











## **ANALISIS SENTIMEN**

**MASYARAKAT INDONESIA PADA MEDIA** SOSIAL TWITTER TERHADAP PROGRAM PRAKERJA DENGAN METODE NAIVE BAYES CLASIFFIER (BERNOULI)

















#### **Data Privacy Infographics**



**Taofik Krisdiyanto** 

Mahasiswa



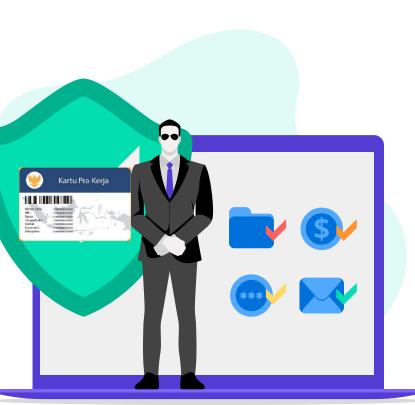
Jemmy Edwing Borong, S.Kom., M.Eng.

Dosen Pembimbing





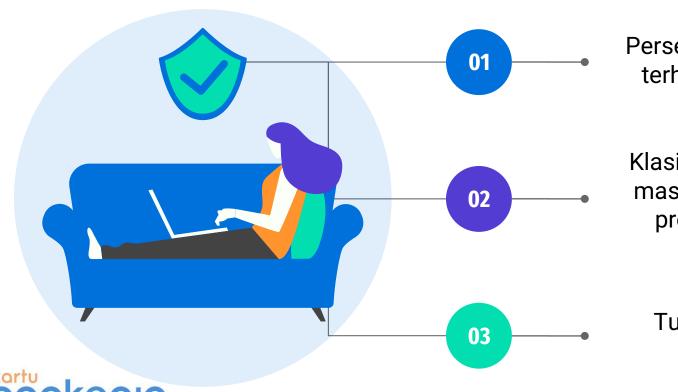
## **Background**



Program pengembangan kompetensi kerja dan kewirausahaan yang ditujukan untuk pencari kerja, pekerja/buruh yang terkena pemutusan hubungan kerja atau pekerja/buruh yang membutuhkan peningkatan kompetensi, termasuk pelaku usaha mikro dan kecil

sejak bulan April 2020 insentif terlambat cair & akses website lemah karena membludak & salah sasaran

## **Purpose**



Persepsi masyarakat terhadap program prakerja

Klasifikasi pandangan masyarakat terhadap program prakerja

Tugas akhir mata kuliah PPL

#### **Application & Package**





















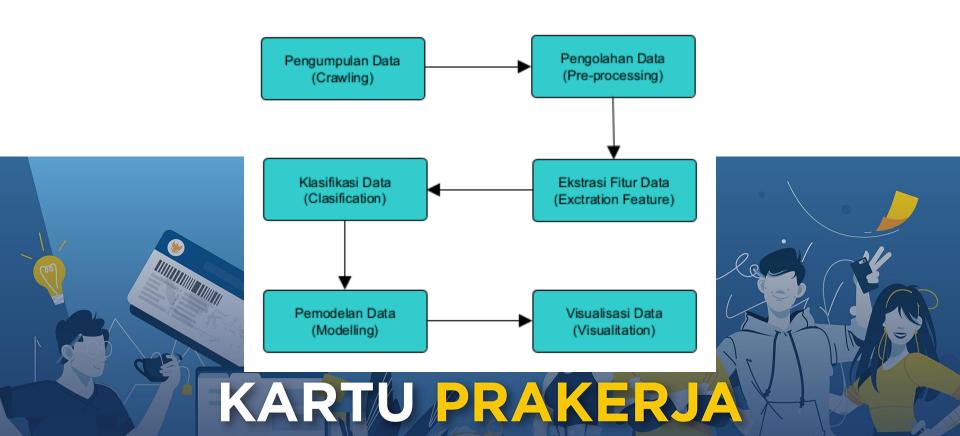








#### **Flowchart System**



## Crawling







#### hashtags

- @izonemfs kemarin dapet prakerja
- @TokopediaCare Halo Kak,\nNama: Muhamad Surya...
- 2 @Chepyzaenal68 Untuk melakukan pendaftaran aku...
- 3 @bukalapak the best lah bukalapak dan prakerja...
- @bukalapak emang bukalapak paling bisa di anda...

# **Preprocessing (Case Folding)**







#mengubah huruf kapital dalam data menjadi huruf kecil
df['lower']=df['hashtags'].str.lower()
df.head(5)

	Unnamed:	0	hashtags	lower
0		0	@izonemfs kemarin dapet prakerja	@izonemfs kemarin dapet prakerja
1		1	@TokopediaCare Halo Kak,\nNama : Muhamad Surya	@tokopediacare halo kak,\nnama : muhamad surya
2		2	@Chepyzaenal68 Untuk melakukan pendaftaran aku	@chepyzaenal68 untuk melakukan pendaftaran aku
3		3	@bukalapak the best lah bukalapak dan prakerja	@bukalapak the best lah bukalapak dan prakerja
4		4	@bukalapak emang bukalapak paling bisa di anda	@bukalapak emang bukalapak paling bisa di anda

# **Preprocessing (Cleansing)**

```
df['clean']=df['lower'].str.replace('(https?://[\w\.\/]*)', '') #hapus mengandung http
df['clean']=df['clean'].str.replace('(http\S+)', '') #hapus mengandung https
df['clean']=df['clean'].str.replace('(?:&(?:lt|nbsp|amp|gt);)', '') #hapus lt,nbsp,gt,amp
df['clean']=df['clean'].str.replace('(@|#)\w+', '') #hapus @ dan #
df['clean']=df['clean'].str.replace('[^A-Za-z0-9\s\-\/]', '') #menghapus selain huruf, spasi dan strip
df['clean']=df['clean'].str.replace('(\-|\/)', ' ') #menghapus -
df['clean']=df['clean'].str.replace('\n', ' ') # hapus enter
df['clean']=df['clean'].str.replace('\t', ' ') # hapus enter
df['clean']=df['clean'].str.replace('\d+', ' ') # hapus angka
df['clean']=df['clean'].str.replace('prakerja', ' ') # hapus kata prakerja
df['clean']=df['clean'].str.replace('\s{2,}', ' ') # hapus spasi lebih dari 2
df['clean']=df['clean'].str.replace('^rt.*', '') #Remove RT Tweets
#df['clean']=df['clean'].str.replace('<.*?>', '') #karakter html
df['clean']=df['clean'].str.replace("\.\.", " ") #hapus .
df['clean']=df['clean'].str.replace('\b[a-zA-Z]\b', '') #huruf mandiri
df.dropna(subset=['clean'], inplace=True) #Remove Empty cell
df = df.drop duplicates() #Remove Duplicate Tweet
df.head(5)
```







clean	lower	hashtags	Unnamed: 0	
kemarin dapet	@izonemfs kemarin dapet prakerja	@izonemfs kemarin dapet prakerja	0	0
halo kak nama muhamad suryadi syahdoa e mail	@tokopediacare halo kak,\nnama : muhamad surya	@TokopediaCare Halo Kak,\nNama : Muhamad Surya	1	1
untuk melakukan pendaftaran akun kartu silaka	@chepyzaenal68 untuk melakukan pendaftaran aku	@Chepyzaenal68 Untuk melakukan pendaftaran aku	2	2
the best lah bukalapak dan dalam membantu mas	@bukalapak the best lah bukalapak dan prakerja	@bukalapak the best lah bukalapak dan prakerja	3	3
emang bukalapak paling bisa di andalkan kalo	@bukalapak emang bukalapak paling bisa di anda	@bukalapak emang bukalapak paling bisa di anda	4	4

## **Preprocessing (Tokenizing)**

hashtags

```
import string
import re #regex library
import nltk
nltk.download('punkt')
# impor word_tokenize & FreqDist dari NLTK
from nltk.tokenize import word_tokenize
from nltk.probability import FreqDist

# NLTK kata tokenize
def word_tokenize_wrapper(text):
    return word_tokenize(text)

df['tokenize'] = df['clean'].apply(word_tokenize_wrapper)

print('Tokenizing Result : \n')
df.head()
```





clean



tokenize

0	0	@izonemfs kemarin dapet prakerja	@izonemfs kemarin dapet prakerja	kemarin dapet	[kemarin, dapet]
1	1 <sub>k</sub>	@TokopediaCare Halo Kak,\nNama : Muhamad Surya	@tokopediacare halo kak,\nnama : muhamad surya	halo kak nama muhamad suryadi syahdoa e mail	[halo, kak, nama, muhamad, suryadi, syahdoa, e
2	2	@Chepyzaenal68 Untuk melakukan pendaftaran aku	@chepyzaenal68 untuk melakukan pendaftaran aku	untuk melakukan pendaftaran akun kartu silaka	[untuk, melakukan, pendaftaran, akun, kartu, s
3	3	@bukalapak the best lah bukalapak dan prakerja	@bukalapak the best lah bukalapak dan prakerja	the best lah bukalapak dan dalam membantu mas	[the, best, lah, bukalapak, dan, dalam, memban
4	4	@bukalapak emang bukalapak paling bisa di anda	@bukalapak emang bukalapak paling bisa di anda	emang bukalapak paling bisa di andalkan kalo	[emang, bukalapak, paling, bisa, di, andalkan,

lower

## **Preprocessing (Filtering)**

```
from nltk.corpus import stopwords
import nltk
nltk.download('stopwords')
# ------ dapatkan stopword dari stopword NLTK ------
# dapatkan stopword indonesia
list stopwords = stopwords.words('indonesian')
print(len(list_stopwords))
# ----- tambahkan stopword secara manual -----
list_stopwords.extend(['yg', 'dg', 'rt', 'dgn', 'ny', 'klo',
                    'kalo', 'amp', 'biar', 'bikin', 'bilang',
                    'gak', 'ga', 'krn', 'nya', 'nih', 'sih',
                    'sih', 'tau', 'tdk', 'tuh', 'utk', 'ya',
                    'id', 'ign', 'sdh', 'aja', 'tp', 'mmg',
                    'nvg', 'hehe', 'pen', 'nan', 'loh', 'nah',
                    'yah','co','ya','aja','the', 'best'])
len(list stopwords)
# ------ tambahkan stopword dari file txt -----
# baca txt stopword menggunakan panda
txt_stopword = pd.read_csv("gdrive/My Drive/PPL/prakerja/ID-Stopwords-master/id.stopwords.02.01.2016.txt",
                       names= ["stopwords"], header = None)
# konversi string stopword ke daftar dan tambahkan stopword tambahan
list_stopwords.extend(txt_stopword["stopwords"][0].split())
len(list stopwords)
# konversi daftar ke kamus
list_stopwords = set(list_stopwords)
#hapus stopword pada token daftar
def stopwords removal(words):
   return [word for word in words if word not in list_stopwords]
df['stopword'] = df['tokenize'].apply(stopwords removal)
df.head(5)
```





stopword	freq_dist_token	tokenize	clean	lower	hashtags	named: 0	Un
[kemarin, dapet]	{'kemarin': 1, 'dapet': 1}	[kemarin, dapet]	kemarin dapet	@izonemfs kemarin dapet prakerja	@izonemfs kemarin dapet prakerja	0	0
[halo, kak, nama, muhamad, suryadi, syahdoa, e	{'halo': 1, 'kak': 1, 'nama': 1, 'muhamad': 1,	[halo, kak, nama, muhamad, suryadi, syahdoa, e	halo kak nama muhamad suryadi syahdoa e mail	@tokopediacare halo kak,\nnama : muhamad surya	@TokopediaCare Halo Kak,\nNama : Muhamad Surya	1	1
[pendaftaran, akun, kartu, silakan, kunjungi,	('untuk': 1, 'melakukan': 1, 'pendaftaran': 1,	[untuk, melakukan, pendaftaran, akun, kartu, s	untuk melakukan pendaftaran akun kartu silaka	@chepyzaenal68 untuk melakukan pendaftaran aku	@Chepyzaenal68 Untuk melakukan pendaftaran aku	2	2
[bukalapak, membantu, masyarakat]	{'the': 1, 'best': 1, 'lah': 1, 'bukalapak': 1	[the, best, lah, bukalapak, dan, dalam, memban	the best lah bukalapak dan dalam membantu mas	@bukalapak the best lah bukalapak dan prakerja	@bukalapak the best lah bukalapak dan prakerja	3	3
[emang, bukalapak, andalkan, penambahan, ilmu,	('emang': 1, 'bukalapak': 1, 'paling': 1, 'bis	[emang, bukalapak, paling, bisa, di, andalkan,	emang bukalapak paling bisa di andalkan kalo	@bukalapak emang bukalapak paling bisa di anda	@bukalapak emang bukalapak paling bisa di anda	4	4

# **Preprocessing (Stemming)**

```
#import paket Sastrawi
!pip install Sastrawi
!pip install swifter
from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
import swifter
#buat stemmer
factory = StemmerFactory()
stemmer = factory.create stemmer()
# stemmer
def stemmed wrapper(term):
   return stemmer.stem(term)
term_dict = {}
for document in df['stopword']:
    for term in document:
       if term not in term dict:
           term dict[term] ="
print(len(term dict))
print("----")
for term in term dict:
    term dict[term] = stemmed wrapper(term)
    print(term,":",term_dict[term])
print(term dict)
print("----")
# terapkan istilah bertangkai ke kerangka data
def get stemmed term(document):
    return [term dict[term] for term in document]
df['stemming'] = df['stopword'].swifter.apply(get stemmed term)
print(df.head())
```







gabung : gabung kompetisi : kompetisi

banyaknya : banyak pengangguran : anggur

akibat : akibat diperbanyak : banyak

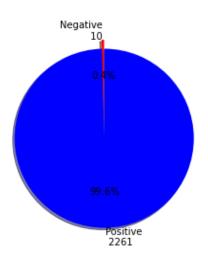
kuotanya : kuota kejelasan : jelas himbara : himbara

stemming	stopword	freq_dist_token	tokenize	clean
[ˈkemarinˈ, ˈdapetˈ]	[ˈkemarinˈ, ˈdapetˈ]	<freqdist 2="" and="" outcomes="" samples="" with=""></freqdist>	[ˈkemarinˈ, ˈdapetˈ]	kemarin dapet
[ˈhaloˈ, ˈkakˈ, ˈnamaˈ,	[ˈhaloˈ, ˈkakˈ, ˈnamaˈ,	<freqdist 22="" and="" outcomes="" samples="" with=""></freqdist>	['halo', 'kak', 'nama',	halo kak nama muhamad
ˈmuhamadˈ, ˈsuryadiˈ,	ˈmuhamadˈ, ˈsuryadiˈ,		'muhamad', 'suryadi',	suryadi syahdoa e mail
['daftar', 'akun', 'kartu', 'sila', 'kunjung',	['pendaftaran', 'akun', 'kartu', 'silakan', 'k	<freqdist 15="" samples<br="" with="">and 16 outcomes&gt;</freqdist>	['untuk', 'melakukan', 'pendaftaran', 'akun',	untuk melakukan pendaftaran akun kartu silaka
[ˈbukalapakˈ, ˈbantuˈ,	[ˈbukalapakˈ, ˈmembantuˈ,	<freqdist 8="" and="" outcomes="" samples="" with=""></freqdist>	['the', 'best', 'lah',	the best lah bukalapak dan
ˈmasyarakat']	ˈmasyarakat']		'bukalapak', 'dan', 'da	dalam membantu mas
['emang', 'bukalapak',	['emang', 'bukalapak',	<freqdist 11="" 12="" and="" outcomes="" samples="" with=""></freqdist>	['emang', 'bukalapak',	emang bukalapak paling bisa
'andal', 'tambah', 'ilm	'andalkan', 'penambahan		'paling', 'bisa', 'di',	di andalkan kalo

## **Data Feature Extraction (LM)**







```
# load positive word
positive = pd.read_csv('gdrive/My Drive/PPL/prakerja/ID-OpinionWords-master/positive.txt', header=None)
positive = positive[0].values.tolist()
positive = '|'.join(positive)
# load negative word
negative = pd.read_csv('gdrive/My Drive/PPL/prakerja/ID-OpinionWords-master/negative.txt', header=None)
negative = negative[0].values.tolist()
negative = '|'.join(negative)

df['positive'] = [len(re.findall(positive, i.lower())) / len(i.split()) for i in df.stemming]
df['negative'] = [len(re.findall(negative, i.lower())) / len(i.split()) for i in df.stemming]
df['sentimen'] = ['positive' if df.iloc[i].positive >= df.iloc[i].negative else 'negative' for i in range(df.shape[0])]
```

	clean	tokenize	freq_dist_token	stopword	stemming	positive	negative	sentimen
	kemarin dapet	[ˈkemarinˈ, ˈdapetˈ]	<freqdist 2="" and="" outcomes="" samples="" with=""></freqdist>	[ˈkemarinˈ, ˈdapetˈ]	[ˈkemarinˈ, ˈdapetˈ]	1.000000	0.000000	positive
	halo kak nama muhamad suryadi syahdoa e mail	['halo', 'kak', 'nama', 'muhamad', 'suryadi', 	<freqdist 22<br="" with="">samples and 22 outcomes&gt;</freqdist>	['halo', 'kak', 'nama', 'muhamad', 'suryadi',	['halo', 'kak', 'nama', 'muhamad', 'suryadi', 	0.800000	0.050000	positive
16	untuk melakukan endaftaran akun kartu silaka	['untuk', 'melakukan', 'pendaftaran', 'akun', 	<freqdist 15<br="" with="">samples and 16 outcomes&gt;</freqdist>	[ˈpendaftaranˈ, ˈakunˈ, ˈkartuˈ, ˈsilakanˈ, ˈk	['daftar', 'akun', 'kartu', 'sila', 'kunjung',	0.900000	0.000000	positive
	he best lah bukalapak dan dalam membantu mas	['the', 'best', 'lah', 'bukalapak', 'dan', 'da	<pre><freqdist 8="" and="" outcomes="" samples="" with=""></freqdist></pre>	[ˈbukalapakˈ, ˈmembantuˈ, ˈmasyarakatˈ]	['bukalapak', 'bantu', 'masyarakat']	2.666667	0.333333	positive
);	emang bukalapak aling bisa di andalkan kalo	['emang', 'bukalapak', 'paling', 'bisa', 'di',	<freqdist 11<br="" with="">samples and 12 outcomes&gt;</freqdist>	['emang', 'bukalapak', 'andalkan', 'penambahan	['emang', 'bukalapak', 'andal', 'tambah', 'ilm	1.833333	0.333333	positive

## **Data Feature Extraction (BOW)**

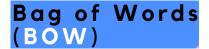


```
bow_model = CountVectorizer(ngram_range=(1,2), min_df=0.0025)
bow_model.fit(df.stemming)
df_vect = pd.DataFrame(bow_model.transform(df.stemming).toarray(), columns=bow_model.get_feature_names())
print(df_vect.shape)
df_vect.head(5)

df_vect['rate_positive_word'] = [i for i in df['positive']]
df_vect['rate_negative_word'] = [i for i in df['negative']]
df_vect.head()
```







	abis	academy	academy bukalapak	academy dapat	academy nilai	ada	ada gelombang	admin	adu	ah	ahli	ahli sdm	ahli sempat	ahli sertifikasi	airlangga	ajar	akademi	akademi klik	akibat	akibat pandemi	akses	aktif	akun	akun kartu	akun linkaja
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	1	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

## **Data Classification**



```
from sklearn.naive_bayes import BernoulliNB
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, classification_report
import seaborn as sn
```

x\_train, x\_test, y\_train, y\_test = train\_test\_split(df\_vect, df.numbered\_sentimen, test\_size=0.2,random\_state=123)

```
test_y = pd.DataFrame(y_test)
test_x.shape[0]

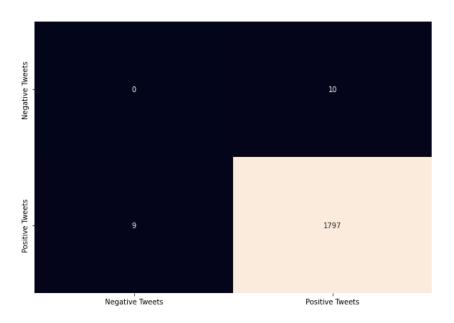
#Buat model Bernoulli
model = BernoulliNB()

#Start Training
classifier = model.fit(x_train,y_train)
```

test x = pd.DataFrame(x test)

		abis	academy	academy bukalapak	academy dapat	academy nilai	ada
	2003	0	0	0	0	0	0
nnamed: 0	976	0	0	0	0	0	0
903 1 976 1	905	0	0	0	0	0	0
95 1 215 1	215	0	0	0	0	0	0
505	4 605	0	1	0	1	0	0

## Data Modeling (CM Train Data)



1816

```
recall f1-score
              precision
                                               support
                   0.00
                             0.00
                                        0.00
                                                    10
                   0.99
                             1.00
                                        0.99
                                                  1806
    accuracy
                                        0.99
                                                  1816
                                        0.50
                                                  1816
  macro avg
                   0.50
                             0.50
weighted avg
                                                  1816
                   0.99
                             0.99
                                        0.99
```

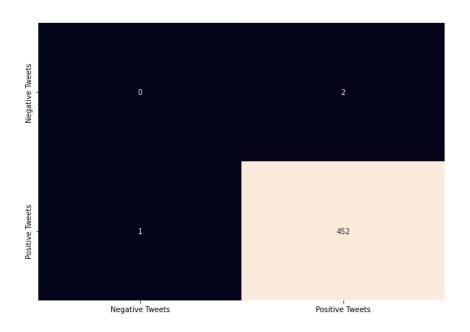
```
#Confusion Matrix Training
cm_train = confusion_matrix(predict,y_train)
print(cm_train)
```

```
[[ 0 10]
[ 9 1797]]
```

```
accuracy = cm_train.trace()/cm_train.sum()
print (accuracy)
print(f"{(accuracy*100):3.2f}%")
```

```
0.9895374449339207
98.95%
```

## Data Modeling (CM Test Data)



455

```
precision
                           recall f1-score
                                              support
           0
                   0.00
                             0.00
                                       0.00
                                                    2
           1
                   1.00
                             1.00
                                       1.00
                                                  453
                                       0.99
                                                  455
    accuracy
  macro avg
                   0.50
                             0.50
                                       0.50
                                                  455
weighted avg
                   0.99
                             0.99
                                       0.99
                                                  455
```

```
#Confusion Matrix Testing
cm = confusion_matrix(predict,y_test)
print(cm)
```

```
[ 1 452]]

predict = classifier.predict(x_test)

classification_report(predict, y_test)
accuracy = cm.trace()/cm.sum()
print (accuracy)
```

0.9934065934065934 99.34%

print(f"{(accuracy\*100):3.2f}%")

## AI Sentimen (NBC-Bernoulli)

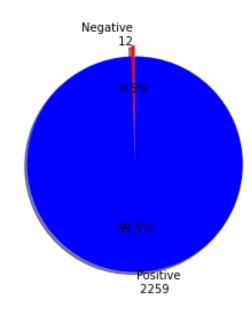


#### matpletlib

```
#Prediksi Bernoulli dengan BOW(df_vect)
result = model.predict(df_vect)

print('jumlah positif : ', len([n for n in result if n == 1]))
print('jumlah negatif : ', len([n for n in result if n == 0]))
```

jumlah positif : 2259
jumlah negatif : 12





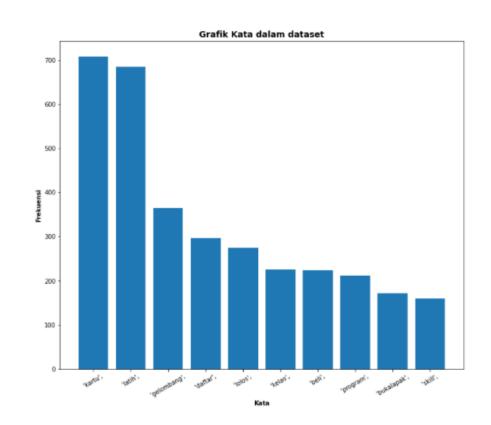
## **Manual vs NBC**



	hashtags	lower	clean	tokenize	freq_dist_token	stopword	stemming	positive	negative	sentimen	numbered_sentimen	ai_predict	sentiment_ai
Unnamed: 0													
0	@izonemfs kemarin dapet prakerja	@izonemfs kemarin dapet prakerja	kemarin dapet	[ˈkemarinˈ, ˈdapetˈ]	<freqdist 2<br="" with="">samples and 2 outcomes&gt;</freqdist>	[ˈkemarinˈ, ˈdapetˈ]	[ˈkemarinˈ, ˈdapetˈ]	1.000000	0.000000	positive	1	1	positive
1	@TokopediaCare Halo Kak,\nNama : Muhamad Surya	@tokopediacare halo kak,\nnama : muhamad surya	halo kak nama muhamad suryadi syahdoa e mail	[ˈhaloʻ, ˈkak', ˈnama', ˈmuhamad', ˈsuryadi',	<freqdist 22<br="" with="">samples and 22 outcomes&gt;</freqdist>	['halo', 'kak', 'nama', 'muhamad', 'suryadi',	['halo', 'kak', 'nama', 'muhamad', 'suryadi',	0.800000	0.050000	positive	1	1	positive
2	@Chepyzaenal68 Untuk melakukan pendaftaran aku	@chepyzaenal68 untuk melakukan pendaftaran aku	untuk melakukan pendaftaran akun kartu silaka	['untuk', 'melakukan', 'pendaftaran', 'akun',	<freqdist 15<br="" with="">samples and 16 outcomes&gt;</freqdist>	['pendaftaran', 'akun', 'kartu', 'silakan', 'k	['daftar', 'akun', 'kartu', 'sila', 'kunjung',	0.900000	0.000000	positive	1	1	positive
3	@bukalapak the best lah bukalapak dan prakerja	@bukalapak the best lah bukalapak dan prakerja	the best lah bukalapak dan dalam membantu mas	['the', 'best', 'lah', 'bukalapak', 'dan', 'da	<freqdist 8<br="" with="">samples and 8 outcomes&gt;</freqdist>	[ˈbukalapakˈ, ˈmembantuˈ, ˈmasyarakat']	('bukalapak', 'bantu', 'masyarakat']	2.666667	0.333333	positive	1	1	positive
4	@bukalapak emang bukalapak paling bisa di anda	@bukalapak emang bukalapak paling bisa di anda	emang bukalapak paling bisa di andalkan kalo 	['emang', 'bukalapak', 'paling', 'bisa', 'di',	<freqdist 11<br="" with="">samples and 12 outcomes&gt;</freqdist>	['emang', 'bukalapak', 'andalkan', 'penambahan	[ˈemangˈ, ˈbukalapakˈ, ˈandalˈ, ˈtambahˈ, ˈilm	1.833333	0.333333	positive	1	1	positive

## **Data Visualization**







## **Data Visualization**







## **Conclusion**

#### **Model Accuracy**

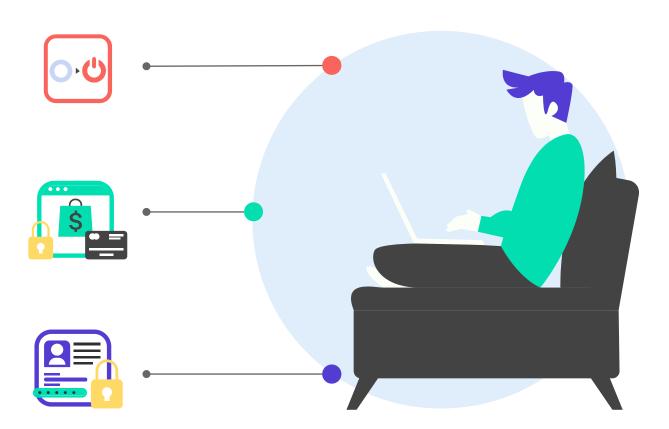
98,95 % Data Train 99, 34 % Data Test

#### **Sentimen**

99,5 % Positive, 0,5 % Negative

#### **Best of Word**

Kartu, Latih, Program, Gelombang



# Maturnaum

