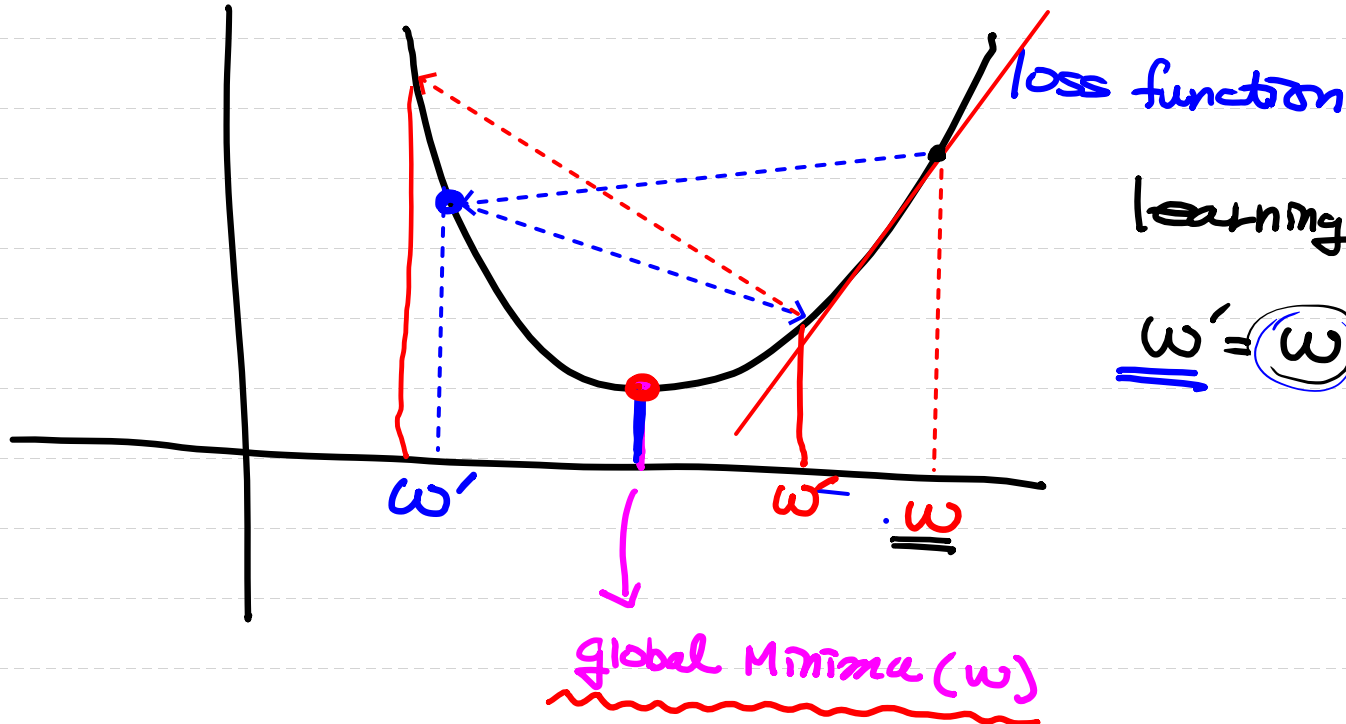


04/05

① learning rate (학습률) $\rightarrow 1e^{-4}$ (크기) loss의 값을 보면서 learning rate를 조절



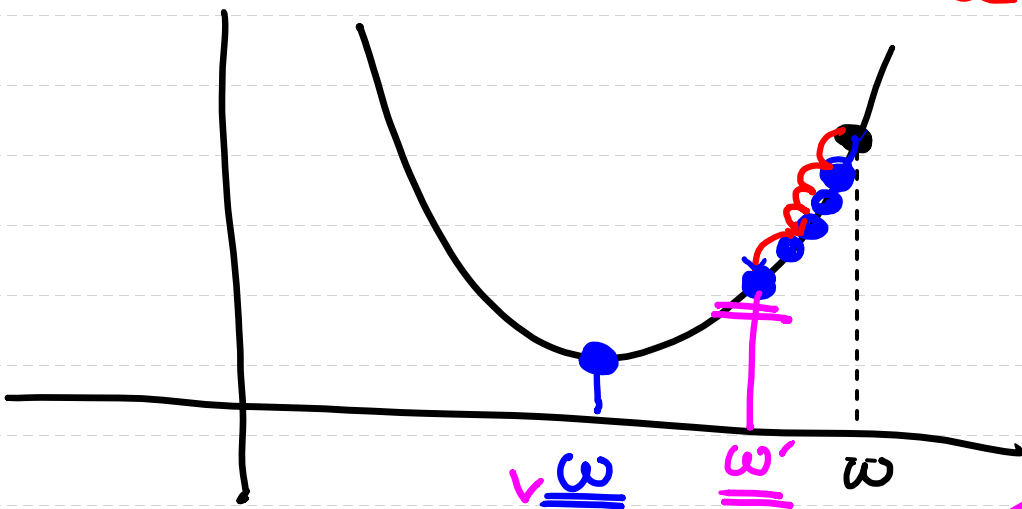
learning rate를 크게 잡으면

$$\underline{w'} = \underline{w} - \alpha \frac{\partial E(w, b)}{\partial w}$$

크게 잡으면

\Rightarrow Overshooting
떨어질

== α 값이 너무 작으면??



local Minima 현상

② Normalization (정규화)

MinMax Scaling (0 ~ 1)

Standardization (표준화)

이상의 값 상대적으로 두배해요
Scale이 동일하지 않아요.

③ Overfitting (과대적합)

✓ Underfitting (과소적합)

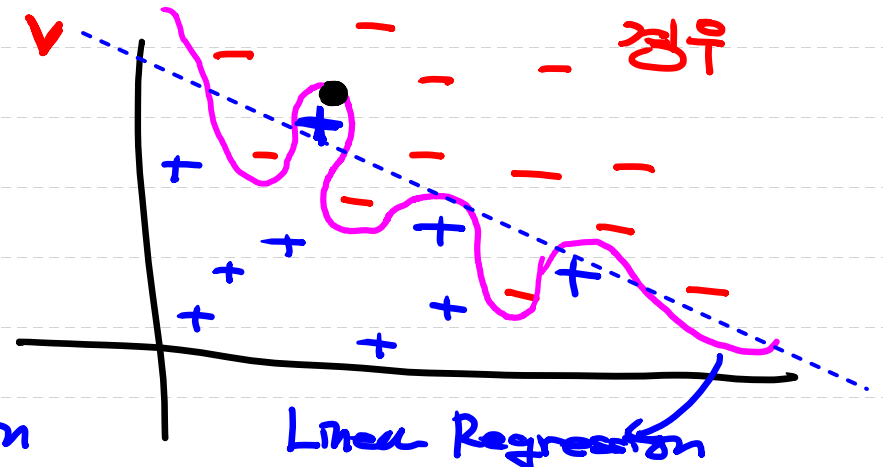
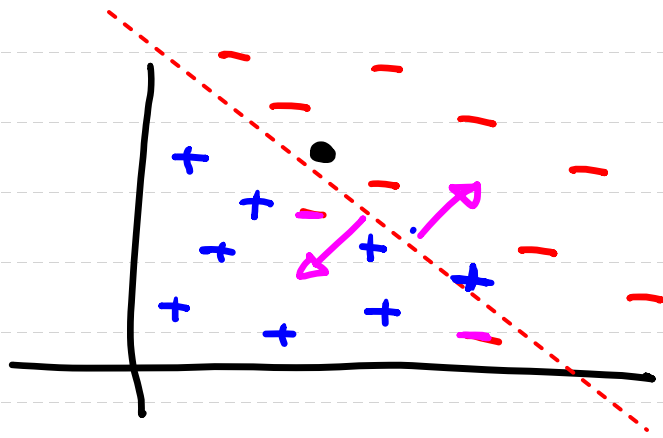
과적합

✓ ✗ ✗ ✗

이렇게 되지
않도록 치기

Overfitting은 Model이 Training Data Set에 너무 잘

맞춰서 오려 Test Data Set에 대해서는 정확도가 떨어지는



* Overfitting을 방지하려면 어떻게 해야 하나요?

① 많은 Training Data를 사용하는 것

내가 가진 Data가 작다. → 데이터 증강 (Augmentation)
하지만 Data가 늘어나면 → resource ↑ 학습시간이 상당히 길어지고
를 많이 사용하고

② Feature의 개수를 줄일 필요가 있어요.
중복된 feature를 삭제
[종속변수인 상관관계가 적은거.
↓
"데이터보정"

③ Weight (w)의 값이 너무 크지 않도록 제한 → 어떻게요??

↳ 값이 크면 그래프가 구불구불하게 그려져요!

↳ loss 함수를 변경

④ Deep Learning 기법

Regularization
(제제)

* loss 함수를 의도적으로
크게 해서 w 값을 작게 제한
시켜서 w 값을 제한

* Dropout이라는
방식을 이용

↳ [L1
L2

* Lasso Regression
수식이 달라요!!
* Ridge Regression

① Logistic Regression (Binary classification)

$$y = \frac{1}{1 + e^{-(wz + b)}}$$

$$\text{Cross Entropy (loss)} = - \sum_{i=1}^n \{ x_i \log y_i + (1 - x_i) \log (1 - y_i) \} +$$

L1
Regularization

$$\frac{1}{2} \propto \sum_{i=1}^n |w_i|$$

Regularization ✓

Strength

(내가 설정하는
parameter)

L2
Regularization

$$\frac{1}{2} \propto \sum_{i=1}^n |w_i|^2$$

code로 표현!

$tf.nn.l2_loss(w)$

④ Over Sampling

데이터의 불균형을 어떻게 해소할 것인가?

↓ (데이터의 비대칭문제)
Class에 속한 데이터의 개수가 크기가 있을때
→ "imbalanced Data problem"

[0 Class가 2개 → Class는 불균형을 지칭
1]

Under Sampling

★ Over Sampling

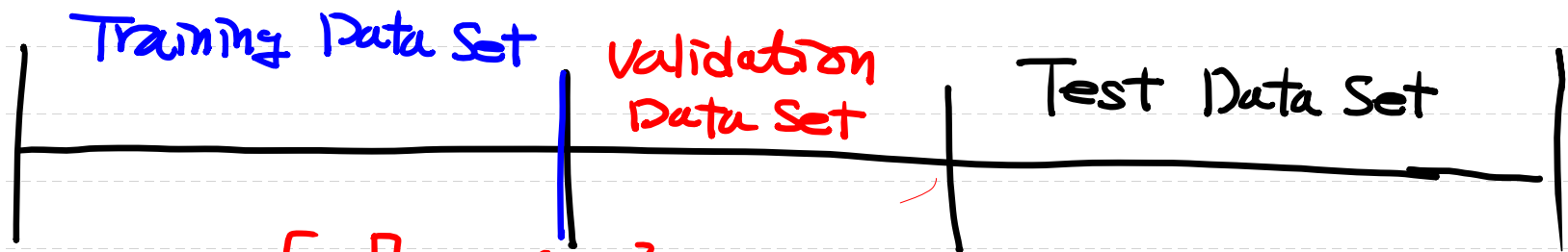
★ 데이터도 부족하고 편향도 있어
★★ 있을때 소수의 데이터를 복제

→ Data가 거리가 충분히 많으면
복제본 (class) [0 1000000 → 50000
1 100000 → 50000]

(데이터가 이렇게 많은 경우가 흔치않고
정확적인 데이터를 얻기버릴수 있어요

→ "SMOTE 기법" ✓
관련된 module도 있어요!!

⑤ k-Fold Cross Validation (vs. Hold-out Validation)



$\begin{bmatrix} 7 & \vdots & 3 \\ 8 & \vdots & 2 \end{bmatrix}$

↳ 보통 3배

$[00000000111] \rightarrow$ 전체 Data

X $[00000000][111]$
 ↓ Training Data Set ↓ Validation Data Set



O $[0000011][001]$

이런 처리는 Sklearn 이 담당!!

★ k-Fold Cross validation ⇒ 시간이 오래 걸린다는
단점
(5-Fold)

★ Sklearn ★

