

KT AIVLE School

4차 미니프로젝트

AI 5반 14조

결측치 제거

```
# 결측치를 제거한 후 확인합니다.  
spam.isna().sum()
```

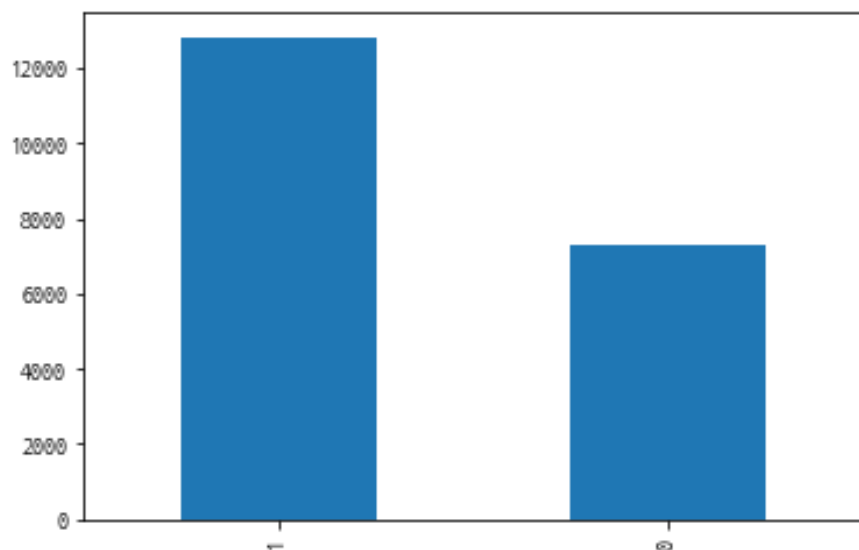
```
text      11  
label      0  
dtype: int64
```

```
spam.dropna(axis=0, inplace=True)|
```

데이터 분석

```
spam['label'].value_counts().plot(kind='bar')
```

<AxesSubplot:>



```
print('정상 메일과 스팸 메일의 개수')
print(spam.groupby('label').size().reset_index(name='count'))

print(f'정상 메일의 비율 = {round(spam["label"].value_counts()[0]/len(spam) * 100,3)}%')
print(f'스팸 메일의 비율 = {round(spam["label"].value_counts()[1]/len(spam) * 100,3)}%')
```

정상 메일과 스팸 메일의 개수

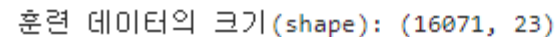
label	count
0	0 7272
1	1 12817

정상 메일의 비율 = 36.199%

스팸 메일의 비율 = 63.801%

▪ Sequence Vectorize

```
[[1, 5299, 9, 877], [2, 1271, 145, 5300, 2546, 122], [1, 878,
등장 빈도가 1번 이하인 희귀 단어의 수: 23920
단어 집합(vocabulary)에서 희귀 단어의 비율: 81.86734204942158
전체 등장 빈도에서 희귀 단어 등장 빈도 비율: 26.65151363216009
단어 집합의 크기: 29219
메일의 최대 길이 : 23
메일의 평균 길이 : 5.584656
```



모델링

▪ RNN

```
from tensorflow.keras.layers import SimpleRNN, Embedding, Dense
from tensorflow.keras.models import Sequential

embedding_dim = 32
hidden_units = 32
|
model = Sequential()
model.add(Embedding(vocab_size, embedding_dim))
model.add(SimpleRNN(hidden_units))
model.add(Dense(1, activation='sigmoid'))

model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['acc'])
model.fit(x_train2, y_train, epochs=4, batch_size=64, validation_split=0.2)
```

```
x_test_encoded2 = tokenizer.texts_to_sequences(x_val)
x_test_padded2 = pad_sequences(x_test_encoded2, maxlen = max_len)
print("\n 테스트 정확도: %.4f" % (model.evaluate(x_test_padded2, y_val)[1]))
```

```
126/126 [=====] - 3s 20ms/step - loss: 0.2969 - acc: 0.8512
```

```
테스트 정확도: 0.8512
```

kt

AIVLE

AIVLE
Let's make it possible