

J : J.H Lee is the criminal

K : the Criminal is Korean

M : the Criminal is male

E : The Criminal is on E

H : The criminal said "Hello"

$$P(J) = \frac{1}{10B}$$

$$P(J|E) = \frac{1}{10B}$$

$$P(J|\neg E) = \text{not defined}$$

① Human Being is  
only on E

② The criminal is H

중요한 힌트 fact (가장 확률)

$$P(J) = \frac{1}{10B}$$

① 2번 2번 ② 2번

$$P(J) = \frac{1}{10B}$$

$$P(J|H) = \frac{1}{10B}$$

$$P(J|\neg H) = \frac{1}{10B} \Rightarrow \text{Independent}$$

## Naïve Bayesian Classifier

A Simple Probabilistic Approach



# Recap: Probability

just assumption

- Conditional Probability, Independence & Bayesian Rule 독립인지
- A and B are independent given C → C is true false면 보장못함

$$\hookrightarrow P(A, B|C) = P(A|C) P(B|C)$$

$$\hookrightarrow P(A|B, C) = P(A|\neg B, C) = P(A|C)$$

Q, "로 쌍을 구분 즉, P(A, B|C)

A와 B의 독립에 대한 조건!

앞에 공식이랑 비슷한 구조 (가정체라고 생각해보기)

- Bayesian Rule given C

$$P(A|B, C) = \frac{P(A, B|C)}{P(B|C)} = \frac{P(B|A, C)P(A|C)}{P(B|C)}$$

조건이 달라졌기 때문에

- Which are True?

without assumption

with assumption

- If A and B are indep., A and B are indep. given C F (C에 따라 다른 수도)
- If A and B are indep. given C, A and B are indep. F
- If A and B are indep. given C, A and B are indep. Given D F

조건이 달라짐 : 참/거짓을 보장하지 못함

# Recap: Probability

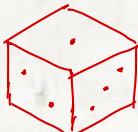
| ~3 page 중요하다!

## • Multinomial Distribution

- Die tossing with

not fair Dice

$$P(1) = \frac{1}{12}, P(2) = \frac{1}{12}, P(3) = \frac{1}{12}, P(4) = \frac{3}{12}, P(5) = \frac{3}{12}, P(6) = \frac{3}{12}$$



What is the probability of “1”=3, “2”=1, “4”=2 ← 6번던짐

순서에 대한 조합

$$\underline{C(6,3)} \times \underline{C(3,1)} \times \underline{C(2,2)} \times \left(\frac{1}{12}\right)^3 \times \left(\frac{1}{12}\right) \times \left(\frac{3}{12}\right)^2$$

- You toss a die several times, and repeat this several times

Die face	1	2	3	4	5	6
Trial 1	1	2	2	2	2	1
Trial 2	2	1	0	1	2	0
Trial 3	2	2	1	1	1	2

10 times  
6 times  
9 times ) total 25t

- Estimate  $P(1), P(2), \dots, P(6)$

$$= \frac{5}{25} \quad = \frac{5}{25} \quad \dots \quad = \frac{3}{25}$$

# Tennis Playing

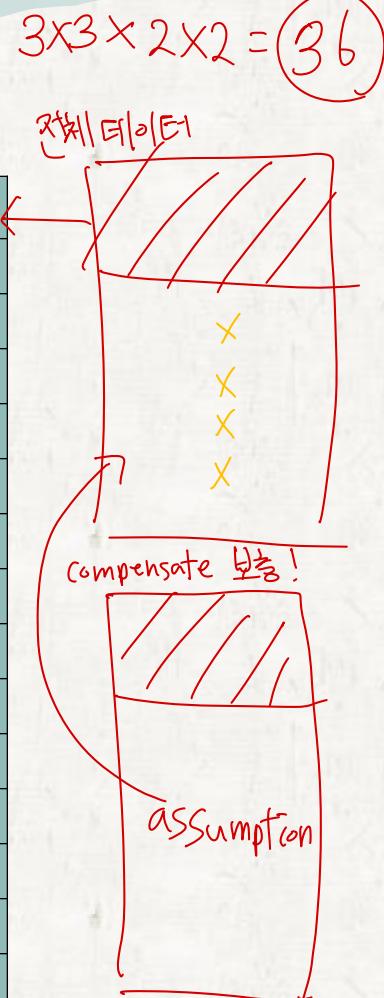
3x3x2x2 = 36

37/36

weather condition

day

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



# Tennis Playing

- Today is Sunny, Mild, High and Strong
  - Will he play tennis today?
- How can you predict his tennis playing? Conditional Probability
  - Let's use Probability
    - $P(\text{yes} | \text{sunny, mild, high, strong})$   $\Rightarrow$  We cannot evaluate!  
Given  
앞서라이드 데이터 36개의 CASE 전부  
나오지도 않았고, 대비터가 완전히 빠져나와  
(partial Data available)
    - $P(\text{no} | \text{sunny, mild, high, strong})$   
 $\Rightarrow$  Estimate! (some assumption 필요)
  - However, we need to observe a lot of “*sunny, mild, high, strong*” cases

# Tennis Playing

## Naïve Bayesian Classifier

- We assume the strongest assumption on independence to compensate shortage of data samples

naïve  
assumption

$\text{e.g. } (\text{sunny}, \text{mild}, \text{high} \dots)$

If class is given,  
Inputs are independent from each other.

assumption!

$$P(A, B | C) = P(A|C) P(B|C)$$

Whole samples

We may not say that  
inputs are independent

Mixed up

“Yes” class samples

We assume that  
inputs are independent

“No” class samples

We assume that  
inputs are independent

“simplest” case  
only!

# Tennis Playing

## Naïve Bayesian Assumption

	A	B	Class
1	O	O	Yes
2	O	X	Yes
3	X	O	Yes
4	X	X	Yes

$$P(A|Yes) = 0.5$$

$$P(B|Yes) = 0.5$$

$$P(A, B|Yes) = 0.25$$

$= P(A) P(B)$  Independent!

$$P(A|No) = 0.75$$

$$P(B|No) = 0.75$$

$$P(A, B|No) = 0.5625$$

Independent!

	A	B	Class
1	O	O	No
2	O	O	No
3	O	O	No
4	O	O	No
5	O	O	No
6	O	O	No
7	O	O	No
8	O	O	No
9	O	O	No
10	O	X	No
11	O	X	No
12	O	X	No
13	X	O	No
14	X	O	No
15	X	O	No
16	X	X	No
17	O	O	Yes
18	O	X	Yes
19	X	O	Yes
20	X	X	Yes

이상과는  
주어지지 X

dependent

$$P(A) = 0.7 \frac{14}{20}$$

$$P(B) = 0.7 \frac{14}{20}$$

$$P(A, B) = 0.5 \neq P(A) P(B)$$

) yes

# Tennis Playing

## Naïve Bayesian Classifier

Bayesian Rule  
switch !!, 0 해석 A (A 해석 0)

$$P(\text{yes}|\text{sunny}, \text{mild}, \text{high}, \text{strong}) = \frac{P(\text{sunny}, \text{mild}, \text{high}, \text{strong}|\text{yes})P(\text{yes})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

도리 가정 :  $P(A, B|C) = P(A|C)P(B|C)$  이용

$$= \frac{P(\text{sunny}|\text{yes})P(\text{mild}|\text{yes})P(\text{high}|\text{yes})P(\text{strong}|\text{yes})P(\text{yes})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

각각을 표에서 구함

$$P(\text{no}|\text{sunny}, \text{mild}, \text{high}, \text{strong}) = \frac{P(\text{sunny}, \text{mild}, \text{high}, \text{strong}|\text{no})P(\text{no})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$= \frac{P(\text{sunny}|\text{no})P(\text{mild}|\text{no})P(\text{high}|\text{no})P(\text{strong}|\text{no})P(\text{no})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

# Tennis Playing

- Naïve Bayesian Classifier

- Then, how to obtain the probabilities?

$$P(\text{yes}) = ?$$

$$P(\text{no}) = ?$$

$$P(\text{sunny}|\text{yes}) = ?$$

$$P(\text{yes}|\text{sunny}, \text{mild}, \text{high}, \text{strong}) =$$

$$P(\text{mild}|\text{yes}) = ?$$

$$\frac{P(\text{sunny} \mid \text{yes})P(\text{mild} \mid \text{yes})P(\text{high} \mid \text{yes})P(\text{strong} \mid \text{yes})P(\text{yes})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$P(\text{high}|\text{yes}) = ?$$

$$P(\text{strong}|\text{yes}) = ?$$

$$P(\text{sunny}|\text{no}) = ?$$

$$P(\text{no}|\text{sunny}, \text{mild}, \text{high}, \text{strong}) =$$

$$P(\text{mild}|\text{no}) = ?$$

$$\frac{P(\text{sunny} \mid \text{no})P(\text{mild} \mid \text{no})P(\text{high} \mid \text{no})P(\text{strong} \mid \text{no})P(\text{no})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$P(\text{high}|\text{no}) = ?$$

$$P(\text{strong}|\text{no}) = ?$$

# Tennis Playing

## Naïve Bayesian Classifier

- Let's ESTIMATE from the table

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(\text{sunny|yes}) = 2/9$$

$$P(\text{mild|yes}) = 4/9$$

$$P(\text{high|yes}) = 3/9$$

$$P(\text{strong|yes}) = 3/9$$

$$P(\text{sunny|no}) = 3/5$$

$$P(\text{mild|no}) = 2/5$$

$$P(\text{high|no}) = 4/5$$

$$P(\text{strong|no}) = 3/5$$

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

# Tennis Playing

## Naïve Bayesian Classifier

$$P(\text{yes} | \text{sunny}, \text{mild}, \text{high}, \text{strong})$$

$$= \frac{P(\text{yes})P(\text{sunny} | \text{yes})P(\text{mild} | \text{yes})P(\text{high} | \text{yes})P(\text{strong} | \text{yes})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$= \frac{1}{\alpha} \left( \frac{9}{14} \times \frac{2}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{3}{9} \right) = \frac{0.007055}{\alpha}$$

가장

$$P(\text{sunny}, \text{mild}, \text{high}, \text{strong}) = \alpha$$

$\Rightarrow$  No, 가능성이 비교할 수 있다.

$$P(\text{no} | \text{sunny}, \text{mild}, \text{high}, \text{strong})$$

$$= \frac{P(\text{no})P(\text{sunny} | \text{no})P(\text{mild} | \text{no})P(\text{high} | \text{no})P(\text{strong} | \text{no})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$= \frac{1}{\alpha} \left( \frac{5}{14} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{3}{5} \right) = \frac{0.04114}{\alpha}$$

둘째 더해서  $\alpha$  같은

구해도 되는데 의미 없다.

Will he play or not?

given!

yes, no 가 주어졌을 때

↓ 독립이지 (전치) 독립

≠

$P(s) P(m) P(h) P(s)$

“Some Special Case”

# Discussion

- Let's do it once more with the following data

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Rain	Mild	High	Weak	Yes
4	Rain	Cool	Normal	Strong	No
5	Sunny	Mild	High	Weak	No
6	Rain	Mild	Normal	Weak	Yes
7	Overcast	Mild	High	Strong	Yes
8	Rain	Mild	High	Strong	No

- Today is Sunny, Mild, High and Strong
  - Will he play tennis today?

# Discussion

## Naïve Bayesian Classifier

- Let's ESTIMATE the followings, but HOW?

$$P(\text{yes}) = 3/8$$

$$P(\text{no}) = 5/8$$

$$P(\text{sunny} | \text{yes}) = 0/3$$

$$P(\text{mild} | \text{yes}) = 3/3$$

$$P(\text{high} | \text{yes}) = 2/3$$

$$P(\text{strong} | \text{yes}) = 1/3$$

$$P(\text{sunny} | \text{no}) = 3/5$$

$$P(\text{mild} | \text{no}) = 2/5$$

$$P(\text{high} | \text{no}) = 4/5$$

$$P(\text{strong} | \text{no}) = 3/5$$

small size  
dataset!

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Rain	Mild	High	Weak	Yes
4	Rain	Cool	Normal	Strong	No
5	Sunny	Mild	High	Weak	No
6	Rain	Mild	Normal	Weak	Yes
7	Overcast	Mild	High	Strong	Yes
8	Rain	Mild	High	Strong	No

Sample 추가 Based on Expert knowledge

??

You!!

어떤 데이터에 smoothing을 적용하는 건  
미친 짓!

# Discussion

## Naïve Bayesian Classifier

$$P(\text{yes} \mid \text{sunny}, \text{mild}, \text{high}, \text{strong})$$

$$= \frac{P(\text{yes})P(\text{sunny} \mid \text{yes})P(\text{mild} \mid \text{yes})P(\text{high} \mid \text{yes})P(\text{strong} \mid \text{yes})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$= \frac{1}{\alpha} \left( \frac{3}{8} \times \frac{0}{3} \times \frac{3}{3} \times \frac{2}{3} \times \frac{1}{3} \right) = 0$$

$$P(\text{no} \mid \text{sunny}, \text{mild}, \text{high}, \text{strong})$$

$$= \frac{P(\text{no})P(\text{sunny} \mid \text{no})P(\text{mild} \mid \text{no})P(\text{high} \mid \text{no})P(\text{strong} \mid \text{no})}{P(\text{sunny}, \text{mild}, \text{high}, \text{strong})}$$

$$= \frac{1}{\alpha} \left( \frac{5}{8} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{3}{5} \right) = \frac{0.072}{\alpha}$$

Is it OK?

# Discussion

## Smoothing

- Adjust the probability

$$P(x_i = t | y = c) = \frac{N_{x_i=t, y=c} + \alpha}{N_{y=c} + \alpha n_i}$$

$$P(\text{sunny} | \text{yes}) = \frac{N_{\text{sunny}} + \alpha \cdot 1}{N_{\text{yes}} + \alpha \cdot 10}$$

↓  
number of

- $N_{x_i=t, y=c}$  : number of training examples for which  $x_i = t, y = c$
- $N_{y=c}$  : number of examples for which  $y = c$
- $n_i$  : number of values  $x_i$  can take
- $\alpha$  : smoothing parameter ( $\alpha \geq 0$ )
  - $\alpha = 1$ : Laplace smoothing
  - $\alpha < 1$ : Lidstone smoothing

추가하는 작업 ; smoothing

↳ avoid overfitting 하는 좋은 방법

내가 추가한 샘플이 자연하다:  $\alpha \uparrow$

과소적합 → 너무 단순한 모델, 학습 복잡한 패턴  
underfitting

특정 데이터에 대한 과도한 신뢰: 훈련 데이터에 없는 데이터를 무시해버림

(No sample)  
too much follow Data

overfitting ; smoothing

하기 어렵다  
과적합 상태임!

<over vs under> 1 dataset  $\Rightarrow p=0$ 으로 만드는 데이터 존재  
2,3 dataset  $\Rightarrow p=0$ 으로 만드는 데이터 없음

1set  $\rightarrow p=0$ 이네! 나머지 set도  $p=0$ 이겠지?: Overfitting: 데이터에 과도하게 적합  
1set, 2set, 3set 모두 성능↓, 데이터 학습 잘안됨: Underfitting: 데이터에 제대로 적합하지 못함  
+ 중요한 패턴 잡지못함  
A, B의 간접상관관계 등...

# Discussion

## Smoothing

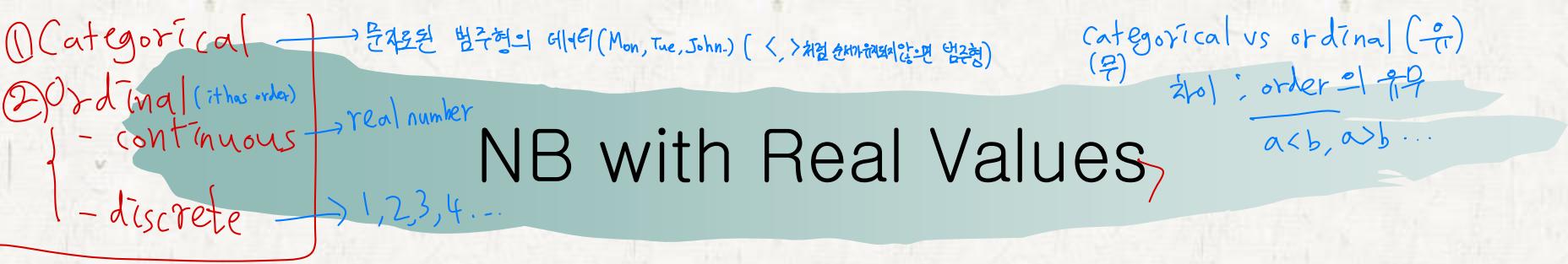
Observed samples  
from given data

$$P(x_i = t | y = c) = \frac{N_{x_i=t, y=c} + \alpha}{N_{y=c} + \alpha n_i}$$

Implicit observations  
by Expert Knowledge

- If there are more  $\alpha n_i$  samples of which  $y$  is  $c$ ,  
we may observe more  $\alpha$  samples of which  $x_i$  is  $t$  and  $y$  is  $c$

$$P(\text{sunny} | \text{yes}) = \frac{0 + 1}{3 + 2} = \frac{1}{5}$$



# NB with Real Values

## Tennis playing record with real values

	Outlook	Temp(°C)	Humidity	Wind	Play
1	Sunny	30	High	Weak	No
2	Sunny	31	High	Strong	No
3	Overcast	28	High	Weak	Yes
4	Rain	23	High	Weak	Yes
5	Rain	10	Normal	Weak	Yes
6	Rain	13	Normal	Strong	No
7	Overcast	14	Normal	Strong	Yes
8	Sunny	25	High	Weak	No
9	Sunny	15	Normal	Weak	Yes
10	Rain	22	Normal	Weak	Yes
11	Sunny	19	Normal	Strong	Yes
12	Overcast	20	High	Strong	Yes
13	Overcast	33	Normal	Weak	Yes
14	Rain	18	High	Strong	No

# NB with Real Values

- Today is
  - Outlook=Sunny, Temp=22, Humidity=High and Wind=Strong
- Will he play tennis today?

$$P(Yes|O = \text{sunny}, T = 22, H = \text{High}, W = \text{Strong})$$

$$= \frac{P(\text{sunny} | \text{yes}) P(T = 22 | \text{yes}) P(\text{high} | \text{yes}) P(\text{strong} | \text{yes}) P(\text{yes})}{P(\text{sunny, mild, high, strong})}$$

$$P(No|O = \text{sunny}, T = 22, H = \text{High}, W = \text{Strong})$$

$$= \frac{P(\text{sunny} | \text{no}) P(T = 22 | \text{no}) P(\text{high} | \text{no}) P(\text{strong} | \text{no}) P(\text{no})}{P(\text{sunny, mild, high, strong})}$$

# NB with Real Values

- If there are continuous values, how can I ESTIMATE the probability?

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(\text{sunny}|\text{yes}) = 2/9$$

$$P(T = 22|\text{yes}) = \boxed{??}$$

$$P(\text{high}|\text{yes}) = 3/9$$

$$P(\text{strong}|\text{yes}) = 3/9$$

$$P(\text{sunny}|\text{no}) = 3/5$$

$$P(T = 22|\text{no}) = \boxed{??}$$

$$P(\text{high}|\text{no}) = 4/5$$

$$P(\text{strong}|\text{no}) = 3/5$$

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	30	High	Weak	No
2	Sunny	31	High	Strong	No
3	Overcast	28	High	Weak	Yes
4	Rain	23	High	Weak	Yes
5	Rain	10	Normal	Weak	Yes
6	Rain	13	Normal	Strong	No
7	Overcast	14	Normal	Strong	Yes
8	Sunny	25	High	Weak	No
9	Sunny	15	Normal	Weak	Yes
10	Rain	22	Normal	Weak	Yes
11	Sunny	19	Normal	Strong	Yes
12	Overcast	20	High	Strong	Yes
13	Overcast	33	Normal	Weak	Yes
14	Rain	18	High	Strong	No

# NB with Real Values

ordinal & Observation

이면 정규분포가 가능할 것임

from observation

- Measurement values  $\rightarrow$  modeled with the Gaussian dist.

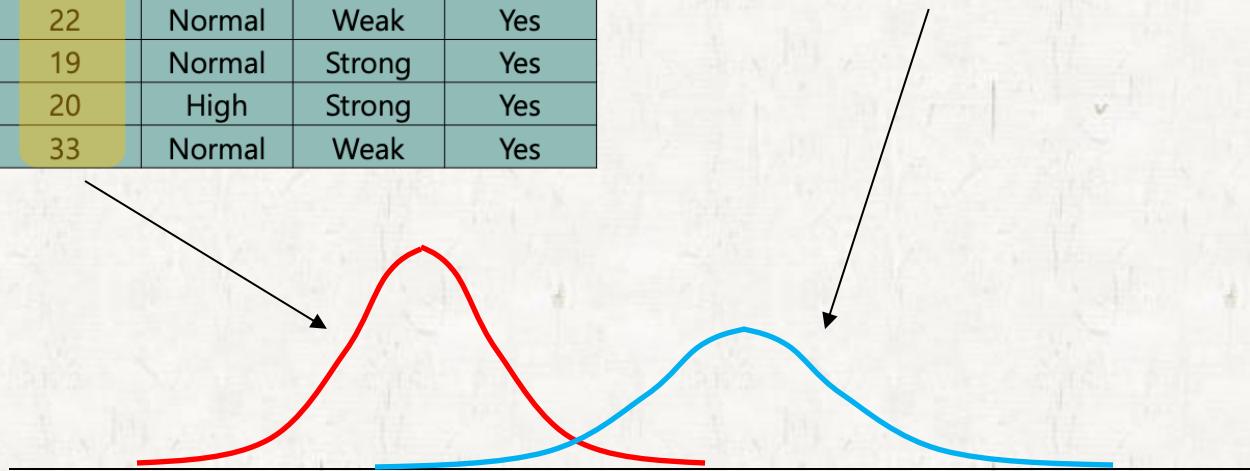
- Assumption: Temperature of each class will follow the Gaussian distribution

Yes 와 No에 대한 가우시안 분포(정규분포) 다르다

ordinal (discrete)

	Outlook	Temp	Humidity	Wind	Play
3	Overcast	28	High	Weak	Yes
4	Rain	23	High	Weak	Yes
5	Rain	10	Normal	Weak	Yes
7	Overcast	14	Normal	Strong	Yes
9	Sunny	15	Normal	Weak	Yes
10	Rain	22	Normal	Weak	Yes
11	Sunny	19	Normal	Strong	Yes
12	Overcast	20	High	Strong	Yes
13	Overcast	33	Normal	Weak	Yes

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	30	High	Weak	No
2	Sunny	31	High	Strong	No
6	Rain	13	Normal	Strong	No
8	Sunny	25	High	Weak	No
14	Rain	18	High	Strong	No



Categorical = {Mon, Tue, ..., Sun}  
= set 집합

XYZ 같은  
다른 string이  
존재할 수 있다

## NB with Real Values

ordinal & observation 둘다 맞지 않으면?

정수로 가정해도 됨

center point에서 확률이  
(평균) 가장 높은 확률을  
통해 확률을 정해줌

observation 값

을 통해 확률을 정해줌

~ 항상 참은 아니다

정수로  
다르는 것은  
아니다

- Measurement values → modeled with the Gaussian dist.
- Assumption: Temperature of each class will follow the Gaussian distribution

	Outlook	Temp	Humidity	Wind	Play
3	Overcast	28	High	Weak	Yes
4	Rain	23	High	Weak	Yes
5	Rain	10	Normal	Weak	Yes
7	Overcast	14	Normal	Strong	Yes
9	Sunny	15	Normal	Weak	Yes
10	Rain	22	Normal	Weak	Yes
11	Sunny	19	Normal	Strong	Yes
12	Overcast	20	High	Strong	Yes
13	Overcast	33	Normal	Weak	Yes

고난률을 통해 구하기

$$\mu = 20.4$$

$$\sigma = 7.1$$

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

ordinal ( $\neq$  continuous)

$$P(T = 22|yes) \propto \frac{1}{7.1 \cdot \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{22-20.4}{7.1}\right)^2} = 0.055$$

# NB with Real Values

- Measurement values → modeled with the Gaussian dist.
  - Assumption: Temperature of each class will follow the Gaussian distribution

	Outlook	Temp	Humidity	Wind	Play
1	Sunny	30	High	Weak	No
2	Sunny	31	High	Strong	No
6	Rain	13	Normal	Strong	No
8	Sunny	25	High	Weak	No
14	Rain	18	High	Strong	No

$$\mu = 23.4 \quad \sigma = 7.8$$

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$P(T = 22|no) \propto \frac{1}{7.8 \cdot \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{22-23.4}{7.8}\right)^2} = 0.050$$

이미지의 예

표준편차가 매우 크다

→ not reliable → 예측하기 어렵다

내의 경험기반으로  
(수학적)

23

# NB with Real Values

From observation

- If there are continuous values, how can I ESTIMATE the probability?

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(\text{sunny}|\text{yes}) = 2/9$$

$$P(T = 22|\text{yes}) = 0.055$$

$$P(\text{high}|\text{yes}) = 3/9$$

$$P(\text{strong}|\text{yes}) = 3/9$$

$$P(\text{sunny}|\text{no}) = 3/5$$

$$P(T = 22|\text{no}) = 0.050$$

$$P(\text{high}|\text{no}) = 4/5$$

$$P(\text{strong}|\text{no}) = 3/5$$

$$\frac{P(\text{sunny} \mid \text{yes})P(T = 22 \mid \text{yes})P(\text{high} \mid \text{yes})P(\text{strong} \mid \text{yes})P(\text{yes})}{P(\text{sunny, mild, high, strong})}$$

$$= \boxed{\frac{1}{\alpha}} \times \frac{2}{9} \times 0.055 \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}$$

$$\frac{P(\text{sunny} \mid \text{no})P(T = 22 \mid \text{no})P(\text{high} \mid \text{no})P(\text{strong} \mid \text{no})P(\text{no})}{P(\text{sunny, mild, high, strong})}$$

$$= \frac{1}{\alpha} \times \frac{3}{5} \times 0.050 \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14}$$

# NB with Real Values

from Count

## Naïve Bayesian Classifier with Counts

- You observe wind blows each day
- We count strong winds and weak winds in a day

one - event  
(한 가지 사건)  
= Wind

	Outlook	Temp	Humidity	S-Wind	W-Wind	Play
1	Sunny	Hot	High	1	3	No
2	Sunny	Hot	High	4	2	No
3	Overcast	Hot	High	2	5	Yes
4	Rain	Mild	High	1	3	Yes
5	Rain	Cool	Normal	2	2	Yes
6	Rain	Cool	Normal	3	2	No
7	Overcast	Cool	Normal	5	1	Yes
8	Sunny	Mild	High	2	3	No
9	Sunny	Cool	Normal	1	2	Yes
10	Rain	Mild	Normal	0	4	Yes
11	Sunny	Mild	Normal	3	1	Yes
12	Overcast	Mild	High	3	2	Yes
13	Overcast	Hot	Normal	1	2	Yes
14	Rain	Mild	High	3	3	No

# NB with Real Values

- Today is

- O=Sunny, T=Mild, H=High, S-Wind=3, and W-Wind=2

- Will he play tennis today?

$$P(Yes|O = \text{Sunny}, T = \text{Mild}, H = \text{High}, SW = 3, WW = 2)$$

$$= \frac{P(\text{Sunny} | yes)P(\text{Mild} | yes)P(\text{High} | yes)P(SW = 3, WW = 2 | yes)P(yes)}{P(\text{sunny, mild, high, 3, 2})}$$

$$P(No|O = \text{sunny}, T = \text{Mild}, H = \text{High}, SW = 3, WW = 2)$$

$$= \frac{P(\text{Sunny} | no)P(\text{Mild} | no)P(\text{High} | no)P(SW = 3, WW = 2 | no)P(no)}{P(\text{sunny, mild, high, 3, 2})}$$

- Gaussian dist. is not good for modeling counts

SW-WW

split x; 강한 상관관계가

있어서 하나로 처리

# NB with Real Values

- Naïve Bayesian Classifier with Counts

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(\text{sunny}|\text{yes}) = 2/9$$

$$P(\text{mild}|\text{yes}) = 4/9$$

$$P(\text{high}|\text{yes}) = 3/9$$

$$P(\text{SW} = 3, \text{WW} = 2|\text{yes}) = ?$$

$$P(\text{sunny}|\text{no}) = 3/5$$

$$P(\text{mild}|\text{no}) = 2/5$$

$$P(\text{high}|\text{no}) = 4/5$$

$$P(\text{SW} = 3, \text{WW} = 2|\text{no}) = ?$$

*# of SW*

	<b>Outlook</b>	<b>Temp</b>	<b>Humidity</b>	<b>S-Wind</b>	<b>W-Wind</b>	<b>Play</b>
<b>1</b>	Sunny	Hot	High	1	3	No
<b>2</b>	Sunny	Hot	High	3	1	No
<b>3</b>	Overcast	Hot	High	2	5	Yes
<b>4</b>	Rain	Mild	High	1	3	Yes
<b>5</b>	Rain	Cool	Normal	1	3	Yes
<b>6</b>	Rain	Cool	Normal	3	1	No
<b>7</b>	Overcast	Cool	Normal	3	3	Yes
<b>8</b>	Sunny	Mild	High	2	2	No
<b>9</b>	Sunny	Cool	Normal	1	2	Yes
<b>10</b>	Rain	Mild	Normal	0	4	Yes
<b>11</b>	Sunny	Mild	Normal	3	1	Yes
<b>12</b>	Overcast	Mild	High	3	2	Yes
<b>13</b>	Overcast	Hot	Normal	1	2	Yes
<b>14</b>	Rain	Mild	High	3	1	No

# NB with Real Values

## Naïve Bayesian Classifier with Counts

$$P(\text{yes}) = 9/14$$

$$P(\text{no}) = 5/14$$

$$P(\text{sunny}|\text{yes}) = 2/9$$

$$P(\text{mild}|\text{yes}) = 4/9$$

$$P(\text{high}|\text{yes}) = 3/9$$

$$\underbrace{P(\text{SW} = 3, \text{WW} = 2|\text{yes})}_{P(\text{sunny}|\text{no}) = 3/5} = \underbrace{C(5,3) \times P(\text{SW}|\text{yes})^3 \times P(\text{WW}|\text{yes})^2}_{\hookrightarrow 5C_3 = \text{SW } 3 \text{번이 들어갈 자리를 선택하기만 하면 됨}}$$

$$P(\text{mild}|\text{no}) = 2/5$$

$$P(\text{high}|\text{no}) = 4/5$$

$$P(\text{SW} = 3, \text{WW} = 2|\text{no}) = C(5,3) \times P(\text{SW}|\text{no})^3 \times P(\text{WW}|\text{no})^2$$

ML: 모든 구성을 simple  
But, 그들이 complex하게 연결되어 있다.

multinomial (동전 던지기 같은 것)

S S S WW

이거  
Data + knowledge experience  $\Rightarrow$  ML의 철학

# NB with Real Values

- We assume that Counts follows Multinomial dist.

	Outlook	Temp	Humidity	S-Wind	W-Wind	Play
3	Overcast	Hot	High	2	5	Yes
4	Rain	Mild	High	1	3	Yes
5	Rain	Cool	Normal	1	3	Yes
7	Overcast	Cool	Normal	3	3	Yes
9	Sunny	Cool	Normal	1	2	Yes
10	Rain	Mild	Normal	0	4	Yes
11	Sunny	Mild	Normal	3	1	Yes
12	Overcast	Mild	High	3	2	Yes
13	Overcast	Hot	Normal	1	2	Yes

- Estimate  $P(SW|yes)$  and  $P(WW|yes)$

$$P(SW=1|yes) = P(SW|yes) = \frac{\# \text{ of Strong winds}}{\# \text{ of All winds}} = \frac{2 + 1 + 1 + \dots}{(2 + 5) + (1 + 3) + (1 + 3) + \dots} = \frac{15}{40}$$

$$P(WW=1|yes) = P(WW|yes) = \frac{\# \text{ of Weak winds}}{\# \text{ of All winds}} = \frac{5 + 3 + 3 + \dots}{(2 + 5) + (1 + 3) + (1 + 3) + \dots} = \frac{25}{40}$$

# NB with Real Values

- We assume that Counts follows Multinomial dist.

	Outlook	Temp	Humidity	S-Wind	W-Wind	Play
1	Sunny	Hot	High	1	3	No
2	Sunny	Hot	High	3	1	No
6	Rain	Cool	Normal	3	1	No
8	Sunny	Mild	High	2	2	No
14	Rain	Mild	High	3	1	No

- Estimate  $P(SW|no)$  and  $P(WW|no)$

$$P(SW|no) = \frac{\# \text{ of Strong winds}}{\# \text{ of All winds}} = \frac{1 + 3 + 3 + 2 + 3}{(1 + 3) + (3 + 1) + (3 + 1) + \dots} = \frac{12}{20}$$

$$P(WW|no) = \frac{\# \text{ of Weak winds}}{\# \text{ of All winds}} = \frac{3 + 1 + 1 + 2 + 1}{(1 + 3) + (3 + 1) + (3 + 1) + \dots} = \frac{8}{20}$$

# Summary

- ➊ Conditional independence assumption is often violated but it is easy to implement and works surprisingly well anyway
- ➋ When to use
  - Moderate or large training set available
  - Attributes that describe instances are conditionally independent given classification
- ➌ Successful applications
  - Diagnosis
  - Classifying text documents: spam or non-spam