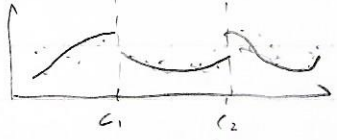


# Spline : 다항식 기반 회귀 함수

단계함수를 기저함수로 사용하며, 다항식을 기반으로 함. (단계함수 + 다항 회귀)

$$\rightarrow y_i = \beta_0 + \beta_1 C_1(x_i) + \dots + \beta_k C_k(x_i) + \epsilon_i$$

단, 여기서  $C_k$  함수는 1과 0을 return하는 함수가 아닌, 내부적으로 다항 회귀를 수행함. (일반적으로 3차함)



→ 불연속적인 변을 방지하기 위해 제약을 줌.

$$C_1(c_1) = C_2(c_1)$$

위와 같이 하더라도 해당포인트에서

비분이 불가능함



→ 추가적인 제약을 줌

$$C_1'(c_1) = C_2'(c_1), C_1''(c_1) = C_2''(c_1)$$

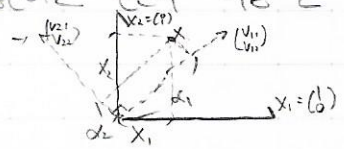
이를 통해 모든 점에서 비분이 가능해짐.



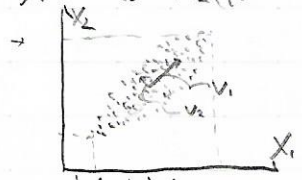
Spline은 구간을 사용하기 때문에 loess 보다 유연성이 적음. (loess는 모든 점에 대해)

## PCR (Principal Component Regression) : 주성분 회귀

주성분이란 분산이 가장 큰 좌표로 변환한 것.



$$X = X_1(c) + X_2(o) = \alpha_1(v_{11}) + \alpha_2(v_{21})$$



$X_1$ 과  $X_2$ 로는 데이터의 최대 분산을 설명할 수 없음. 새로운 좌표계  $V_1$ 을 통해라도 최대 분산을 설명할 수 있음

$V_1$ 은 첫번째 주성분으로 가장 큰 데이터의 분산을 설명하고,  $V_2$ 는 두번째 주성분으로  $V_1$ 과 직교함.  $\beta$ 개의 변수에 대해서  $P$ 개의 주성분이 있음. 데이터의 정보는 곧 데이터의 분산임. 주성분을 순서대로 데이터의 분산 정보를 많이 포함하고 있게 됨.



이를 통해 적은 분산의 값으로 회귀식으로 예측할 수 있음. PC를 찾는 과정은 PCA와 같은 방법으로 수행되며, (Y를 고려하지 않음) 이후 원하는 회귀분석을 수행하면 되겠음

→ 많은 경우에 기본적인 방향이 홀리이 관계가 있으나, 모든 경우에 최상은 아님.

PCR에서 사용하는 PC에 결측치에 따라 유연성이 감소함.

## PLS (Partial Least Square) : 최소 제곱 회귀

$$\text{좌표 변환} : Z_m = \sum_{j=1}^p \omega_{mj} X_j$$

새로운 좌표 군은 원래 좌표 X의 가중합임.

PCR을 X의 최대 분산만 군을 찾는 방식임.

PLS는 X와 Y의 최대 분산을 통해 군을 찾고자

하는 지도학습의 방법으로 아예짐.

→  $X_1, \dots, X_p$ 의 분산이 있을 때 각 분산을 기준으로 선형회귀를 수행하여 각 회귀계수인  $\beta_1, \dots, \beta_p$  ( $\omega_1, \dots, \omega_p$ )를 얻음.

→ 단 이때 모든 분산은 분산이 1이므로

중요한 변 선형으로 잘라내어 주어져야 함.

표준편차를 나눔으로써 각  $X$ 와 Y의 상관계수를 얻고, 새 좌표계는 상관계수에 따라 얻어짐.

$$Z_1 = \sum_{j=1}^p \omega_{1j} X_j, Z_2$$

를 기준으로 회귀식을 구하면,  $Y \sim \beta_0 + \beta_1 Z_1$  이 됨

회귀식이 설명하지 못하는 잔차  $r = Y - \beta_0 - \beta_1 Z_1$  이 됨. 이후 이 잔차 r을 설명하기 위해 또 다른 좌표계를 구할 수 있음