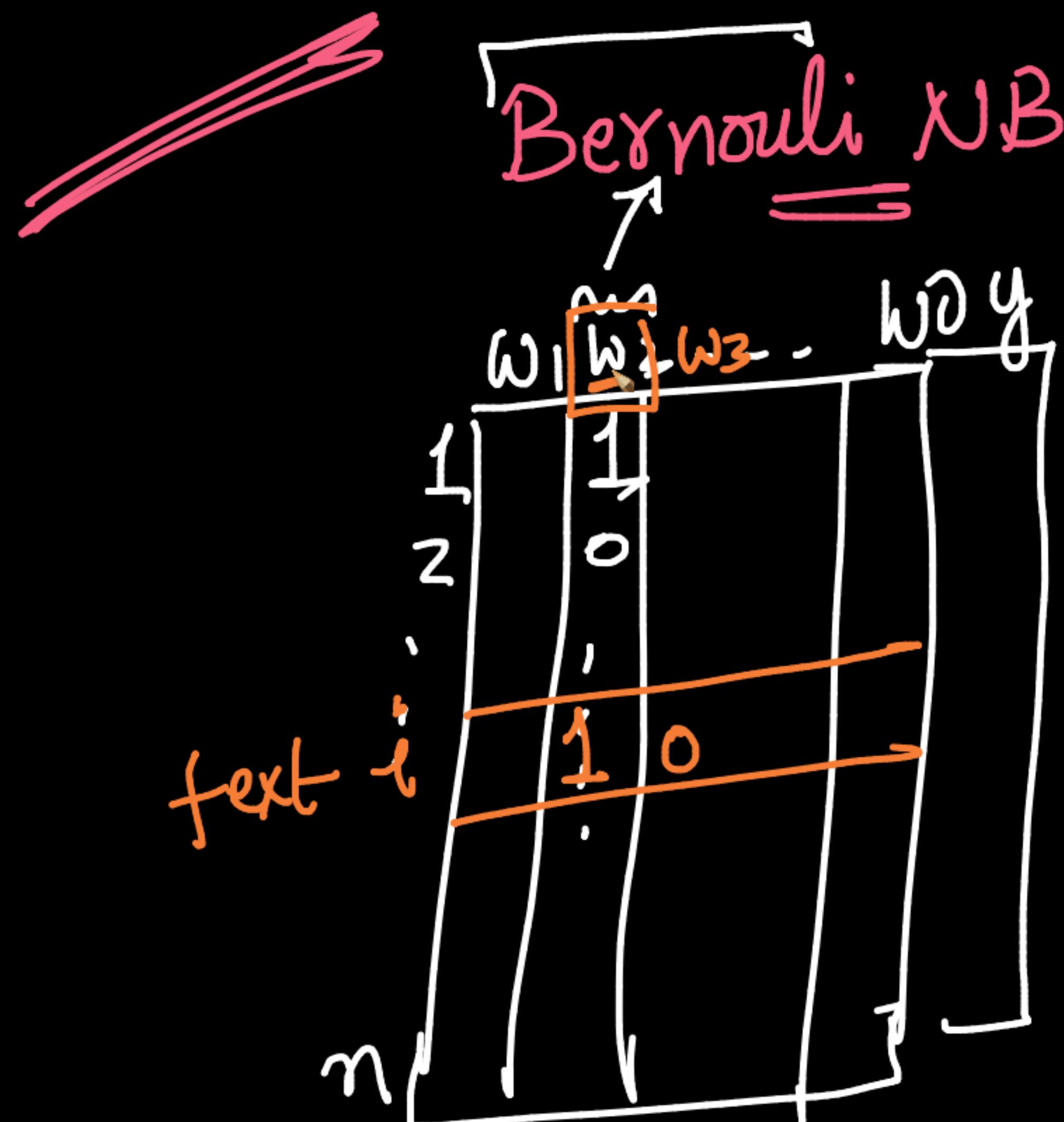


Theoretically: one α per word

Practical → we avoid (one α)

#woots: 10,000

$\alpha_1 \alpha_2 \dots \alpha_{10,000} \rightarrow$ v. hard to
HPT



Bernoulli(p)

$$0 - 1 - p = q$$

$$1 - p$$

Multinomial
(k -values)

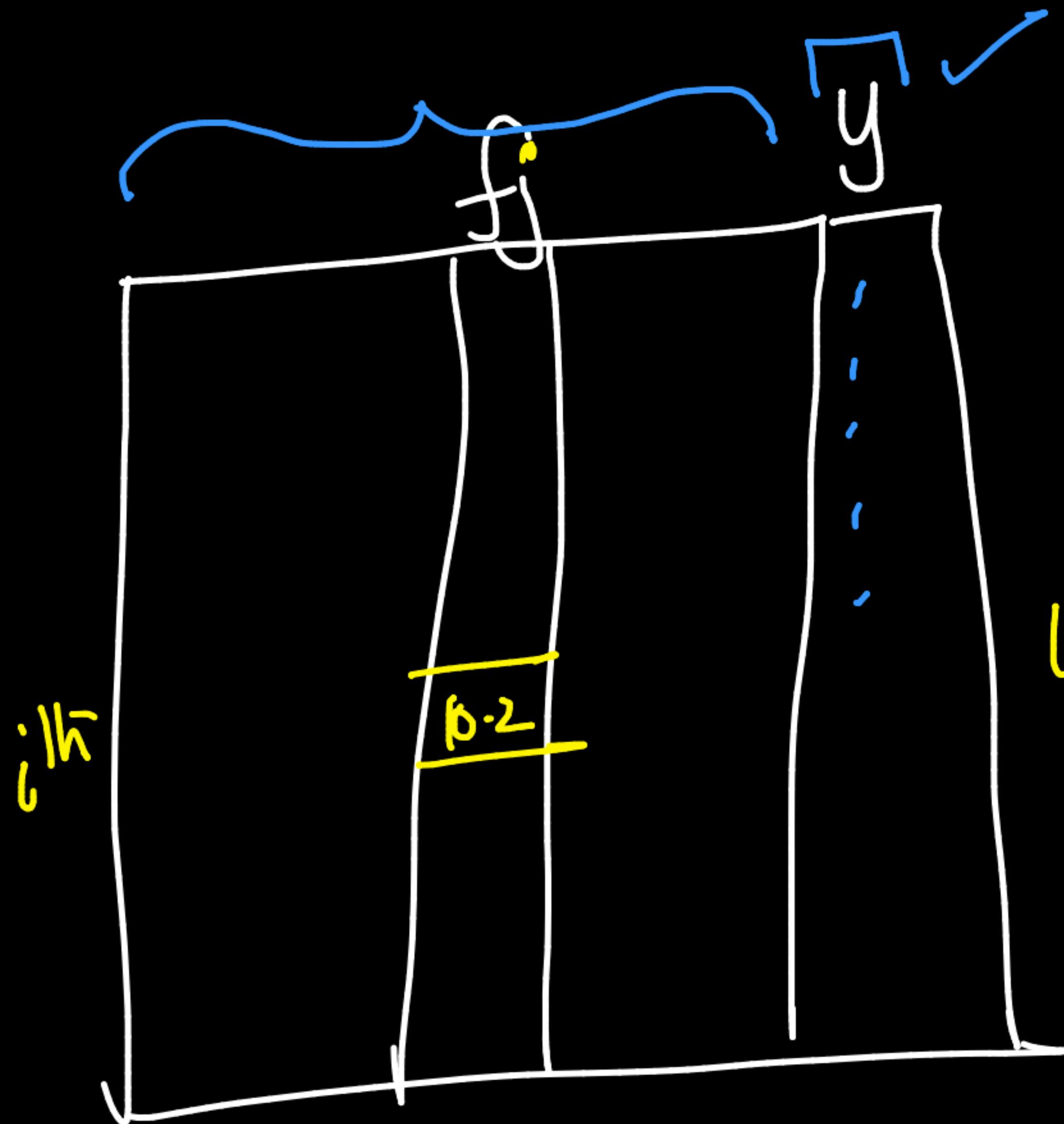
$$p_1, p_2, \dots, p_{k-1}$$

NB $\xrightarrow{\text{Class priors}} p(y=1) \quad p(y=\bar{0})$

$\xrightarrow{\text{Likelihoods}} p(w_j=c | y=1)$

$p(w_j=c | y=0)$

Non
text



f_j : real-valued

$$P(y=1) \& P(y=0)$$

Likelihoods

$$P(f_j = 0.2 | y=1)$$

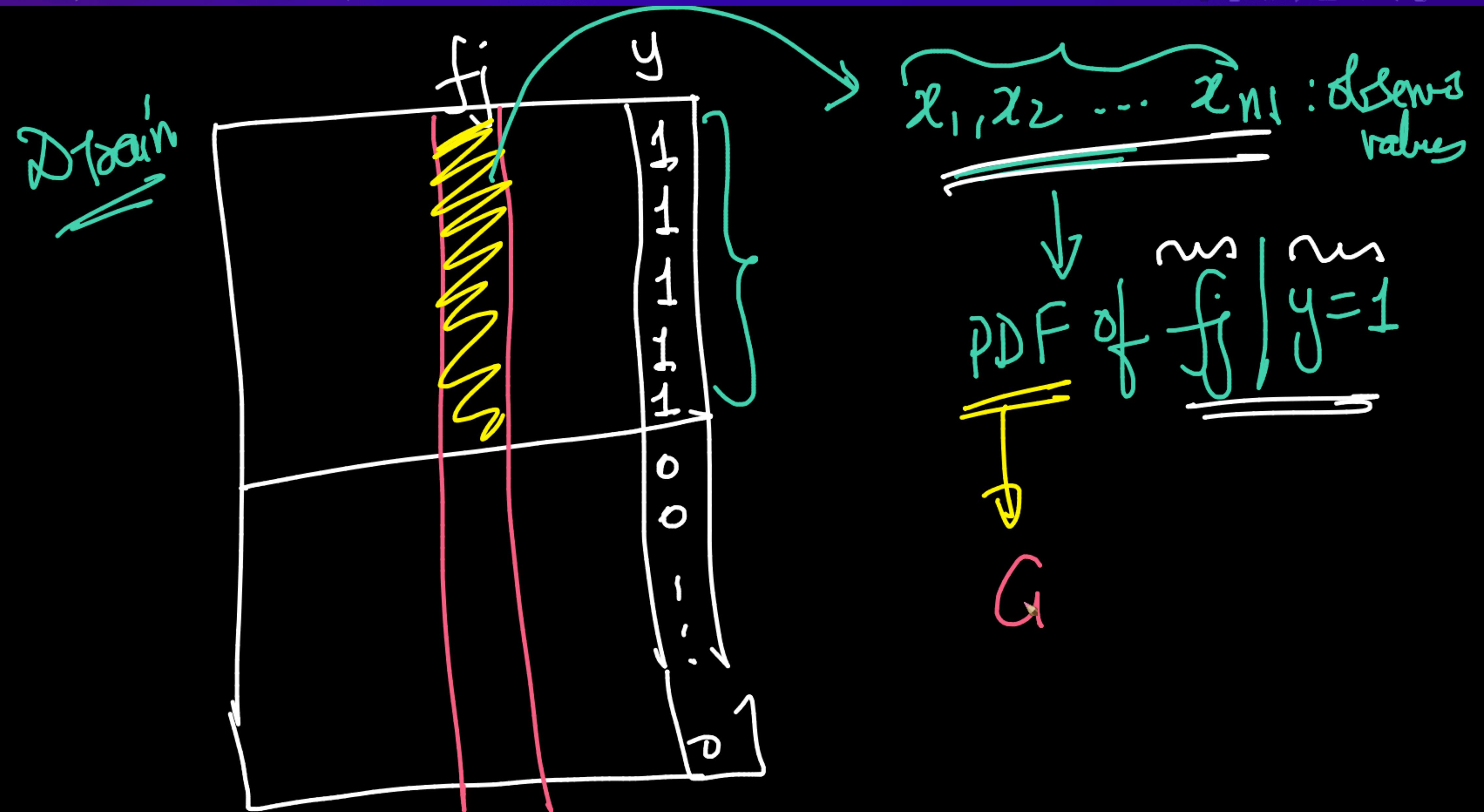
$$P(f_j = 0.2 | y=0)$$

continuous r.v

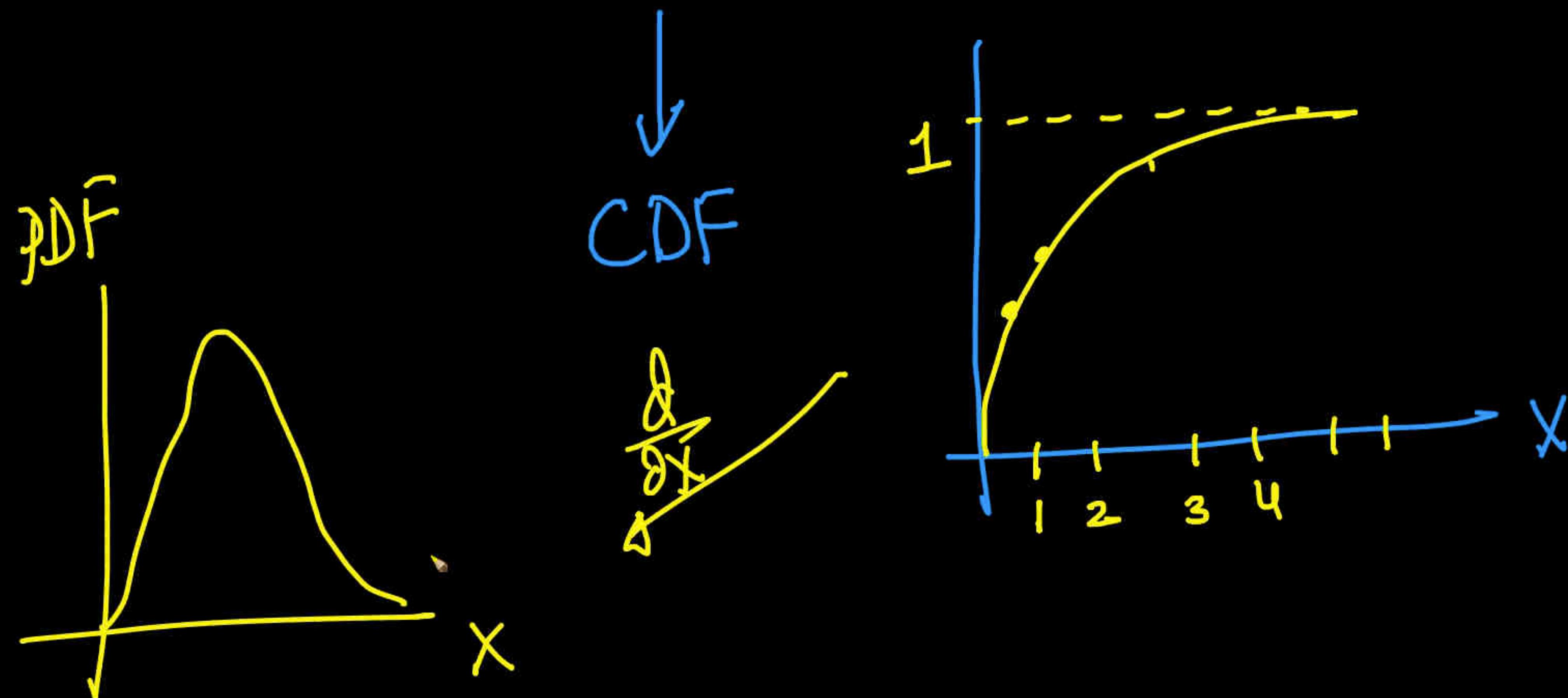
$$P(f_j = 0.2 | y = 1)$$

conditional PDF
pdf of w_j
 $| y = 1$

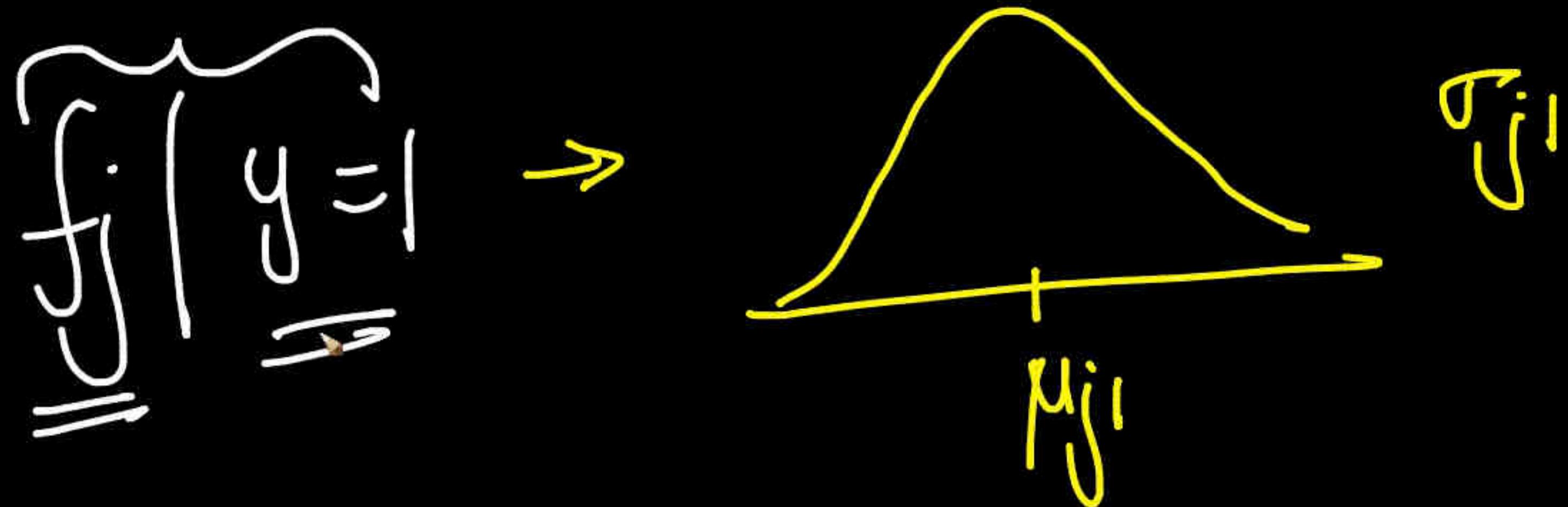
✓ $P(f_j) = \text{PDF of } f_j$



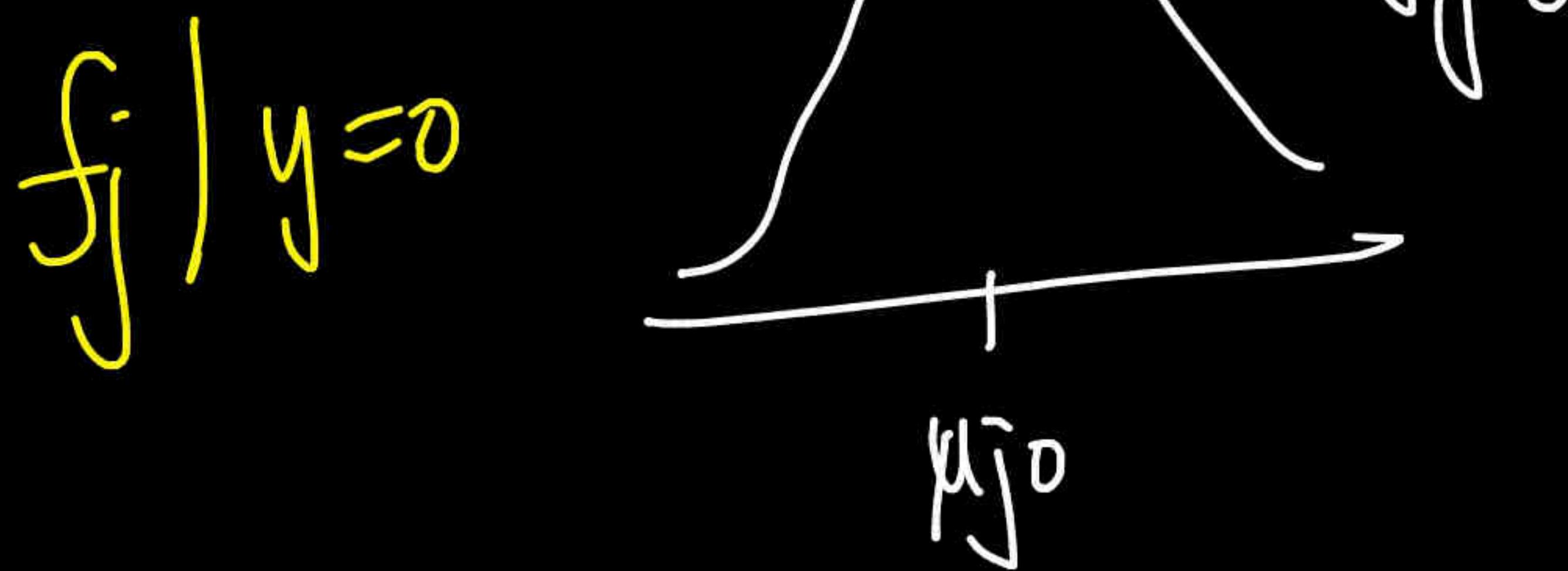
$x_1, x_2, \dots, x_n \rightarrow$ empirical PDF



f_j^-



f_j^-



if all conditional PDFs are gaussian

Gaussian NB

Gaussian NB

Class
points

Likelihoods!

$$P(y=1) \text{ & } P(y=0)$$

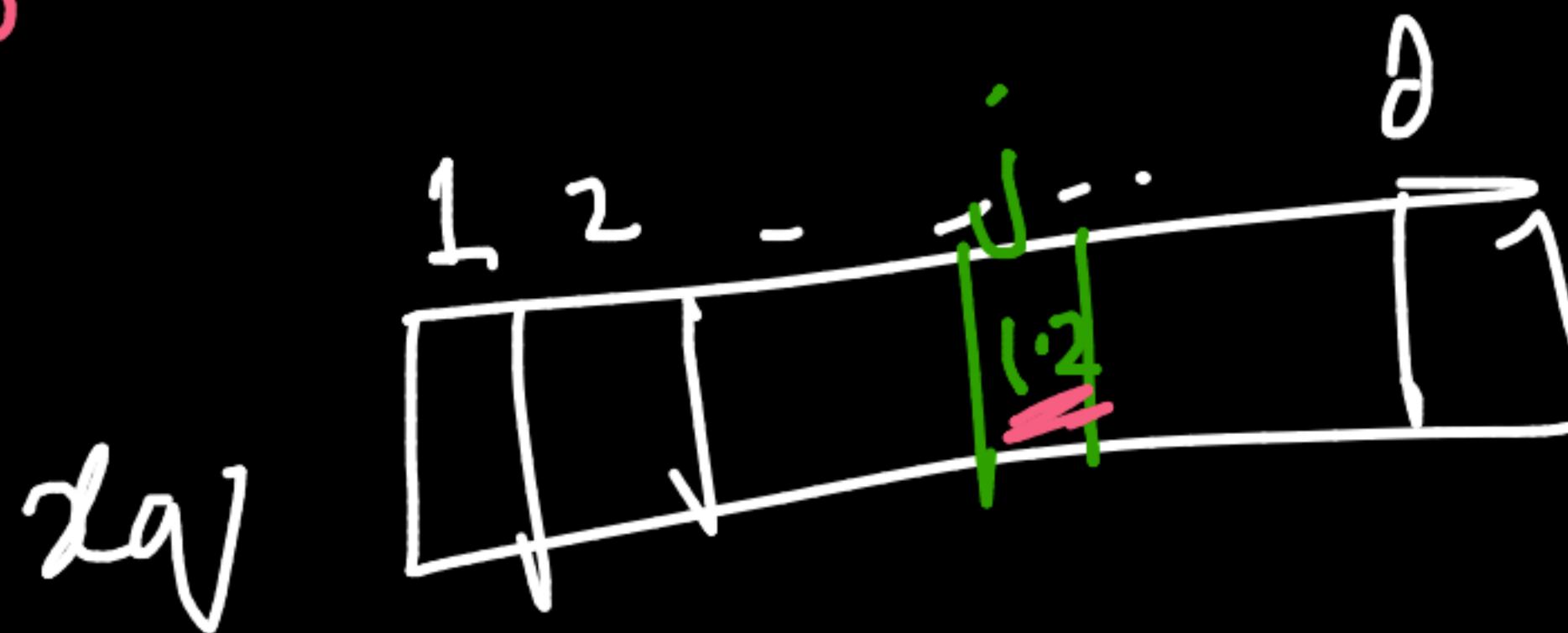
$$f_j | y=1 \rightarrow \underline{\underline{N(\mu_{j1}, \sigma_{j1}^2)}}$$

$$f_j | y=0 \rightarrow \underline{\underline{N(\mu_{j0}, \sigma_{j0}^2)}}$$

Toán:

QS-PLA / SW / AD / FS tests

Test time

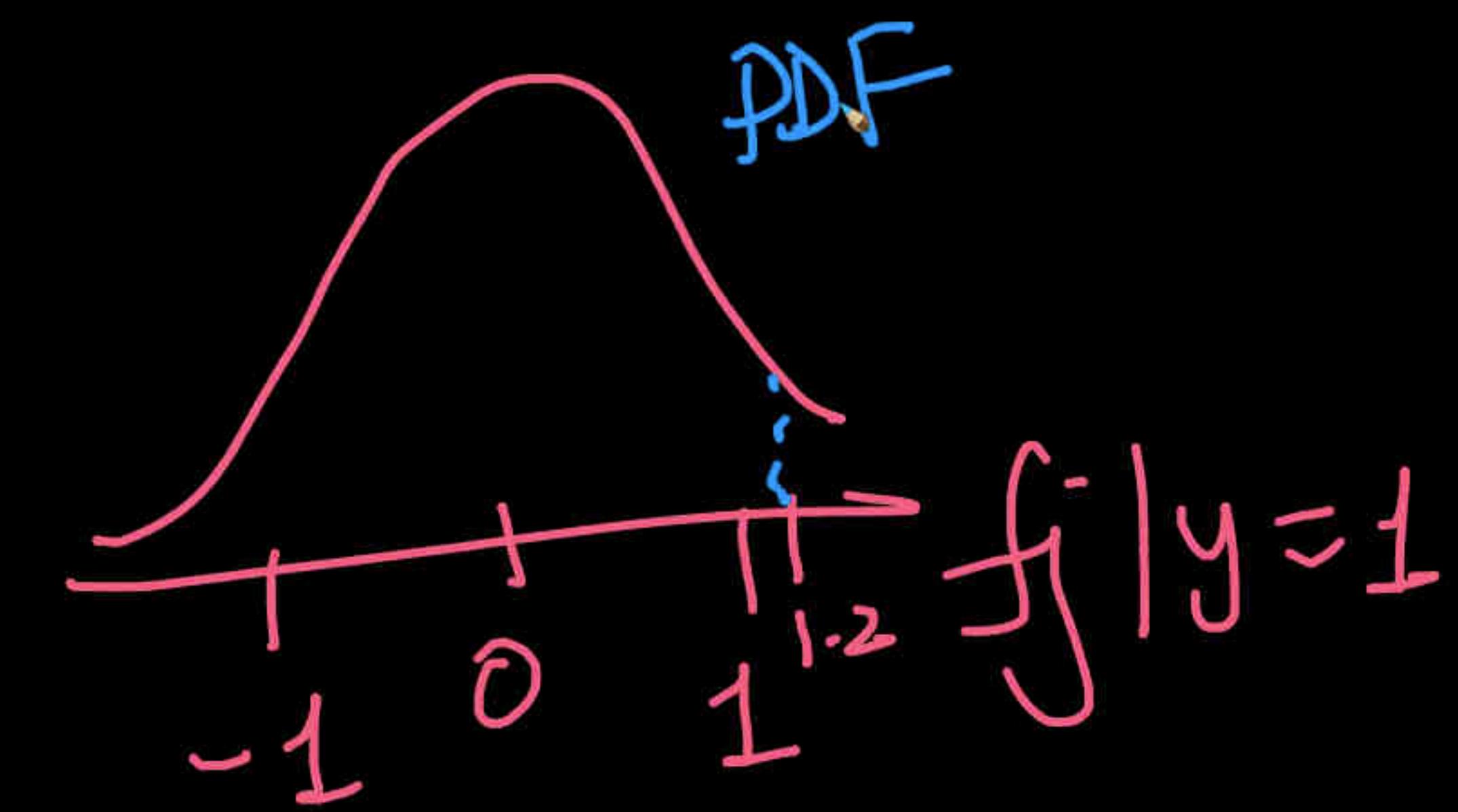


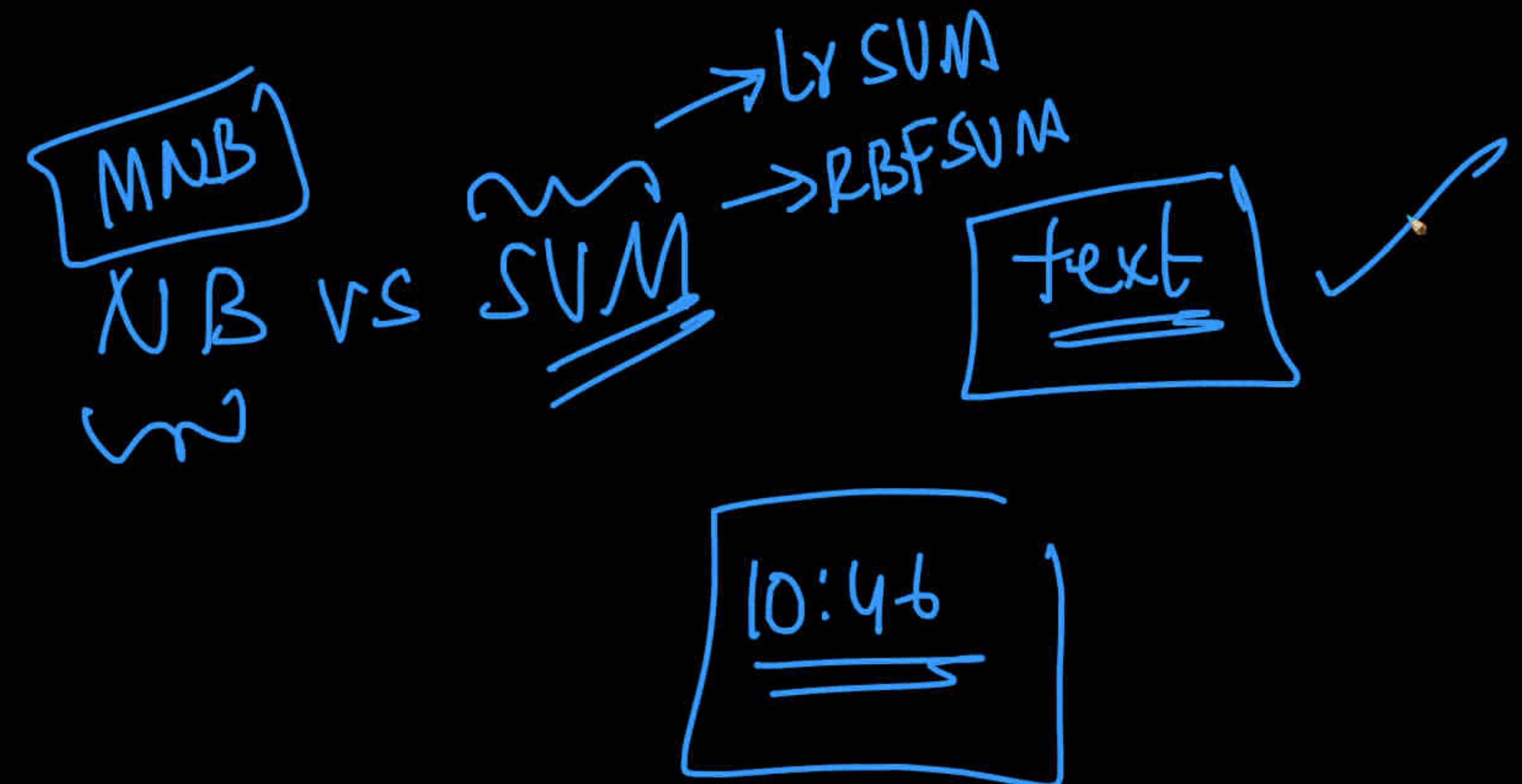
$$P(y=1 | x_w) \propto P(y=1) \prod_{j=1}^d P(w_j | y=1)$$

$$P(y=0 | x_w) \propto P(y=0) \prod_{j=1}^d P(w_j | y=0)$$

$$\begin{aligned} & \xrightarrow{\text{C} \sim N(\mu_j, \sigma_j)} \\ & P(w_j = 1.2 | y=1) \end{aligned}$$

$$\prod_{j=1}^d P(w_j | y=1)$$





+ Code + Text

```
[ ] # Importing libraries

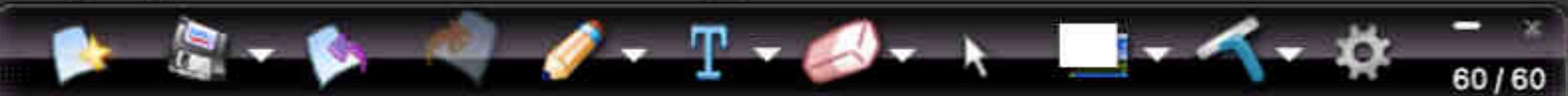
{x}
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from collections import Counter 
from sklearn import feature_extraction, model_selection, naive_bayes, metrics, svm
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler

from sklearn.svm import SVC 
# Libraries for text processing
import re, nltk
nltk.download('punkt')
nltk.download('stopwords')
from nltk import word_tokenize, sent_tokenize
from nltk.corpus import stopwords

import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

> [nltk_data] Downloading package punkt to /root/nltk_data...

> [nltk_data] Package punkt



SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business C x +

colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=XJPNgzCaFvX_



+ Code + Text

Reconnect



```
[ ] from sklearn.model_selection import train_test_split, cross_val_score
[ ] from sklearn.preprocessing import StandardScaler

from sklearn.svm import SVC

# Libraries for text processing
import re, nltk
{nltk.download('punkt')} ←
{nltk.download('stopwords')} ←
from nltk import word_tokenize, sent_tokenize
from nltk.corpus import stopwords

import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
[ ] #Data: https://drive.google.com/file/d/1UEc9IY2HgIAZOsm4qpXdqU05a1ZyeyTV/view?usp=sharing
id = "1UEc9IY2HgIAZOsm4qpXdqU05a1ZyeyTV"
path = "https://docs.google.com/uc?export=download&id=" + id
print(path)
```





+ Code + Text

[] 2022-05-27 14:32:11 (121 MB/s) - 'spam_clean.csv' saved [483640/483640]

{x} [] df = pd.read_csv('../spam_clean.csv', encoding='latin-1')
df.head()

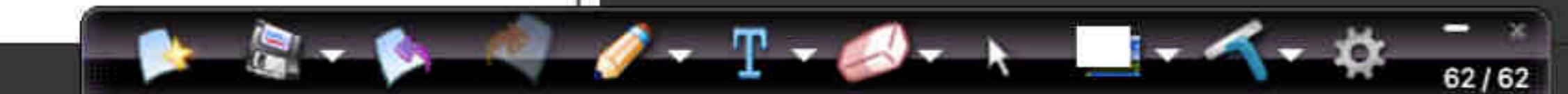
	type	message
0	ham ✓	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

SMS | WhatsApp

< > [] freq = pd.value_counts(df["type"], sort= True)
freq.plot(kind= 'bar', color= ["red", "green"])
plt.title('Distribution of Dependent feature')
plt.show()

Distribution of Dependent feature

5000



SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=XJPNgzCaFvX_ Reconnect

+ Code + Text [] Can say so early now... already done say...

4 ham Nah I don't think he goes to usf, he lives aro...

```
{x} [ ] freq = pd.value_counts(df["type"], sort= True)
freq.plot(kind= 'bar', color= ["red", "green"])
plt.title('Distribution of Dependent feature')
plt.show()
```

Distribution of Dependent feature

5:1

imbalanced data

ham

spam



Text Cleaning & Preprocessing

w₁ w₂ w₃

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

```
▶ df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)
```



type

message

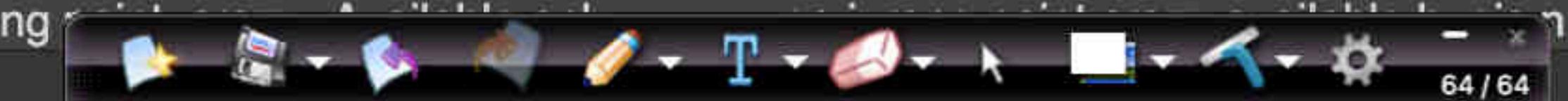
cleaned_message



0 ham

Go until jurong

great ...



+ Code + Text

Reconnect



Text Cleaning & Preprocessing

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

```
▶ df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)
```



type



0

ham

Go until jurong

message

cleaned_message



SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=CN15bHB4xmkp

+ Code + Text Reconnect

Text Cleaning & Preprocessing

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong	great ...

66 / 66

+ Code + Text

Reconnect



Text Cleaning & Preprocessing

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

```
▶ df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)
```



type



0

ham

Go until jurong

message

cleaned_message



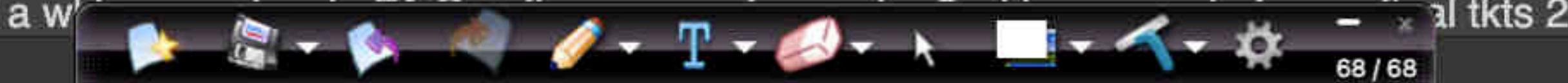


+ Code + Text

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

▶ df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugs n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...



+ Code + Text

Reconnect ▾



```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower()
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word)
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip())
```

```
▶ df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)
```

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugs n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...



colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPi0Wz

+ Code + Text Reconnect

[] def clean_tokenized_sentence(s):
 """Performs basic cleaning of a tokenized sentence"""
 cleaned_s = "" # Create empty string to store processed sentence.
 words = nltk.word_tokenize(s)
 for word in words:
 # Convert to lowercase #
 c_word = word.lower()
 # Remove punctuations #
 c_word = re.sub(r'[^w\s]', ' ', c_word)
 # Remove stopwords #
 if c_word != '' and c_word not in stopwords.words('english'):
 cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
 return(cleaned_s.strip())

Go go

df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugs n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...

70 / 70

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPii0Wz

+ Code + Text Reconnect

[] def clean_tokenized_sentence(s):
 """Performs basic cleaning of a tokenized sentence"""
 cleaned_s = "" # Create empty string to store processed sentence.
 words = nltk.word_tokenize(s)
 for word in words:
 # Convert to lowercase #
 c_word = word.lower()
 # Remove punctuations #
 c_word = re.sub(r'[^w\s]', '', c_word)
 # Remove stopwords #
 if c_word != '' and c_word not in stopwords.words('english'): ✓
 cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
 return(cleaned_s.strip())

{ is the a
of for an
- - until

df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugis n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...

71/71

+ Code + Text

Reconnect ▾

```
# Importing libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from collections import Counter
from sklearn import feature_extraction, model_selection, naive_bayes, metrics, svm
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler

from sklearn.svm import SVC

# Libraries for text processing
import re, nltk
nltk.download('punkt')
nltk.download('stopwords')
from nltk import word_tokenize, sent_tokenize
from nltk.corpus import stopwords
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading pa...



colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPii0Wz

+ Code + Text Reconnect

```
[ ] def clean_tokenized_sentence(s):
    """Performs basic cleaning of a tokenized sentence"""
    cleaned_s = "" # Create empty string to store processed sentence.
    words = nltk.word_tokenize(s)
    for word in words:
        # Convert to lowercase #
        c_word = word.lower() # 2
        # Remove punctuations #
        c_word = re.sub(r'[^w\s]', '', c_word) # 3
        # Remove stopwords #
        if c_word != '' and c_word not in stopwords.words('english'):
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.
    return(cleaned_s.strip()) # 4
```

df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugis n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...

73 / 73

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPii0Wz

+ Code + Text Reconnect

```
[ ] def clean_tokenized_sentence(s):  
    """Performs basic cleaning of a tokenized sentence"""  
    cleaned_s = "" # Create empty string to store processed sentence.  
    words = nltk.word_tokenize(s)  
    for word in words:  
        # Convert to lowercase  
        c_word = word.lower()  
        # Remove punctuations  
        c_word = re.sub(r'[^w\s]', '', c_word)  
        # Remove stopwords  
        if c_word != '' and c_word not in stopwords.words('english'):  
            cleaned_s = cleaned_s + " " + c_word # Append processed words to new list.  
    return(cleaned_s.strip())  
  
df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)  
df.head(10)
```

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugs n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a w...	al tkts 2...

Reconnect

Code Text

Search

File

Code

Text

Reconnect

Profile

Settings

Help

10:57

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPii0Wz

+ Code + Text Reconnect

```
[ ] for word in words:  
    # Convert to lowercase #  
    c_word = word.lower()  
    # Remove punctuations #  
    c_word = re.sub(r'[^\\w\\s]', '', c_word)  
    # Remove stopwords #  
    if c_word != '' and c_word not in stopwords.words('english'):  
        cleaned_s = cleaned_s + " " + c_word    # Append processed words to new list.  
return(cleaned_s.strip())
```

df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugis n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	free entry 2 wkly comp win fa cup final tkts 2...
3	ham	U dun say so early hor... U c already then say...	u dun say early hor u c already say
4	ham	Nah I don't think he goes to usf, he lives aro...	nah nt think goes usf lives around though
5	spam	FreeMsg Hey there darling it's been 3 week's n...	freemsg hey darling 3 week word back like fun ...

75 / 75

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=64NrpQPii0Wz

+ Code + Text Reconnect

```
[ ] c_word = re.sub(r'^\w\s]', '', c_word)
# Remove stopwords #
if c_word != '' and c_word not in stopwords.words('english'):
    cleaned_s = cleaned_s + " " + c_word    # Append processed words to new list.
return(cleaned_s.strip())
```

```
{x} [ ] df["cleaned_message"] = df["message"].apply(clean_tokenized_sentence)
df.head(10)
```

	type	message	cleaned_message
0	ham	Go until jurong point, crazy.. Available only ...	go jurong point crazy available bugis n great ...
1	ham	Ok lar... Joking wif u oni...	ok lar joking wif u oni
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	free entry 2 wkly comp win fa cup final tkts 2...
3	ham	U dun say so early hor... U c already then say...	u dun say early hor u c already say
4	ham	Nah I don't think he goes to usf, he lives aro...	nah nt think goes usf lives around though
5	spam	FreeMsg Hey there darling it's been 3 week's n...	freemsg hey darling 3 week word back like fun ...
6	ham	Even my brother is not like to speak with me. ...	even brother like speak treat like aids patent
7	ham	As per your request 'Melle Melle (Oru Minnamin...	per request melle melle oru minnaminunginte nu...
8	spam	WINNER!! As a valued network customer you have	winner valued network customer selected receiv...

76 / 76

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=FQtdmZxTxt90

+ Code + Text Reconnect

Most common words in Spam and Ham

```
[ ] Counter(" ".join(df[df['type']=='ham']["cleaned_message"]).split())
```

'luv': 30,
'kettoda': 4,
'manda': 4,
'ups': 1,
'3days': 1,
'shipping': 3,
'takes': 12,
'2wks': 1,
'usps': 1,
'lag': 2,
'bribe': 1,
'nipost': 1,
'lem': 8,
'necessarily': 2,
'expect': 4,
'headin': 2,
'mmm': 4,
'jolt': 2,
'suzy': 1,

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=FQtdmZxTxt90

+ Code + Text Reconnect

Reconnect

Most common words in Spam and Ham

```
[ ] Counter(".join(df[df['type']=='ham']['cleaned_message']).split())")
```

luv: 30
kettoda: 4
manda: 4
ups: 1
3days: 1
shipping: 3
takes: 12
2wks: 1
usps: 1
lag: 2
bribe: 1
nipost: 1
lem: 8
necessarily: 2
expect: 4
headin: 2
mmm: 4
jolt: 2
suzy: 1

- non-spam words

- cleaned message

- splitting

$w_1: 10$

$w_2: 20$

$w_3: 26$

⋮

78 / 78

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=FQtdmZxTxt90

+ Code + Text Reconnect

Code:

```
gata : 1,  
'b': 56,  
-----  
counter_ham = Counter(" ".join(df[df['type'] == 'ham']["cleaned_message"]).split()).most_common(20)  
df_ham = pd.DataFrame.from_dict(counter_ham)  
df_ham = df_ham.rename(columns={0:"words in non-spam", 1:"count"})  
  
df_ham.plot.bar(legend = False)  
y_pos = np.arange(len(df_ham["words in non-spam"]))  
plt.xticks(y_pos, df_ham["words in non-spam"])  
plt.title('More frequent words in non-spam messages')  
plt.xlabel('words')  
plt.ylabel('number')  
plt.show()
```

Output:

The chart displays the frequency distribution of words in non-spam messages. The most frequent word appears to be 'the', followed by other common stop words like 'and', 'a', 'in', etc.



+ Code + Text

```
[ ] plt.ylabel('number')  
plt.show()
```

✓



EDA

```
[ ] counter_spam = Counter(" ".join(df[df['type']=='spam']["cleaned_message"])).split()).most_common(20)  
df_spam = pd.DataFrame.from_dict(counter_spam)  
df_spam = df_spam.rename(columns={0:"words in spam", 1:"count_"})  
  
df_spam.plot.bar(legend = False, color = 'orange')
```



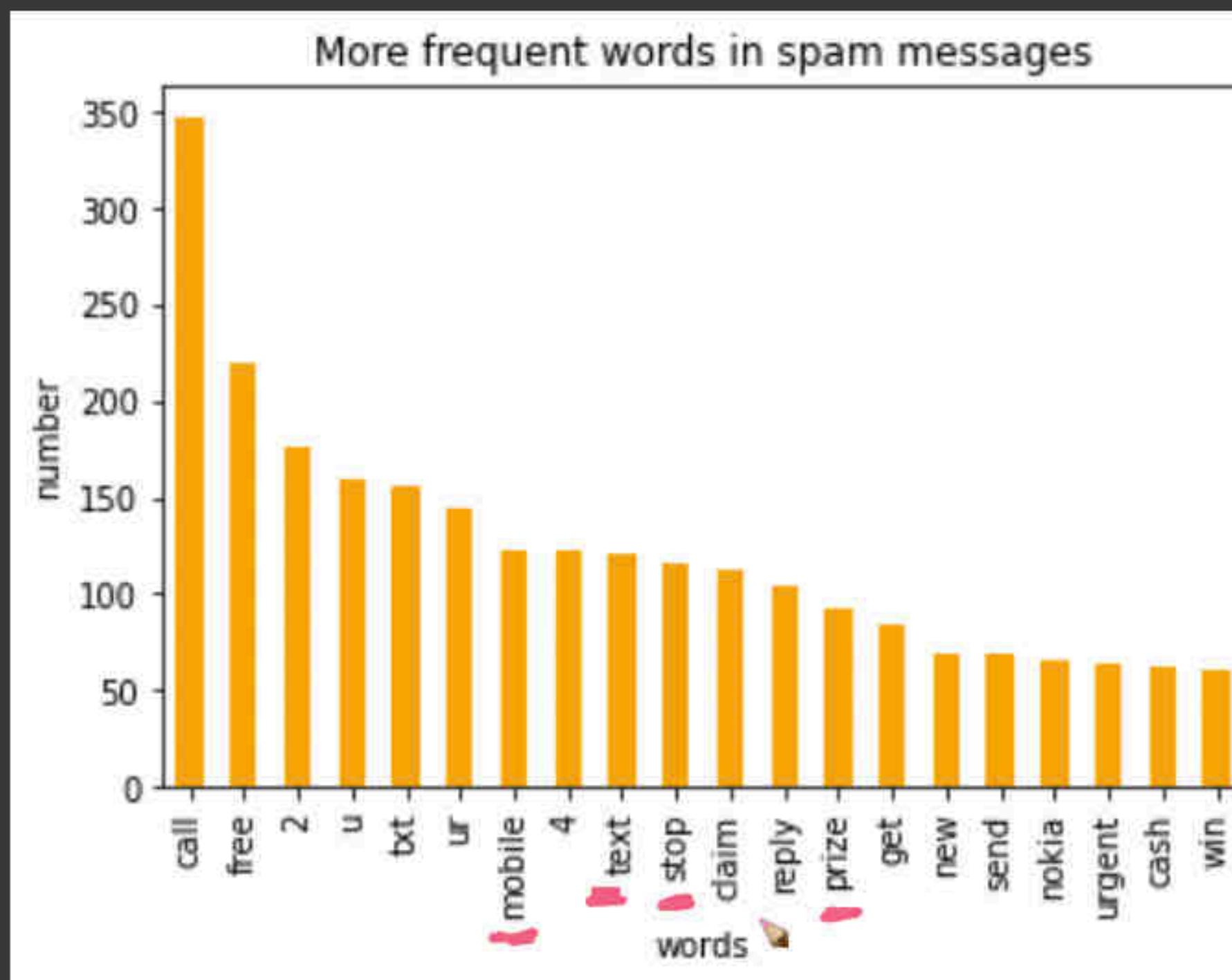
+ Code + Text

```
[ ] df_spam.plot.bar(legend = False, color = 'orange')
y_pos = np.arange(len(df_spam["words in spam"]))
plt.xticks(y_pos, df_spam["words in spam"])
plt.title('More frequent words in spam messages')
plt.xlabel('words')
plt.ylabel('number')
plt.show()
```



+ Code + Text

```
[ ] df_spam.plot.bar(legend = False, color = 'orange')
y_pos = np.arange(len(df_spam["words in spam"]))
plt.xticks(y_pos, df_spam["words in spam"])
plt.title('More frequent words in spam messages')
plt.xlabel('words')
plt.ylabel('number')
plt.show()
```



SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=BO_5WZ12x01X

+ Code + Text Reconnect  

[]



{x}

Generate Datasets

[] df["type"] = df["type"].map({'spam':1,'ham':0})

[] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

[] #Count Vectorizer
f = feature_extraction.text.CountVectorizer()

<> X_train = f.fit_transform(df_X_train)
[] X_test = f.transform(df_X_test)

83 / 83

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=BO_5WZ12x01X

+ Code + Text Reconnect

[]

{x}

Generate Datasets

K-fold CV

75% Train

25% Test

df["type"] = df["type"].map({'spam':1,'ham':0})

[] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

[] #Count Vectorizer
f = feature_extraction.text.CountVectorizer()
X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC - scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes - scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=BO_5WZ12x01X

+ Code + Text Reconnect  

[]



{x}

Generate Datasets

df["type"] = df["type"].map({'spam':1,'ham':0})

[] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

[] #Count Vectorizer
f = feature_extraction.text.CountVectorizer()

X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

85 / 85

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

{x} [] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

```
#Count Vectorizer
f = feature_extraction.text.CountVectorizer()
```

```
x_train = f.fit_transform(df_X_train)
x_test = f.transform(df_X_test)
```

```
# Standard Scaler
scaler = StandardScaler()
x_train = scaler.fit_transform(x_train)
x_test = scaler.transform(x_test)
```

ValueError

Traceback (most recent call last)

<ipython-input-46-ae0768e59d91> in <module>()

86 / 86

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

{x} [] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

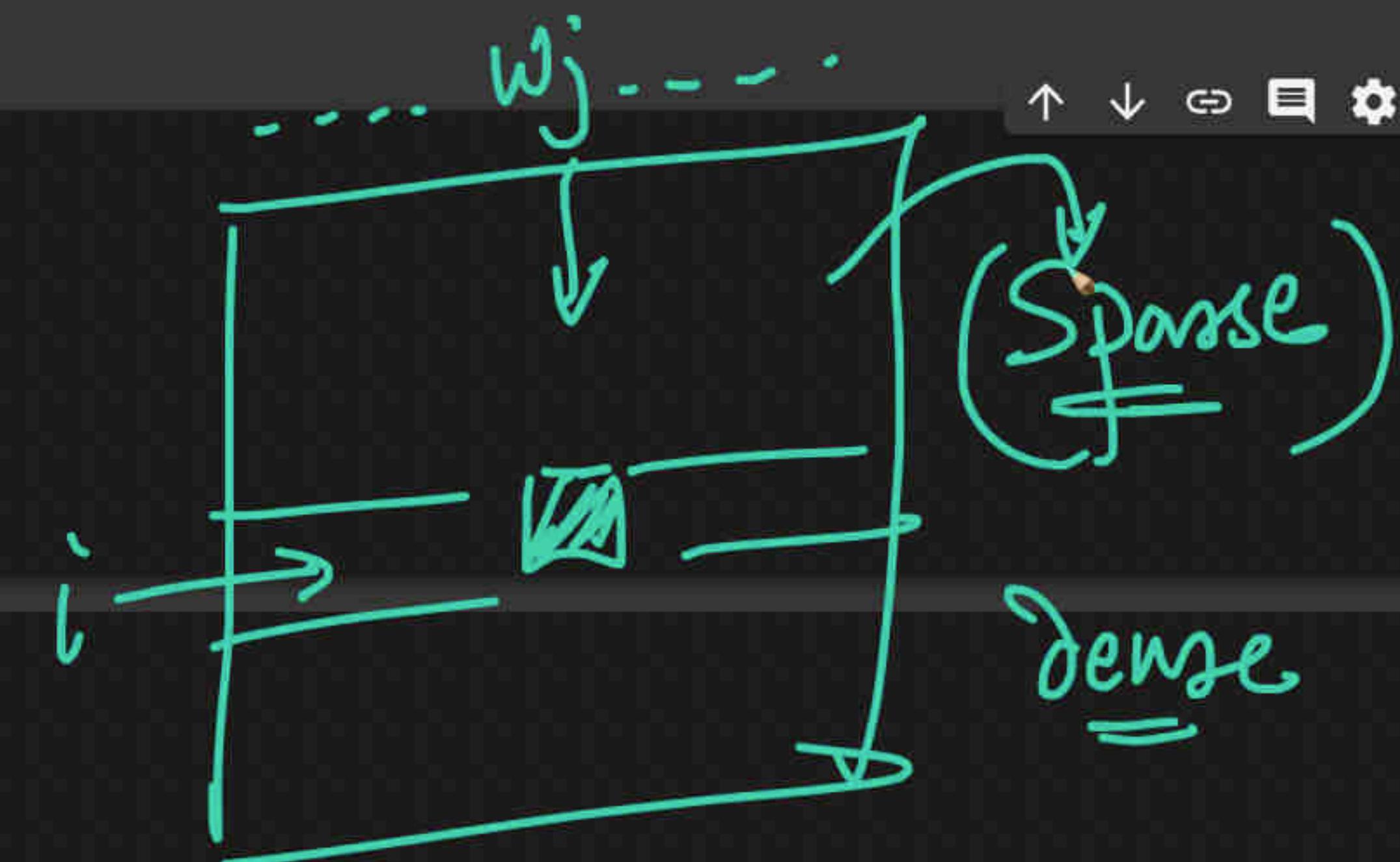
#Count Vectorizer
f = feature_extraction.text.CountVectorizer()

X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

Standard Scaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

ValueError

Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()



87 / 87

SVM_NB.ipynb - Colaboratory | sklearn.preprocessing.StandardScaler | sklearn.svm.SVC — scikit-learn | sklearn.naive_bayes.BernoulliNB | sklearn.naive_bayes.MultinomialNB | 1.9. Naive Bayes — scikit-learn | dsml-course/SVM Business Case | +

+ Code + Text

Reconnect ▾



```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

```
[ ] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])
```

$[(4179,), (1393,)]$

↑ ↓ ⌂ ☰ ⚙ ⌛ 🗑 ⋮

#Count Vectorizer

```
f = feature_extraction.text.CountVectorizer()
```

```
X_train = f.fit_transform(df_X_train)
```

```
X test = f.transform(df X test)
```

▶ # Standard Scaler

```
scaler = StandardScaler()
```

```
x_train = scaler.fit_transform(x_train)
```

```
X test = scaler.transform(X test)
```

→

Traceback (most recent call last)

```
<ipython-input-46-ae0768e59d91> in <module>()
```

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect  

```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

{x} [] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

 #Count Vectorizer 
f = feature_extraction.text.CountVectorizer()

X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

 # Standard Scaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()

89 / 89

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC -- scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes -- scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect  

```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

{x} [] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

 #Count Vectorizer
f = feature_extraction.text.CountVectorizer()
 X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

 # Standard Scaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

<>  ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()

90 / 90

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business Case x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

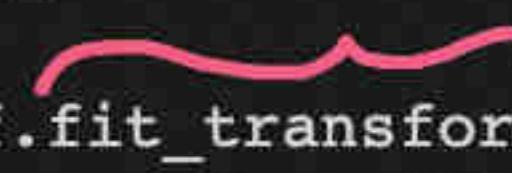
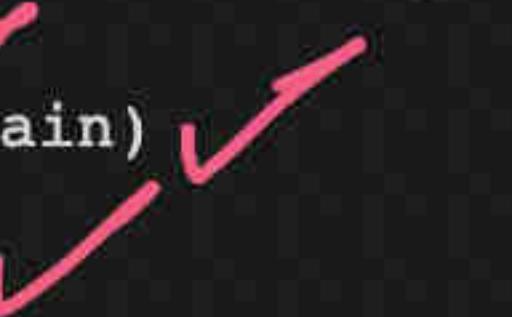
+ Code + Text Reconnect  

```
[ ] df["type"] = df["type"].map({'spam':1,'ham':0})
```

{x} [] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_state=42)
print([np.shape(df_X_train), np.shape(df_X_test)])

[(4179,), (1393,)]

 #Count Vectorizer
f = feature_extraction.text.CountVectorizer()

X_train = f.fit_transform(df_X_train) 
X_test = f.transform(df_X_test) 

 # Standard Scaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

 ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()

91/91

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text

```
[ ] print([np.shape(df_X_train), np.shape(df_X_test)])
```

[(4179,), (1393,)]

{x}

Count Vectorizer
f = feature_extraction.text.CountVectorizer()

X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

Standard Scaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()
 1 # Standard Scaler
 2 scaler = StandardScaler()
----> 3 X_train = scaler.fit_transform(X_train)
 4 X_test = scaler.transform(X_test)

SUM

SH

MNB

NB

$P(\omega_j = \text{?} | y = \text{?})$

Reconnect

92 / 92

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

#Count Vectorizer
f = feature_extraction.text.CountVectorizer()

{x} X_train = f.fit_transform(df_X_train)
X_test = f.transform(df_X_test)

[] # Standard Scaler
{scaler = StandardScaler()
X_train = scaler.fit_transform(X_train) ✓
X_test = scaler.transform(X_test) ✓

fit → (el-means)
& Shi-de ✓

Shi → train
transform → ...

ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()
 1 # Standard Scaler
 2 scaler = StandardScaler()
----> 3 X_train = scaler.fit_transform(X_train)
 4 X_test = scaler.transform(X_test)

2 frames /usr/local/lib/python3.7/dist-packages/sklearn/preprocessing/_data.py in partial_fit(self, X, y, sample_weight)
 869 if self.with_mean:
 870 raise ValueError(
 871 "Cannot center sparse matrices. Consider using with_mean=False or use" 93 / 93

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

```
[ ] X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

{x} {x} ValueError Traceback (most recent call last)
<ipython-input-46-ae0768e59d91> in <module>()
 1 # Standard Scaler
 2 scaler = StandardScaler()
--> 3 X_train = scaler.fit_transform(X_train)
 4 X_test = scaler.transform(X_test)

2 frames /usr/local/lib/python3.7/dist-packages/sklearn/preprocessing/_data.py in partial_fit(self, X, y, sample_weight)
 869 if self.with_mean:
 870 raise ValueError(
--> 871 "Cannot center sparse matrices: pass `with_mean=False` "
 872 "instead. See docstring for motivation and alternatives."
 873)

ValueError: Cannot center **sparse matrices**. pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

```
[ ] # Standard Scaler
scaler = StandardScaler(with_mean=False) # Problems with dense matrix
```

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text

Reconnect

Code Editor

870 raise ValueError(
871 "Cannot center sparse matrices: pass `with_mean=False` "
872 "instead. See docstring for motivation and alternatives."
873)

{x} **ValueError: Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.**

SEARCH STACK OVERFLOW

[] # Standard Scaler
[] {scaler = StandardScaler(with_mean=False) # Problems with dense matrix
[] X_train = scaler.fit_transform(X_train)
[] X_test = scaler.transform(X_test.todense())
[] print([np.shape(X_train), np.shape(X_test)])
[] ((4179, 7692), (1393, 7692))
[] type(X_train)
[] **scipy.sparse.csr.csr_matrix**

95 / 95

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text

Reconnect

Code Editor

870 raise ValueError(
871 "Cannot center sparse matrices: pass `with_mean=False` "
872 "instead. See docstring for motivation and alternatives."
873)

{x} ValueError: Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

[] # Standard Scaler
scaler = StandardScaler(with_mean=False) # Problems with dense matrix
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test.todense())

[] print([np.shape(X_train), np.shape(X_test)])
[(4179, 7692), (1393, 7692)]

[] type(X_train)

scipy.sparse.csr.csr_matrix

bts of memory

Sparse → Memory

96 / 96

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

[] 870 raise ValueError(
[] --> 871 "Cannot center sparse matrices: pass `with_mean=False` "
[] 872 "instead. See docstring for motivation and alternatives."
[] 873)

{x} **ValueError:** Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

```
[ ] # Standard Scaler  
scaler = StandardScaler(with_mean=False) # Problems with dense matrix  
X_train = scaler.fit_transform(X_train)  
X_test = scaler.transform(X_test.todense())
```

```
[ ] print([np.shape(X_train), np.shape(X_test)])
```

```
[ ] [(4179, 7692), (1393, 7692)]
```

```
[ ] type(X_train)
```

```
[ ] <__main__.X_train>: scipy.sparse.csr.csr_matrix
```

scipy.sparse.csr.csr_matrix

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text

[] raise ValueError(
 "Cannot center sparse matrices: pass `with_mean=False` "
 "instead. See docstring for motivation and alternatives."
)

{x} ValueError: Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

[] # Standard Scaler
scaler = StandardScaler(with_mean=False) # Problems with dense matrix
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test.todense())

[] print([np.shape(X_train), np.shape(X_test)])

[(4179, 7692), (1393, 7692)]

[] type(X_train)

scipy.sparse.csr.csr_matrix

BernoulliNB
MultinomialNB

NoScaling

Scaling

Gaussian NB

f N

SVM_NB.ipynb - Colaboratory | sklearn.preprocessing.StandardScaler | sklearn.svm.SVC — scikit-learn | sklearn.naive_bayes.BernoulliNB | sklearn.naive_bayes.MultinomialNB | 1.9. Naive Bayes — scikit-learn | dsml-course/SVM Business Case | +

colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text

Reconnect

Code Editor

Cell Output

File Cell Kernel Help

870 raise ValueError(
871 "Cannot center sparse matrices: pass `with_mean=False` "
872 "instead. See docstring for motivation and alternatives."
873)

{x} **ValueError:** Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

[] # Standard Scaler
scaler = StandardScaler(with_mean=False) # Problems with dense matrix
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test.todense())

[] print([np.shape(X_train), np.shape(X_test)])

[(4179, 7692), (1393, 7692)]

[] type(X_train)

scipy.sparse.csr.csr_matrix

GNB:-
 $f_j | y=j \sim N(\mu_j, \sigma_j^2)$

99 / 99

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=o3yHdbvxyYeY

+ Code + Text Reconnect

[] --> 871 "Cannot center sparse matrices: pass `with_mean=False`"
872 "instead. See docstring for motivation and alternatives."
873)

{x} ValueError: Cannot center sparse matrices: pass `with_mean=False` instead. See docstring for motivation and alternatives.

SEARCH STACK OVERFLOW

```
[ ] # Standard Scaler
scaler = StandardScaler(with_mean=False) # Problems with dense matrix
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test.todense())
```

```
[ ] print([np.shape(X_train), np.shape(X_test)])
[(4179, 7692), (1393, 7692)]
```

```
[ ] type(X_train)
scipy.sparse.csr.csr_matrix
```

Linear SVM Models

colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=4kfeld_7ziZq

+ Code + Text

Reconnect

scipy.sparse.csr.csr_matrix

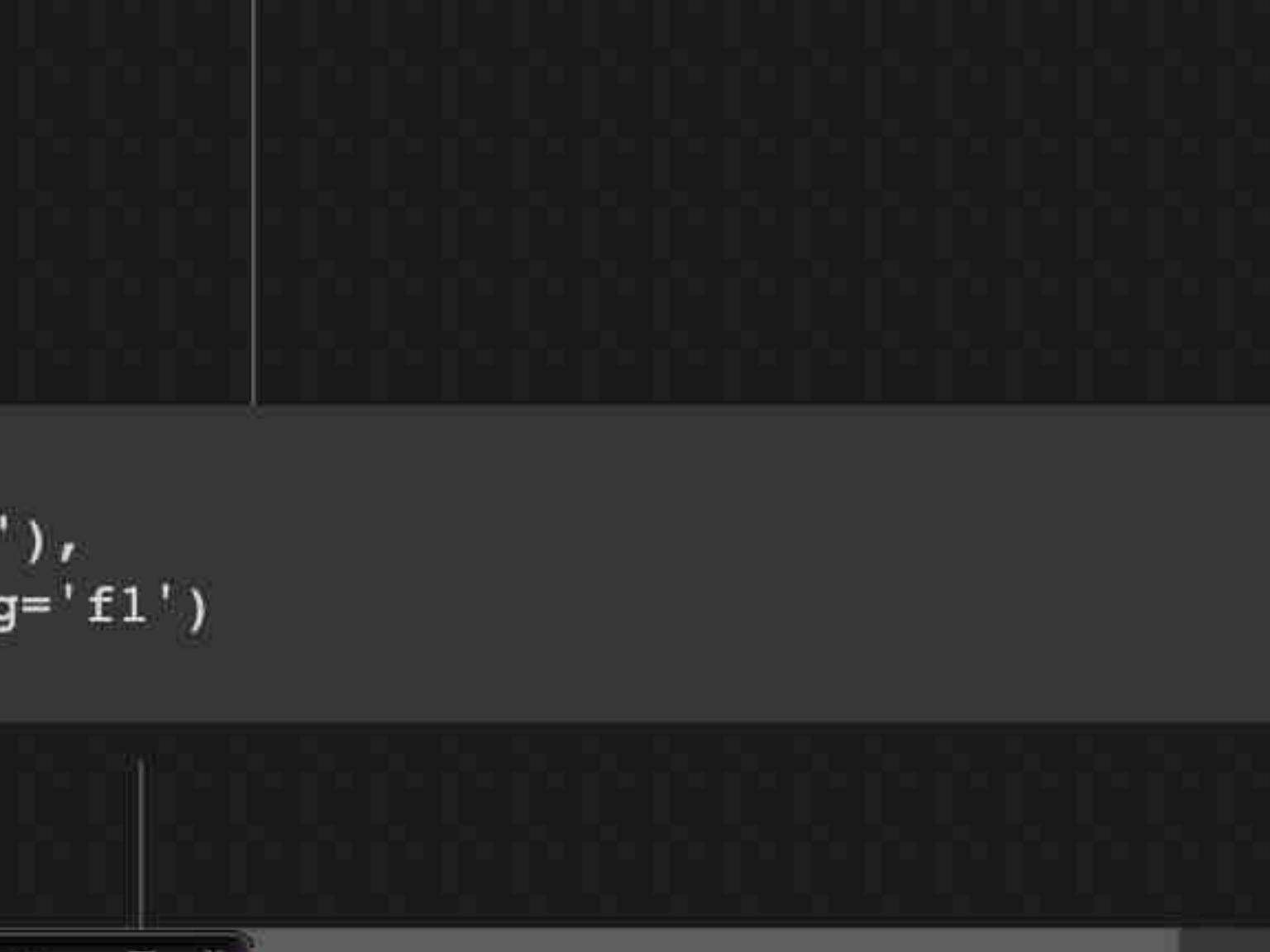
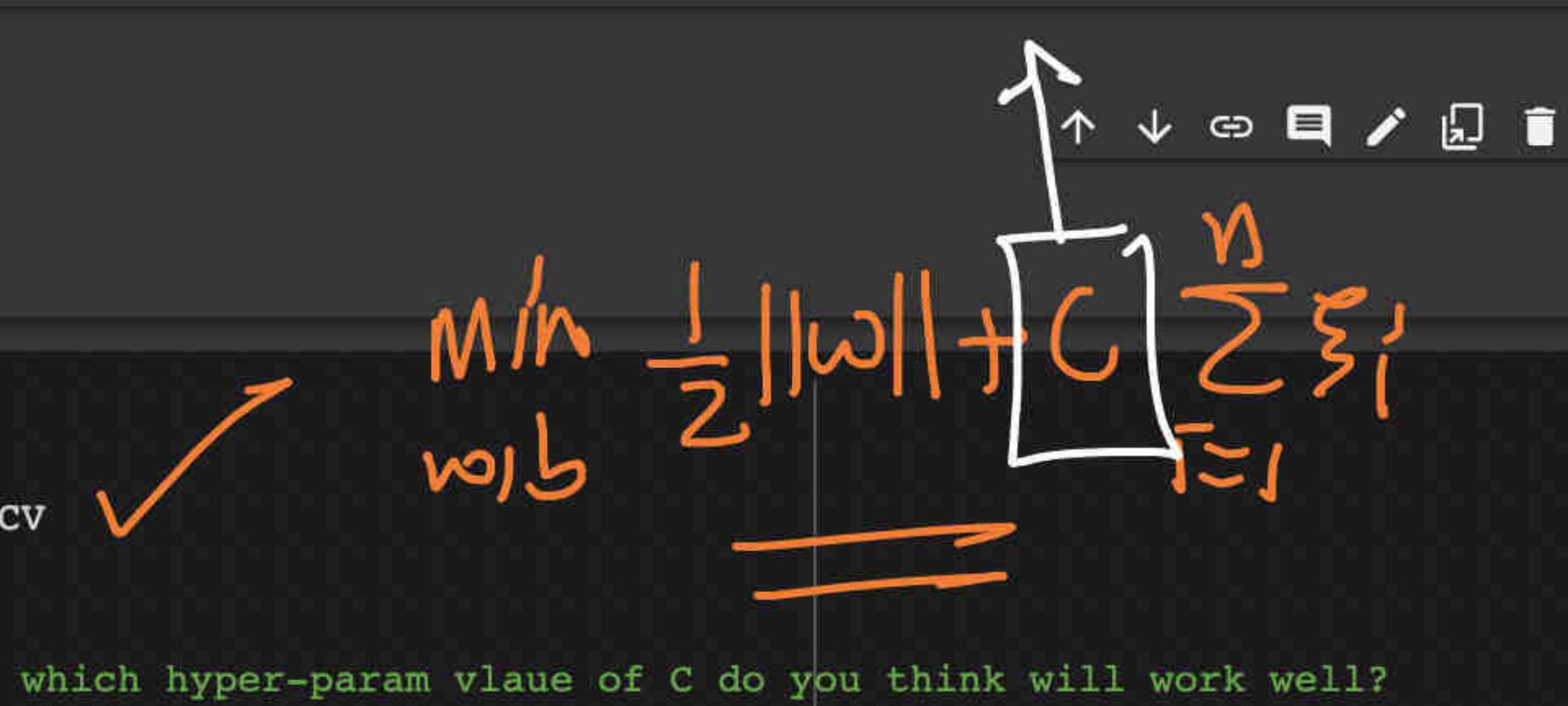
Linear SVM Models

```
[ ] # SVC
from sklearn.model_selection import GridSearchCV
params = {
    'C': [1e-4, 0.001, 0.01, 0.1, 1, 10] # which hyper-param value of C do you think will work well?
}
svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='linear')
clf = GridSearchCV(svc, params, scoring = "f1", cv=3)
clf.fit(X_train, y_train)
```

GridSearchCV(cv=3,
estimator=SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear'),
param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10]}, scoring='f1')

```
[ ] = clf.cv_results_
```

Min $\frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$



SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=4kfeld_7ziZq

+ Code + Text Reconnect

scipy.sparse.csr.csr_matrix

{x} Linear SVM Models

[] # SVC

```
from sklearn.model_selection import GridSearchCV

params = {
    'C': [1e-4, 0.001, 0.01, 0.1, 1, 10] # which hyper-param value of C do you think will work well?
}
svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='linear')
clf = GridSearchCV(svc, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

GridSearchCV(cv=3,
            estimator=SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear'),
            param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10]}, scoring='f1')
```

[] = clf.cv_results_



+ Code + Text

```
SVC = SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear', random_state=42)
[ ] clf = GridSearchCV(svc, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

{x}
GridSearchCV(cv=3,
            estimator=SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear'),
            param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10]}, scoring='f1')
```



```
res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")
```

Parameters:{'C': 0.0001} Mean_score: 0.6497817092458725 Rank: 6
Parameters:{'C': 0.001} Mean_score: 0.7839609357761174 Rank: 1 0.7839
Parameters:{'C': 0.01} Mean_score: 0.7814333814333815 Rank: 2
Parameters:{'C': 0.1} Mean_score: 0.7746267554214574 Rank: 3
Parameters:{'C': 1} Mean_score: 0.7746267554214574 Rank: 3
Parameters:{'C': 10} Mean_score: 0.7746267554214574 Rank: 3

RBF SVM

[] ↳ 1 cell hidden



+ Code + Text

```
'C': [1e-4, 0.001, 0.01, 0.1, 1, 10] # which hyper-param value of C do you think will work well?  
[ ]  
    }  
    svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='linear')  
    clf = GridSearchCV(svc, params, scoring = "f1", cv=3)  
  
    clf.fit(X_train, y_train)  
  
GridSearchCV(cv=3,  
            estimator=SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear'),  
            param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10]}, scoring='f1')
```



```
res = clf.cv_results_  
  
for i in range(len(res["params"])):   
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")  
  
Parameters:{'C': 0.0001} Mean_score: 0.6497817092458725 Rank: 6  
✓ Parameters:{'C': 0.001} Mean_score: 0.7839609357761174 Rank: 1  
Parameters:{'C': 0.01} Mean_score: 0.7814333814333815 Rank: 2  
Parameters:{'C': 0.1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 10} Mean_score: 0.7746267554214574 Rank: 3
```

RBF SVM



+ Code + Text

```
'C': [1e-4, 0.001, 0.01, 0.1, 1, 10] # which hyper-param value of C do you think will work well?  
[ ]  
    }  
    svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='linear')  
    clf = GridSearchCV(svc, params, scoring = "f1", cv=3)  
  
    clf.fit(X_train, y_train)  
  
GridSearchCV(cv=3,  
            estimator=SVC(class_weight={0: 0.1, 1: 0.5}, kernel='linear'),  
            param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10]}, scoring='f1')
```

```
res = clf.cv_results_  
  
for i in range(len(res["params"])):   
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")  
  
Parameters:{'C': 0.0001} Mean_score: 0.6497817092458725 Rank: 6  
Parameters:{'C': 0.001} Mean_score: 0.7839609357761174 Rank: 1  
Parameters:{'C': 0.01} Mean_score: 0.7814333814333815 Rank: 2  
Parameters:{'C': 0.1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 10} Mean_score: 0.7746267554214574 Rank: 3
```

dim is high

RBF SVM

+ Code + Text

```
Parameters:{'C': 0.01} Mean_score: 0.7814333814333815 Rank: 2  
Parameters:{'C': 0.1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 1} Mean_score: 0.7746267554214574 Rank: 3  
Parameters:{'C': 10} Mean_score: 0.7746267554214574 Rank: 3
```



COSINE SIM

d = 7692

▶ RBF SVM

[] ↳ 1 cell hidden

▶ Multinomial NB

[] ↳ 1 cell hidden

$$K(x_i, x_j) = \frac{x_i^T \cdot x_j}{\text{euc-dist}}$$

$$K_{RBF}(x_i, x_j) = \exp\left\{-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right\}$$

+ Code + Text

Reconnect         

```
[ ] # RBF-SVC

{x}

params = {
    'C': [0.1, 1, 10], ✓
    'gamma':[1e-3, 0.01, 0.1, 1] ✓
}
svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='rbf')
clf = GridSearchCV(svc, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")

# Euclidean distance in high dimensions in RBF calculation.
```

$$\gamma = \frac{1}{\sigma}$$

```
<> Parameters:{'C': 0.1, 'gamma': 0.001} Mean_score: 0.0 Rank: 9
```

```
= Parameters:{'C': 0.1, 'gamma': 0.01} Mean_score: 0.0 Rank: 9
```

```
>- Parameters:{'C': 0.1, 'gamma': 0.1} Mean_score: 0.0 Rank: 9
```

```
Parameters:{'C': 0.1, 'gamma': 1} Mean_score: 0.0 Rank: 9
```

```
Parameters:{'C': 1, 'gamma': 0.001} Mean_score: 0.3400195328811167 Rank: 1
```



+ Code + Text

```
[ ] clf = GridSearchCV(svc, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

{x}

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")

# Euclidean distance in high dimensions in RBF calculation.
```

```
Parameters:{'C': 0.1, 'gamma': 0.001} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 0.01} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 0.1} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 1} Mean_score: 0.0 Rank: 9
Parameters:{'C': 1, 'gamma': 0.001} Mean_score: 0.3400185328811167 Rank: 1
Parameters:{'C': 1, 'gamma': 0.01} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 1, 'gamma': 0.1} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 1, 'gamma': 1} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 10, 'gamma': 0.001} Mean_score: 0.26494585373090046 Rank: 2
Parameters:{'C': 10, 'gamma': 0.01} Mean_score: 0.25141937096101824 Rank: 3
Parameters:{'C': 10, 'gamma': 0.1} Mean_score: 0.2486880986310215 Rank: 4
Parameters:{'C': 10, 'gamma': 1} Mean_score: 0.2486880986310215 Rank: 4
```

f1:
f1: 0.34



+ Code + Text

```
        }
[ ] svc = SVC(class_weight={ 0:0.1, 1:0.5 }, kernel='rbf')
clf = GridSearchCV(svc, params, scoring = "f1", cv=3)

{x}
clf.fit(X_train, y_train)

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")
# Euclidean distance in high dimensions in RBF calculation.
```

Parameters:{'C': 0.1, 'gamma': 0.001} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 0.01} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 0.1} Mean_score: 0.0 Rank: 9
Parameters:{'C': 0.1, 'gamma': 1} Mean_score: 0.0 Rank: 9
Parameters:{'C': 1, 'gamma': 0.001} Mean_score: 0.3400185328811167 Rank: 1
Parameters:{'C': 1, 'gamma': 0.01} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 1, 'gamma': 0.1} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 1, 'gamma': 1} Mean_score: 0.01780655274942489 Rank: 6
Parameters:{'C': 10, 'gamma': 0.001} Mean_score: 0.26494585373090046 Rank: 2
Parameters:{'C': 10, 'gamma': 0.01} Mean_score: 0.25141937096101824 Rank: 3
Parameters:{'C': 10, 'gamma': 0.1} Mean_score: 0.2486880986310215 Rank: 4
Parameters:{'C': 10, 'gamma': 1} Mean_score: 0.2486880986310215 Rank: 4

SVM_NB.ipynb - Colaboratory x sklearn.preprocessing.StandardScaler x sklearn.svm.SVC — scikit-learn x sklearn.naive_bayes.BernoulliNB x sklearn.naive_bayes.MultinomialNB x 1.9. Naive Bayes — scikit-learn x dsml-course/SVM Business C x + colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=4kfeld_7ziZq

+ Code + Text Reconnect

Reconnect

2 cells hidden

{x}

RBF SVM

[] ↳ 1 cell hidden

Multinomial NB

[] ↳ 1 cell hidden

RBF-SVM

MNB

by SVM

110 / 110



+ Code

Multinomial NB

```
{x}
[ 1] df_X_train, df_X_test, y_train, y_test = train_test_split(df['cleaned_message'], df['type'], test_size=0.25, random_
print([np.shape(df_X_train), np.shape(df_X_test)])
```

#Count Vectorizer

```
{ 2} f = feature_extraction.text.CountVectorizer()
```

```
 3] X_train = f.fit_transform(df_X_train)
```

```
 4] X_test = f.transform(df_X_test)
```

No need of scaling

~~ZZZZ~~

Multinomial NB

```
from sklearn.model_selection import GridSearchCV
```

```
<> params = {
```

```
    'alpha':[0.01, 0.1, 1, 10]
```

```
    }
```

```
mnb = naive_bayes.MultinomialNB()
```

```
clf = GridSearchCV(mnb, params, scoring = "f1", cv=3)
```



+ Code + Text

```
[ ] X_train = f.fit_transform(df_X_train)
[ ] X_test = f.transform(df_X_test)

# No need of scaling

{x}

# Multinomial NB

from sklearn.model_selection import GridSearchCV

params = {
    'alpha':[0.01, 0.1, 1, 10]
}
mnb = naive_bayes.MultinomialNB()
clf = GridSearchCV(mnb, params, scoring = "f1", cv=3)
clf.fit(X_train, y_train)

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")
```



+ Code

```
[ ] X_train = f.fit_transform(df_X_train)
[ ] X_test = f.transform(df_X_test)

# No need of scaling

{x}

# Multinomial NB

from sklearn.model_selection import GridSearchCV

params = {
    'alpha':[0.01, 0.1, 1, 10]
}
mnb = naive_bayes.MultinomialNB()
clf = GridSearchCV(mnb, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")
```

colab.research.google.com/drive/1_V0nqZq54rk_XVtIC_j1ez7vPaKDaZVM#scrollTo=-RtxxwDA90TK

+ Code + Text Reconnect

```
[ ]     'alpha':[0.01, 0.1, 1, 10]
}
mnb = naive_bayes.MultinomialNB()
clf = GridSearchCV(mnb, params, scoring = "f1", cv=3)

clf.fit(X_train, y_train)

res = clf.cv_results_

for i in range(len(res["params"])):
    print(f"Parameters:{res['params'][i]} Mean_score: {res['mean_test_score'][i]} Rank: {res['rank_test_score'][i]}")
```

MNB > Y SUM >>
RBF SVM

[(4179,), (1393,)]
Parameters:{'alpha': 0.01} Mean_score: 0.8921625905385584 Rank: 2
Parameters:{'alpha': 0.1} Mean_score: 0.8883633779156167 Rank: 3
Parameters:{'alpha': 1} Mean_score: 0.9006541826507674 Rank: 1
Parameters:{'alpha': 10} Mean_score: 0.8580328947757173 Rank: 4

 SVM_NB.ipynb 

File Edit View Insert Runtime Tools Help All changes saved

 Comment

 Share



+ Code + Text

Reconnect ▾

Editing

1

▼ Data

[x]

↑ ↓ ↻ ⏷ ⏵ ⏴ ⏵ ⏴

```
# Importing libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from collections import Counter
from sklearn import feature_extraction, model_selection, naive_bayes, metrics, svm
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler

from sklearn.svm import SVC

# Libraries for text processing
import re, nltk
nltk.download('punkt')
nltk.download('stopwords')
```

