

# 데이터 분석을 통한 유방암 진단

제작 : 문경환

# 01 CONTENTS

## Overview

### Purpose

- AI를 기반으로 한 유방암 진단

### Summary

- 데이터 분석을 통한 주요 feature 추출로 ML 학습 Score 향상(91% -> 94% & Overfitting 방지)  
-> barplot, violinplot, heatmap, correlation 사용
- 이미지 데이터를 통한 최적의 DL 모델 구축과 고찰(최대 91%)  
-> tensorflow, keras 사용



# CONTENTS

01

Machine Learning

02

Deep Learning

03

Analysis

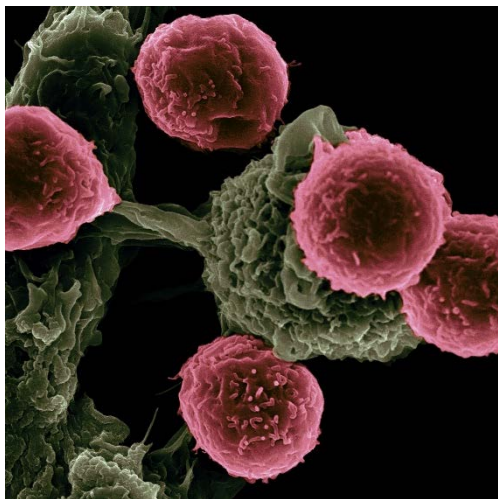


# 01

Breast Cancer diagnosis with ML

# 01 CONTENTS

## Inspection & Target



### Breast Cancer

- 세계에서 가장 흔한 여성 암
- 생활습관이 세계적으로 바뀔에 따른 큰 증가 폭

```
1 len(data.columns)
2 >> 33
```



```
1 len(data.columns)
2 >> 9
```

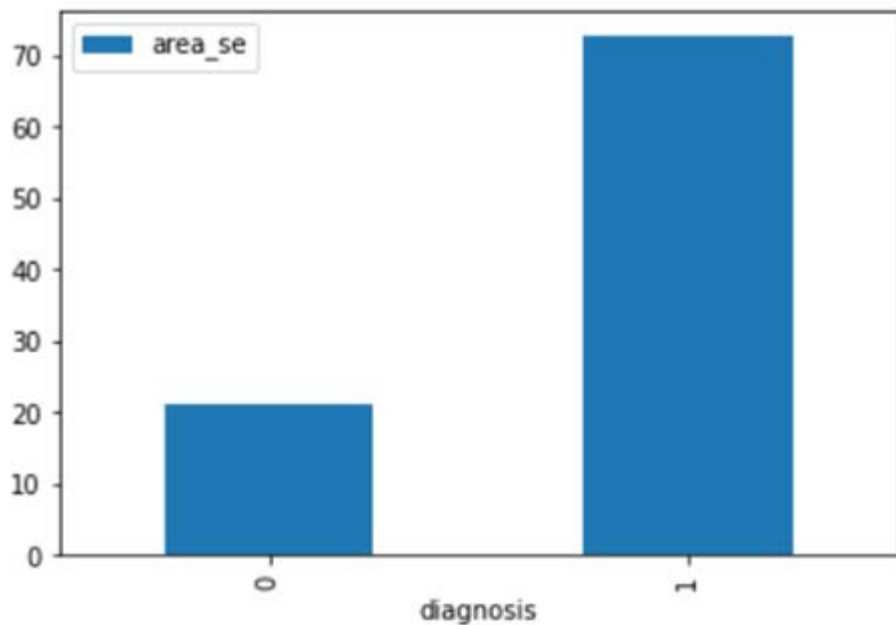
- Target : feature reduction을 통한 머신러닝의 성능 최적화

# 01 CONTENTS

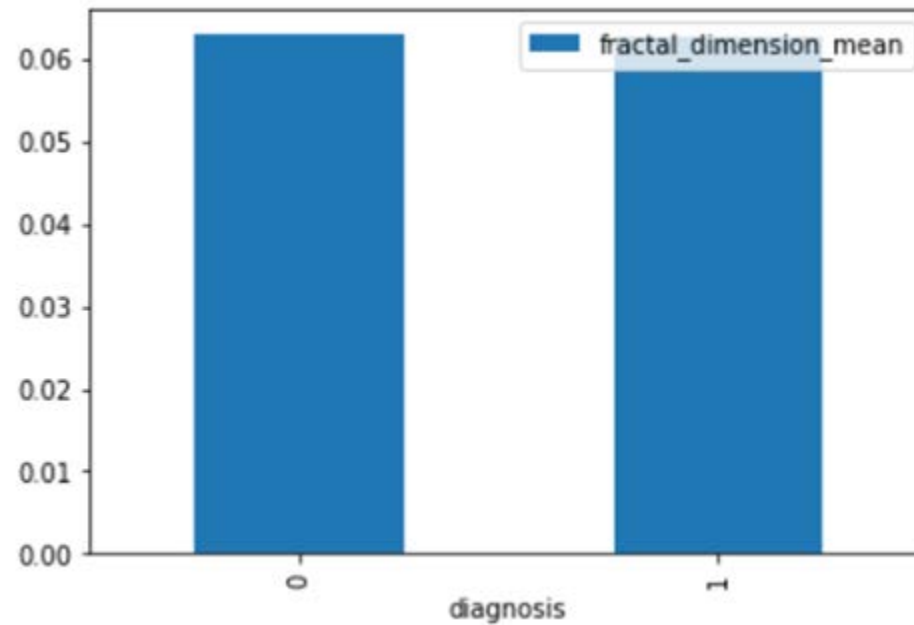
## Mean by target

- 데이터의 target과 각 feature의 평균값을 구하였다.

Best fit



Worst fit

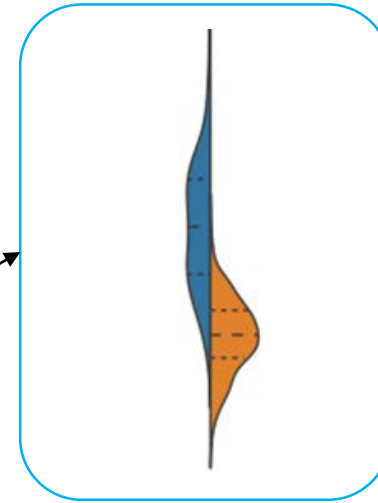
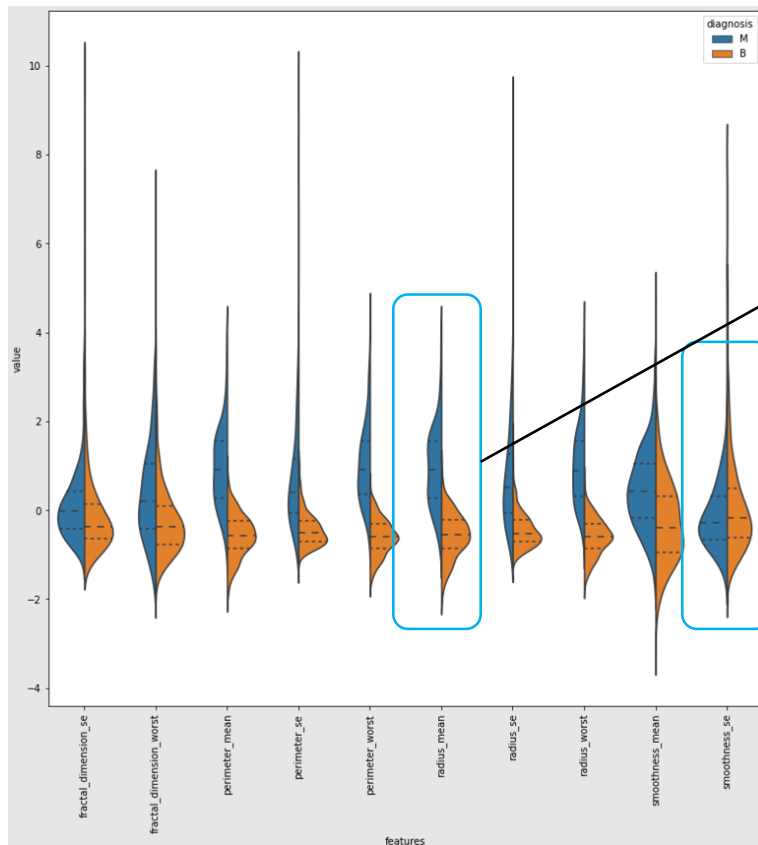


- 데이터의 label과 각 feature의 평균값이 비슷할수록 학습에 방해가 될 수 있다 -> feature 삭제

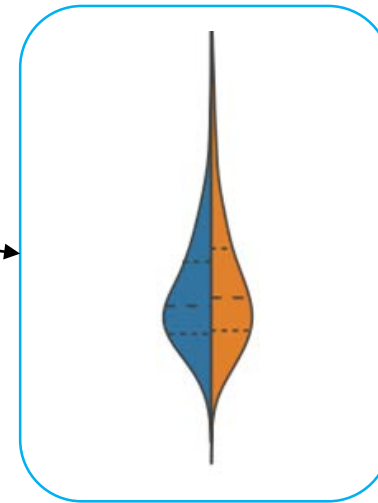
# 01 CONTENTS

## Distribution by target

- 데이터의 target에 따른 각 feature 값의 분포를 구하였다.



Best fit

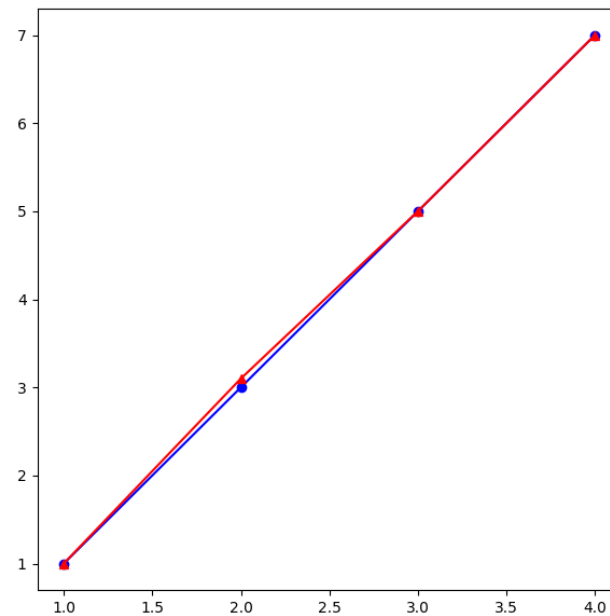


Worst fit

# 01 CONTENTS

Why distribution important?

```
1 import matplotlib.pyplot as plt
2
3 feature_1 = [1, 3, 5, 7]
4 feature_2 = [1, 3.1, 5, 7]
5
6 label = [1, 2, 3, 4]
7
8 f, ax = plt.subplots(figsize=(7,7))
9 ax.plot(label, feature_1, 'bo-')
10 ax.plot(label, feature_2, 'r^-')
11 plt.show()
```



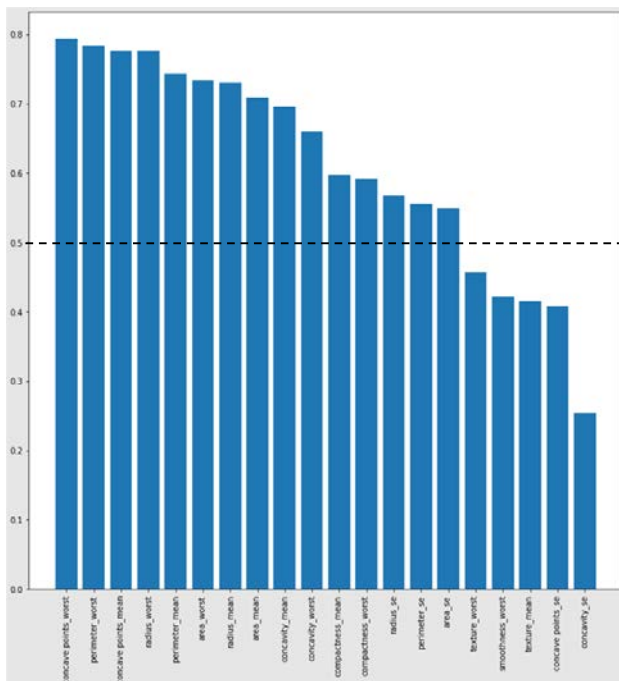
- 데이터의 feature의 값 분포가 비슷하니 그려진 선이 비슷할 수 있음을 알 수 있다 -> 따라서 삭제가 필요하다.



# 01 CONTENTS

## Correlation with target

- 데이터의 target과 각 feature간의 correlation을 구하였다.
- Correlation이 낮다는 것은 비교 대상의 관계가 크지 않다고 볼 수 있다.

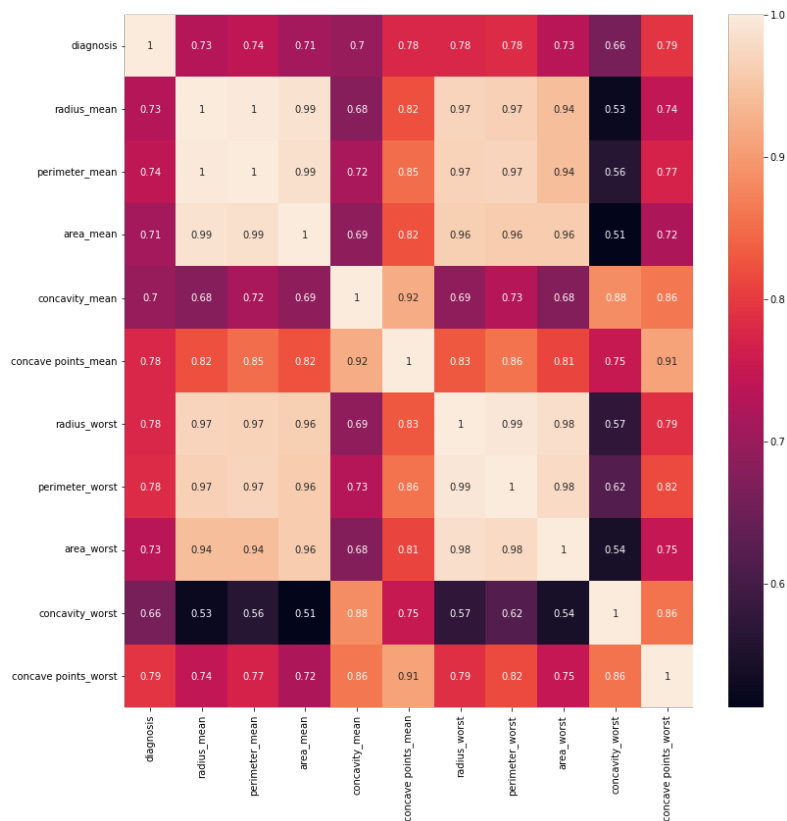


Correlation이 0.5 내외인 feature은  
학습에 방해가 될지도 모르니 삭제하겠다.

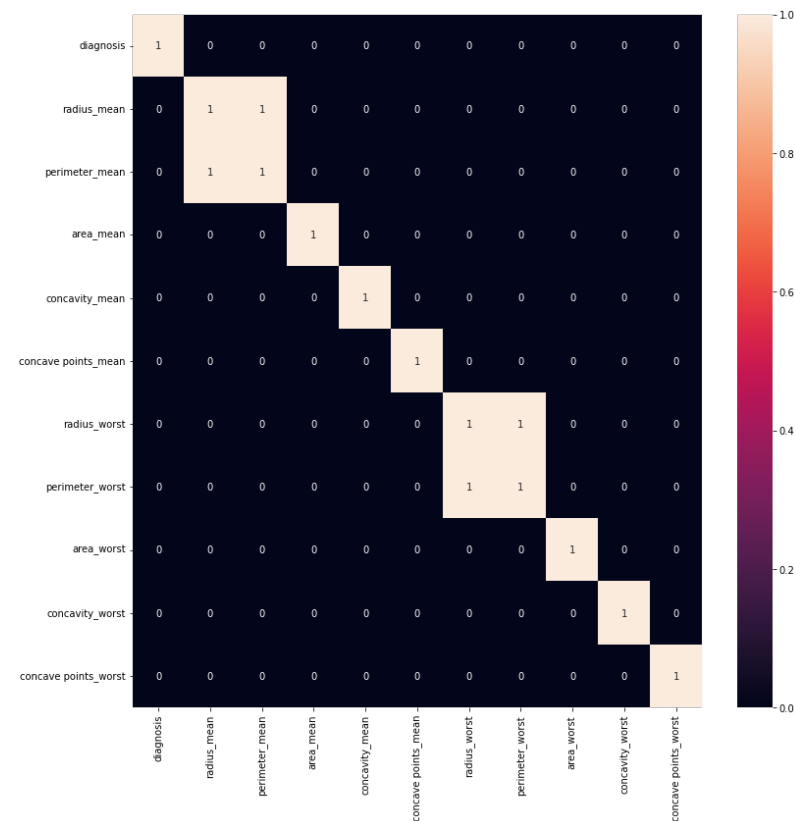
# 01 CONTENTS

## Draw heatmap

- 전체 데이터의 상관관계를 heatmap으로 나타내었다.



Correlation > 0.99



# 01 CONTENTS

How to create new feature with original data

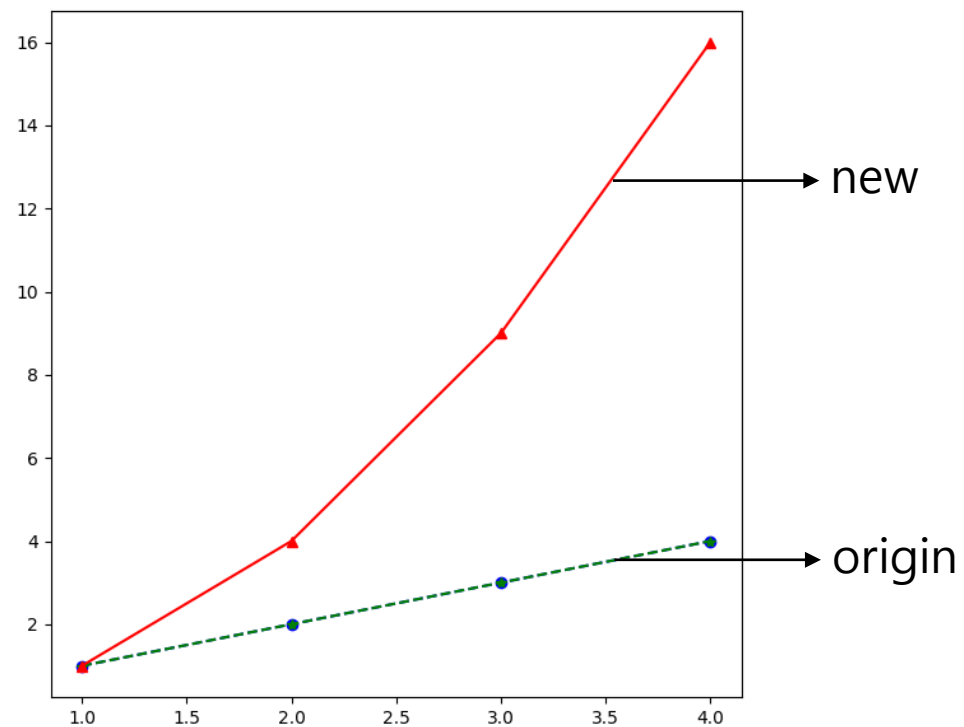
```
1  ->creat new feature
2  data['radXper'] = data['radius_mean'] * data['perimeter_mean']
3  data['radXper']
4
5  ->creat new feature
6  data['radXperW'] = data['radius_worst'] * data['perimeter_worst']
7  data['radXperW']
8
9  ->drop used data
10 data.drop(columns=['radius_mean', 'perimeter_mean', 'radius_worst', 'perimeter_worst'], inplace=True)
```

- 데이터의 correlation이 1에 가깝다는 것은 두 feature의 관계가 두텁다는 것이다.
  - Correlation이 비슷한 데이터 끼리 곱 연산을 하여 새 feature을 창조한다.

# 01 CONTENTS

Is new feature useful?

```
1 import matplotlib.pyplot as plt
2
3 origin1 = [1, 2, 3, 4]
4 origin2 = [1, 2, 3, 4]
5 new = [1*1, 2*2, 3*3, 4*4]
6
7 f, ax = plt.subplots(figsize=(7,7))
8 ax.plot([1,2,3,4], origin1, 'bo--')
9 ax.plot([1,2,3,4], origin2, 'g*--')
10 ax.plot([1,2,3,4], new, 'r^--')
11 plt.show()
```



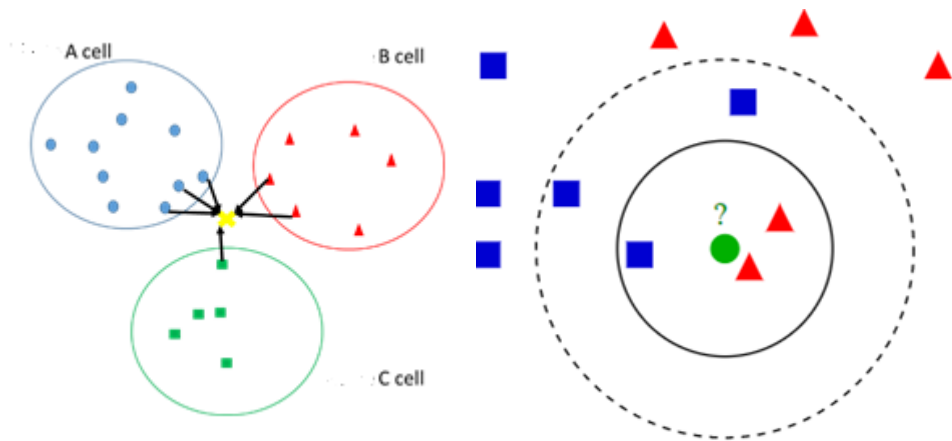
- 새 데이터의 x축이 증가할수록 기존의 데이터에 비해 큰 폭으로 차이가 남을 알 수 있다  
-> 분류에 도움이 될 수 있다.

# 01 CONTENTS

## Machine Learning

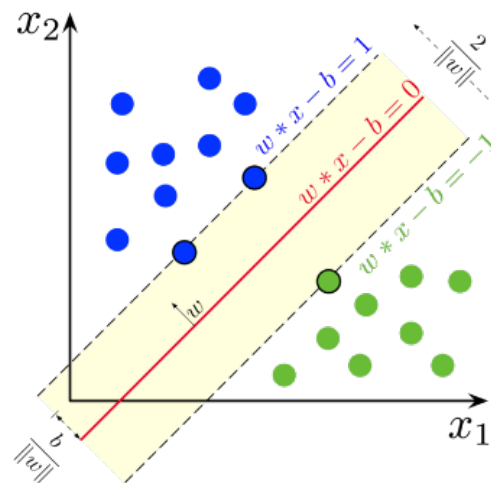
- 머신러닝의 KNN과 SVM을 이용하여 데이터를 분류하려고 한다.

K - Nearest Neighbor



출처: Wikimedia commons

Support Vector Machine



출처: Wikimedia commons

# 01 CONTENTS

## Score



```
1 KNN
2 knn = KNeighborsClassifier()
3 knn.fit(X_train, y_train)
4 knn.score(X_test, y_test)
5 >> 0.9440559440559441
6
7 SVM
8 svmc = SVC(kernel='linear')
9 svmc.fit(X_train, y_train)
10 svmc.score(X_test, y_test)
11 >> 0.951048951048951
```

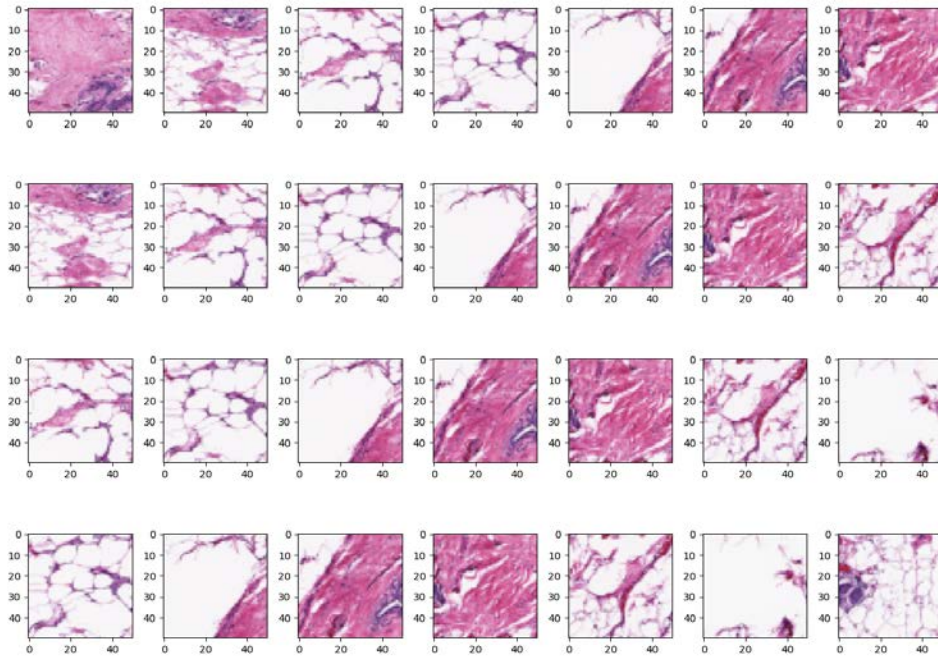
- KNN과 SVM을 사용한 결과 약 94%, 95%의 score를 달성하게 되었다.
  - 따라서 위에서 진행한 process가 적절했다고 생각할 수 있다.
- 예측에 사용된 feature을 돌아보면 모두 크기와 형태에 관련된 것만 남았음을 알 수 있다.
  - 따라서 도메인 지식이 없어도 암을 판별하는 기준을 예상 할수있다.

# 02

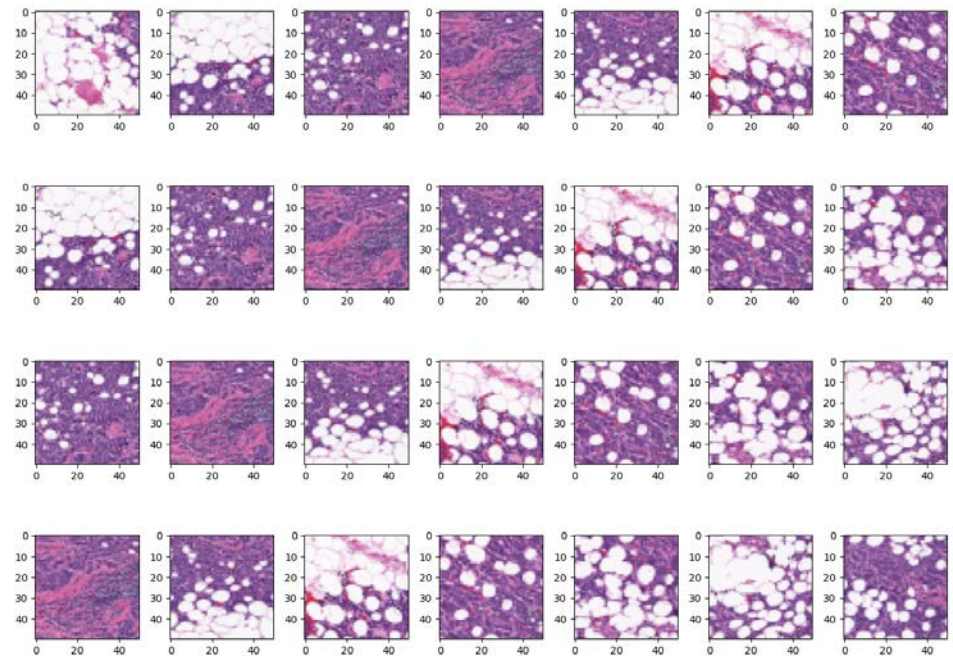
Breast Cancer diagnosis with DL

## 02 CONTENTS

### Inspection & Target



- Normal Image



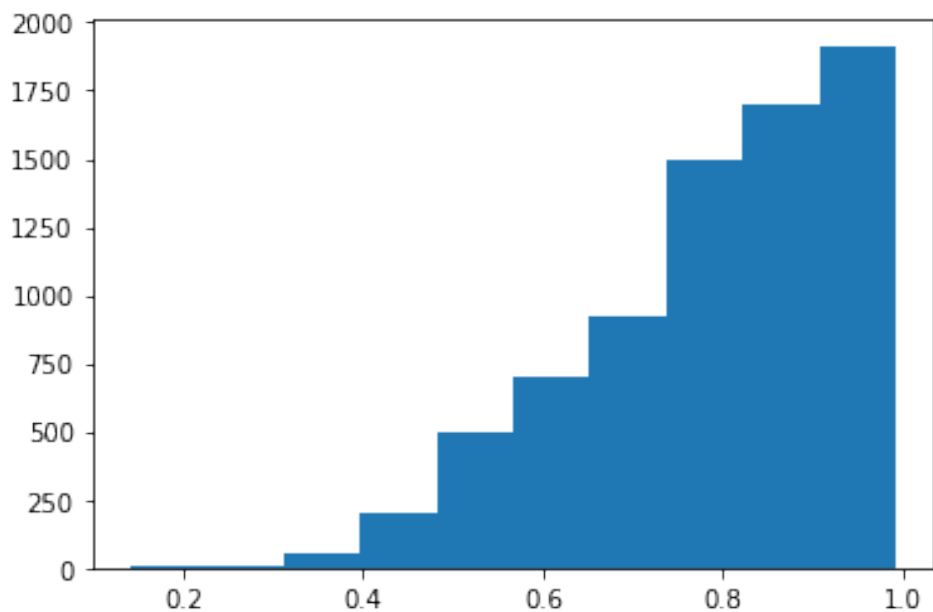
- Cancer Image

- Target : 딥러닝을 이용하여 암 환자 진단하기

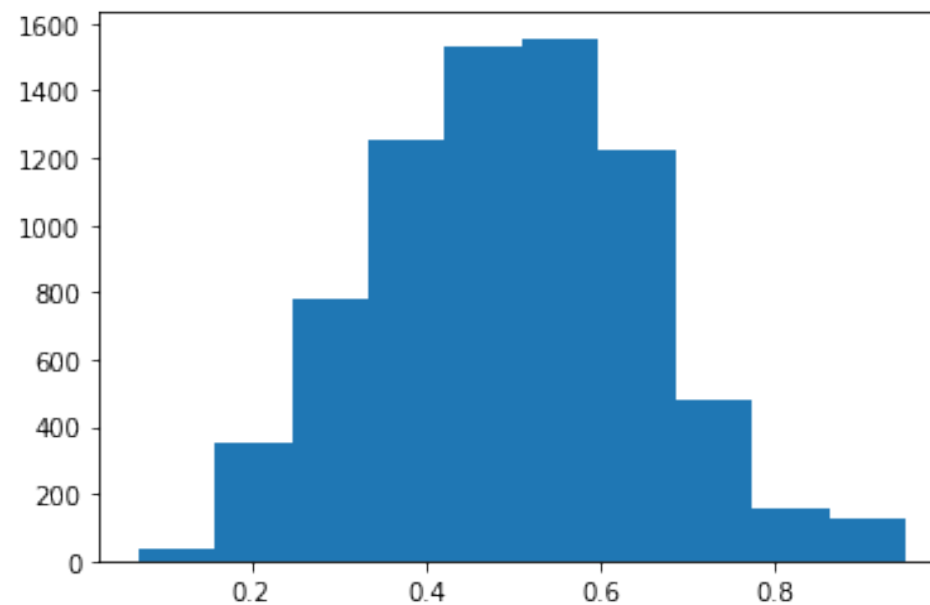


## 02 CONTENTS

### Brightness



• Normal Image

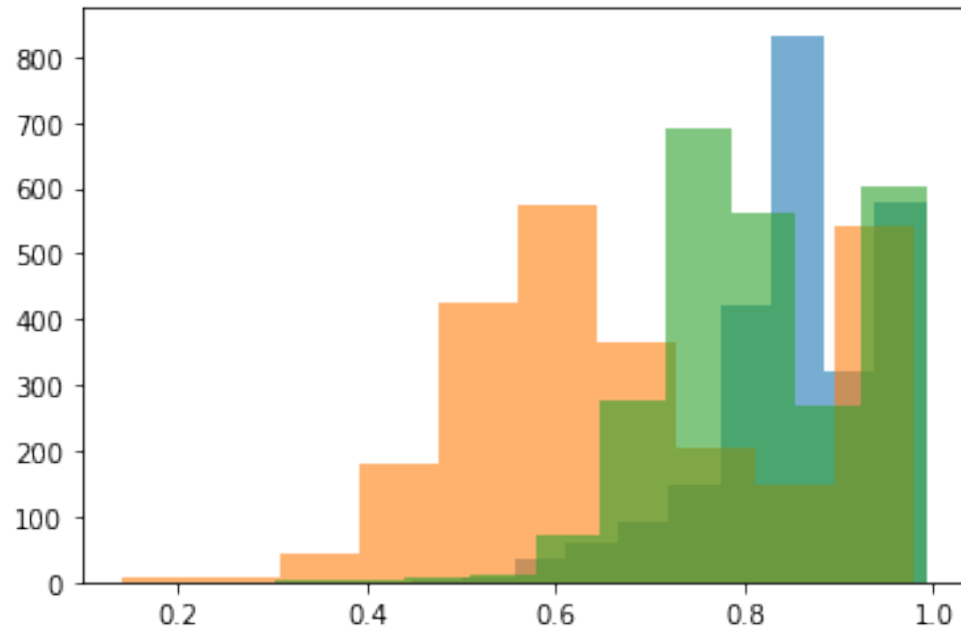


• Cancer Image

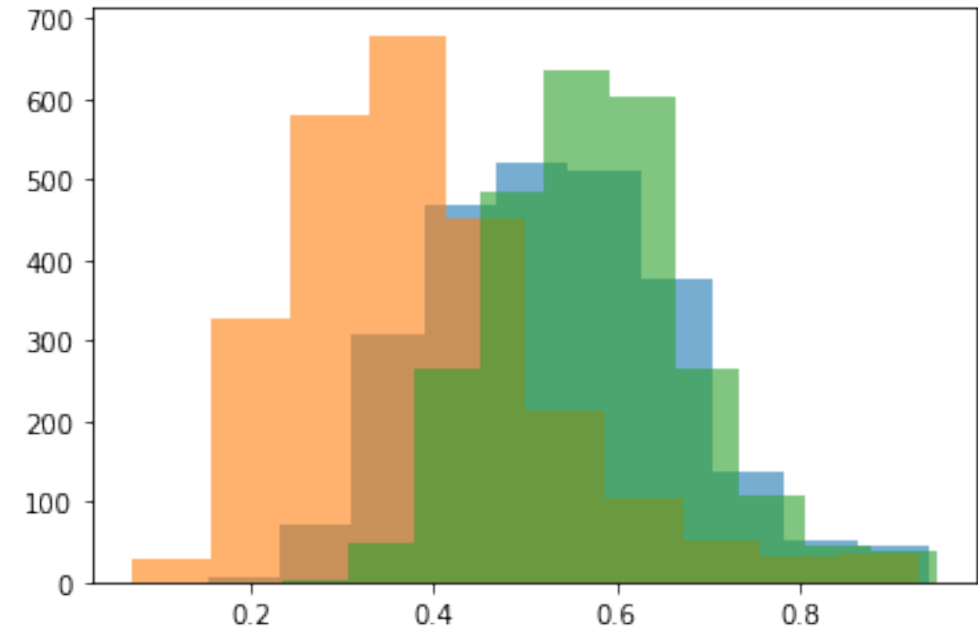
- 암 환자의 세포가 정상인에 비해 어두운 부분이 많다는 것을 알 수 있다.

## 02 CONTENTS

### Density



- Normal Image

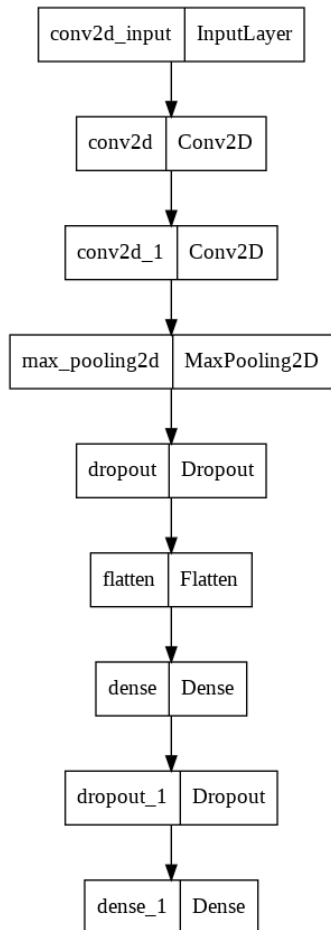


- Cancer Image

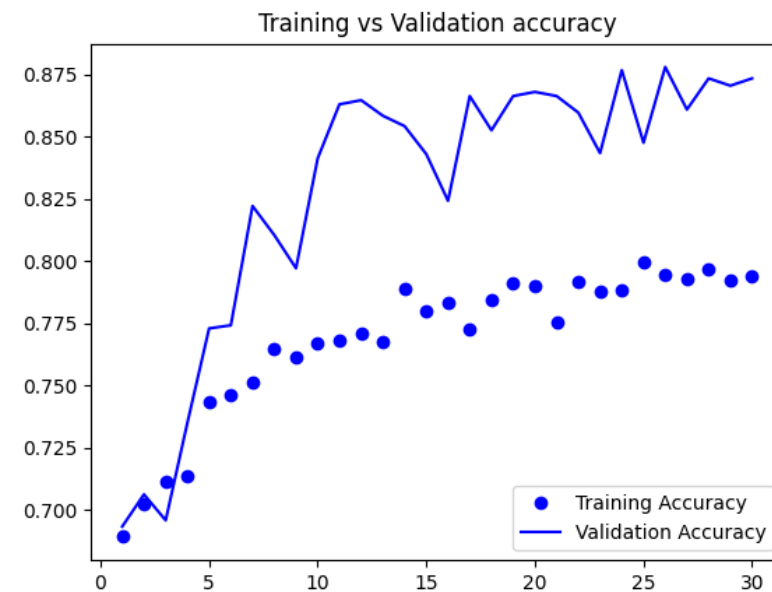
- 암 환자의 세포가 정상인에 비해 세포가 밀집하게 분포한 것을 알 수 있다.

## 02 CONTENTS

Model : CNN

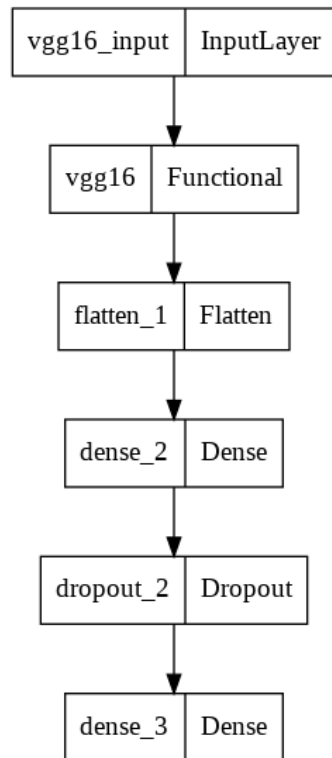


```
1 model = models.Sequential()
2 model.add(layers.Conv2D(32, 3, activation='relu', input_shape=(50,50,3)))
3 model.add(layers.Conv2D(64, 3, activation='relu'))
4 model.add(layers.MaxPooling2D((2, 2)))
5 model.add(layers.Dropout(0.25))
6 model.add(layers.Flatten())
7 model.add(layers.Dense(128, activation='relu'))
8 model.add(layers.Dropout(0.5))
9 model.add(layers.Dense(1, activation='sigmoid'))
10 model.compile(loss=tf.keras.losses.BinaryCrossentropy(),
11               optimizer=tf.keras.optimizers.RMSprop(learning_rate=1e-5),
12               metrics=['accuracy'])
```

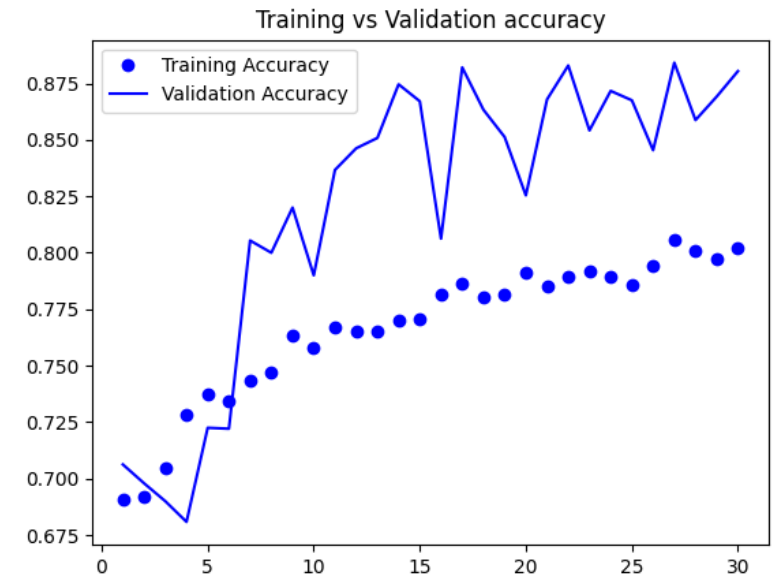


## 02 CONTENTS

Model : VGG16

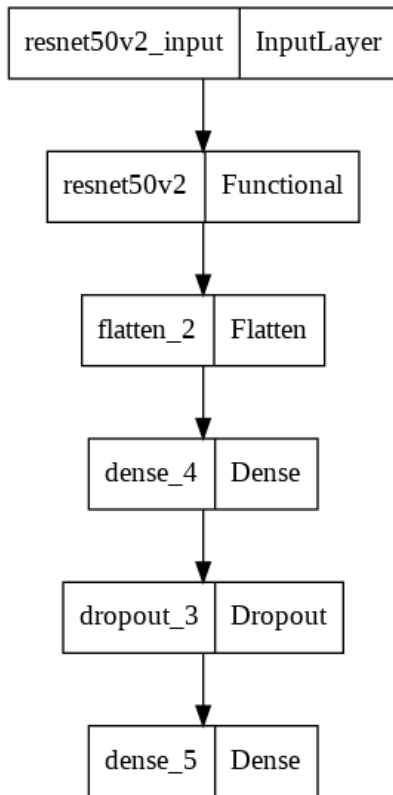


```
1 model = models.Sequential()
2 model.add(VGG16(weights='imagenet', include_top=False, input_shape=(50,50,3)))
3 model.add(layers.Flatten())
4 model.add(layers.Dense(128, activation='relu'))
5 model.add(layers.Dropout(0.5))
6 model.add(layers.Dense(1, activation='sigmoid'))
7 model.compile(loss=tf.keras.losses.BinaryCrossentropy(),
8               optimizer=tf.keras.optimizers.RMSprop(learning_rate=1e-5),
9               metrics=['accuracy'])
```

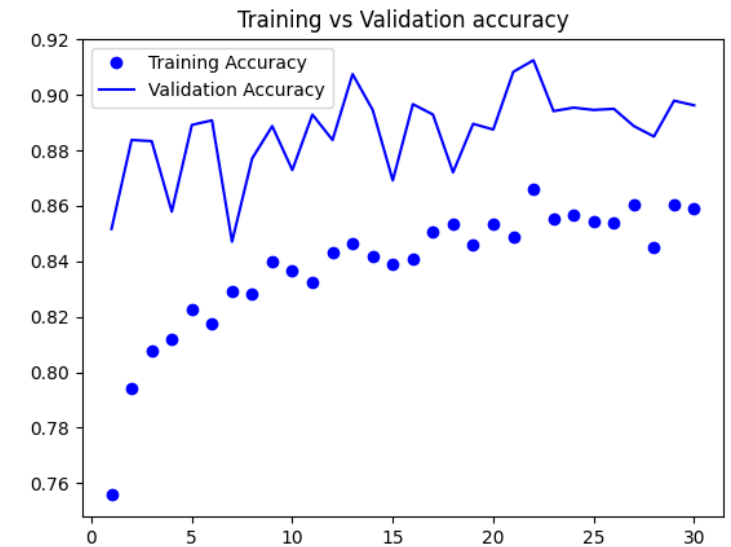


## 02 CONTENTS

Model : ResNet



```
1 model = models.Sequential()
2 model.add(ResNet50V2(weights='imagenet', include_top=False, input_shape=(50,50,3)))
3 model.add(layers.Flatten())
4 model.add(layers.Dense(128, activation='relu'))
5 model.add(layers.Dropout(0.5))
6 model.add(layers.Dense(1, activation='sigmoid'))
7 model.compile(loss=tf.keras.losses.BinaryCrossentropy(),
8               optimizer=tf.keras.optimizers.RMSprop(learning_rate=1e-5),
9               metrics=['accuracy'])
```



03

Analysis

## 03 CONTENTS

### Diagnosis Function

```
1 def diagnosis(data, image):
2     # 데이터 전처리
3     data = pd.DataFrame(data).T
4     data.drop(columns=['id', 'fractal_dimension_mean', 'texture_se', 'symmetry_se', 'Unnamed: 32',
5                        'compactness_se', 'fractal_dimension_se', 'fractal_dimension_worst', 'smoothness_mean',
6                        'smoothness_se', 'symmetry_mean', 'symmetry_worst', 'compactness_mean',
7                        'compactness_worst', 'radius_se', 'perimeter_se', 'area_se', 'texture_worst', 'smoothness_worst',
8                        'texture_mean', 'concave points_se', 'concavity_se'], inplace=True)
9     data['radXper'] = data['radius_mean'] * data['perimeter_mean']
10    data['radXperW'] = data['radius_worst'] * data['perimeter_worst']
11    data.drop(columns=['radius_mean', 'perimeter_mean', 'radius_worst', 'perimeter_worst'], inplace=True)
12
13    # 모델 로드
14    svm_load = joblib.load('/content/drive/MyDrive/svm.pkl')
15    resnet_load = models.load_model('/content/drive/MyDrive/ResNetTrain.h5')
16
17    # 예측
18    data_pred = svm_load.predict(data.T.to_numpy()[1:].reshape(-1, 8))[0]
19    image_pred = resnet_load.predict(image.reshape(-1, 50, 50, 3))[0][0]
20
21    if not data_pred:
22        print(f'유방암일 확률은 {int(image_pred*100)}% 입니다.' if image_pred > 0.5 else '의사의 확인이 필요합니다.')
23    else:
24        print(f'유방암일 확률은 {int(image_pred*100)}% 입니다.' if image_pred <= 0.5 else '의사의 확인이 필요합니다.')
```

```
1 ->정상인 진단
2 diagnosis(data_none, img_none)
3 >> 유방암일 확률은 1% 입니다.
4
5 ->암환자 진단
6 diagnosis(data_cancer, img_cancer)
7 >> 유방암일 확률은 98% 입니다.
8
9 ->정상인과 암환자의 데이터를 섞었을 때
10 diagnosis(data_none, img_cancer)
11 >> 의사의 확인이 필요합니다.
```

- 실제로 바로 사용 가능한 함수를 구성하였다. 데이터에 따른 예측도 오류 없이 도출되었다.

## 03 CONTENTS

### Limit of Model

```
1 base = ResNet50V2(weights='imagenet', include_top=False, input_shape=(50,50,3))
2
3 base.trainable = True
4 set_train = False
5
6 for layer in base.layers:
7     if layer.name == 'conv5_block3_3_conv':
8         set_train = True
9         if set_train:
10             layer.trainable = True
11         else:
12             layer.trainable = False
13
14 model = models.Sequential()
15 model.add(base)
16 model.add(layers.Flatten())
17 model.add(layers.Dense(128, activation='relu'))
18 model.add(layers.Dropout(0.5))
19 model.add(layers.Dense(1, activation='sigmoid'))
20 model.compile(loss=tf.keras.losses.BinaryCrossentropy(),
21               optimizer=tf.keras.optimizers.RMSprop(learning_rate=1e-5),
22               metrics=['accuracy'])
```

- Imagenet의 모델은 동결 해제시킨 학습이나 일반 학습이나 큰 차이가 없었다.

```
1 knn.predict(test_data)
2 >> breast cancer
3
4 svm.predict(test_data)
5 >> breast cancer
6
7 model.predict(test_data)
8 >> breast cancer
```

- 모든 모델이 92% 내외로 예측을 하고 있다.



## 03 CONTENTS

### Limit & Apply

#### Limit

1. 적은 데이터의 양
2. validation accuracy > accuracy  
->underfitting이나 validation data가 예측하기 쉬운 구성일 수 있음.
3. Accuracy를 높이려면 앙상블이나, 최적화된 layer 구현 필요 예상.

#### Apply

1. 데이터를 입력하면 진단결과를 출력해주는 웹사이트 제작
2. 의료현장의 유방암 진단 보조기구로 사용

Thank  
You

---