# HMM for Bioinformatics

**Paweł Błażej**

Department of Genomics, Faculty of Biotechnology,

blazej@smorfland.uni.wroc.pl

6 marca 2019

1. $\mathcal{X}, \mathcal{Y}$ – alphabet;

1. $\mathcal{X}, \mathcal{Y}$ – alphabet;
2. $|\mathcal{X}| = |\mathcal{Y}| = L$ – size of an alphabet;

1. $\mathcal{X}, \mathcal{Y}$ – alphabet;
2. $|\mathcal{X}| = |\mathcal{Y}| = L$ – size of an alphabet;
3. $X, Y$ discrete random variables ($X : \Omega \to \mathcal{X}$);

1. $\mathcal{X}, \mathcal{Y}$ – alphabet;
2. $|\mathcal{X}| = |\mathcal{Y}| = L$ – size of an alphabet;
3. $X, Y$ discrete random variables ($X : \Omega \to \mathcal{X}$);
4. $P_X(x_i)$ probability.

We can think of entropy as the level of uncertainty associated with a random variable.

We define the entropy of the random variable $X \sim P$ as:

$$H(P) = H(X) = -\sum_{i=1}^{L} P(x_i) log_2(P(x_i))$$

where $H(0) = 0$ (convention $0 \cdot log_2 0 = 0$).

$$H(X, Y) = -\sum_x \sum_y P(x, y) log P(x, y)$$

$$H(X|Y) = -\sum_x \sum_y P(x,y) log P(x|y).$$

1. $H(X) \geqslant 0$

1. $H(X) \geqslant 0$
2. $H(X) \leqslant \log(L)$ with equality iff $X$ is uniformly distributed i.e. $P(x_i) = 1/L$ for every $i$;

1. $H(X) \geqslant 0$
2. $H(X) \leqslant log(L)$ with equality iff $X$ is uniformly distributed i.e. $P(x_i) = 1/L$ for every $i$;
3. $H(X, Y) \leqslant H(X) + H(Y)$ with equality iff $X$ and $Y$ are independent;

1. $H(X) \geqslant 0$
2. $H(X) \leqslant log(L)$ with equality iff $X$ is uniformly distributed i.e. $P(x_i) = 1/L$ for every $i$;
3. $H(X, Y) \leqslant H(X) + H(Y)$ with equality iff $X$ and $Y$ are independent;
4. $H(X)$ is concave in $X$ .

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y).$$

$$M(X, Y) = \sum_x \sum_y P(x, y) log \frac{P(x, y)}{P(x)P(y)}$$

$$M(X, Y) = \sum_x \sum_y P(x, y) log \frac{P(x, y)}{P(x)P(y)}$$

### Remark

$$M(X, Y) = H(H) + H(Y) - H(X, Y).$$

$$\mathcal{X} = \{x_1, x_2\}, \ P(x_1) = p, \ P(x_2) = 1 - p$$

$$H(p) = -p\log_2 p - (1 - p)\log_2(1 - p)$$

1. $\mathcal{X} = \{A, T, G, C\}$;

1. $\mathcal{X} = \{A, T, G, C\}$;
2. $p(X = A) = p(X = T) = \ldots = p(X = C) = 0.25$;

1. $\mathcal{X} = \{A, T, G, C\}$;
2. $p(X = A) = p(X = T) = \ldots = p(X = C) = 0.25$;
3. $H(X) = (\frac{1}{4} log \frac{1}{4} + \ldots + \frac{1}{4} log \frac{1}{4}) = 2$;

1. $\mathcal{X} = \{A, T, G, C\}$;
2. $p(X = A) = p(X = T) = \ldots = p(X = C) = 0.25$;
3. $H(X) = (\frac{1}{4} log \frac{1}{4} + \ldots + \frac{1}{4} log \frac{1}{4}) = 2$;
4. We can think in this case of the entropy as a number of yes/no questions needed to identify an autcome.

### Definition

A code is called uniquely decodable if any string composed of finite number of symbols from $\mathcal{X}$ gets a unique codestring.

## Definition

A code is called uniquely decodable if any string composed of finite number of symbols from $\mathcal{X}$ gets a unique codestring.

## Kraft inequality

A necessary and sufficient condition on the lengths of codestring $l_i$ in a uniquely decodable code is:

$$\sum_{i=1}^{L} 2^{-l_i} \leqslant 1.$$

The expected length of the optimal uniquely decodable code belongs to the interval $[H(X), H(X) + 1)$

$$D(P_1|P_2) = \sum_{i=1}^{L} P_1(x_i) log \frac{P_1(x_i)}{P_2(x_i)}$$

$$D(P_1|P_2) = \sum_{i=1}^{L} P_1(x_i) log \frac{P_1(x_i)}{P_2(x_i)}$$

### Convention

1. $0 \cdot log \frac{0}{P_2(x_i)} = 0$;
2. $P_2(x_i) log \frac{P_1(x_i)}{0} = \infty$

If X is a random variable with the probability $P_1$ on an alphabet of $L$ symbols, and $P_2$ is the uniform distribution on this alphabet then

$$D(P_1|P_2) = log(L) - H(X).$$

1. T. Inglot, „Teoria informacji a statystyka matematyczne" wykład wygłoszony na XXXVIII Konferencji Statystyka Matematyczna Wisła 2012;

2. L. Dębowski, "Information Theory and Statistics", Institute of Computer Science Polish Academy of Science, Warsaw, 2013.