

AI 스터디

조 상 구
2025.12.28



1. Getting Started with scikit learn API

- Cross validation with hyperparameter tuning
- RandomGridsearch with cross validation
- Pipeline (data preprocessing and estimators)
- Model selection and evaluation with performance index
- Loss function
- Bias and variance, overfitting & under fitting

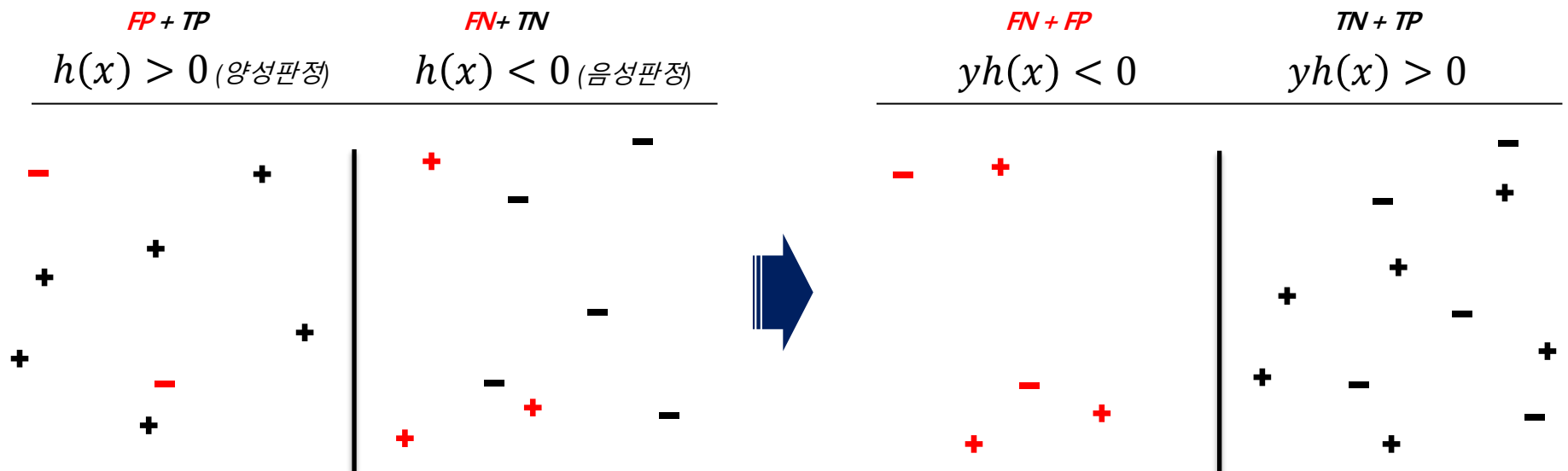


Loss function

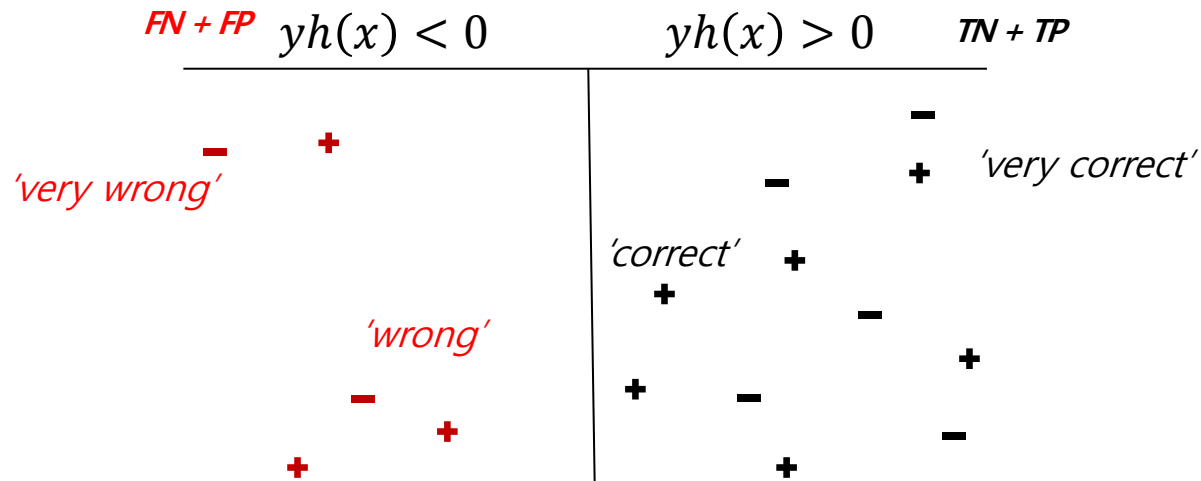
실제 데이터를 정확하게 예측 판단하지 못하는 비율

- 허위음성율(FNR, False Negative Ratio)과 허위양성율(FPR, False Positive Ratio)을 최소화하는 것이 목적

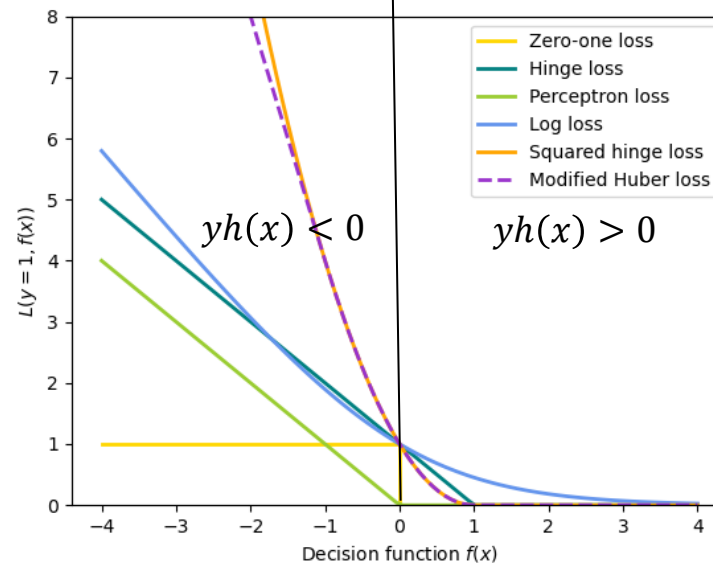
$$\text{Error rate} = \text{FN} + \text{FP} = \frac{1}{n} \sum_{i=1}^n [y_i \neq \text{sign}(h(x))]$$



Loss functions



- ✓ 예측이 맞다면 손실은 '0' 예측이 틀리면 '1' 손실의 penalty를 부여 (zero, one loss function)
- ✓ 알고리즘별로 loss function이 다름



SGD: convex loss functions

https://scikit-learn.org/stable/auto_examples/linear_model/plot_sgd_loss_functions.html

Loss functions

평균 제곱 오차(mean squared error, MAE)

머신러닝 모형이 예측한 확률벡터와 실제 라벨(원핫코딩)을 고차원 공간의 점으로 이해하여 유클리디언(피카고라스정리)로 거리(Distance)를 측정한 값이 Error

$$\text{Error} = \frac{1}{2} \sum_k (y_k - t_k)^2$$

t_k
실제 원핫코딩 라벨 k 번째 좌표

y_k
예측 모형이 k라고 예측한 확률 값

Error
평균제곱오차

0 = [1,0,0,0,0,0,0,0,0,0]

1 = [0,1,0,0,0,0,0,0,0,0]

⋮

9 = [0,0,0,0,0,0,0,0,0,9]

$f_{\text{model A}}$



[0.9,0.1,0,0,0,0,0,0,0,0]

$f_{\text{model B}}$



[0.1,0.9,0,0,0,0,0,0,0,0]

$$E(A) = \frac{1}{2} ((1 - 0.9)^2 + (1 - 0.1)^2 + 0^2 \dots + 0^2) = 0.01$$

$$E(B) = \frac{1}{2} ((1 - 0.1)^2 + (1 - 0.9)^2 + 0^2 \dots + 0^2) = 0.81$$

Loss functions

교차 엔트로피 오차(cross entropy error)

교차 엔트로피는 정보이론에서 확률분포사이의 거리를 재는 방법

$$\text{Error} = - \sum_k t_k \times \log y_k$$

t_k
실제 원핫코딩 라벨 k 번째 좌표

y_k
예측 모형이 k라고 예측한 확률 값

Error
교차 엔트로피 오차

0 = [1,0,0,0,0,0,0,0,0,0]

1 = [0,1,0,0,0,0,0,0,0,0]

⋮

9 = [0,0,0,0,0,0,0,0,0,9]

$f_{\text{model A}}$



[0.9,0.1,0,0,0,0,0,0,0,0]

$$E(A) = -(1 \log 0.9 + 0.1 \log 0.1 + \dots + 0) = -\log 0.9 = 0.11$$

Error = $-\log y_k$
만약 라벨이 y_k 이면

$f_{\text{model B}}$



[0.1,0.9,0,0,0,0,0,0,0,0]

$$E(B) = -(0 \log 0.9 + 1 \log 0.1 + \dots + 0) = -\log 0.1 = 2.3$$

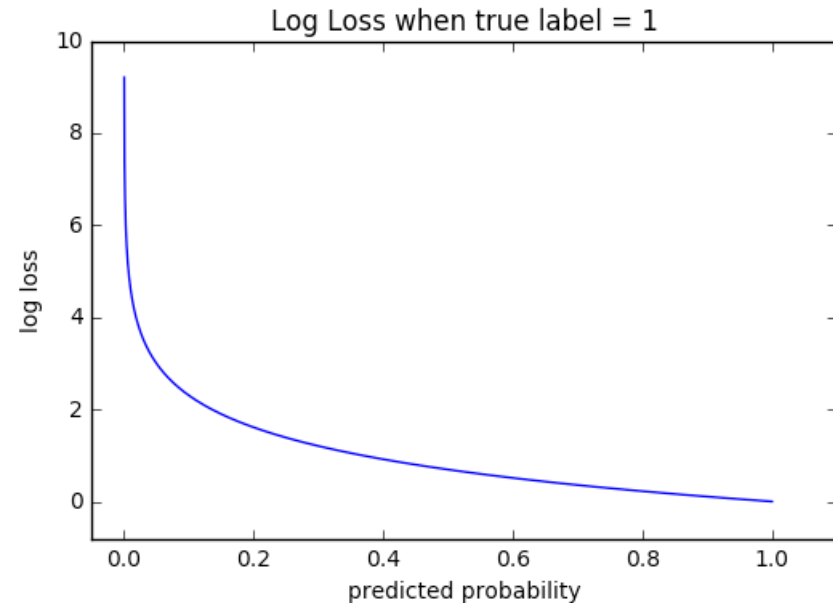
Loss functions

이진분류 교차 엔트로피 오차(binary cross entropy error)

이진분류 교차엔트로피는 log loss 함수

$$\text{Log loss} = -y\log p - (1 - y)(1 - p)\log(1 - p)$$

$$= \begin{cases} -y\log p, & \text{if } y = 1 \\ -(1 - y)(1 - p)\log(1 - p), & \text{if } y = 0 \end{cases}$$



예) 실제 양성($y=1$)인데 예측 모형 A와 B가 각각 '0.9', '0.7'의 확률(p)로 양성(1)이라고 판정할 경우 각각 로그손실

$$Loss_{(p=0.9)} = -1\log^{0.9} - (1 - 1)(1 - 0.9)\log^{(1-0.9)} = \log^{0.9} = 0.10536 \quad : \text{예측모형 A}$$

$$Loss_{(p=0.7)} = -1\log^{0.7} - (1 - 1)(1 - 0.7)\log^{(1-0.7)} = \log^{0.7} = 0.35667 \quad : \text{예측모형 B}$$