



P. Vitanyi, M&X Li, B. Ma

“Similarity Distance and Phylogeny”
aka “The Similarity Metric”

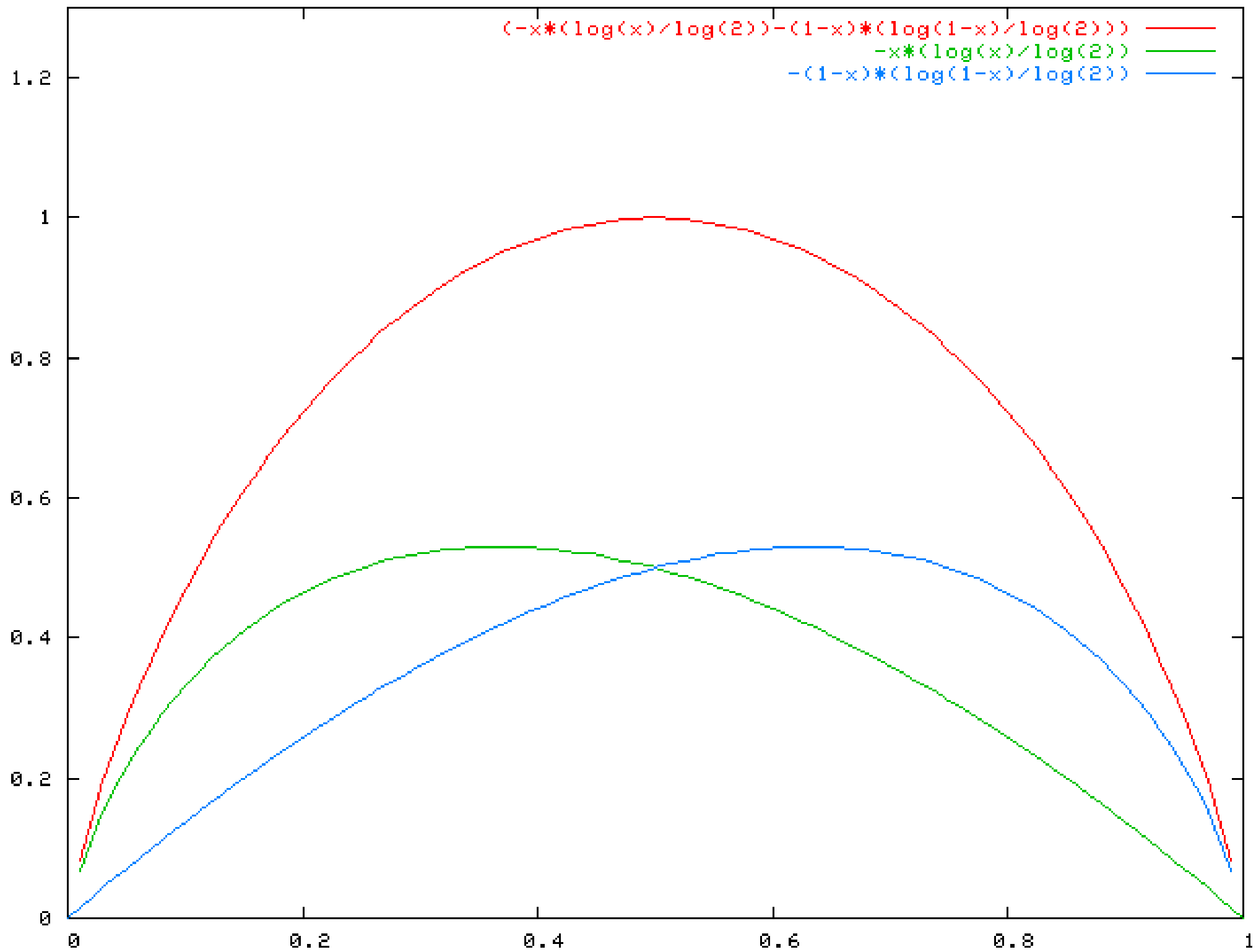
SIG-SIAM Discrete Algorithms '03

Sisukord

- ♦ Informatsioonisisaldus
- ♦ Kolmogorovi keerukus
- ♦ Tingimuslik keerukus
- ♦ Sarnasuse mõõtmine pakkimisega
- ♦ Näide: fülogeneesipuu
- ♦ Näide: keelte puu

Informatsioonisisaldus

- ♦ Infosisaldus
 - Sagedaste (tõenäoliste) sündmuste infosisaldus on väike
 - $I(s) = -\log(\text{Pr}(s))$
 - $I(s) \sim$ bittide arv sõnumi kodeerimiseks
- ♦ Entroopia
 - Vähetõenäolised sündmused esinevad harva, kuid nende infosisaldus on suur
 - $H = \text{SUM}[-\text{Pr}(s) * \log(\text{Pr}(s))]$



Kolmogorovi keerukus (1965)

- ♦ 123456789101112131415161718192..
 - Iga sümbol võrdse tõenäosusega
 - Kui palju on selles arvus informatsiooni?
 - Milline on tema vähim esitus?
 - For $i=1$ to n do print i ;
- ♦ Definiitsioon (Kolmogorovi keerukus)
 - Objekti (teksti, arvu, pildi, bitistringi) keerukus on lühima programmi (Turingi masina) pikkus, mis väljastab algse objekti.

K(x) arvutamine

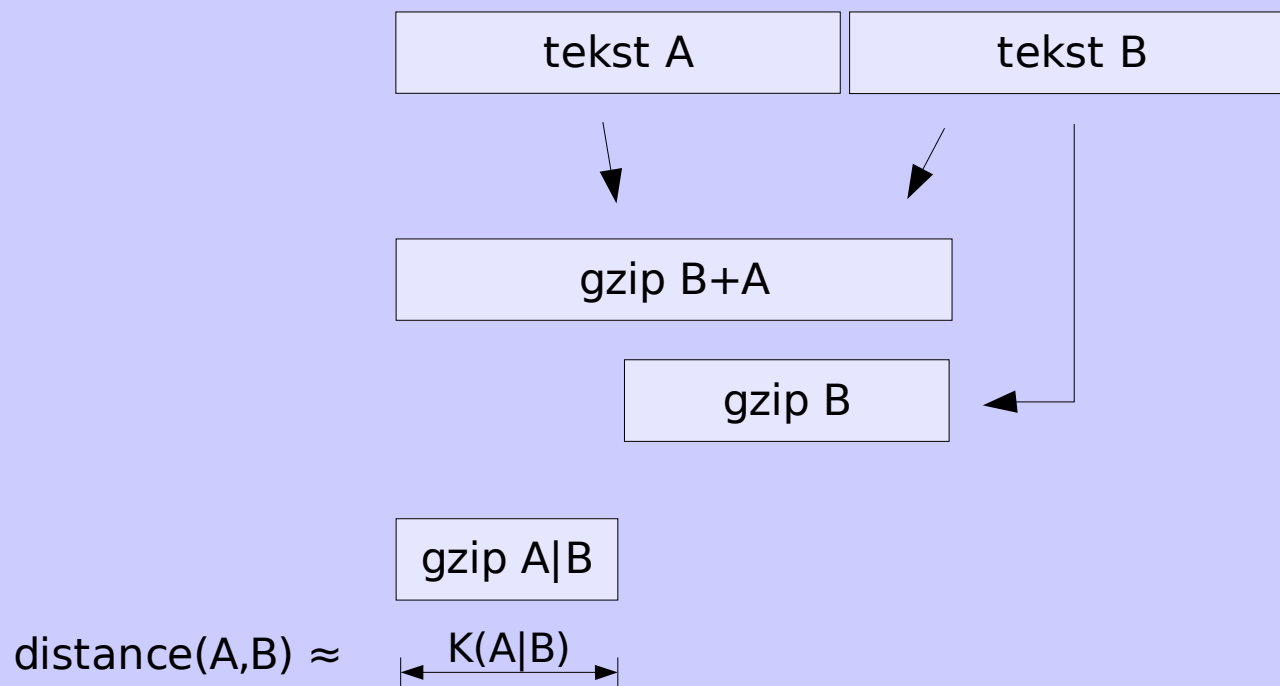
- ♦ $K(x)$ ei ole arvutatav !!!
 - s.t ei leidu algoritmi, mis suudab iga x jaoks leida selle Kolmogorovi keerukuse (lühima programmi, mis väljastab x).
- ♦ $K(x)$ on hinnatav
 - $K(x) \leq x$ kahendesitus
 - $K(x+1) \leq K(x) + c$
- ♦ $K(x)$ on lähendatav
 - $K(x) \leq \text{length}(\text{gzip}) + \text{length}(x.\text{zip})$

Tinglik keerukus

- ♦ $K(x|y)$ = lühima programmi pikkus, mis väljastab x , kui y on teada.
 - Lühim programm, mis transformeerib $y \rightarrow x$
 - Mõiste: transformatsiooni keerukus.
 - $K(x,y) \leq K(y) + K(x|y)$
 - $K(x|y) \leq K(x,y) - K(y)$
- ♦ Sarnasuse mõõtmine
 - Idee: y on parajasti nii sarnane x -le, kui lihtne on x -st saada y .
 - $distance(x,y) \sim K(x|y)$

Sarnasus pakkimise kaudu

- $K(x|y) \approx K(x,y) - K(y)$
- $length(\text{gzip } "y+x") - length(\text{gzip } "y")$



Evolutsioonipuu

- ♦ Levinud meetodid
 - K-gon (K-mer) statistik
 - Geenide järjekord genoomis
 - Genoomi ühisosa (ühiste geenide arv)
 - Ümberjärjestuskaugus
 - Transformatsioonikaugus
- ♦ Pakkimiskaugus (Vitanyi jt.)
 - $d(A,B) \approx [K(x|y) + K(y|x)] / [K(xy)]$

