# Naive Bayes

Random Variables and Probability

A: Random variable

A: the next patient you examine has a cough

$P(A)$ : fraction of possible world which $A$ is true

↓
probability

$$0 \leq P(A) \leq 1 \quad P(true) = 1$$
$$P(false) = 0$$

S: Sample space

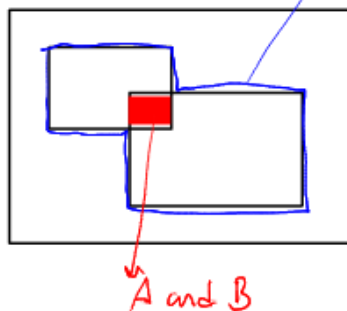$$P(A) = \frac{\text{\# of A's}}{\text{\# of elements in } S}$$

$$\sum P(A) = 1$$

## Dependence of Variables

Independence of Variables

$$P(A \text{ or } B) = \frac{P(A) + P(B)}{}$$



A or B

A and B

$$P(A) = P(A \text{ and } B) + P(A \text{ and not } B)$$

## Conditional Probability

$P(A|B)$ = Fraction of worlds in which B is TRUE that also have A TRUE

$P(H) = 1/10$
↓
headache

$P(F) = 1/40$
↓
flu

$P(H|F) = 1/2$

$P(F|H) = ?$

$$P(F/H) = \frac{P(F \text{ and } H)}{P(H)} =$$

$$P(F \text{ and } H) = P(H|F) * P(F) = \frac{1}{2} \cdot \frac{1}{40} = \frac{1}{80}$$

$$= P(F|H) * P(H)$$

$$P(F|H) = \frac{1/80}{1/10} = \frac{1}{8}$$

$$\boxed{\text{Chain Rule: } P(A|B) = \frac{P(A \text{ and } B)}{P(B)}} \searrow \neq 0$$

$$\boxed{\text{Bayes Rule: } P(A|B) = \frac{P(B|A) * P(A)}{P(B)}}$$

likelihood (Conditional probability)
prior probability
class · evidence (attribute)
evidence

$$P\left(elma \,\Big|\, \begin{matrix} ince \\ yaprak \end{matrix}\right) = \frac{P(incey. | elma) * P(elma)}{P(incey.)}$$

**Example**

$$P(\oplus|cancer) = 0{,}98$$
$$P(\theta | \overline{cancer}) = 0{,}97$$
$$P(cancer) = 0{,}08$$

Test yapılıyor, hasta kanser mi?

$$P(cancer | \theta) = ?$$

$$\left.\begin{matrix} P(cancer | \theta) \\ P(\overline{cancer} | \theta) \end{matrix}\right\} > \begin{matrix} ? \\ maximum \\ a posteriori \end{matrix} \Big\} \rightarrow MAP$$

$$h_{MAP} = \underset{h \in H}{argmax} (P(h|A))$$

$$\downarrow \begin{matrix} cancer \\ \overline{cancer} \end{matrix}$$

$$P(cancer | \theta) = \frac{P(\theta | cancer) * P(cancer)}{P(\theta)}$$

$$P(\overline{cancer} | \theta) = \frac{P(\theta | \overline{cancer}) * P(\overline{cancer})}{P(\theta)}$$

$$h_{MAP} = \{ P(\theta | cancer) * P(cancer), P(\theta | \overline{cancer}) * P(\overline{cancer}) \}$$

$$= \{ 0{,}98 * 0{,}008 , \; 0{,}03 * 0{,}992 \}$$

$$= \{ 0{,}07, \; 0{,}029 \} \rightarrow \overline{cancer}$$

$a_1, a_2, a_3, \ldots a_n$; attributes (bağımsız)

$$P(a_1 \cap a_2 \cap a_3 \ldots \cap a_n) = \prod_i P(a_i | v)$$

a özellikleri independent ↓ target value (class)

$$V_{MAP} = \underset{v_j \in V}{argmax} \; \frac{P(a_1, a_2, \ldots a_n | v_j) * P(v_j)}{\underbrace{P(a_1, \ldots a_n)}} = P(a_1, \ldots a_n | v_j) * P(v_j) \Rightarrow$$

↗ class ↘ constant

$$\underset{v_j \in V}{argmax} \; P(v_j | a_1, \ldots a_n) = ?$$

---

**NAIVE BAYES CLASSIFIER**

$$\Rightarrow V_{MAP} = \underset{v_j \in V}{argmax} \; \prod_i P(a_i | v_j) * P(v_j)$$

---

## Training

① Her sınıf için sınıflara ait örnek dökümanlar ayrı ayrı tek dök yapılır ($docs_j$)

② stop words temizlenir ($a, in, the \ldots$)

③ $P(v_j) = \frac{|docs_j|}{|All \; docs|} = \frac{n_j}{N}$

④ Kelimelerden sözlük oluşturulur

|  | class A | class B |
|---|---|---|
| fantastic | 0,2 | 0,5 |

$n_k$: $w_k$ kelimesinin o sınıfta kaç defa geçtiği

$$P(w_k | v_j) = \frac{n_k}{n_j}$$

↗ k kelimesinin o sınıftaki sayısı
↘ j. sınıftaki kelime sayısı

$$P(w_k | v_j) = \frac{n_k + m}{n_j + |vocabulary|}$$

→ m-estimate (sabit)
↳ sözlükteki toplam farklı kelime sayısı

---

**Example**

class A: " the cat crabs the crolls off the stairs "

class B: " It is raining cats and dogs "

Vocabulary = { cat, crab, croll, stair, rain, dog }

x = " cats eat mice and dogs bury bones "

$$V_{MAP} = \underset{\{c_a, c_b\}}{argmax} \; \prod_i P(a_i | v_j) * P(v_j)$$

$$P(v_j) = \frac{1}{2} \qquad P(w_k | v_j) = \frac{n_{kj} + 1}{n_j + 6}$$

a.) $P(A) * P(cat|A) * P(dog|A) = \frac{1}{2} * \frac{2}{4+6} * \frac{1}{10} = 0{,}01$

b.) $P(B) * P(cat|B) * P(dog|B) = \frac{1}{2} * \frac{2}{3+6} * \frac{2}{3+6} = 0{,}024$

$\left.\begin{array}{c}\end{array}\right\}$ $V_{MAP} \Rightarrow V_B$

log likelihood
↓
değerler çok küçükse
değerlerin log'ları alınabilir