

Classification

$y \in C$ 인 경우, $f(x) \in C$ 인 \hat{y} 를 찾는 것이 목적이다.

$f(x)$ 가 0, 1로 classification이 가능할지도 모르지만 애매한 경우도 존재.

그런 경우 $f(x)$ 가 특정한 C의 확률을 최대화 ↑ 비보자.

방법 1) Just Linear model에 넣어서 class는 $\in \{0, 1\}$ 로 설정하면

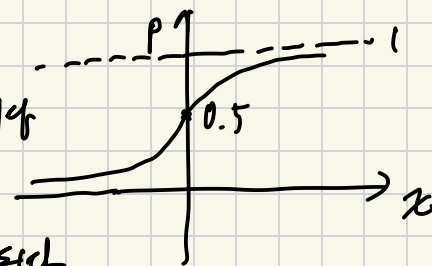
되나? NO, Linear model의 result는 $-\infty \sim \infty$ 사이 값이라 \hat{y} 설정 x
& class가 3개 이상일 경우 $\in \{0, 1, 2\}$ 시 1과 2의 2배 차이가 생긴다.

< Logistic Regression > Linear Regression + Sigmoid.
(이진분류는 연속값을 따온다는 가정)

$x \in \mathbb{R}^d$ 이고 $y \in \{-1, +1\}$ 인 Binary classification 가정

이때 $f(x)$ 가 1 일 확률이 높다면 1이 가깝게
" +1 이 정답 애매한 경우 0.5 " } 아주 Mapping 하고 싶다.
" -1 이 될 확률이 높다면 0 이 " }

그래서 4분 function $\text{Sigmoid} = \frac{1}{1 + e^{-x}}$ 이다
Linear Regression 결과 ∞ 라면 1
" 0 라면 0.5 } mapping 된다.
" $-\infty$ 라면 -1



확률의 성질 : 모든 outcome들이 $0 \sim 1$ 사이이 존재)
조건 = 을 만족해야 한다.

그러면 $p(y=1|x) + p(y=-1|x) = 1$ 을 만족해야 한다.

$$\frac{1}{1 + e^{\beta^T x}} + \frac{1}{1 + e^{\beta^T x}} = 1 \quad \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n \text{ 이다.}$$

2.2.2 $p(y=1|x) = \frac{1}{1+e^{-\beta x}}$ 을 식을 정리해보면

$$\log\left(\frac{p(x;\beta)}{1-p(x;\beta)}\right) = \beta_0 + \beta_1 x, \text{ 이다.}$$

$\rightarrow y=1$ 일 확률
 $\rightarrow y=0$ 일 확률.

$p(x;\beta)$ 가 ↑일수록 값이 ↑ 이고 $p(x;\beta)$ 가 ↓일수록 값이 ↓ 된다.

이렇게 $\log\left(\frac{p(x;\beta)}{1-p(x;\beta)}\right)$ 를 "log odds" 라고 부르고

Logistic Regression은 log odds가 선형적이라는 가정을 하고 있다.

odds는 항상 0을 넘지 않는다. y 가 양성일 확률 / 음성일 확률의 비가 되기 때문

$-\infty < \log \text{ odds} < \infty$ 이 범위를 갖는다.

그래서 $\log \text{ odds} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$ 처럼 선형적 해석가능.

Linear Regression의 경우 데이터들의 관계 자체가 선형적이라는 가정을 한 모델이라면, Logistic Regression의 경우 log odds가 선형적이라는 가정을 한 모델이다.

< MLE of Logistic Regression >

$$L(\beta) = \prod_{i, y_i=1} p(y_i|x_i) \cdot \prod_{i, y_i=0} (1-p(y_i|x_i))$$

$$\log L(\beta) = \sum \log\left(\frac{1}{1+e^{-\beta x_i}}\right) + \sum \log\left(\frac{1}{1+e^{\beta x_i}}\right)$$

$$= \sum_{i=1}^n \log\left(\frac{1}{1+e^{-y_i \beta^T x_i}}\right) = - \sum_{i=1}^n \log(1+e^{-y_i \beta^T x_i})$$

$$\Rightarrow \sum_{i=1}^n \frac{x_i y_i}{1+e^{-y_i \beta^T x_i}} = 0 \quad \neq \text{not closed solution.}$$

β 를 구하는 공식이 존재하지 않는다.

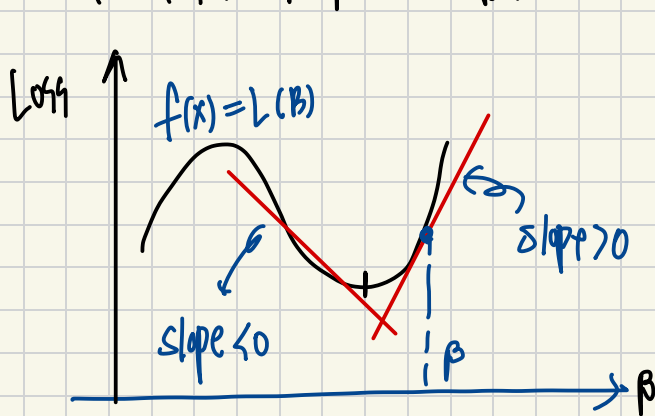
Closed-form 해와 convex 문제는 완전히 별개의 개념입니다. 우선 closed-form 해의 존재 여부는 MLE나 최소 제곱으로 세운 손실 함수 $L(\beta)$ 에 대해 $\nabla_{\beta} L = 0$ 을 풀었을 때, 이 식이 $-A\beta + b = 0$ 처럼 β 에 대해 2차식 또는 선형 방정식 형태로 정리되어 $\beta = A^{-1}b$ 와 같은 해로 바로 구해질 수 있는지를 의미합니다. 반면 목적 함수가 convex라는 것은, 그 함수가 볼록 형태여서 모든 국소 최소값이 곧 전역 최소값이 된다는 뜻이며, 이 경우 경사 하강법 같은 수치 최적화 기법을 사용해 전역 최적해를 안정적으로 찾을 수 있다는 보장을 제공합니다. 따라서 "손실이 2차식이면 closed-form 해가 가능"하고 "손실이 convex이면 경사법으로 전역 최적해 보장"이라는 두 판단 기준은 서로 다른 개념임을 명확히 구분해야 합니다.

⇒ MLE의 결과로 나온 방정식이 선형식이 아니라면 closed-form이 존재하지 않는다.

< gradient Descent >

optimization : 특정 set에서 조건 만족하는 최선의 element를 구해보자.

그래서 "최적해" 하나를 찾는 방법을 구하자 라는 β 에 대해서 미분을 계속해보자.



$$\beta_{\text{new}} = \beta_{\text{old}} - \underbrace{\gamma \cdot \frac{2}{2\beta_{\text{old}}}}_{\text{미분값이}} L(\beta)$$

↑
기울기라 반대방향.
↓
차라 안바뀌어 이동량지

different 1) Not convex problem인 경우 ⇒ Local optim에 빠질수 있다.

different 2) 미분이 가능하면 cost function 사용 가능

different 3) 수렴이 느려진다 (거기 수렴했으면 경우) ⇒ SGD로 개선.

< Stochastic Gradient Descent >

data를 특정 개수의 Batch 만큼만 활용해서 최적화 진행 (mini-batch)

diminishing returns 효과 (Batch크기 ↑ 학습률 얻는 정도 ↓)

지금까지는 1개의 변수와 이진분류문제를 풀었다.

But 변수가 N개가 된다면 아래와 같이 표현이 가능해진다.

$$\log\left(\frac{p(x;\beta)}{1-p(x;\beta)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n \text{ 이 된다 (}\beta \text{이 선형값이기때문에)}$$

$$p(x;\beta) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n}} \text{ 이 된다.}$$

(class가 2개 이상인 경우)

이진분류에서 $p(y=1|x) = \frac{1}{1+e^{-\beta^T x}}$ & $p(y=-1|x) = \frac{1}{1+e^{\beta^T x}}$ 이다.

이제 부호를 중립기에 해주면

$$p(y=1|x) = \frac{e^{\beta^T x}}{e^{-\beta^T x} + e^{\beta^T x}}, \quad p(y=-1|x) = \frac{e^{-\beta^T x}}{e^{-\beta^T x} + e^{\beta^T x}} \text{ 가 됨}$$

이를 확장해서 N개의 class인 경우

$$p(y=k|x) = \frac{e^{x^k}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}} \text{ 로 설정이 가능}$$

↳ 모든 class에 대한 합을 sum.

이걸 softmax function 이라.