# 지도학습: 분류 (실습)

홍익 대학교
Hyun-Sun Ryu

# 분류 프로세스

문제정의 → 데이터 수집 → 탐색적 데이터분석 → 모델 생성 → 학습/예측 → 평가

# 분류 프로세스

| 문제정의 | 데이터수집 | 탐색적 데이터분석 | 모델생성 | 학습/예측 | 평가 |
|---|---|---|---|---|---|

- **Titanic 생존자 예측 모델 만들기**

< titanic 탑승자 데이터를 바탕으로 생존자와 사망자 분류하기 >

# 분류 프로세스

- 캐글 누리집 – compete항목에서 다운로드

  https://www.kaggle.com/c/titanic/data

- 모델 훈련을 위하여 train 데이터가 필요함

# 분류 프로세스

# 분류 프로세스

# 분류 프로세스

# 분류 프로세스

# 분류 프로세스

# 분류 프로세스

**사망 ROC**



**생존 ROC**

# 분류 프로세스

# 분류 프로세스

| 문제정의 | 데이터수집 | 탐색적데이터분석 | 모델생성 | 학습/예측 | **평가** |
| --- | --- | --- | --- | --- | --- |

Test and Score, Confusion matrix, ROC analysis 위젯으로

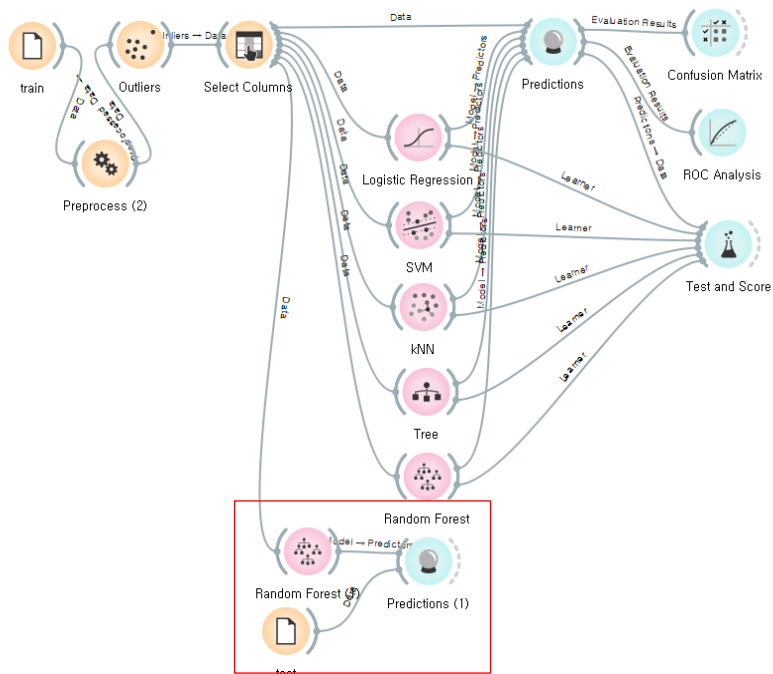분류 모델을 평가한 결과 Random forest의 분류 결과가 가장

신뢰도가 높다고 판단

# 분류 프로세스

# 질문 있나요?

hsryu13@hongik.ac.kr