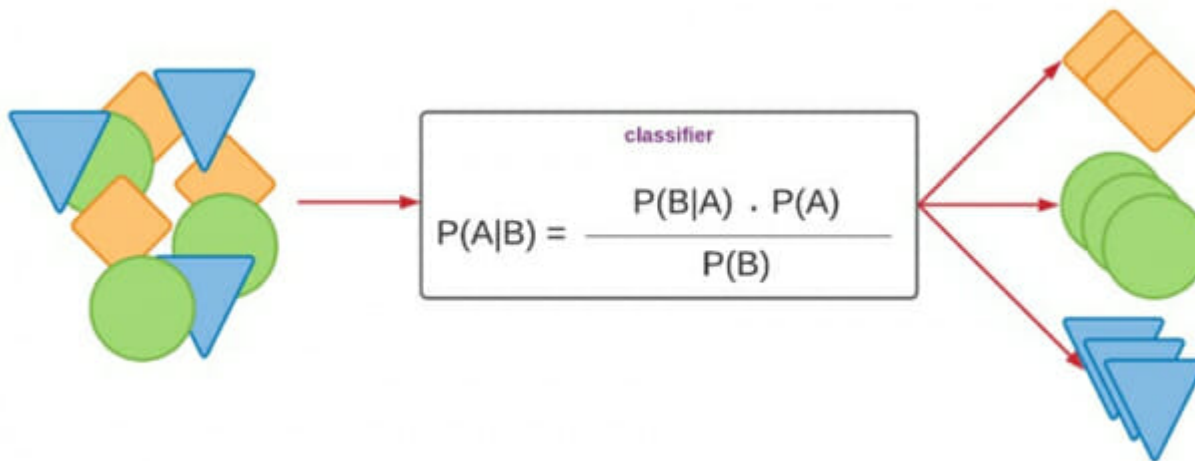
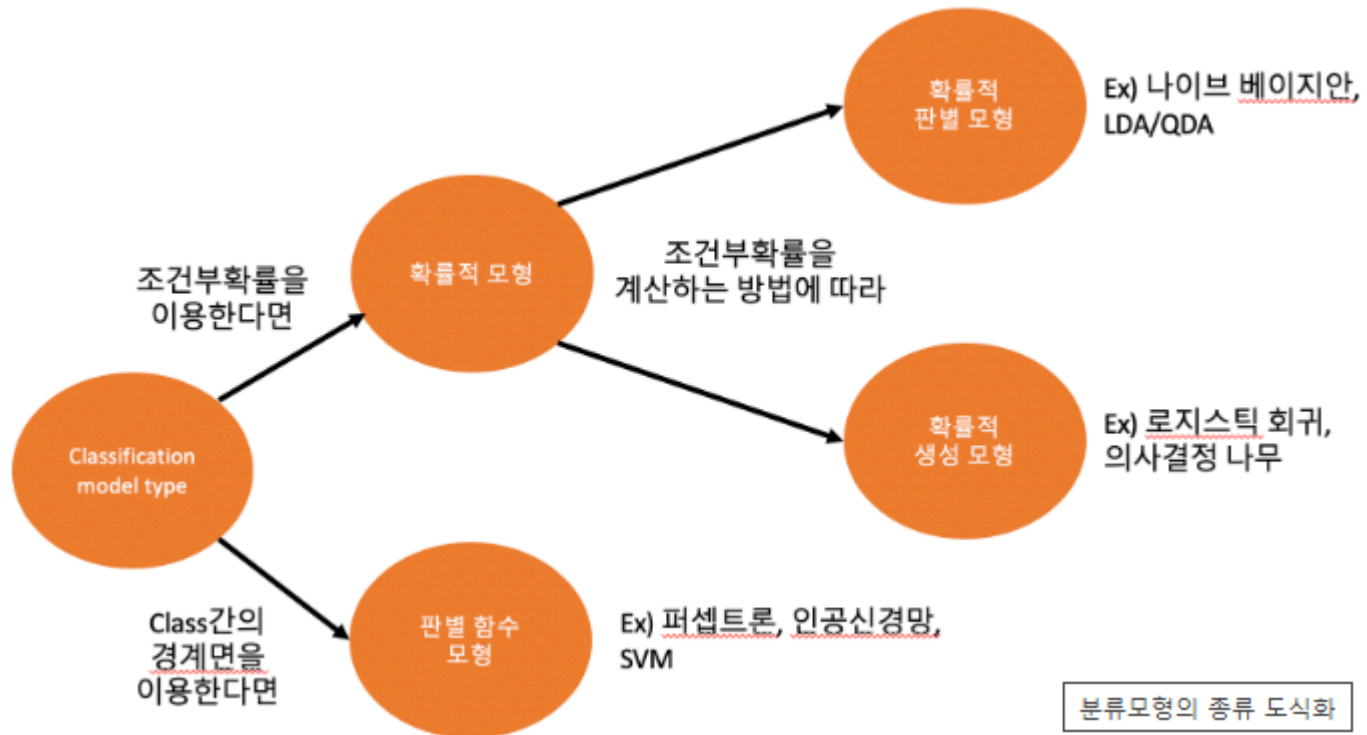


✓ 나이브 베이즈 정리(Naive Bayes Theorem)-추론

✓ 나이브 베이즈 분류(Naive Bayesian Classification)란?

- 데이터의 특징을 가지고 각 클래스(레이블)에 속할 확률을 계산하는 조건부 확률 기반의 분류 방법
- 데이터의 특징이 모두 상호 독립적이라는 가정하에 확률 계산을 단순화 \Rightarrow 나이브(naive)하다
- Bayes Theorem에 의해 데이터의 특징을 통해 클래스 전체의 확률 분포 대비 특정 클래스에 속할 확률을 구하는 것
- 나이브 베이즈 분류를 통해 데이터 특징이 하나 이상일 때 나이브 베이즈 공식으로 해당 데이터가 어떤 레이블에 속할 확률이 가장 높은지를 알 수 있음





✓ 1) 결합 확률

$$P(X, Y) = P(X | Y) \cdot P(Y)$$

두 가지 이상의 사건이 동시에 발생하는 확률

$$P(X, Y) = P(X | Y) \cdot P(Y) = \frac{P(X \cap Y)}{P(Y)} \cdot P(Y) = P(X \cap Y) = P(X) \cdot P(Y)$$



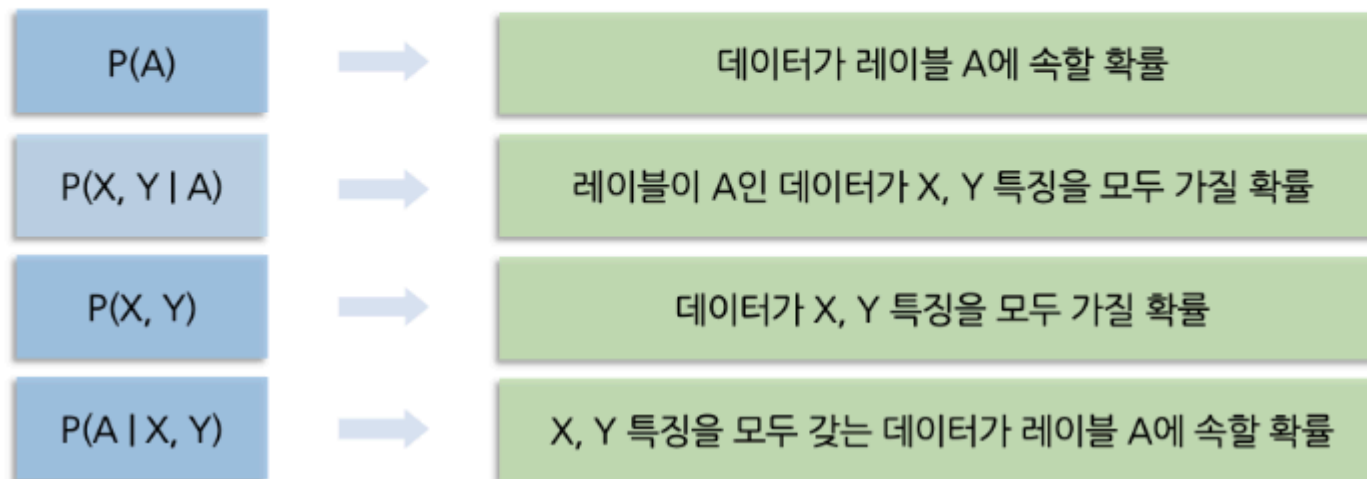
나이브 베이즈 알고리즘에서는 X, Y 사건이 독립이라고 가정한다.

✓ 2) 결합확률과 베이즈 정리

데이터의
특징: X, Y

전체 레이블
중 하나: A

$$P(A | X, Y) = \frac{P(X, Y | A) \cdot P(A)}{P(X, Y)} \quad (\text{by Bayes' Theorem})$$



$$P(X, Y) = P(X) \cdot P(Y)$$

$$P(X, Y | A) = P(X | A) \cdot P(Y | A)$$

어떤 한 데이터가 각 특징을 갖는 사건끼리는 서로 독립이라고 베이즈 정리에서 전제한다.

레이블이 A인 데이터가 특징 X를 가질 사건과 특징 Y를 가질 사건은 서로 독립이다.



$$P(A | X, Y) = \frac{P(X | A) \cdot P(Y | A) \cdot P(A)}{P(X) \cdot P(Y)}$$



데이터가 각 특징을 갖는 사건끼리 독립인 것이지, 특징과 레이블이 서로 독립이라는 의미가 아니다!

✓ 3) 나이브 베이즈 분류기의 원리

- 나이브 베이즈 분류기는 베이즈 정리를 사용하여 주어진 데이터를 특정 클래스에 속할 확률을 계산하고, 가장 높은 확률을 가진 클래스를 선택하는 방식
 - 이때, 모든 특성(feature)이 서로 독립적이라고 가정하며, 이러한 가정이 '나이브'라는 이름에 기인함
- 분류 과정
 - 각 클래스의 사전 확률 $P(C)$ 를 계산
 - 각 특성별로 조건부 확률 $P(x_i | C)$ 를 계산

- 테스트 데이터의 모든 특성에 대해 조건부 확률을 곱하여 $P(x|C)$ 를 구함
- $P(x|C)$ 와 $P(C)$ 를 곱하여 사후 확률 $P(C|x)$ 를 구함
- 가장 높은 확률을 가진 클래스를 선택

✓ 4) 나이브 베이즈 분류 적용 예시

💡 어떻게 쓰일 수 있을까? 💡



받은 메일이 스팸 메일인지 아닌지를 판단하는 분류기를 만들 수 있다!



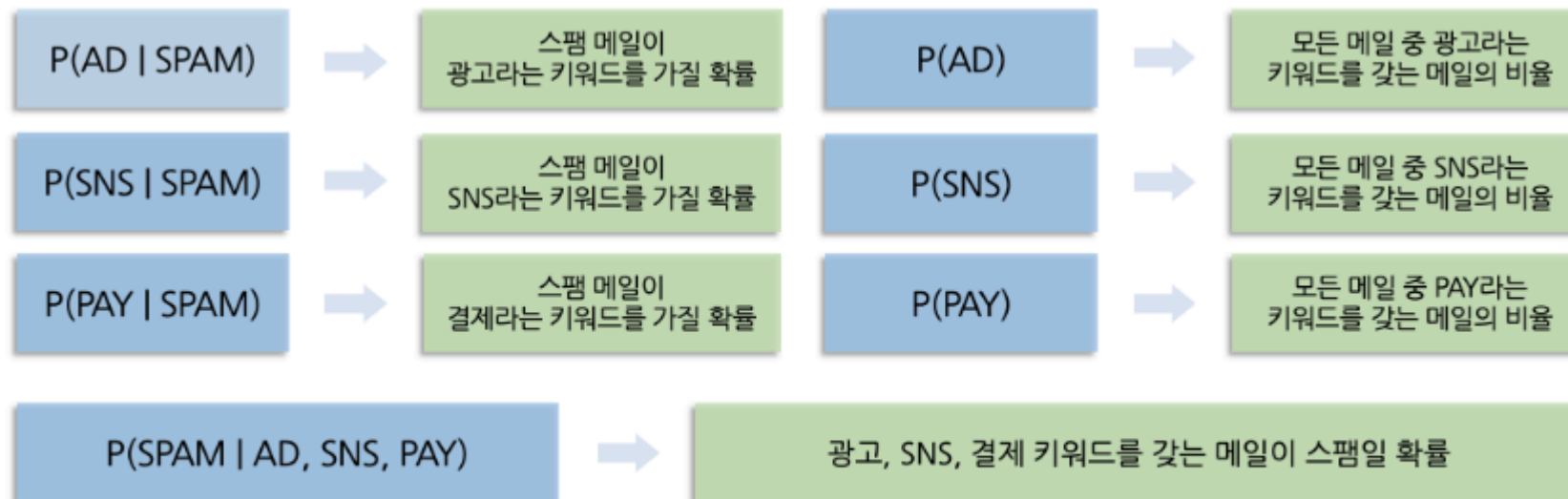
수신한 어떤 한 메일의 내용을 키워드로 추출했을 때
그 결과가 '광고', 'SNS', '결제'라고 하자.
이 키워드를 가지고 해당 메일이 스팸 메일인지 아닌지를 판단해 보자!

광고(AD), SNS(SNS), 결제(PAY) 키워드를 갖는 메일이 스팸(SPAM)일 확률을 구한다.



단, 메일이 광고, SNS, 결제 라는 각 키워드를 가질 사건은 서로 독립이라고 가정한다.

$$P(SPAM | AD, SNS, PAY) = \frac{P(AD | SPAM) \cdot P(SNS | SPAM) \cdot P(PAY | SPAM) \cdot P(SPAM)}{P(AD) \cdot P(SNS) \cdot P(PAY)}$$



스팸으로 분류될 확률이 임계치 이상이면 해당 메일을 스팸으로 분류한다!

✓ 스팸 메일 분류(이진 분류)에 적용

학습 세트



- 어떤 문장이 주어졌을 때 스팸 메일일 확률은 $P(\text{스팸 메일}|\text{문장})$ 로, 정상 메일일 확률은 $P(\text{정상 메일}|\text{문장})$ 로 표기할 수 있음
- 메일의 모든 문장이 2개의 단어로 구성되었다고 가정했을 때 나이브 베이즈 분류에서

- 스팸 메일일 확률

- $P(\text{스팸 메일}|\text{문장}) = P(\text{단어1}|\text{스팸 메일})P(\text{단어2}|\text{스팸 메일})P(\text{스팸 메일})$

- 정상 메일일 확률

- $P(\text{정상 메일}|\text{문장}) = P(\text{단어1}|\text{정상 메일})P(\text{단어2}|\text{정상 메일})P(\text{정상 메일})$


```

1 #라이브러리 불러오기
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.naive_bayes import MultinomialNB

```

```

1 #사이킷런에서 문서 데이터 세트를 가져옴
2 news = datasets.fetch_20newsgroups()

```

```

1 #입력 데이터와 타겟을 준비
2 X, y = news.data, news.target
3
4 #데이터 세트를 학습 세트와 테스트 세트로 분할
5 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)
6
7 #학습 세트의 첫 번째 입력 데이터를 확인
8 print(X_train[0])

```

↩ From: clay@rsd.dl.nec.com (Clay Finley)
 Subject: Re: Carrying Arms
 Nntp-Posting-Host: rsd21.rsd.dl.nec.com
 Organization: NEC America, Radio Software Dept
 Distribution: usa
 Lines: 22

```

|> In article <1993Apr5.220457.6800@spdc.ti.com> dwhite@epcot.spdc.ti.com (Dan White) writes:
|>
|> >However, haven't we already lost our right to bear arms?
|>
|> > It seems that in most states, like Texas, a citizen may own a
|> >gun and carry while at his home or business. But a citizen is severely
|> >restricted from bearing outside these areas. Here in Texas you cannot
|> >carry in your car except when "traveling" which is usually defined as
|> >"traveling across a county line." How did this come about? Are there
|> >any court rulings on the legality of restricting the carrying of a
|> >weapon outside the home?
|>

```

In Texas, it is legal to carry handguns while "traveling", and also to and from sporting activities. ^^^^^^

Chapter 46 of the Texas State Penal Code does NOT restrict long guns.
Therefore, it is legal to carry and transport long guns any place in Texas.

Regards,
Clay

```
1 #텍스트 데이터를 TF-IDF 벡터 값으로 변환
2 vectorizer = TfidfVectorizer()
3 X_train_vec = vectorizer.fit_transform(X_train)
4 X_test_vec = vectorizer.transform(X_test)
```

- TF-IDF에서 TF(Term Frequency)는 단어 빈도, 즉 단어의 출현 빈도를 말함

◦ 예를 들어 스티브 잡스가 말한 "the journey is the reward"라는 문장의 단어 빈도는 다음 같이 나타낼 수 있음

단어	the	hello	journey	world	is	reward
출현빈도	2	0	1	0	1	1

```
1 #다중 분류 나이브 베이즈 분류 모델 객체를 생성
2 model = MultinomialNB(alpha=0.01)
```