

Machine Learning

Linear Model



START

- 선형 분류모델을 이해하고 사용 할 수 있다.
- 다양한 분류평가 지표를 이해 할 수 있다.



Linear Model

(Classification)

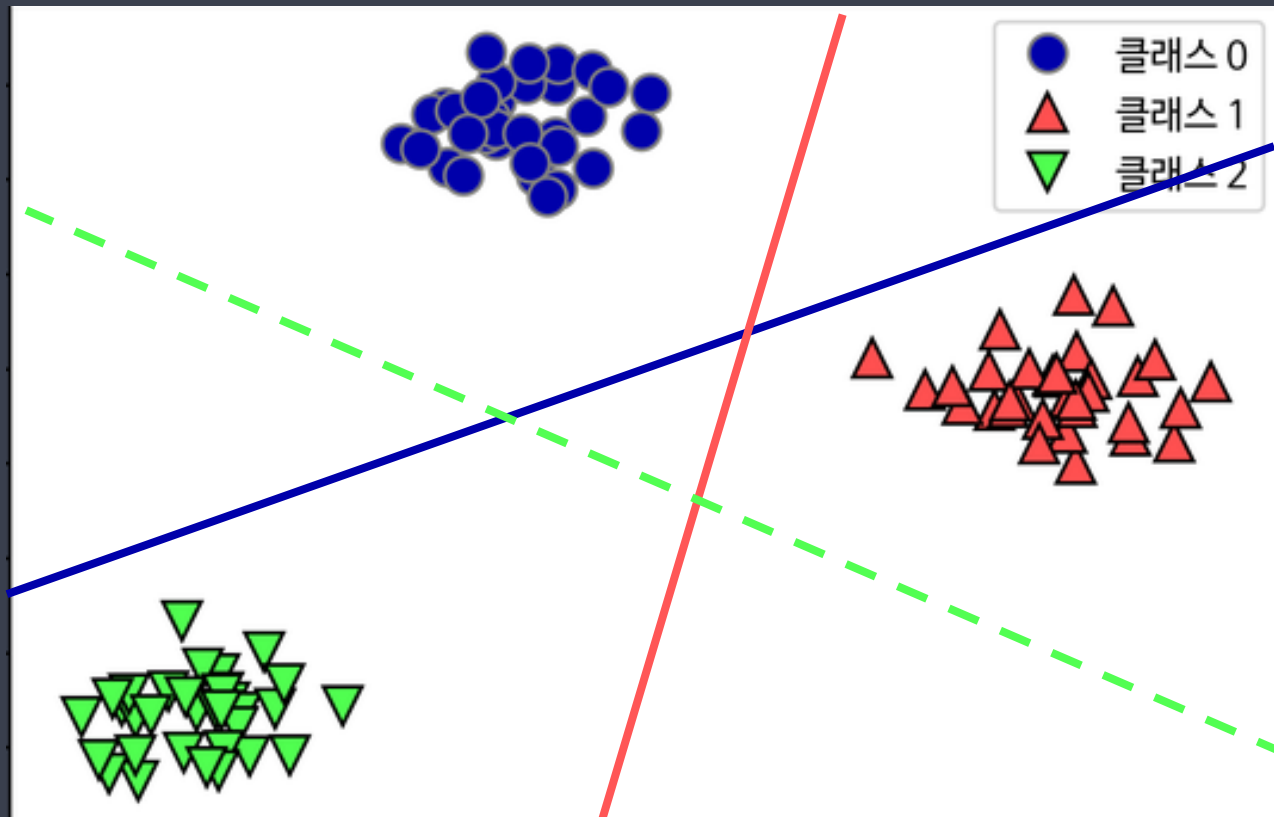
- 선형 분류모델을 이해하고 사용 할 수 있다.
- Logistic Regression 모델의 필요성을 이해할 수 있다.
- Logistic Regression 모델의 비용함수를 이해할 수 있다.

분류용 선형 모델

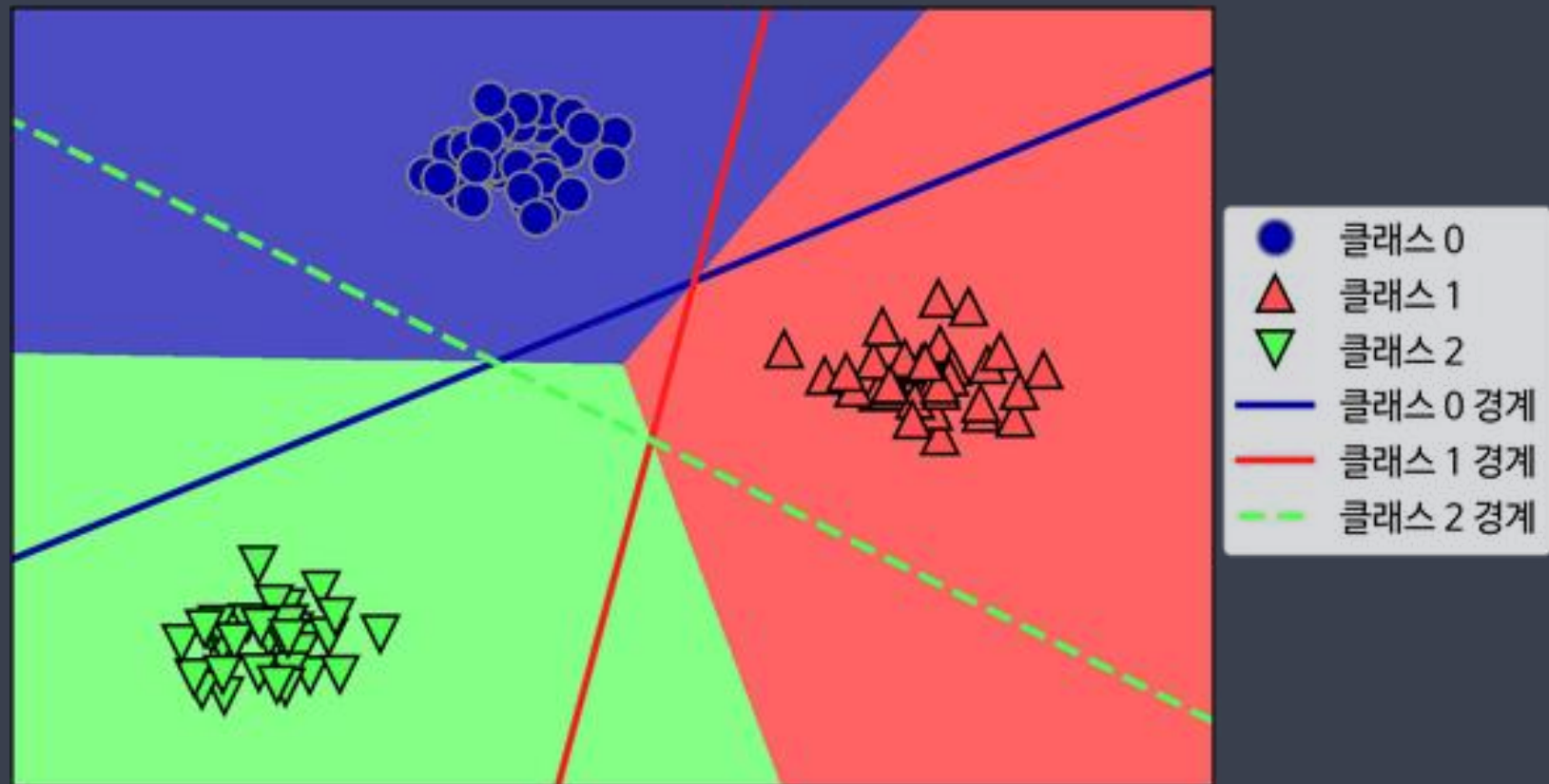
$$y = w_1x_1 + w_2x_2 + w_3x_3 + \cdots + w_px_p + b > 0$$

- 특성들의 가중치 합이 0보다 크면 class를 +1(양성클래스)
0보다 작다면 클래스를 -1(음성클래스)로 분류한다.
- 분류용 선형모델은 결정 경계가 입력의 선형함수
- 일대다 방법을 통해 다중 클래스 분류

분류용 선형 모델



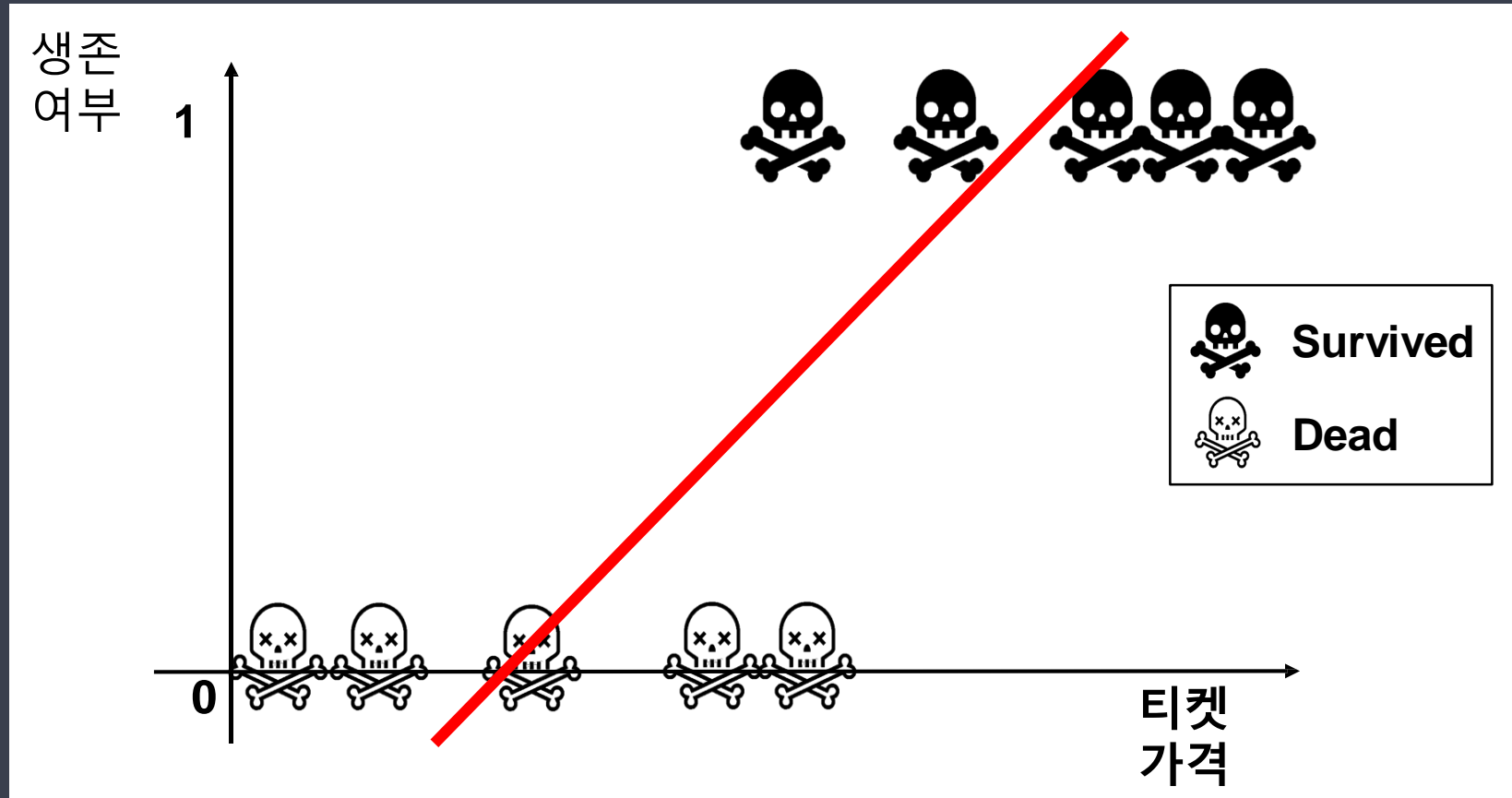
분류용 선형 모델



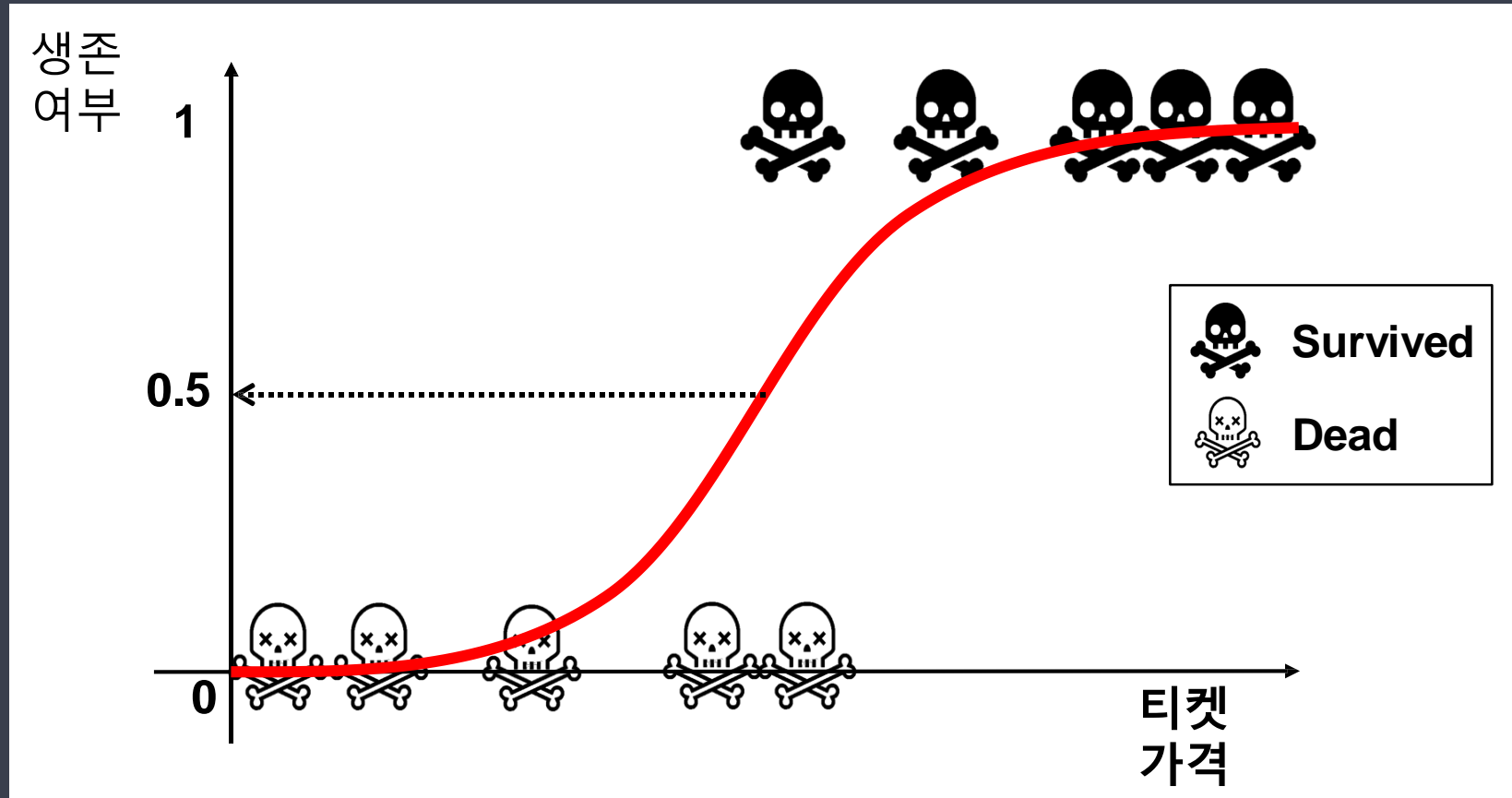
Logistic Regression

- 회귀공식을 사용해서 Regression이라는 이름이 붙음
- 결정 경계가 선형이기 때문에 선형 모델
- 시그모이드 함수의 최적선을 찾고 반환값을 확률로 간주

Linear Model – Classification

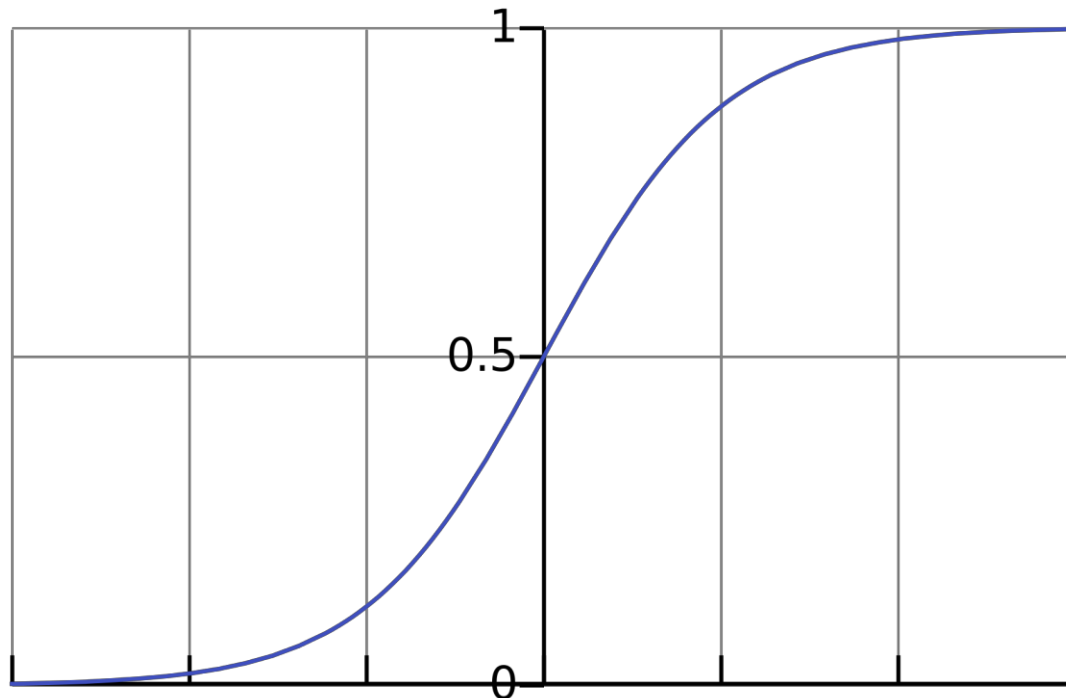


Linear Model – Classification



Logistic Regression

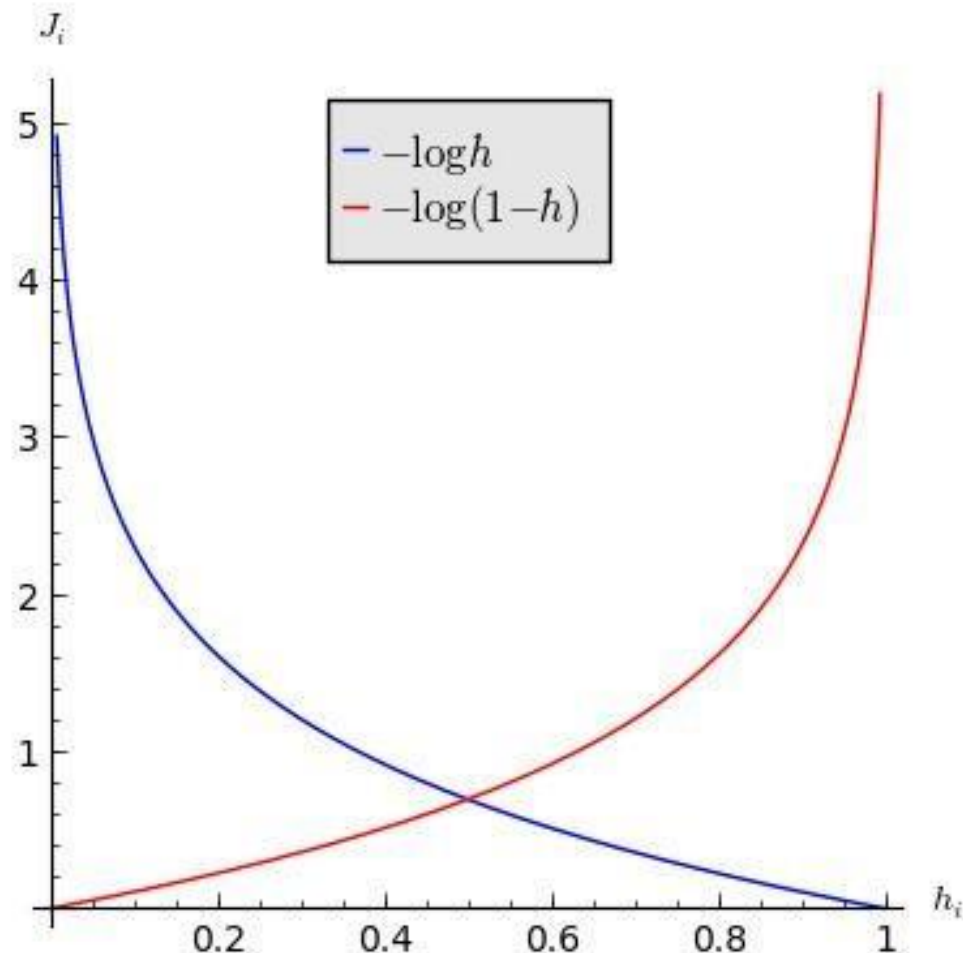
- 선형 함수의 결과 값을 Sigmoid Function(Logistic Function)을 이용해 0과 1로 변환한다.



$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

Linear Model – Classification

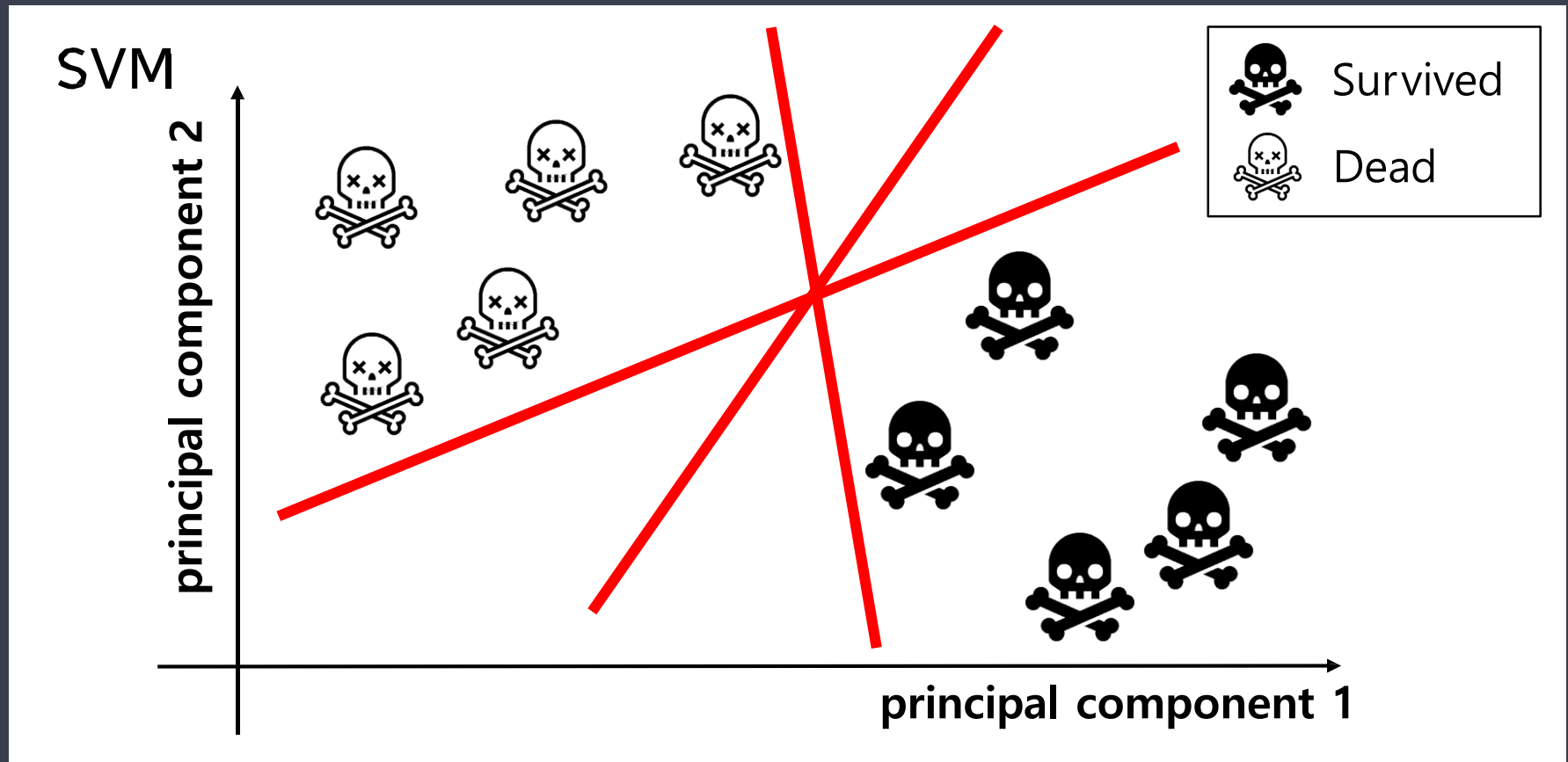
$$\begin{aligned}\text{Cost}(h_{\theta}(x), y) &= -\log(h_{\theta}(x)) && \text{if } y = 1 \\ \text{Cost}(h_{\theta}(x), y) &= -\log(1 - h_{\theta}(x)) && \text{if } y = 0\end{aligned}$$



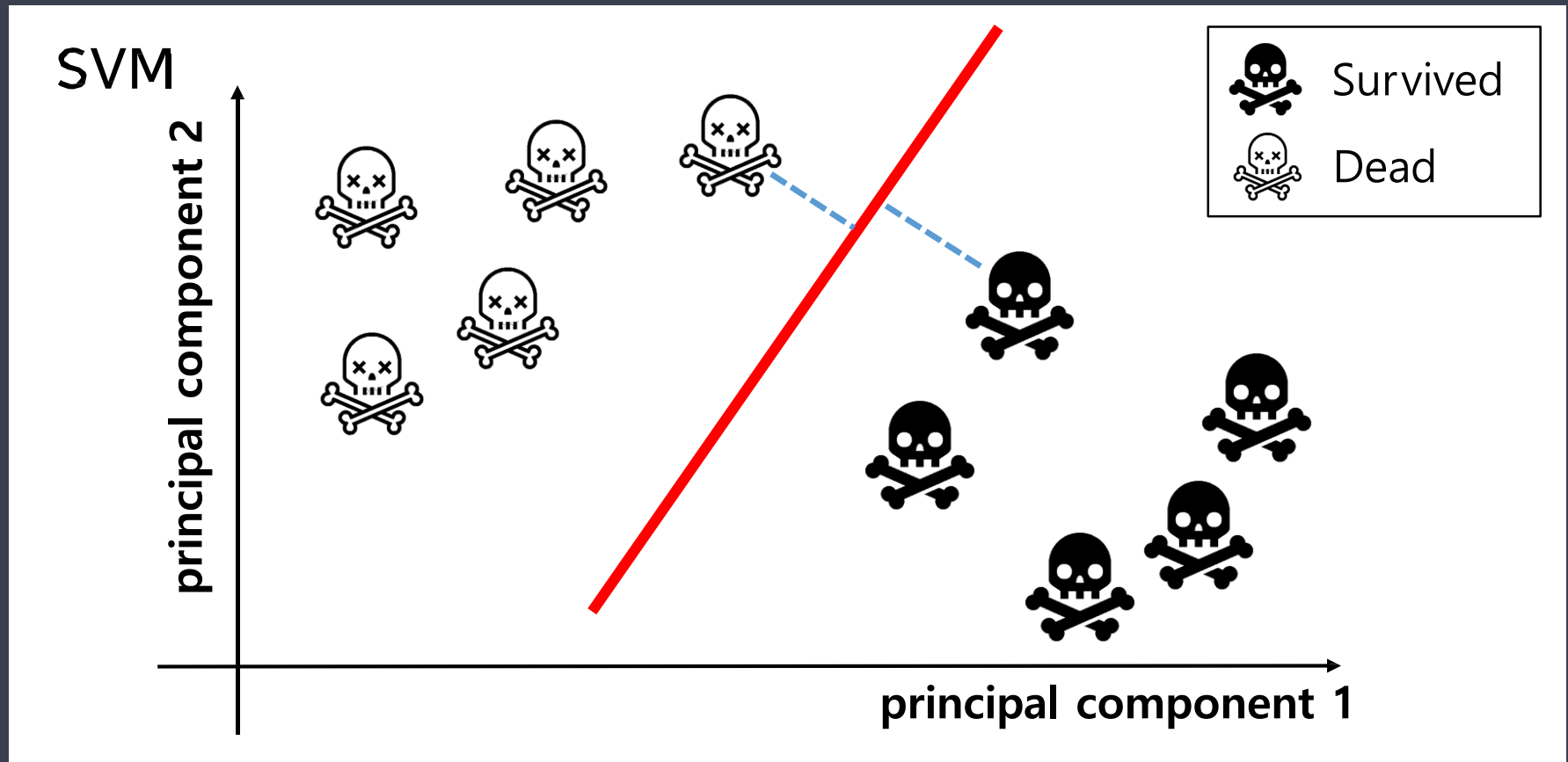
주요 매개변수(Hyperparameter)

- 선형 분류 모델 : C
(값이 클수록 규제가 약해진다.)
- 기본적으로 $L2$ 규제를 사용, 하지만 중요한 특성이 몇 개 없다면 $L1$ 규제를 사용해도 무방
(주요 특성을 알고 싶을 때 $L1$ 규제를 사용하기도 한다.)

Support Vector Machine (SVM)



Support Vector Machine (SVM)



장단점 및 주요 매개변수(Hyperparameter)

- 회귀 선형 모델 : α
(값이 클수록 규제가 강해진다.)
- 선형 분류 모델 : C
(값이 클수록 규제가 약해진다.)
- 기본적으로 L2규제를 사용, 하지만 중요한 특성이 몇 개 없다면
L1규제를 사용해도 무방
(주요 특성을 알고 싶을 때 L1 규제를 사용하기도 한다.)

장단점 및 주요 매개변수(Hyperparameter)

- 선형 모델은 학습 속도가 빠르고 예측도 빠르다.
- 매우 큰 데이터 세트와 희소 (sparse)한 데이터 세트에서도 잘 동작한다.
- 특성이 많을 수록 더욱 잘 동작한다.
- 저차원(특성이 적은)데이터에서는 다른 모델이 더 좋은 경우가 많다.

손글씨 데이터 분류 실습



분류 평가 지표

Confusion_matrix

negative class	TN	FP
positive class	FN	TP
	predicted negative	predicted positive

**정확도
(Accuracy)**
전체 중에 정확히
맞춘 비율

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Confusion_matrix

100명 중 암 환자는 5명

실제 암 X	95	0	
	TN	FP	
실제 암 O	0	5	
	FN	TP	
예측 암 X			
예측 암 O			

100

100

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Confusion_matrix

100명 중 암 환자는 5명

실제 암 X	95	0
	TN	FP
	FN	TP
실제 암 O	5	0
예측 암 X		
예측 암 O		

95

100

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Confusion_matrix

negative class	TN	FP
positive class	FN	TP
	predicted negative	predicted positive

$$\text{Recall} = \frac{TP}{TP + FN}$$

**재현율
(Recall)**
실제 양성 중에
예측 양성 비율

Confusion_matrix

100명 중 암 환자는 5명

실제 암 X	95	0
	TN	FP
실제 암 O	5	0
	FN	TP
예측 암 X		
예측 암 O		

$$\frac{0}{5}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

Confusion_matrix

100명 중 암 환자는 5명

실제 암 X	0	95
	TN	FP
실제 암 O	0	5
	FN	TP
예측 암 X		
예측 암 O		

$$\frac{5}{5}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

Confusion_matrix

negative class	TN	FP
positive class	FN	TP
	predicted negative	predicted positive

$$\text{Precision} = \frac{TP}{TP + FP}$$

정밀도
(Precision)

예측 양성 중에
실제 양성 비율

Confusion_matrix

100명 중 암 환자는 5명

실제 암 X	0	95
	TN	FP
실제 암 O	0	5
	FN	TP
예측 암 X		
예측 암 O		

$$\frac{5}{100}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

Confusion_matrix

negative class	TN	FP
positive class	FN	TP
	predicted negative	predicted positive

F1 - score

정밀도와 재현율의
조화평균

$$F = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

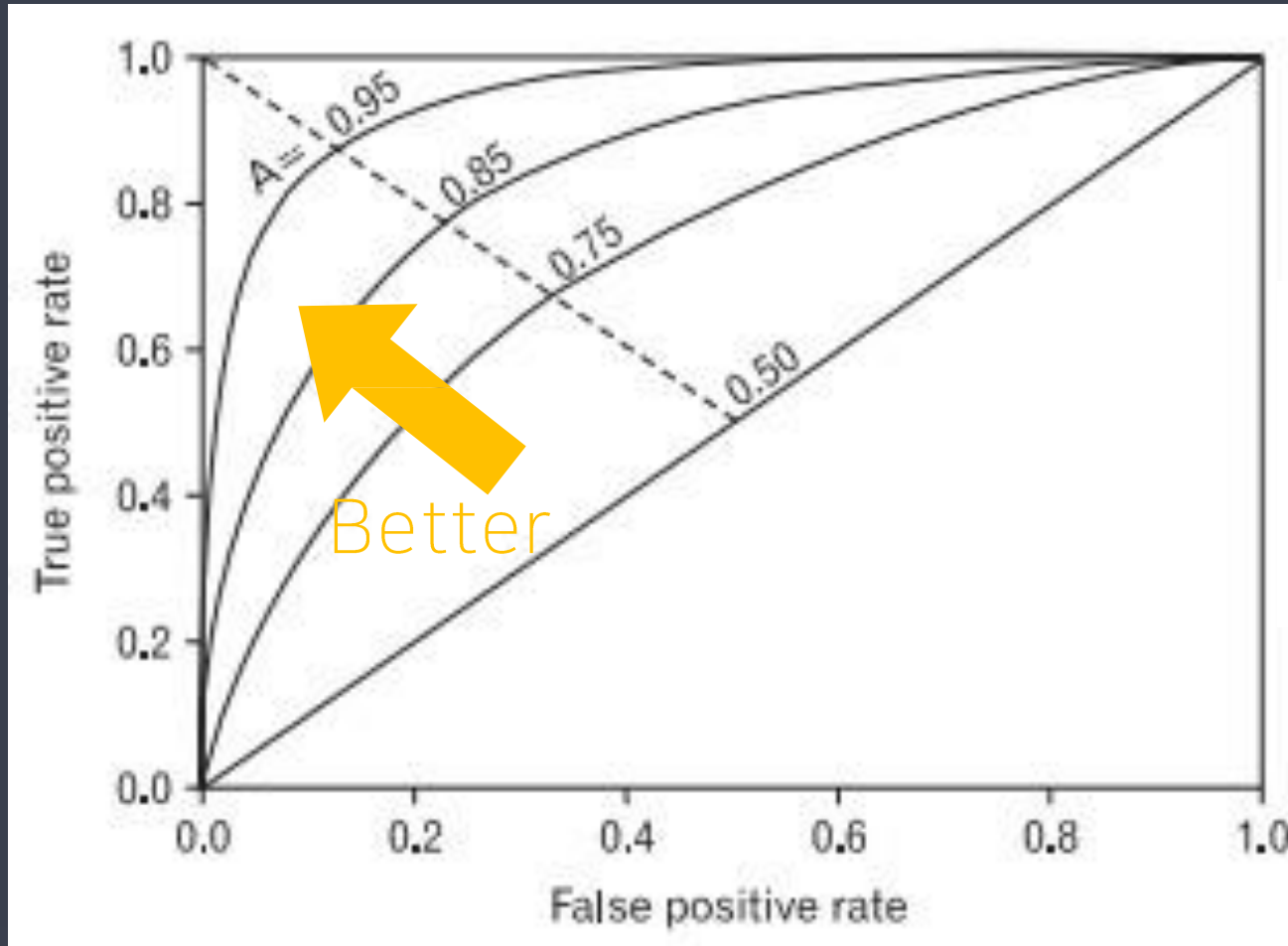
낮은 재현율보다 높은 정밀도를 선호하는 경우

어린이에게 안전한 동영상(양성)을 걸러내는 분류기를 훈련시킬 경우 좋은 동영상이 많이 제외되더라도(낮은 재현율) 안전한 것들만 노출시키는(높은 정밀도) 분류기가 더 좋다.

낮은 정밀도보다 높은 재현율을 선호하는 경우

감시 카메라로 줌도둑(양성)을 잡아내는 분류기를 훈련시킬 경우 경비원이 잘못된 호출을 종종 받지만(낮은 정밀도) 거의 모든 줌도둑을 잡는(높은 재현율) 분류기가 더 좋다.

ROC curve



$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \text{recall}$$

**분류 결과의 불확실성 및 분류결과에 대한
평가지표를 확인해보자.**