# Lesson 4
# Entropy

## 4-3: Algorithmic Entropy: Kolmogorov Complexity

### Jan Reimann

Math 574, Topics in Logic

Penn State, Spring 2014

# Information Content of Strings

In the previous lectures, we defined an information measure *I* for random variables/probability distributions.

Then we defined entropy as expected information.

Q: Can we define the information content of an *individual string* $\sigma \in 2^{<\mathbb{N}}$?

# Information Content of Strings

We could try: Put a uniform probability distribution $\lambda$ on $\{0, 1\}^n$, and define

$$I(\sigma) = -\log \lambda(\sigma).$$

$$\lambda(\sigma) = 2^{-n}$$

Problem: all strings have the same probability, hence the same information content.

But we would expect the string

$$0000000000 \ldots 0000$$

to have low information content, while the

outcome of a coin toss

has high information content.

Q: Can we find a probability distribution on strings such that

simple strings have high probability,  $-\log P$

complex strings have low probability?

*We will see later that this indeed possible.*

# Kolmogorov Complexity

Idea: low information content = high compressibility

Kolmogorov complexity makes this idea rigorous.

Let $M$ be a Turing machine, $\sigma \in 2^{<\mathbb{N}}$.

$$C_M(\sigma) = \min\{|p| : M(\underbrace{p}_{\text{code}}) = \underline{\sigma}\},$$

where we let $\min \emptyset = \infty$.

$\underline{M}$-complexity = length of shortest $\underline{M}$-description (code)

*decoder*

Problem: arbitrariness in the choice of $M$. Different machines can assign the same string drastically different complexities.

*Solution:*     *Use universal Turing machines.*

$\sigma$   Complex

$M_0$   long pgt

$M_1$   "$\sigma$ burned in"

We define a pairing function for strings as $\langle \sigma, \tau \rangle = 0^{|\sigma|} 1 \sigma \tau$.

We also identify natural numbers with their binary representation.

Note: $|m| = \log(m)$.

Recall a universal Turing machine $U$ emulates all other TM's:

$$U(\langle e, \sigma \rangle) = M_e(\sigma).$$

Fix any universal TM $U$ and define $C(\sigma) = C_U(\sigma)$.

# Kolmogorov Complexity

**THM:** [Invariance Theorem]

*Kolmogorov, Solomonoff*

For any TM $M$ there exists a constant $c_M$ such that

$$\forall \sigma \; C(\sigma) \leqslant C_M(\sigma) + c_M.$$

$M = M_e$

▶ Proof: If $p$ is a shortest $M$-program for $\sigma$, and $e$ is the Gödel number of $M$, then $\langle e, p \rangle$ is a $U$-program for $\sigma$, and hence

$c_M$

$$C(\sigma) \leqslant |\langle e, p \rangle| = |0^{|e|}1ep| = 2\log(e) + |p| + 1 = C_M(\sigma) + 2\log(e) + 1.$$

$$|e| + 1 + \log(e) + |p|$$
$$\text{\tiny "} \log(e)$$

$U(\langle e, p \rangle)$

$= M(p)$

*Notation: $C(\sigma) \leqslant^+ f(\sigma)$ means: There exists $c$ s.t.*

$$\forall \sigma \; C(\sigma) \leqslant f(\sigma) + c$$

$$\leqslant f(\sigma) + O(1)$$

.

1.  $C(\sigma) \leqslant^{+} |\sigma|$.

    Consider the copy machine $M(p) = p$. We have $C_M(\sigma) \leqslant |\sigma|$ and hence by the invariance theorem $C(\sigma) \leqslant^{+} |\sigma|$.

    *incompressib[le]*

2.  For any $n$, there exists a string $\sigma$ of length $n$ with $C(\sigma) \geqslant |\sigma| = n$.

    A simple counting argument: There are $2^n$ strings of length $n$, but only $1 + 2^1 + 2^2 + \cdots + 2^{n-1} = 2^n - 1$ programs of length $< n$.

3.  If $h : 2^{<\mathbb{N}} \to 2^{<\mathbb{N}}$ is computable then $C(h(\sigma)) \leqslant^{+} C(\sigma)$.

    Let $M$ be a Turing machine given by $M(\sigma) = h(U(\sigma))$. Then by the invariance theorem,

    $$C(h(\sigma)) \leqslant^{+} C_M(h(\sigma)) \leqslant C(\sigma).$$
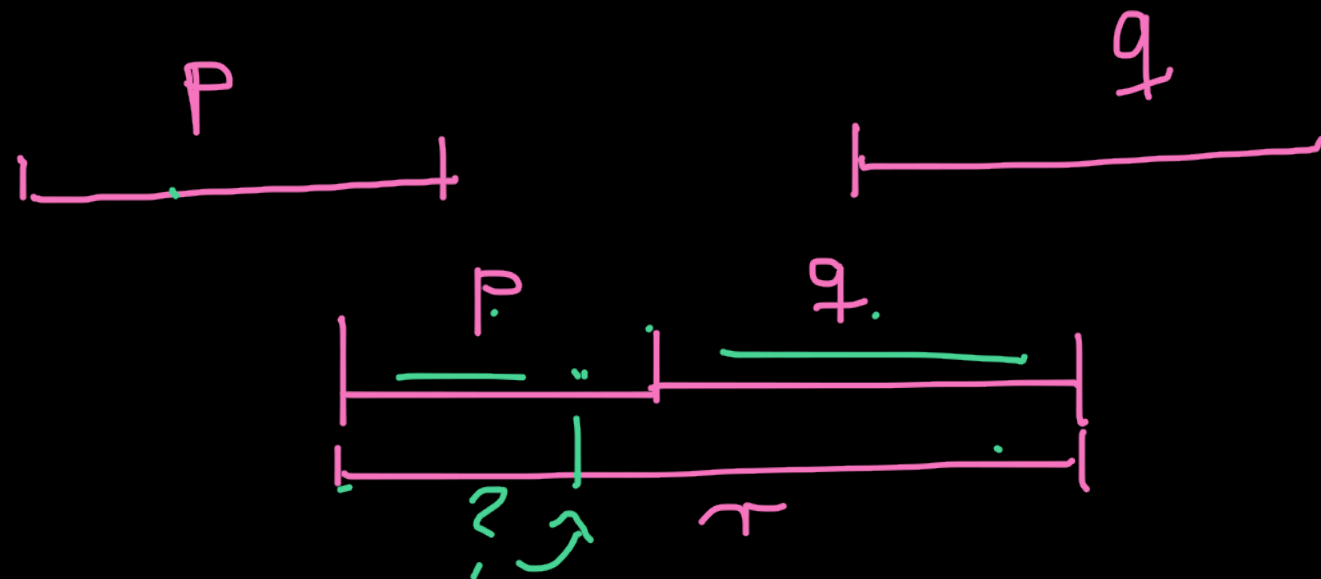
    $U(p) = \sigma$

    $M(p) = h(\sigma)$

How robust is *C* as an information measure?

► Do we have $C(\sigma, \tau) \leqslant^+ C(\sigma) + C(\tau)$?

$$H(X, Y)$$

$$C(\sigma, \tau) = C(\langle \sigma, \tau \rangle)$$

$$M(r) = \cdots$$

THM: [Martin-Löf]

Suppose $k$ is fixed. For any sufficiently long $\tau$ there exists $\sigma \sqsubseteq \tau$ such that $C(\sigma) < |\sigma| - k$.

Proof:

▶ Order all finite strings length-lexicographically, i.e.

$$\langle\rangle < 0 < 1 < 00 < 01 < 10 < 11 < 000 < 001 < \dots$$

and let $n(\sigma)$ be the position of $\sigma$ in this ordering.

$\sigma = 00$
$n(\sigma) = 3$

▶ Suppose $\vartheta \sqsubseteq \tau$. Let $n = n(\vartheta)$, and let $\rho$ be the next $n$ bits of $\tau$.

▶ Put $\alpha = \vartheta ^\frown \rho$. Then $C(\alpha) \leqslant |\rho| + c$ for some constant $c$.

▶ If we choose $|\vartheta| > k + c$, then

$$C(\alpha) \leqslant |\rho| + c = (|\alpha| - |\vartheta|) + c < |\alpha| - k.$$

$\vartheta$ $\rho$ $\tau$

$n$ bits

COR: For any $d$ there exists $\tau = \vartheta \frown \sigma$ such that

$$C(\tau) = C(\vartheta \frown \sigma) > C(\vartheta) + C(\sigma) + d.$$

$$C(\vartheta, \sigma) \geqslant C(\vartheta \frown \sigma)$$

Proof:

- ▶ Pick $c$ such that $C(\alpha) \leqslant |\alpha| + c$ for all $\alpha$.

- ▶ Choose a sufficiently long string $\tau$ with $C(\tau) \geqslant |\tau|$ and $C(\vartheta) < |\vartheta| - (c + d)$ for some $\vartheta \sqsubset \tau$ (by THM).

- ▶ Let $\sigma$ be such that $\tau = \vartheta \frown \sigma$.

- ▶ Then

$$C(\vartheta) + C(\sigma) < |\vartheta| - (c + d) + |\sigma| + c = |\tau| - d \leqslant C(\tau) - d.$$

11

What went wrong?

*We exploited that fact that a string not only provides information through its bits, but also through its length.*

This fact is not captured by *C*.

*Question: Can we alter the definition of complexity to take this into account?*

⟶ *Prefix-free complexity*

$K$