

Supervised Learning의 종류

(know the answer)

- Regression (회귀)

연속적인 값을 예측. (차이가 있다)

- Classification (분류)

이산적인 값 (맞거나 틀리거나)

- **UNSupervised Learning**. : Don't know Answer.
PCA.

- **Reinforcement Learning**.

알려진 프로그램에 따라가 좋은 결과를 얻게 된다.

Supervised Learning의 종류

- 예측하려는 변수의 종류에 따라
- Regression (회귀): 연속 변수를 예측
- Classification (분류): 이산 변수를 예측

회귀 vs. 분류

- 예측하는 변수가 다름(연속 vs. 이산)
- 오류의 형태가 다름
- 회귀: 예측과 실제의 거리
 - 예: 3.84로 예측했는데 4.28 오차를 수치로 산수한다.
- 분류: 예측과 실제의 차이
 - 예: 고릴라로 예측했는데 판다

평가지표의 종류

- 크로스 엔트로피 ~~오류~~ ^{loss}를 fitting 시킬 때 중요.

{ 아예 90% 정확률로 A라고 했는게 B인 경우 큰 페널티를 줌.
49% vs 51% 일 경우는 페널티 ↓.

정확도에 따라
페널티 강도가
달라진다.

$$H(p, q) = - \sum p(x) \log q(x)$$

- 분류에서 확률로 예측할 때 로그-우도

분류에서 나올 수 있는 경우

| | | 예측 | |
|----|-------------|-----------------------|-----------------------|
| | | 양성 Positive | 음성 Negative |
| 실제 | 양성 Positive | 진양성 True Positive | 위음성 False Negative |
| | 음성 Negative | 위양성 False Positive | 진음성 True Negative |

정확도 Accuracy

진양성, 진음성이라 한 것만
주게는 정확도

| | | 예측 | |
|----|-------------|-----------------------|-----------------------|
| | | 양성 Positive | 음성 Negative |
| 실제 | 양성 Positive | 진양성 True Positive | 위음성 False Negative |
| | 음성 Negative | 위양성 False Positive | 진음성 True Negative |

정밀도 Precision

해답이 positive
→ 이치상 맞는
해답.

ex) 회사에서도 음성인식률 파악 ⇒ 관습x ⇒ 신중하게 다루기
ex) 병원에서도 환자 음성인식률 파악 ⇒ 관습x. 음성인식률 파악
아니고x.

| | | 예측 | |
|----|-------------|-----------------------|-----------------------|
| | | 양성 Positive | 음성 Negative |
| 실제 | 양성 Positive | 진양성 True Positive | 위음성 False Negative |
| | 음성 Negative | 위양성 False Positive | 진음성 True Negative |

→ 최대한 음성을 바꾸어 버리면 된다. } 정밀도 ↑
ex) 회사 신임직원 ⇒ 최대한 파악

실제 양성 중에 양성 맞춘 횟수

재현율 Recall

정답률 : 양성으로
보아야 성공

→ 정답률은 높여야 하겠다.

· 양성으로 많이 넣어 버린다.

| | | 예측 | |
|----|-------------|-----------------------|-----------------------|
| | | 양성 Positive | 음성 Negative |
| 실제 | 양성 Positive | 진양성 True Positive | 위음성 False Negative |
| | 음성 Negative | 위양성 False Positive | 진음성 True Negative |

· 정답률 : 예측 양성 중에 맞춘 횟수

· 재현율 : 실제 양성 중에 양성 맞춘 횟수

Kappa

$$Kappa = \frac{O - E}{1 - E}$$

O : observed accuracy, E : Expected Accuracy

F1 score

ex) 10km ↑
90km ↓
[평균 80km X]

- 정밀도와 재현율의 조화평균

높은 속도, 높은 정확도를

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

LM의 정규화

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (x_{1i} - x_{2i})^2}{n}}$$

- RMSE만 최소화하는 대신

평균 제곱 편차 (Root Mean Square Deviation)

- RMSE + (w의 크기)를 동시에 최적화

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

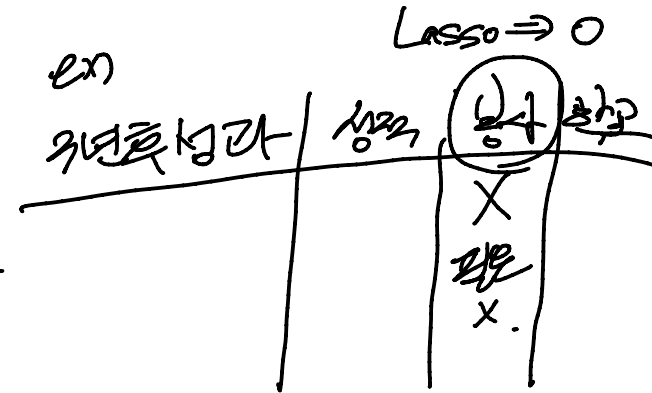
Cost function.

(정규화)
모든 w의 크기를
제한하는 것
모든 w를
0으로 만드는 것

y = (ax + b) (채워)

모든 w가 0이 되면 (over fitting) ⇒ cross validation으로
해결되지 않는 것은 맞지 않음. 해를 ↑
값수있지
제방X.

Lasso



- $q = 1$ summary : 자살을 일으키는 요인 w 를 0으로

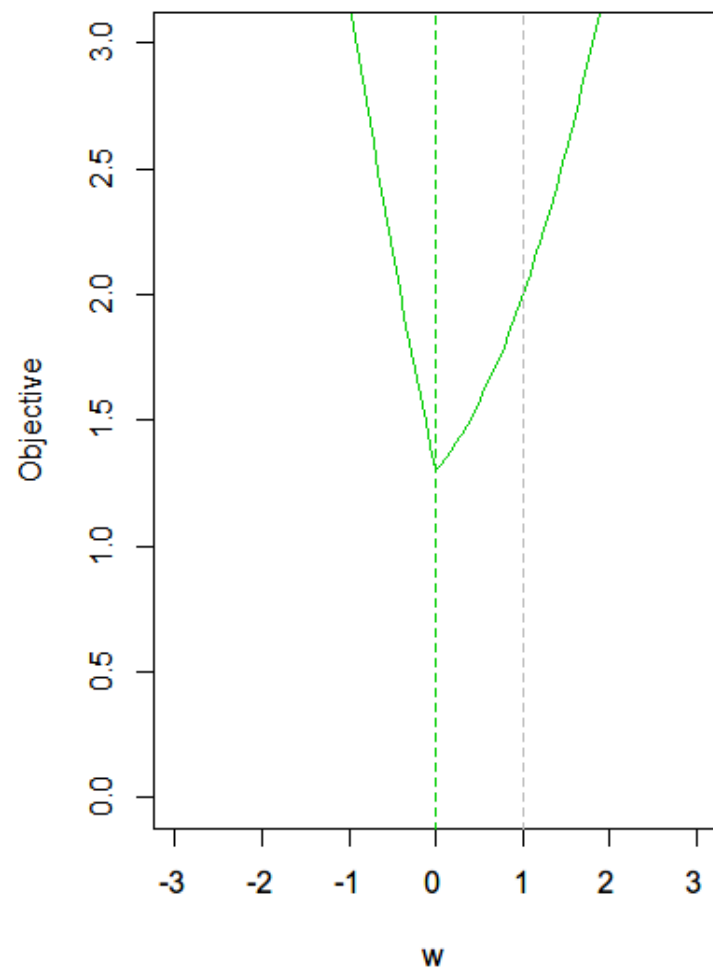
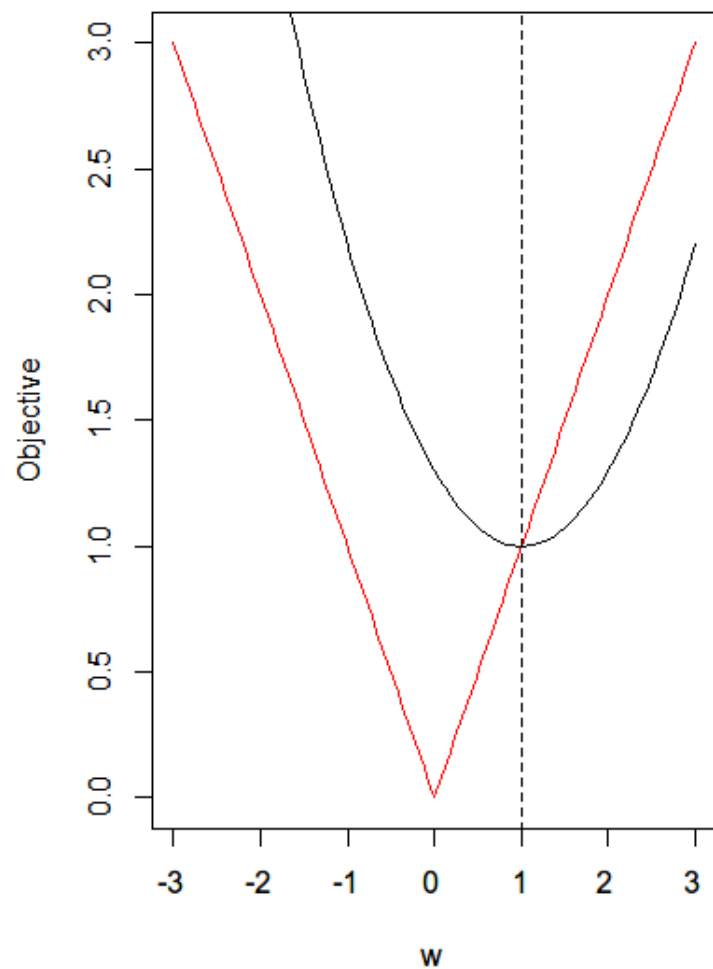
- w 의 절대값의 합도 함께 최소화

- w 를 0으로 만드는 경향이 있음 $y = ax + b$
 여기서 $y = 0ax + b$
 x 의 값에

- 변수 선택의 기능

가장 낮은 값을 가지고 LASSO를
 넣으면. 최솟값은 변수들을 제거해
 있다.

Lasso

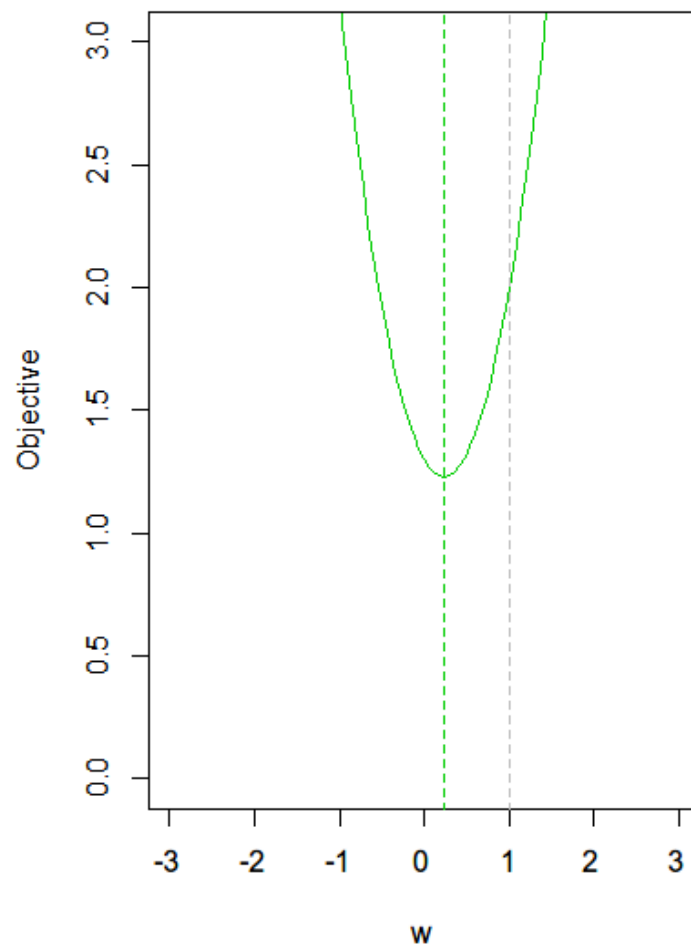
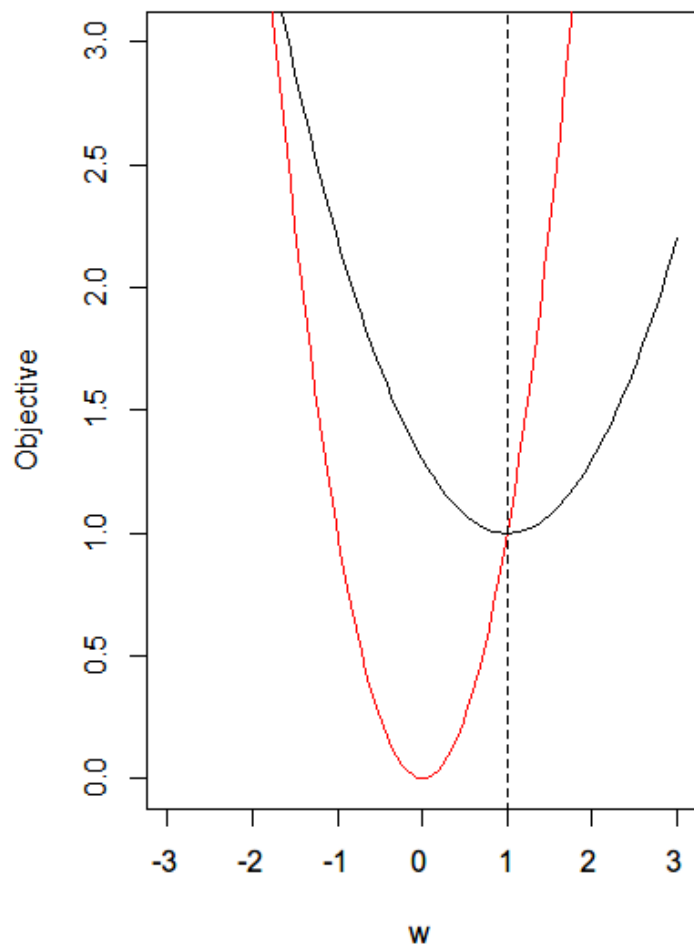


Ridge

(일반적으로 Lasso 보다 overfitting
예방능력이 좋다.)

- $q = 2$
- w 의 제곱의 합도 함께 최소화
- 대체로 Lasso에 비해 예측력이 좋음
- 변수 선택 X

Ridge



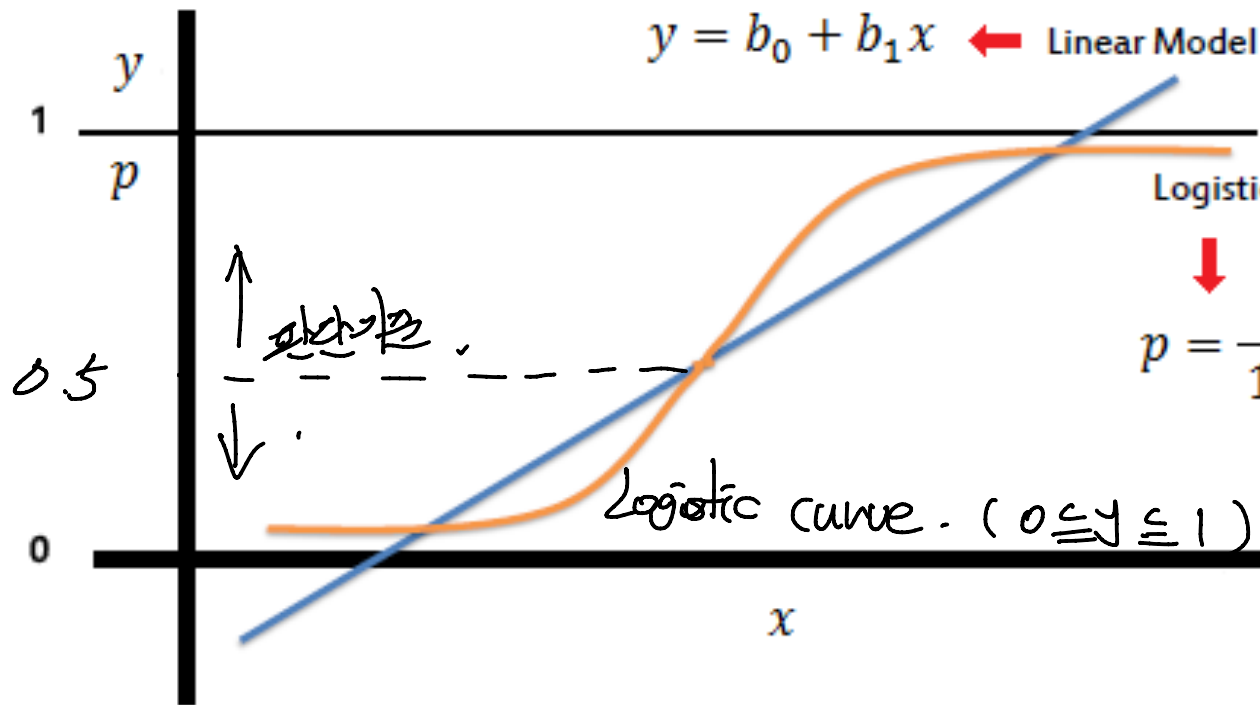
Elastic Net

- RMSE + Lasso + Ridge
 - Lambda: 정규화항의 가중치
 - Alpha: 정규화항에서 Lasso의 비중
- CV로 결정

Logistic Regression

- Linear Model for Classification

단순의 선형 회귀 모형이
잘못되어있다.



결과 $\rightarrow (0, 1)$
 Linear \rightarrow Can be
 over 1
 or under 0
 \rightarrow Logistic func
 $\Rightarrow 0 \sim 1$

판단기준을 확률값으로
해석가능.