

< 간단한 머신러닝 >

Date.

No.

Chapter 3. 선형모델

* 3.1 기본 개념

$$x = (x_1, x_2, \dots, x_d)$$

선형모델 ; 속성들의 linear combination으로 예측하는 함수.

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$$

$$= w^T x + b \quad \text{where } w = (w_1, w_2, \dots, w_d)$$

* non-linear model \rightarrow 선형모델을 기반으로 좀더 좋게 하거나
고차원으로 투영하여 만들어짐.

* 해석능력 comprehensibility (이해가능성 understandability)

\rightarrow w 는 예측에 있어 각 속성의 중요성을 직관적으로 드러냄

* 3.2 Linear Regression.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} = \{(x_i, y_i)\}_{i=1}^m$$

$$x_i = (x_{i1}, x_{i2}, \dots, x_{id}), y_i \in \mathbb{R}$$

first, $d=1$.

If discrete feature,

case 1. ordered \Rightarrow Label Encoding

case 2. else \Rightarrow One Hot Encoding

\rightarrow 이를 선택해 해결하면 정확히 높은
순서 관계를 가지고 여러 계층 등에
영향을 미침!

(0 평균 제곱차 (square loss))

\rightarrow Euclidean Distance 사용.

$$J(w, b) = \sum_{i=1}^m (y_i - (x_i w + b))^2 ; w \text{와 } b \text{에 대한 convex function.}$$

$[a, b]$ 에서 정의된 함수 f 에 대해, 만약 구간 내 두 점 x_1, x_2 가

$$f\left(\frac{x_1 + x_2}{2}\right) \leq \frac{f(x_1) + f(x_2)}{2}$$

를 만족하면 f 는 구간 $[a, b]$ 상의 convex function 이라 함.

Second, $d \geq 2$, **Multivariate Linear Regression**

$$\text{ex: } f(x_0) = w^T x_0 + b \rightarrow f(x_0) \approx y_0$$

$$\text{Let } \hat{w} = (w, b), \quad D \text{ is } X \text{ set. } X \in \mathbb{R}^{m \times (d+1)}$$

$$\hat{w}^* = \underset{\hat{w}}{\operatorname{argmin}} \underbrace{(y - X\hat{w})^T (y - X\hat{w})}_{E\hat{w}}$$

$$\frac{\partial E\hat{w}}{\partial \hat{w}} = 2X^T(X\hat{w} - y)$$

* Caution! $X^T X$: full-rank or positive definite,

$$\hat{w}^* = (X^T X)^{-1} X^T y$$

→ 현실에서 full-rank 일 경우는 드물!!

ex) 샘플 정보량, 유전자 데이터, 전 안개나 밤의 특성
(몇십 ~ 몇백의 샘플)

* Reference, Connection between rank and positive-definite / math-stackexchange

✓ 한편 $X^T X \in M_n(\mathbb{R})$, full-rank가 아니라면,

$\lambda_0 = 0$ 일 수가 존재, 즉, PD X.

✓ 위 조건에서, nontrivial kernel 이 가짐.

$\Rightarrow \exists v \neq 0$ s.t. $Mv = 0$.

$v^T M v = 0$, Hence, PD X.

✓

$$\ln y = w^T x + b$$

log-linear function

$$y = g^{-1}(w^T x + b)$$

where g is \rightarrow link function
monotone differentiable function
generalize linear model

* 3.3 Logistic Regression.

문제는 \rightarrow "GLM", ^상 선형 함수 z 를 찾아 z 를 y 로
 실제 출력 레이블 y 와 비교하는 것!!

Heaviside function (단계-제한 함수 ^{mit-} step function)

$$y = \begin{cases} 0, & z < 0 \\ 0.5, & z = 0 \\ 1, & z > 0 \end{cases} \rightarrow \text{불연속!!}$$

그러나, \hookrightarrow Heaviside function과 유사한 sigmoid function (연속함수)
 \hookrightarrow monotone differentiable

Sigmoid function $\rightarrow x \rightarrow -\infty, f(x) = 0, x \rightarrow \infty, f(x) = 1$
 \hookrightarrow logistic function

$$y = \frac{1}{1 + e^{-z}}$$

GLM처럼! $y = \frac{1}{1 + e^{-(w^T x + b)}}$

$$\ln \frac{y}{1-y} = w^T x + b$$

odds

$\log \text{ odds (logit)}$

\rightarrow 실제 데이터 y 와 z 에 근사!!

logistic regression! ^{최적화 방법!}

① 사전 데이터 x 에 대한 가중치 필요 X.

② 근사 확률에 대한 예측 가능.

③ solution의 목적 함수가 convex function!

\rightarrow 수렴 속도 굉장히 Good!!

Date.

No.

let $y = p(y=1|x)$. then,

$$\ln \frac{p(y=1|x)}{p(y=0|x)} = w^T x + b$$

odds

$$p(y=1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}$$

$$p(y=0|x) = \frac{1}{1 + e^{w^T x + b}}$$

w 와 b 를 추정하기 위해, 최대우도법 maximum likelihood method 사용!

$$l(w, b) = \sum_{i=1}^m \ln p(y_i | x_i; w, b)$$

(22+3) 최대치 찾기!

why?

$$p(y_i | x_i; w, b) = y_i p_1(\hat{x}_i; \beta) + (1 - y_i) p_0(\hat{x}_i; \beta)$$

$$l(\beta) = \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i}))$$

minimize!

where $\beta = (w; b)$, $\hat{x} = (x; 1)$

$$(\because y_i = \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} + (1 - y_i) \cdot \frac{1}{1 + e^{\beta^T \hat{x}_i}})$$

$$= \frac{y_i}{1 + e^{\beta^T \hat{x}_i}} (e^{\beta^T \hat{x}_i} - 1) + \frac{1}{1 + e^{\beta^T \hat{x}_i}} = \frac{1}{1 + e^{\beta^T \hat{x}_i}} (y_i e^{\beta^T \hat{x}_i} - y_i + 1)$$

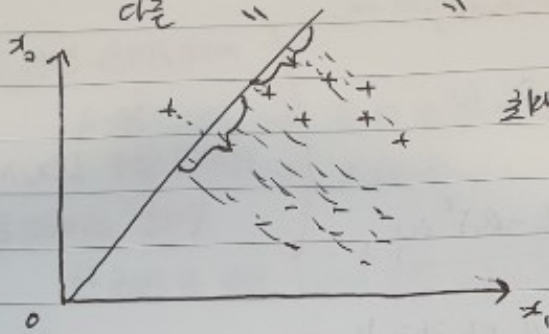
* 3.4 Linear Discriminant Analysis

↳ 각 클래스 샘플의 공분산 행렬이 같고 full rank.

↳ Fisher 판별 분석이라고도 함 (이름 뒤의 가정을 따르지 X)

Idea: 훈련 데이터셋을 전하고 샘플은 하나의 직선에 투영.

같은 class 에 속하는 sample을 가능한 가까운 선에 투영
다른 " " " " 최대한 먼 위치에!



크게한 가깝게 vs 크게한 멀게!

$$\text{Let } D = \left\{ (x_0, y_0) \right\}_{x=1}^m, \quad y_0 \in \{0, 1\}$$

X_0, μ_0, Σ_0 ; n 번째 클래스 평균, 공분산 행렬

① 같은 클래스 샘플 최대한 가까이

$$\text{minimize } \omega^T \Sigma_0 \omega + \omega^T \Sigma_1 \omega$$

(\because 직선은 1-d,
모두 실수!)

② 다른 클래스 샘플 충분히 멀리

$$\text{maximize } \| \omega^T \mu_0 - \omega^T \mu_1 \|^2$$

$$\text{maximize } J = \frac{\| \omega^T \mu_0 - \omega^T \mu_1 \|^2}{\omega^T \Sigma_0 \omega + \omega^T \Sigma_1 \omega} = \frac{\omega^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T \omega}{\omega^T \Sigma_0 \omega + \omega^T \Sigma_1 \omega}$$

집안 내 산포행렬 within-class scatter matrix

$$S_w = \Sigma_0 + \Sigma_1 = \sum_{x \in X_0} (x - \mu_0)(x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1)(x - \mu_1)^T$$

집안 간 산포행렬 between-class scatter matrix

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$$

$$J = \frac{\omega^T S_b \omega}{\omega^T S_w \omega}$$

S_b, S_w 의 일반화된 레일리 몫 generalized Rayleigh quotient

Date.

No.

In LDA (Latent Discriminator Analysis),

어떻게 w 를 계산하는가?

$$\text{WLOG, } w^T S_w w = 1,$$

$$\min_w -w^T S_b w \quad \text{s.t.} \quad w^T S_w w = 1$$

by Lagrange multiplier methods,

$$S_b w = \lambda S_w w$$

$$\text{Since } S_b w = (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w,$$

scalar

 $S_b w$ 의 방향은 $\mu_0 - \mu_1$ 과 같음 !!

$$S_b w = \lambda (\mu_0 - \mu_1), \quad \boxed{w = S_w^{-1} (\mu_0 - \mu_1)}$$

$$S_w = U \Sigma V^T \rightarrow S_w^{-1} = V \Sigma^{-1} U^T$$

singular value decomposition!

제약 최적화 문제

↓

라그랑주 함수 $L(x, \lambda)$ 의

문제와 최적화 문제

* 3.1 라쿤 분류 학습

WLOG, Let C_1, C_2, \dots, C_N classes.

기본 Idea = 분해법

$$\begin{cases} O_vO, \text{ One vs One} \\ O_vR, \text{ One vs Rest} \\ M_vM, \text{ Many vs Many} \end{cases}$$

① O_vO

N 개의 클래스 \rightarrow 둘씩 묶어서 분류기 구축

$\rightarrow N(N-1)/2$ 개의 이진 분류기 생성.

② O_vR

한 클래스만 양성, 나머지 음성 $\rightarrow N$ 개의 분류기!

③ M_vM

대부분 몇 개의 클래스를 양성값에, 나머지 클래스를 음성값으로 분해.

음? 어떻게?

✓ ECC, 오류 수정 코드 Error Correcting Output Codes

• 코드 ; N 개 클래스를 M 개로 나눈 ~~$(N > M)$~~ $(N < M)$
+ 훈련

• 코드 해독 ; M 개의 분류기는 각각 test sample의 예측을 진행
해당 예측 레이블의 하나의 코드가 됨!

코딩 matrix \rightarrow 이진코드 + 상위코드

	J_1	J_2	J_3	J_4	J_5	해당 예측
C_1	-1	+1	-1	+1	+1	$\rightarrow 3 \quad 2\sqrt{3}$
C_2	+1	-1	-1	+1	-1	$\rightarrow 4 \quad 4$
C_3	-1	+1	+1	-1	+1	$\rightarrow 1 \quad 2$
C_4	-1	-1	+1	+1	-1	$\rightarrow 2 \quad 2\sqrt{2}$
C_5	-1	-1	+1	-1	+1	$\uparrow \quad \uparrow$

(DAG directed Acyclic Graph)도 있음!

SVM에 대한 연구도!

※ 3.6 클래스 불균형 문제

※ 우선 모든 분류학습기의 ~~문제~~ 공통 가능한 것!

(\Rightarrow 서로 다른 클래스의 훈련 샘플들의 수가 같음!

\downarrow

이와 다르게 $PRF=2$ 이면 시야하면?

odds $\left(\frac{y}{1-y} \right)$: 양성값일 가능성과 음성값일 가능성의 비!

$\frac{y}{1-y} > \frac{m^+}{m^-}$ \rightarrow 관측 오즈!! (과연 맞는지 살펴보기)

$$\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^-}{m^+} \quad \text{rescaling, rebalancing} \quad \dots (1)$$

※ 관측 오즈 실제 오즈에 대해 효과적인 쿼리 X! \leftarrow

① under sampling

② over sampling (ex SMOTE, Interpolation 사용)

③ CD는 결점과 장점에 사용, 임계값 이동 threshold-moving.