

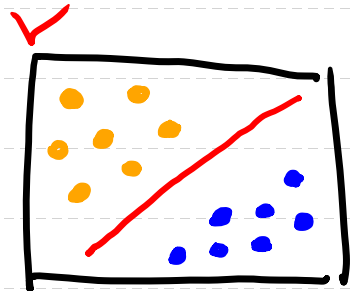
10/20 (월)

Sklearn으로 하기.
Sklearn. SVM

- SVM (support vector machine)
 - regression (회귀) → SUR ✓
 - classification (어떤것) → SUC ✓

→ Deep Learning이 나오니까 위상이 쿨한 딥러닝 그대로 좋은 성능의
모델 자체도 가능하도록 함.

SVM → Decision Boundaries (결정경계) ⇒ 데이터를
분류하기 위한 기법.



새로운 데이터가 들어왔을때 기법을 이용해서 어느 쪽에 포함되는
가를 예측

- 분류 (classification)
- 분류안의 label의 평균 (regression)

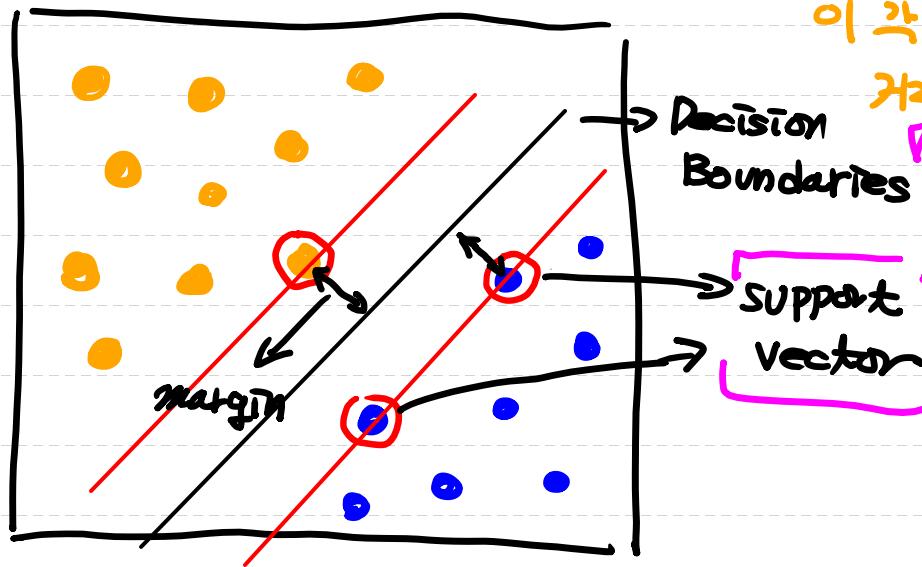
SVM → 데이터를 선형으로 분리하는 최적의 선형 결정경계를 찾는 알고리즘 !!

☆ ((선형으로 직선으로 포함하는 개념 직선 → 직선으로 포함 ⇒ 선형은 1차 함수로 표현
되는 건가요?? ✗
평면이나 곡선도 포함하는 개념.

• feature 가 3개이면 Decision Boundaries는 평면으로 표현

feature의 개수가 증가하면 Decision Boundaries도 "고차원" → 분리 초평면
✓ hyperplane ✓

Decision Boundaries를 구하기 위해 ⇒ Support vector



↪ 각 영역에 들어있는 데이터 포인트
이 각 group의 support vector들간의
거리가 가장 먼 support vector로 구해요!!

참고로 → SVM은 support vector만
이용하는 기법이기에

⇓
✓ 시간을 줄일수 있어요!!

· 주의점

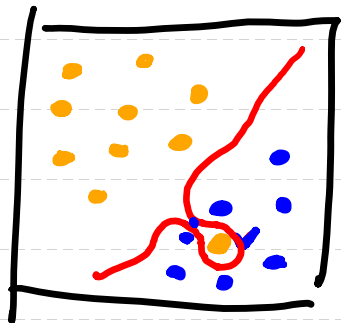
→ 만약에 이상치(outlier)가 존재하면.
이 이상치가 support vector가 될 수 있고

"C" 절충해야
하는 값!!

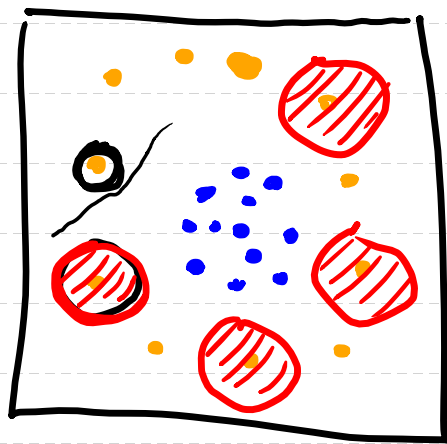
이런 경우 overfitting의 문제가 발생 \Rightarrow margin이 작아지게 되며
↓
Hard margin

(\Rightarrow 아니 그러면 이상치를 제거, 변경하면 되는거 아냐?)
 \Rightarrow 항상 그런건 아니예요!

그중기 때문에 모델로 만들때 \rightarrow 약간의 오차를 허용하는 방식을 고려.



↓ sklearn에서는 cost라고 표현되는 "C" hyperparameter 이용
"C" 기본값 1 C의 값을 크게 하면. 다른 클래스에 있는 데이터 포인트를 적게 허용. \rightarrow overfitting
C의 값은 ↓



→ "kernel 기법"

→ 주어진 data를 고차원의 공간으로 projection (투영)

sklearn [kernel 호지킹]

- linear (직선으로 Decision Boundaries 가 구해짐)
- poly (2 → 3)
- rbf (가장 좋음) →

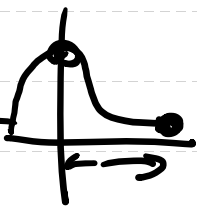
• "gamma"

→ 데이터들이 영향력을 행사하는

거리

→ gamma 값이 클수록 영향력 행사 거리 ↓

gamma 값이 작을수록 " " " ↑



"gamma" hyperparameter.

→ Decision Boundaries를 얼마나 유연하게 그릴지를 조절하는 값!!

이 값이 크면 구불구불 → overfitting

" " " 작으면 직선 → underfitting

구불구불 → Overfitting

· SVM ["C" , "gamma"] \Rightarrow 잘 선택할 수 있을까요?
"hyperparameter"

① 수동으로 적당히 "C", "gamma" 값을 찾는거예요!! \rightarrow X (노동집약적)

② Sklearn 최적의 parameter를 찾아주는 방법 제공 (자동인)

\rightarrow Grid SearchCV \rightarrow "Cross validation"

★ 코드
가리고 찾아보아요! \Rightarrow 사용할 parameter를
여러개 지정해 놓고 CV를 실행해서 최적의 accuracy를
도출하는 hyperparameter를 찾아주세요! \odot