# Information Theory and Statistics
# Lecture 8: Complexity and entropy

Łukasz Dębowski
ldebowsk@ipipan.waw.pl

Ph. D. Programme 2013/2014

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Kolmogorov complexity and entropy

- Prefix-free complexity is the length of a prefix-free code.
- Hence we may expect that it can be related to entropy.

# Distribution of random variables

- Let $(\mathbf{X_i})_{i \in \mathbb{N}}$ be a sequence of random variables, $\mathbf{X_i : \Omega \to \Gamma}$.
- We denote its a probability distribution

$$\mathbf{P(x_1^m) = P(X_1^m = x_1^m)}.$$

- We will consider Kolmogorov complexity $\mathbf{K(x_1^m | P)}$, which is the prefix-free Kolmogorov complexity of $\mathbf{x_1^m}$ given the definition of distribution of $\mathbf{P}$ on the infinite tape.
- If $\mathbf{P}$ is computable then

$$\mathbf{K(x_1^m | P) \overset{+}{<} K(x_1^m) \overset{+}{<} K(x_1^m | P) + K(P)}.$$

# Shannon-Fano coding

### Theorem

*For any distribution* $\mathbf{P}$,

$$K(x_1^m|P) \overset{+}{<} - \log P(x_1^m) + 2 \log m.$$

### Proof

The inequality follows from the fact that a certain program that computes $\mathbf{x_1^m}$ has form "having the definition of $\mathbf{P}$ and the length of string $\mathbf{x_1^m}$, take the Shannon-Fano code word for $\mathbf{x_1^m}$ with respect to $\mathbf{P}$ and compute $\mathbf{x_1^m}$ from it."

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Source coding inequality

> **Theorem (source coding inequality)**
>
> Let $\mathbf{B} : \mathbf{\Gamma}^m \to \mathbf{\Gamma}^*$ be a prefix-free code. For any distribution $\mathbf{P}$,
> $$\sum_{x_1^m} P(x_1^m) \left[ |B(x_1^m)| + \log P(x_1^m) \right] \geq 0.$$

Since prefix-free Kolmogorov complexity is the length of a prefix-free code, we obtain the following result:

> **Theorem**
>
> $$0 \leq \sum_{x_1^m} P(x_1^m) \left[ K(x_1^m|P) + \log P(x_1^m) \right] \overset{+}{<} 2 \log m.$$

# Barron theorem

> **Theorem (Barron theorem)**
>
> Let $\mathbf{B} : \mathbf{\Gamma}^* \to \mathbf{\Gamma}^*$ be a prefix-free code. For any distribution $\mathbf{P}$,
>
> $$\lim_{\mathbf{m} \to \infty} [|B(X_1^m)| + \log P(X_1^m)] = \infty$$
>
> holds with $\mathbf{P}$-probability $\mathbf{1}$.

Since prefix-free Kolmogorov complexity is the length of a prefix-free code, we obtain the following result:

> **Theorem**
>
> $$0 \leq [K(X_1^m|P) + \log P(X_1^m)] \overset{+}{<} 2\log m$$
>
> holds for sufficiently large $\mathbf{m}$ with $\mathbf{P}$-probability $\mathbf{1}$.

# Markov inequality

## Theorem (Markov inequality)

*Let $\epsilon > 0$ be a fixed constant and let $\mathbf{Y}$ be a random variable such that $\mathbf{Y} \geq \mathbf{0}$. We have*

$$\mathbf{P}(\mathbf{Y} \geq \epsilon) \leq \frac{\mathbf{E}\,\mathbf{Y}}{\epsilon}.$$

## Proof

Consider random variable $\mathbf{Z} = \mathbf{Y}/\epsilon$. We have

$$\mathbf{P}(\mathbf{Y} \geq \epsilon) = \int_{\mathbf{Z} \geq 1} d\mathbf{P} \leq \int_{\mathbf{Z} \geq 1} \mathbf{Z}\,d\mathbf{P} \leq \int \mathbf{Z}\,d\mathbf{P} = \frac{\mathbf{E}\,\mathbf{Y}}{\epsilon}.$$

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Borel-Cantelli lemma

Denote

$$\limsup_{n \to \infty} \mathbf{A_n} = \{\omega : \omega \in \mathbf{A_m} \text{ for infinitely many } \mathbf{m}\}.$$

We have

$$\left(\limsup_{n \to \infty} \mathbf{A_n}\right)^{\mathbf{c}} = \{\omega : \omega \notin \mathbf{A_m} \text{ for sufficiently large } \mathbf{m}\}.$$

For proving that some events hold with probability $\mathbf{1}$, the following proposition is particularly useful.

---

**Theorem (Borel-Cantelli lemma)**

*If $\sum_{m=1}^{\infty} \mathbf{P(A_m)} < \infty$ for a family of events $\mathbf{A_1, A_2, A_3, ...}$ then*

$$\mathbf{P}\left(\limsup_{n \to \infty} \mathbf{A_n}\right) = \mathbf{0}.$$

---

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Proof of the Borel-Cantelli lemma

Notice that $\sum_{m=1}^{\infty} \mathbf{P}(\mathbf{A_m}) < \infty$ implies

$$\lim_{m \to \infty} \sum_{k=m}^{\infty} P(A_k) = 0.$$

Hence we obtain

$$\mathbf{P}(\{\omega : \omega \in \mathbf{A_m} \text{ for infinitely many } \mathbf{m}\})$$
$$= \mathbf{P}(\{\omega : \forall_{m \geq 1} \exists_{k \geq m} \omega \in \mathbf{A_k}\})$$
$$= \mathbf{P}\left(\bigcap_{m=1}^{\infty} \bigcup_{k=m}^{\infty} \mathbf{A_k}\right)$$
$$\leq \inf_{m \geq 1} \mathbf{P}\left(\bigcup_{k=m}^{\infty} \mathbf{A_k}\right) \leq \inf_{m \geq 1} \sum_{k=m}^{\infty} \mathbf{P}(\mathbf{A_k}) = \mathbf{0}.$$

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Proof of Barron theorem

Let us write

$$W(x_1^m) = \frac{2^{-|B(x_1^m)|}}{P(x_1^m)2^{-n}}.$$

By the Markov inequality we obtain

$$\sum_{m=1}^{\infty} P\left(|B(X_1^m)| + \log P(X_1^m) \le n\right)$$

$$= \sum_{m=1}^{\infty} P\left(W(X_1^m) \ge 1\right)$$

$$\le \sum_{m=1}^{\infty} \sum_{x_1^m} P(x_1^m)W(x_1^m) = \sum_{m=1}^{\infty} \sum_{x_1^m} 2^{-|B(x_1^m)|+n}$$

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Proof (continued)

Continuing, by the Kraft inequality we obtain,

$$\sum_{m=1}^{\infty} P\left(|B(X_1^m)| + \log P(X_1^m) \le n\right)$$

$$\le \sum_{m=1}^{\infty} \sum_{x_1^m} 2^{-|B(x_1^m)|+n} \le 2^n < \infty.$$

Hence from the Borel-Cantelli lemma we obtain that

$$|B(X_1^m)| + \log P(X_1^m) > n \text{ for sufficiently large } m$$

holds with $P$-probability $1$.

The $n$ in this statement is arbitrary so the claim follows.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# An analogue of the chain rule

The parallels between prefix-free complexity and entropy can be drawn further. The following theorem is an analogue of the chain rule
$H(X, Y) = H(X) + H(Y|X)$.

> **Theorem**
>
> $$K(\langle u, w \rangle) \stackrel{+}{=} K(u) + K(w | \langle u, K(u) \rangle).$$

In the proposition above, it is easy to show that the left hand side is smaller than the right hand side. The proof of the converse inequality is harder.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Partial proof of the chain rule

We will demonstrate that

$$K(\langle u, w \rangle) \overset{+}{<} K(u) + K(w| \langle u, K(u) \rangle).$$

Let $p$ be the shortest program that satisfies $V(p) = u$ and let $p'$ be the shortest program that satisfies $V(p'| \langle u, K(u) \rangle) = w$. Then there exists a prefix-free machine $S$ that satisfies $S(pp') = \langle u, w \rangle$. Hence we obtain the claim.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

Ph.D
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Incomplete analogy

In the algorithmic chain rule there appears term $K(w| \langle u, K(u)\rangle)$ rather than $K(w|u)$. Although $K(w| \langle u, K(u)\rangle)$ differs from $K(w|u)$, we can see that $K(\langle u, K(u)\rangle)$ and $K(u)$ are approximately equal.

### Theorem

$$K(\langle w, K(w)\rangle) \overset{\pm}{=} K(w).$$

### Proof

From the shortest program that computes $w$, we may reconstruct both $w$ and $K(w)$. Hence $K(\langle w, K(w)\rangle) \overset{+}{<} K(w)$. On the hand, we have $K(\langle w, K(w)\rangle) \overset{+}{>} K(w)$ from a previous theorem.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Algorithmic information

### Definition

We define *algorithmic information* between strings **u** and **w** as

$$\mathsf{I(u; w) = K(w) - K(w|\,\langle u, K(u)\rangle)}.$$

# Symmetry of algorithmic information

**Theorem**

$$I(u; w) \overset{+}{=} I(w; u).$$

**Proof**

Observe that

$$I(u; w) = K(w) - K(w| \langle u, K(u) \rangle)$$
$$\overset{+}{=} K(w) + K(u) - K(\langle u, w \rangle)$$
$$\overset{+}{=} K(u) - K(u| \langle w, K(w) \rangle) = I(w; u).$$

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY