

Naive Bayes Learning

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

$$P(\text{Rain} | \text{Wet}) = \frac{\overset{\text{fact}}{P(W|R)} \cdot P(R)}{\underset{P(W)}{P(W)}}$$

การบ่งชี้ฝนเปียก
ฝนตกด้วยค่าความน่าจะเป็นเท่าไร

เปลี่ยนตามสถานการณ์
เช่น อุณหภูมิ คำน้อยมาก
เพราะตกวันเดียว

$$P(\text{Covid} | +) = \frac{P(+ | \text{Covid}) P(\text{Covid})}{P(+)}$$

$$P(- | \text{Covid} | +) = \frac{P(+ | - \text{Covid}) P(- \text{Covid})}{P(+)}$$

100 คน เป็น + 99 คน

100 คนที่ไม่เป็น ติด 3 คน

Posterior Hypothesis

จับ 2 ขั้นตอน

$$P(+ | \text{Covid}) P(\text{Covid})$$

$$(0.99)(0.01)$$

$$= 0.0099$$

$$P(+ | - \text{Covid}) P(- \text{Covid})$$

$$(0.03)(0.99)$$

$$= 0.0297$$

ลบกับ $P(\text{Covid})$ ได้

ผล 16 ผิดน้อยมาก แต่ $P(\text{Covid})$ ที่เปลี่ยนตามสถานการณ์แตกต่างกัน

ถ้ามีการเลือกโกนกัน จะทราบ

$$P(\text{Chest} | \text{Hell})$$

$$\frac{P(H|c) P(c)}{P(H)}$$

$$(0.95)(0.9)$$

$$= 0.855$$

$$P(\neg \text{Chest} | \text{Hell})$$

โอกาสที่ทราบเพราะสาเหตุอื่น

$$\frac{P(H|\neg c) P(\neg c)}{P(H)}$$

$$(0.2)(0.1)$$

$$= 0.02$$

$$(0.95)(x)$$

=

$$(0.2)(1-x)$$

$$0.95x$$

=

$$0.2 - 0.2x$$

$$1.15x$$

=

$$0.2$$

$$x$$

=

$$\frac{0.2}{1.15}$$

ต้องไม่ถึง 20% ถึงจะเชื่อได้ว่า,
ไม่ได้ทำผิด เมื่อไม่นาน

$$0.7 | 0.3$$

$$0.8 | 0.2$$

ถูกทั้งคู่

$$0.56$$

$$0.14$$

ผิดทั้งคู่

$$0.01$$

$$0.24$$

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

D15 Overcast Mild Normal Weak ?

ถ้า 2 เหตุการณ์ที่อิสระต่อกันก็หาพร้อมกัน

ค่าความน่าจะเป็นคือ 2 เหตุการณ์คูณกัน

$$P(Y | \text{Overcast, Mild, Normal, Weak})$$

$$P(N | \text{Overcast, Mild, Normal, Weak})$$

$$P(Y | \text{Overcast, Mild, Normal, Weak})$$

$$P(N | \text{Overcast, Mild, Normal, Weak})$$

$$P(\text{Overcast, mild, Normal, weak} | Y) P(Y)$$

$$\cancel{P(\text{Overcast, Mild, Normal, Weak})} \rightarrow$$

$$P(\text{Overcast, mild, Normal, weak} | N) P(N)$$

$$\cancel{P(\text{Overcast, Mild, Normal, Weak})} \rightarrow$$

$$P(\text{Overcast} | Y) P(\text{Mild} | Y) P(\text{Normal} | Y) P(\text{Weak} | Y) P(Y)$$

$$= \left(\frac{4+1}{9+1} \right) \left(\frac{4+1}{9+1} \right) \left(\frac{4+1}{9+1} \right) \left(\frac{6+1}{9+1} \right) \left(\frac{9+1}{14+1} \right)$$

$$P(\text{Overcast} | N) P(\text{Mild} | N) P(\text{Normal} | N) P(\text{Weak} | N) P(N)$$

$$= \left(\frac{0+1}{5+1} \right) \left(\frac{2+1}{5+1} \right) \left(\frac{1+1}{5+1} \right) \left(\frac{2+1}{5+1} \right) \left(\frac{5+1}{14+1} \right)$$

$$\frac{0 + 0.05}{5 + 0.05}$$

TF - IDF

↓
ความถี่ไม่คง

Andhuan Temp			Humidity			Line
w_1	w_2	w_3	w_4	w_5	...	w_{500}
0	0	0	1	1		0
0	0	0	1	2

$\frac{1}{\text{max}}$
 $\frac{1}{\text{max}}$
TF-IDF

```

1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction.text import TfidfTransformer
3 from sklearn.naive_bayes import MultinomialNB
4 from sklearn.metrics import accuracy_score
5
6 from sklearn.model_selection import train_test_split
7
8 vectorizer=CountVectorizer(tokenizer=lambda x: x.split())
9 tfidf=TfidfTransformer()
10 X=vectorizer.fit_transform(pos_data)
11 X=tfidf.fit_transform(X)
12 X_train, X_test, Y_train, Y_test = train_test_split(X, pos_target, test_size = 0.10, random_state = 42)
13 clf=MultinomialNB().fit(X_train,Y_train)
14 Y_predict = clf.predict(X_test)
15 accuracy_score(Y_test, Y_predict)

```

```

1 X.todense()

matrix([[1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1],
        [0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0],
        [1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 2, 0, 0]])

```

Not TF - IDF

```

1 Y.todense()

matrix([[0.25117442, 0., 0.33026418, 0.33026418, 0.25117442,
        0.25117442, 0.19505935, 0.33026418, 0.25117442, 0.25117442,
        0., 0.25117442, 0., 0.19505935, 0.33026418,
        0.33026418],
        [0., 0.40786601, 0., 0., 0.31019261,
        0., 0.24089223, 0., 0.31019261, 0.31019261,
        0.40786601, 0.31019261, 0.40786601, 0.24089223, 0.,
        ],
        [0.44652407, 0., 0., 0., 0.,
        0.44652407, 0.34676577, 0., 0., 0.,

```

↓
ลดความซ้ำซ้อน

Use TF - IDF