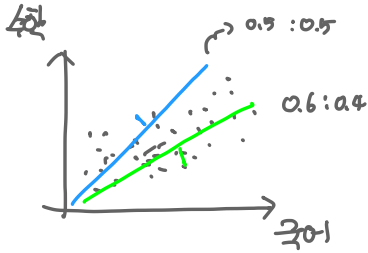


PCA (Principal Component Analysis)

⇒ 데이터 분포에서 어떤 선으로 잘 설명 시킬까? 데이터 분포 잘 표현 되었나? (각각 어떤 공할까? 변형 비율 (비율) 잘)



⇒ 이렇게 변형 비율 다

5:5로 변형할지 6:4로 변형할지
데이터 분포 따라 우리가 잘 변형할지
다르다! 최적의 projection 할 line 찾기

↳ 2차원 데이터를 2개의 변수로 찾아야 한다

principle component (주성분) = data의 분산이 가장 큰 방향 벡터 (잘 퍼짐)

ex) 얼굴 인식, 1000개 사진 받아 1000개 주성분 벡터들 분산이 큰 20개 받아서 사용 (앞면부터 얼굴 전반적 형태, 뒤면부터 세부적)

▣ K-mean clustering

⇒ 키와 몸무게에 따라 옷을 생선할 때, 모든 사이즈 고쳐할 수 없어
S.M.L로 사이즈 나눌 때 이용 가능 하다.

* Centroid = Cluster의 중심

데이터 분할, K 결정. Centroid 설정

while (변함 있음)

모든 데이터 Centroid와 유클리드 거리로 분류

Centroid를 cluster된 데이터의 공간으로 다시 설정

다시 새로운 Centroid로 모든 데이터 유클리드 거리로 분류

⇒ K 최적 값은 무라 단점

Soft max regression

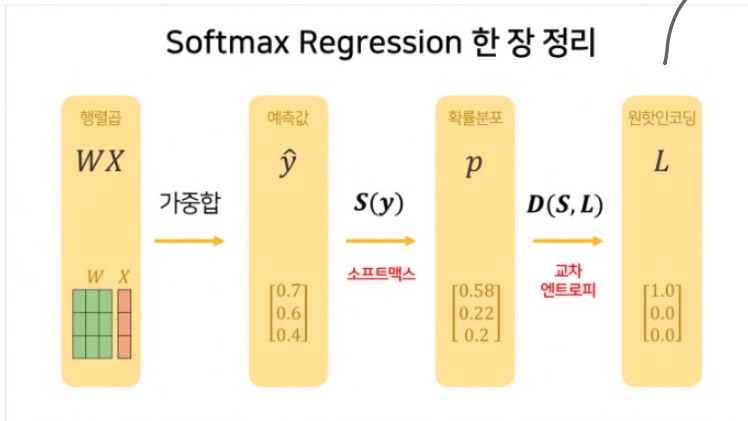
⇒ 로지스틱 회귀를 multi classification O/C.

(로지스틱은 성공/실패 2개 분류, soft max는 A,B,C 처럼 3개 이상 classification)

- 회색을 포함한 1의 확률 분포로 바꾸기 위해 이를 $S(y)$ 로 표현 되고 이를 거쳐 확률 분포 p 로 바뀐다

$$S(b_i) = \frac{e^{b_i}}{\sum_k e^{b_k}}$$

전체 그림



1 2 3. 가중치 따르니
0과 1로 가중치 없애고
표준화 효율적임!

SVM (support vector machine)

=> margin 을 최대화 (class 잘 분류) 하는 선을 찾는 것

- * Support Vector = 선과 가장 가까운 point
- * margin = support vector 와 선의 distance
- * decision boundary = 두 data 구분하는 선
- * robust 하다 = outlier의 영향을 덜 받는다.

		mean	median
(,)	1 2 3 4 5	3	3
	1 2 3 4 100	22	3
		↓	↓
		not robust	robust
		(outlier에 취약. 강건)	

-> 무조건 margin이 큰 선이 분류 잘하리 않을까?

과해서 나온 데이터를 분류하는 법이 많지 않을까?

그 법들 안에서 margin 최대화 하는 선 찾는 것

but 무조건 outlier 있는 경우 어느정도 outlier 무시라고 선찾을

-> Gamma 과다 decision boundary 휘어질, 작다. 직선에 가깝다.