

① bias  $\rightarrow$  예측값의 표준과 실제값과의 차이.  
(정답에 얼마나 가까워졌는지.)

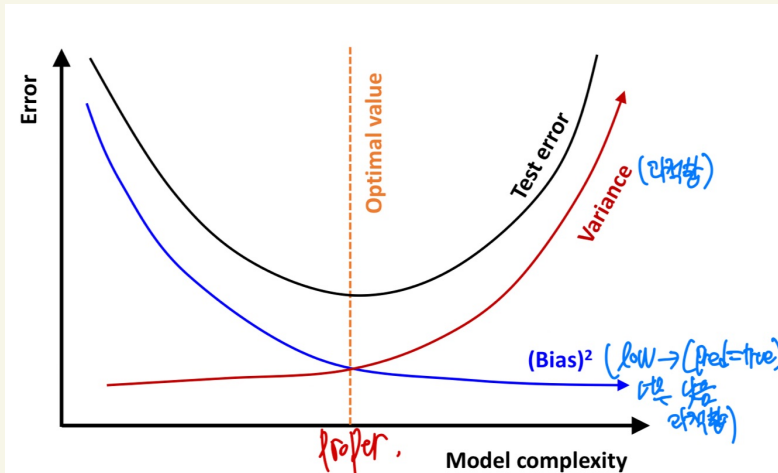
bias  $\uparrow \rightarrow$  정답과 차이  $\uparrow \rightarrow$  정답 값 로 맞음.

② Variance  $\rightarrow$  예측한 값을 실제의 다른 정도.

Variance  $\uparrow \rightarrow$  각각의 예측값 차이  $\uparrow \rightarrow$  데이터셋 민감.

따라서 bias 높다는 것은 성능 안좋은 Underfitting.

Variance 높다는 것은 특정 데이터셋에 overfitting.



Bias-variance tradeoff.

이런 Bias와 variance를 줄이기 위한 방법.



앙상블

① Bagging.

→ 복원추출하여 다른 데이터셋 만들.  
(bootstrap)

Aggregating 방법 ① majority voting

② weighted voting

③ stacking.

↳ 예측값을 한번더 input으로  
넣음.

Bagging → 독립-4 개행하여 variance가  $1/m$ 으로  
줄어듦에 따라 분산 감소는 부분 앙상블 효과 ↓

# Random forest

Bagging 한 종류로 변수 랜덤하게 선택.

∴ 더 각각의 bootstrap 적용.

① 변수 중요도까지 산출

## ② Boosting

Ada boost → 약한 학습기.

미분가능한 데이터 셋 찾고 그것에

분류도에 weight 순차.

$$w^{(1)} \leftarrow w^{(2)} \cdot \exp(-\alpha e^{y_i} h_e(x_i))$$

→  $\alpha$  -1  
강한 |  
↓  
오류면  $w \uparrow$   
정답  $w \downarrow$

# GBM

↳ 정사각형 + Boosting.

잔차를 예측하고 그것을 이용하여  
보완하여 다음 모델에 사용.

loss function의 negative gradient. = 잔차

## XGBoost

↳ GBM의 상위 버전

속도 + 성능 향상 (병렬 처리)