

Fast evaluation and root finding for polynomials with floating-point coefficients

Guillaume Moroz, Rémi Imbach

9 mars 2023

Floating-point arithmetic

Representation

Light-year: 9 460 730 472 580 800 m

$$\underbrace{9.460}_{\text{mantissa}} \cdot 10^{\underbrace{15}_{\text{exponent}}} \text{ m}$$

mantissa **exponent**

number of digits $< m$

absolute value $< \tau$

Polynomial evaluation

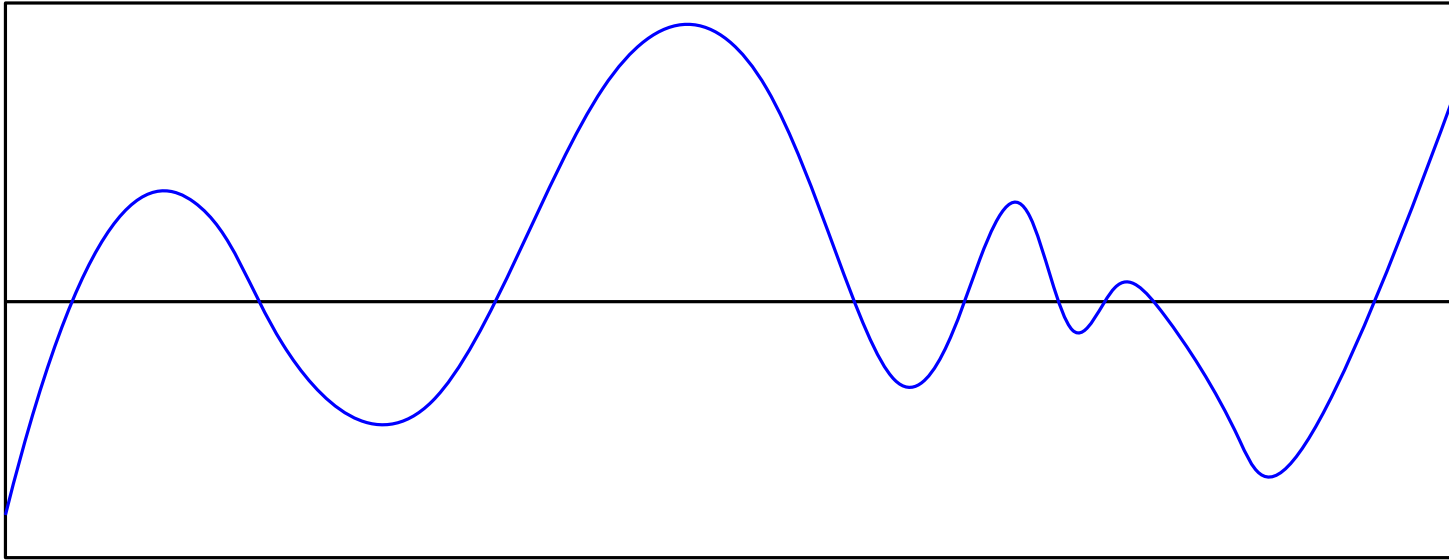
$$f(z) = f_0 + \cdots + f_d z^d$$

- Bit complexity : $\tilde{O}(d(m + \log \tau))$

- Error: $O(d2^{-m})\tilde{f}(|z|)$ where $\tilde{f}(|z|) = \sum |f_j||z|^j$

Main result

$$f(z) = f_0 + \cdots + f_d z^d$$



Theorem (Piecewise polynomial approximation)

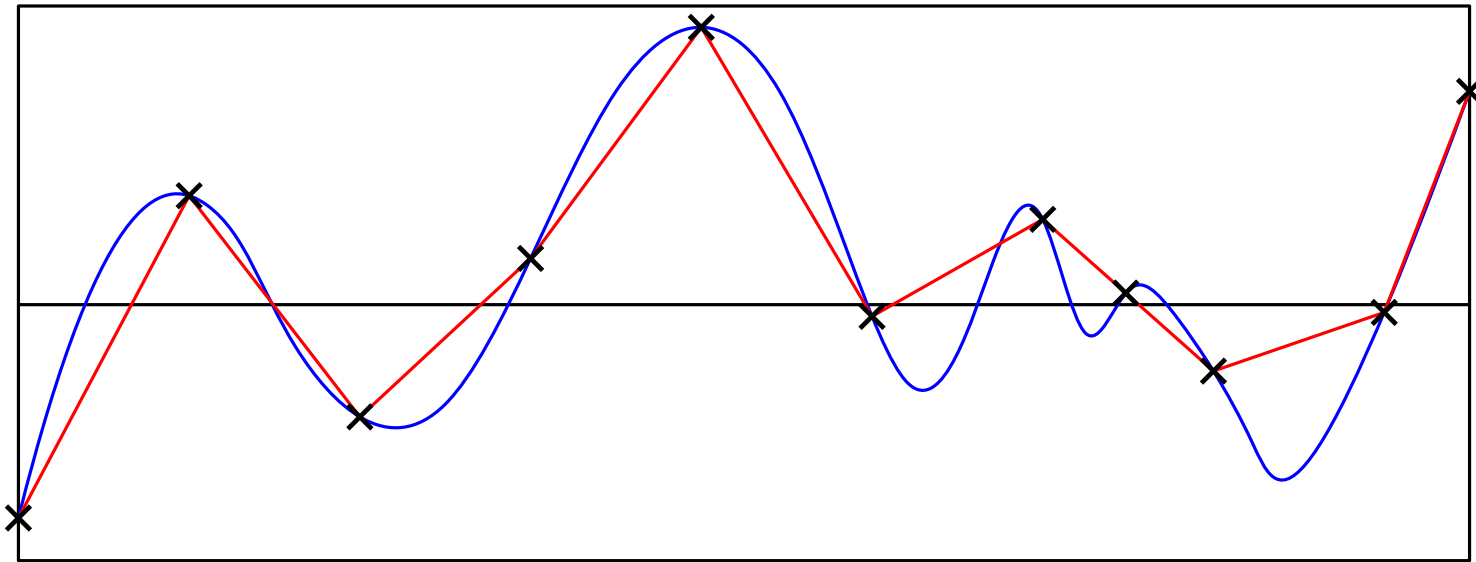
It is possible to compute all g_n of degree m satisfying

$$|f(z) - g_n(z)| < 2^{-m} \tilde{f}(|z|)$$

in $\tilde{O}(d(m + \log \tau))$ bit operations

Main result

$$f(z) = f_0 + \cdots + f_d z^d$$



Theorem (Piecewise polynomial approximation)

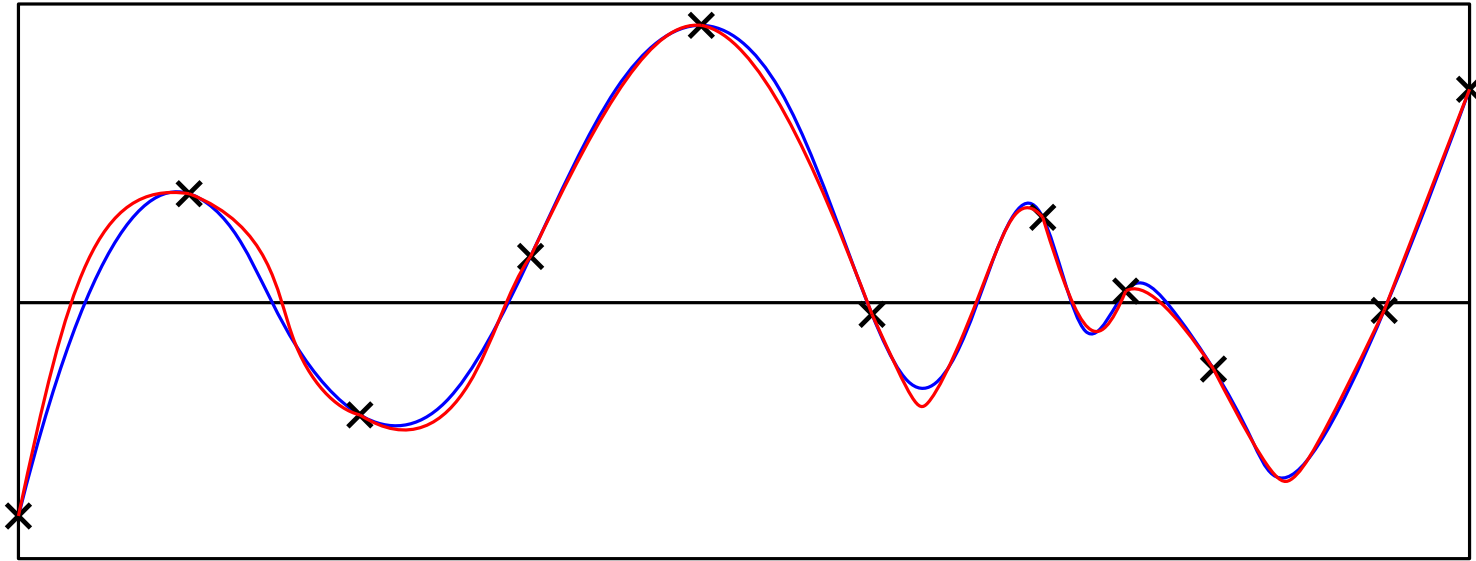
It is possible to compute all g_n of degree m satisfying

$$|f(z) - g_n(z)| < 2^{-m} \tilde{f}(|z|)$$

in $\tilde{O}(d(m + \log \tau))$ bit operations

Main result

$$f(z) = f_0 + \cdots + f_d z^d$$



Theorem (Piecewise polynomial approximation)

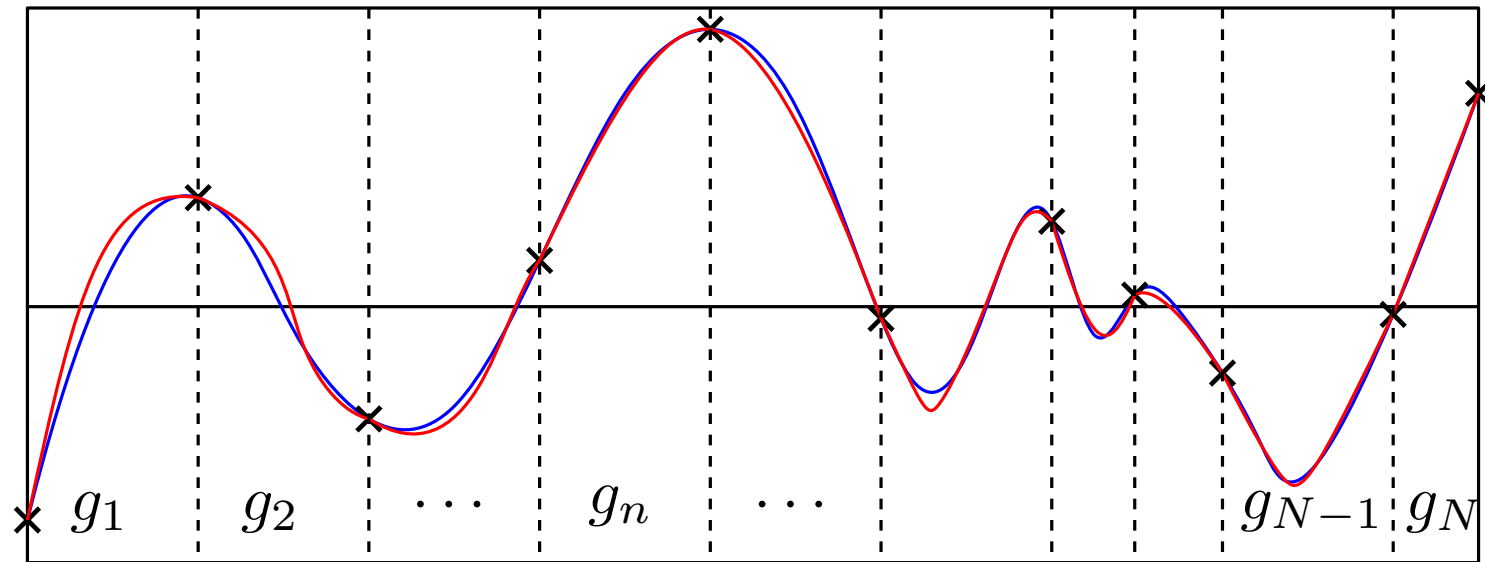
It is possible to compute all g_n of degree m satisfying

$$|f(z) - g_n(z)| < 2^{-m} \tilde{f}(|z|)$$

in $\tilde{O}(d(m + \log \tau))$ bit operations

Main result

$$f(z) = f_0 + \cdots + f_d z^d$$



Theorem (Piecewise polynomial approximation)

It is possible to compute all g_n of degree m satisfying

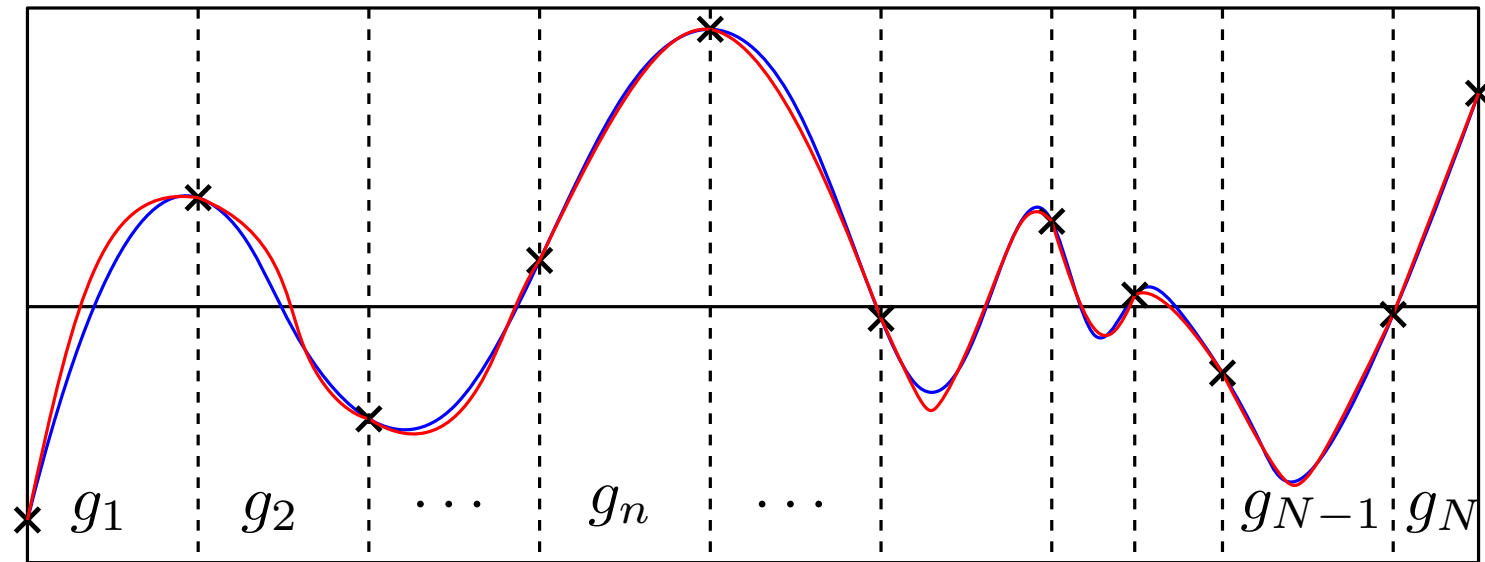
$$|f(z) - g_n(z)| < 2^{-m} \tilde{f}(|z|)$$

in $\tilde{O}(d(m + \log \tau))$ bit operations

Corollary

$$f(z) = f_0 + \cdots + f_d z^d$$

$$m > \log d$$



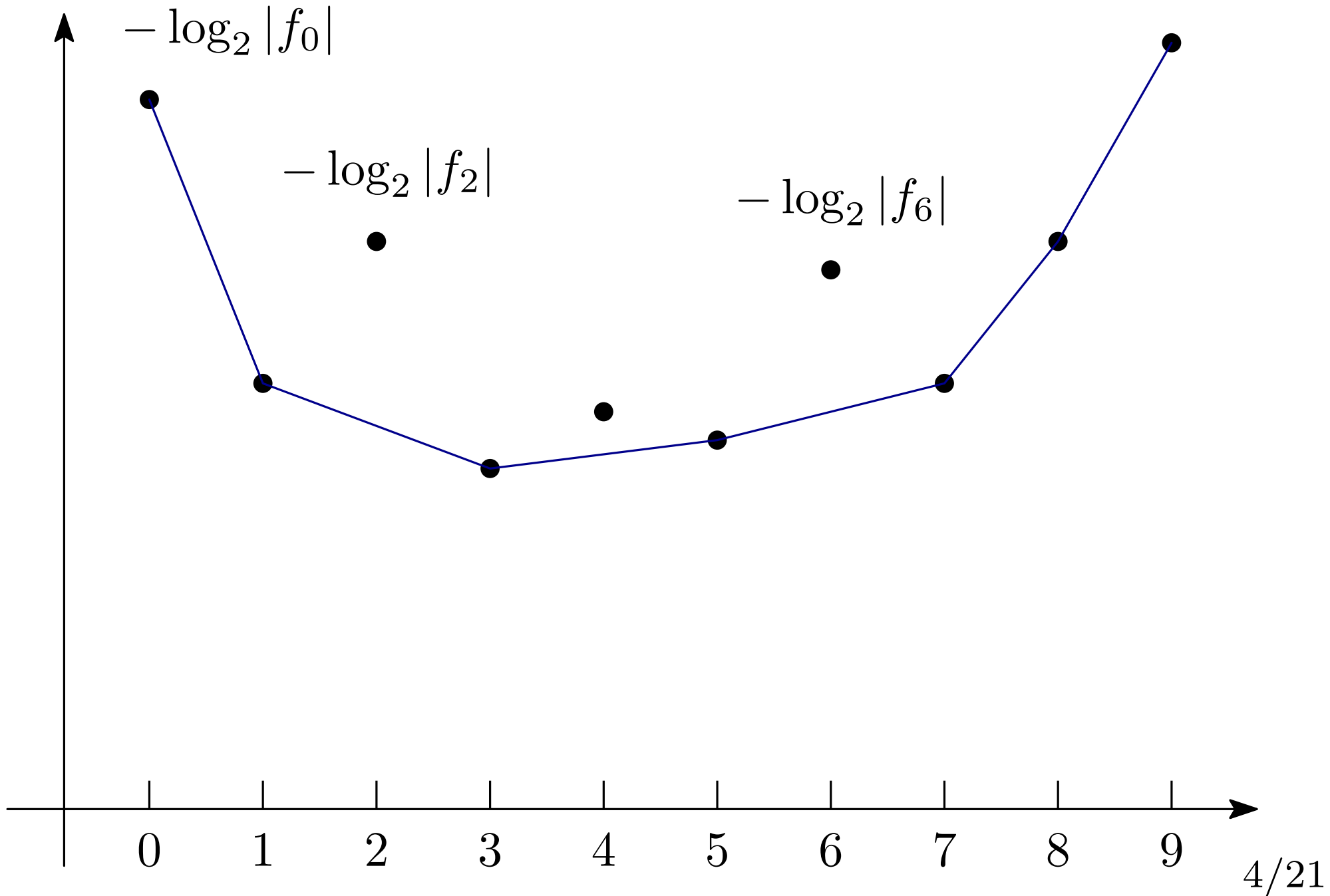
Fast evaluation

With error $2^{-m} \tilde{f}(|z|)$
in
 $\tilde{O}(m(m + \log \tau))$

Fast root finding

All roots with m bits
in
 $\tilde{O}(d(m + \log \tau + \log \kappa))$

Newton polygon



Newton polygon

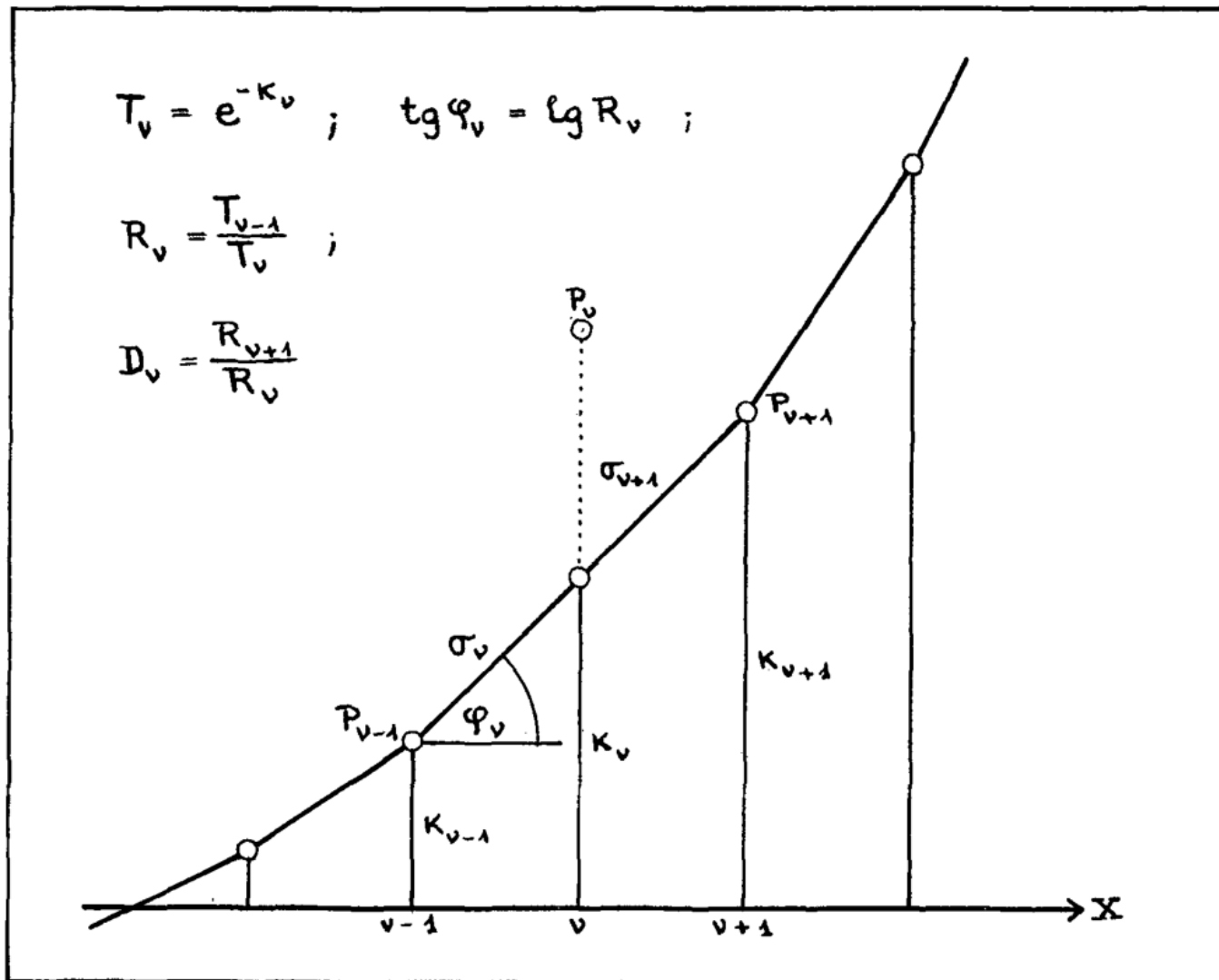
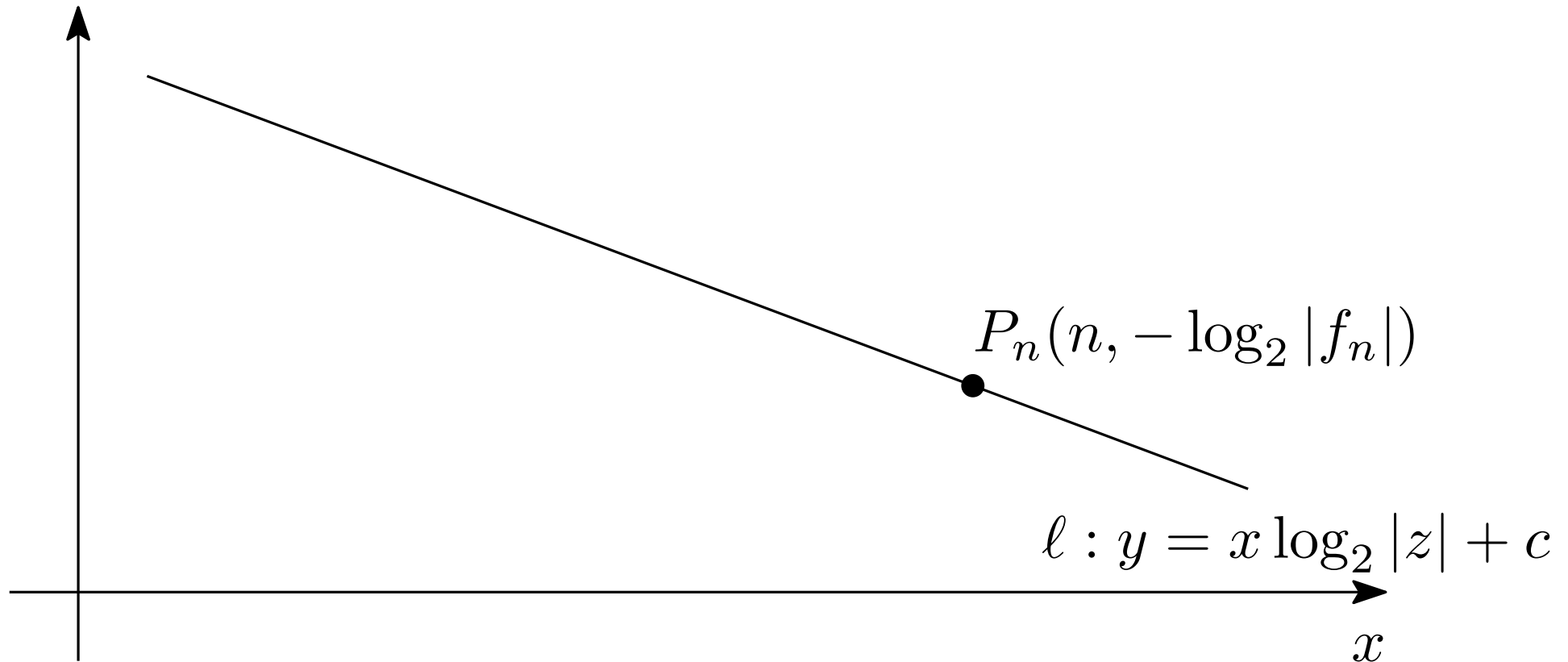


Fig. 1.

Alexandre Ostrowski. *Recherches sur la méthode de Graeffe et les zéros des polynômes et des séries de Laurent*. Acta Mathematica, 1940

Newton polygon

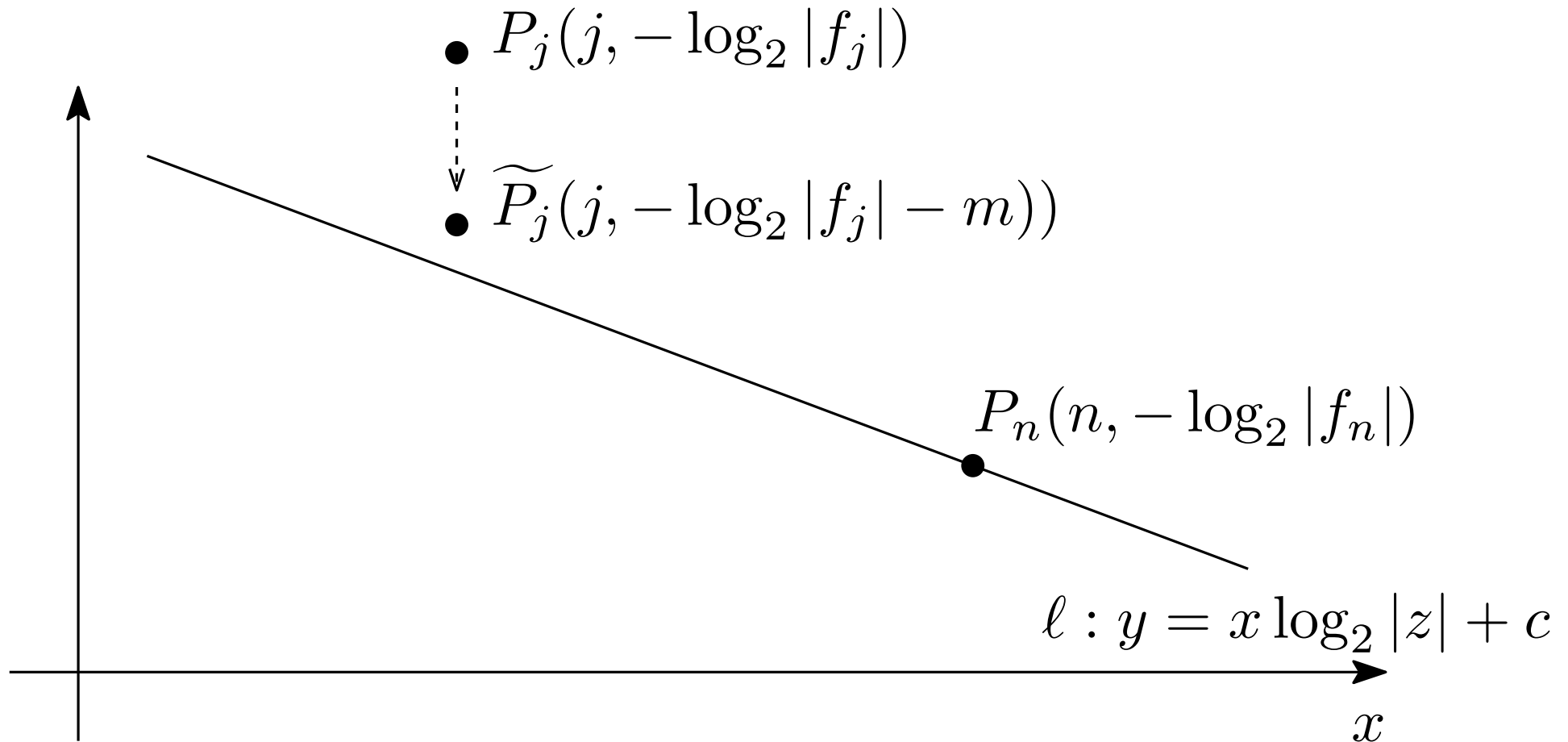
- $P_j(j, -\log_2 |f_j|)$



Dominant monomial

$$P_j \text{ above } \ell \iff |f_j| |z|^j < |f_n| |z|^n$$

Newton polygon

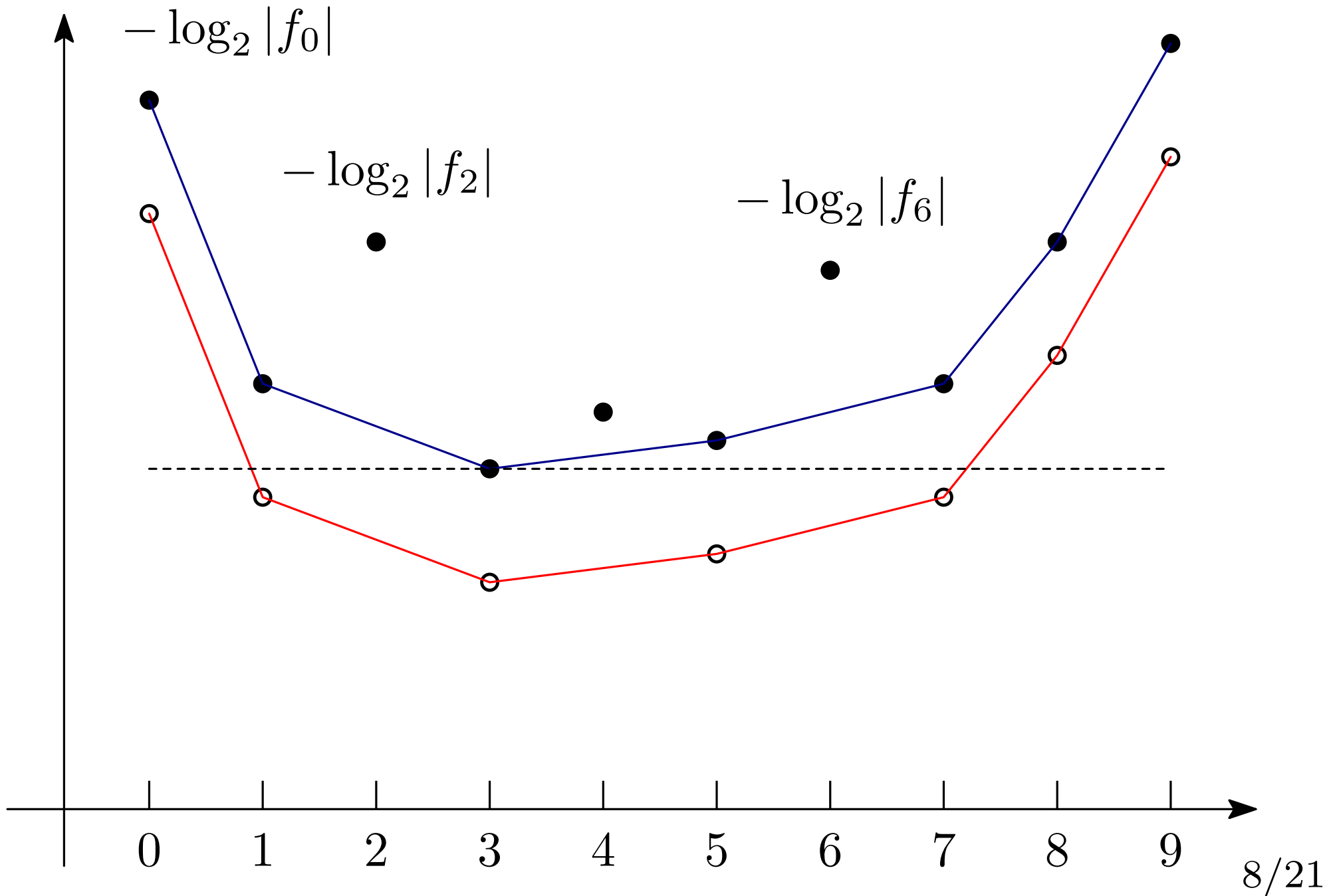


Dominant monomial

$$P_j \text{ above } \ell \iff |f_j| |z|^j < |f_n| |z|^n$$

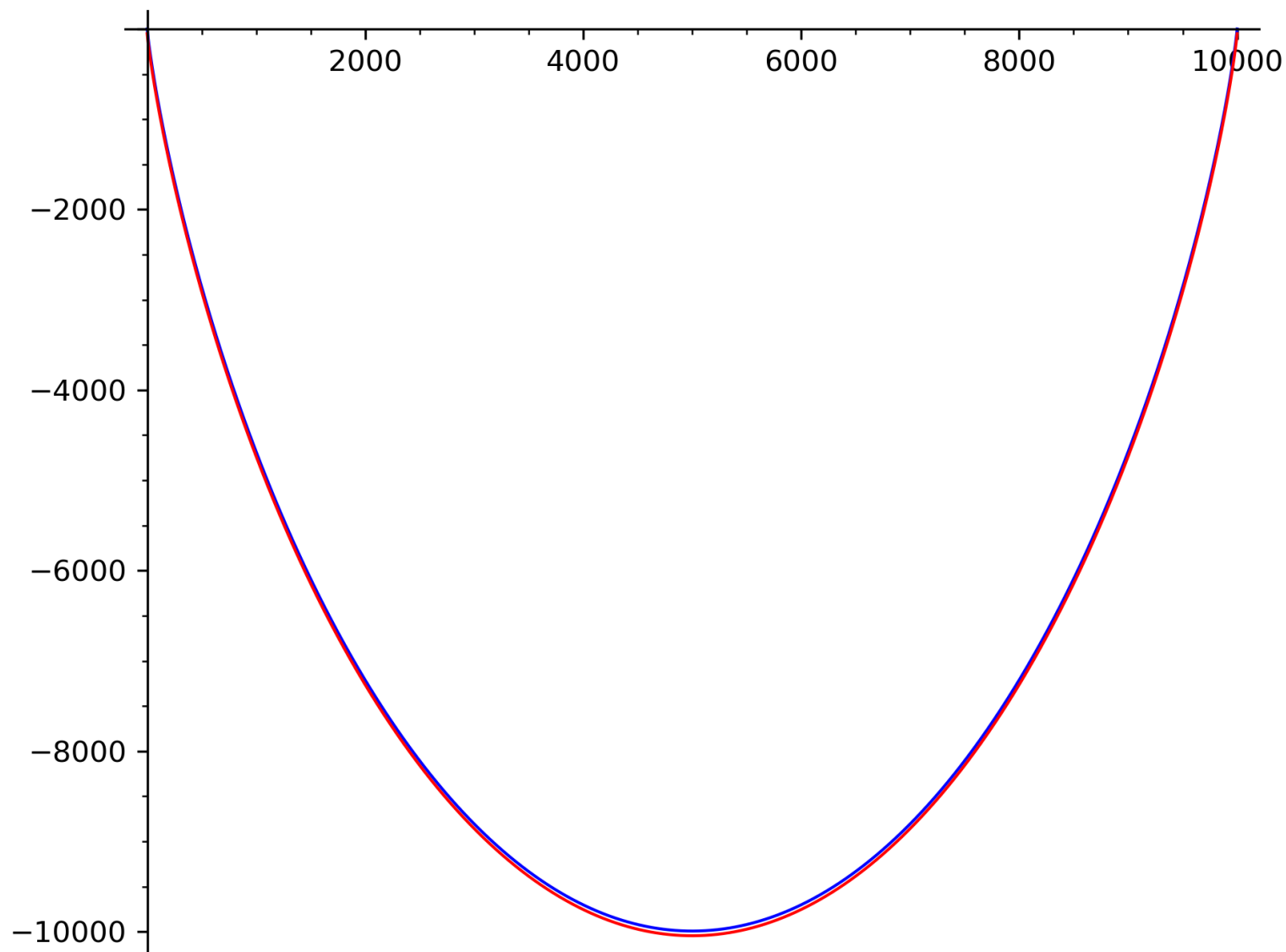
$$\widetilde{P}_j \text{ above } \ell \iff |f_j| |z|^j < 2^{-m} |f_n| |z|^n$$

Newton polygon



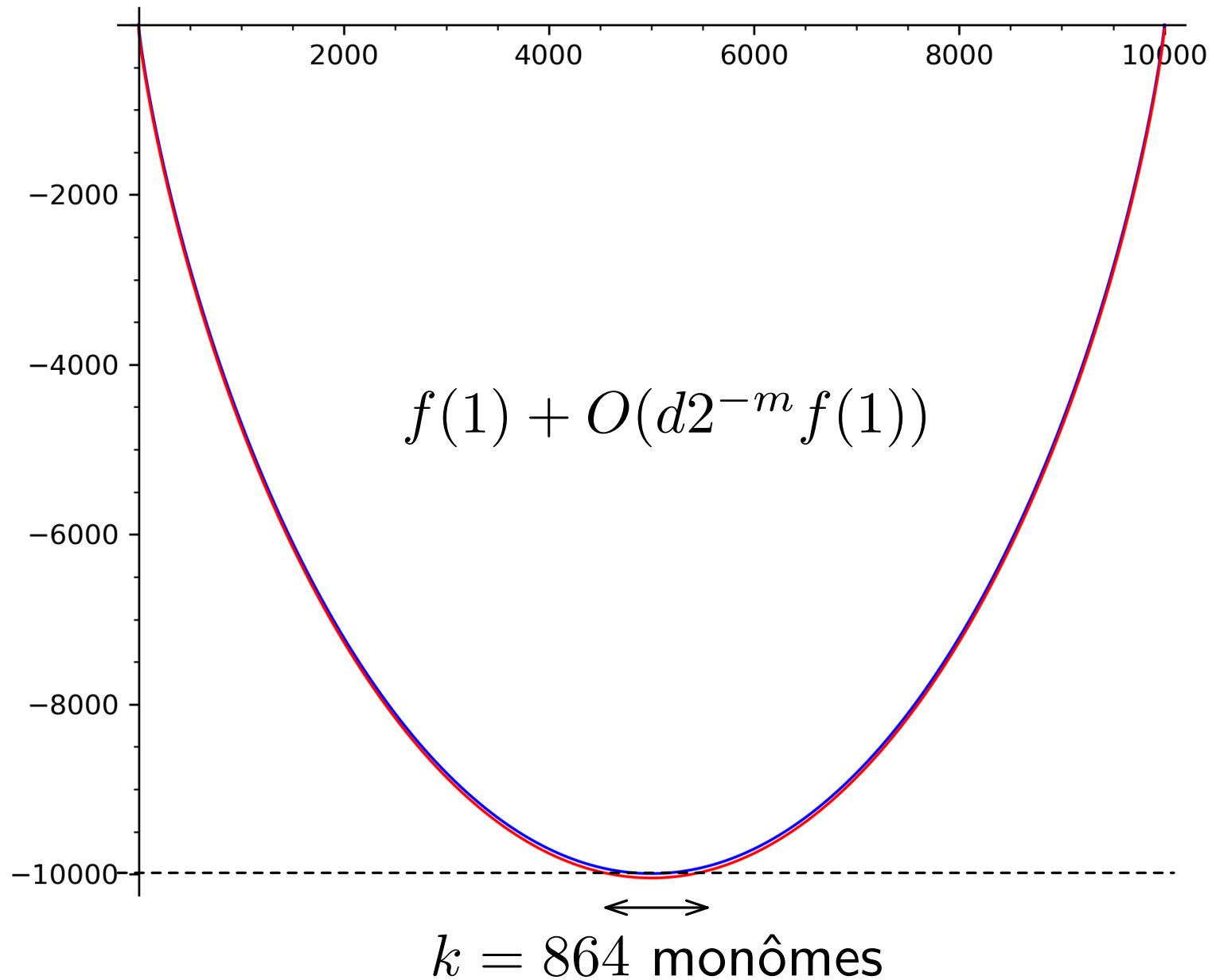
Example: selecting monomials

$$f(z) = (1 + z)^{10000} = \sum \binom{10000}{j} z^j$$



Example: selecting monomials

$$f(z) = (1+z)^{10000} = \sum \binom{10000}{j} z^j$$



Example: Taylor expansion

1

10 000z

⋮

$\binom{10\,000}{\ell} z^\ell$

$\binom{10\,000}{j} z^j$

$\binom{10\,000}{u} z^u$

⋮

$z^{10\,000}$

Example: Taylor expansion

$$1 \quad (1 + \varepsilon t)^k = 1 + \dots + \binom{k}{m} \varepsilon^m t^m + \mathcal{O}\left(\frac{k}{m} \varepsilon\right)^m$$

10 000z

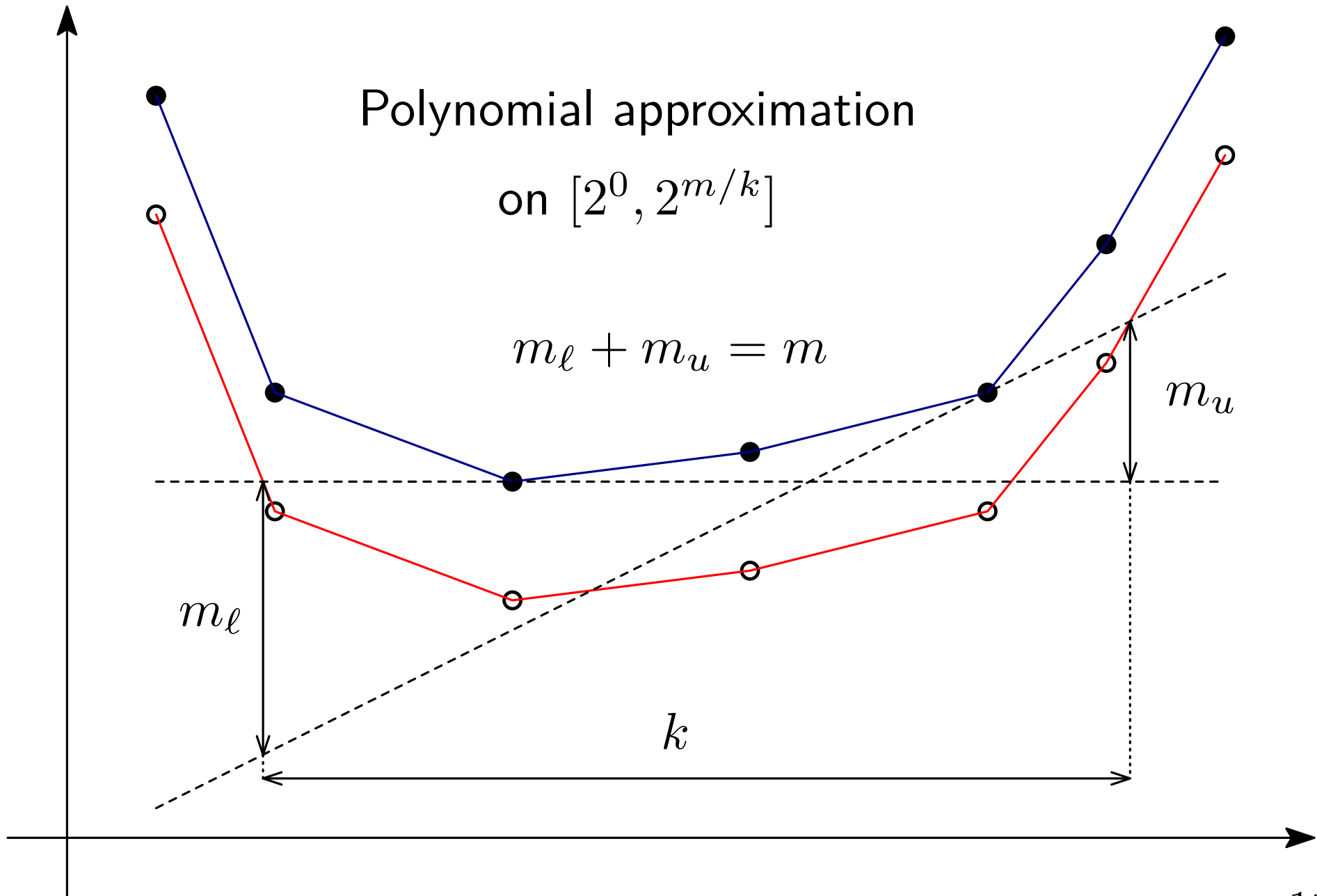
⋮

$$k \left\{ \begin{array}{l} \binom{10\,000}{\ell} z^\ell = \binom{10\,000}{\ell} z^\ell \\ \binom{10\,000}{j} z^j = \binom{10\,000}{j} z^\ell (1 + \varepsilon t)^{j-\ell} \\ \binom{10\,000}{u} z^u = \binom{10\,000}{u} z^\ell (1 + \varepsilon t)^{u-\ell} \end{array} \right.$$

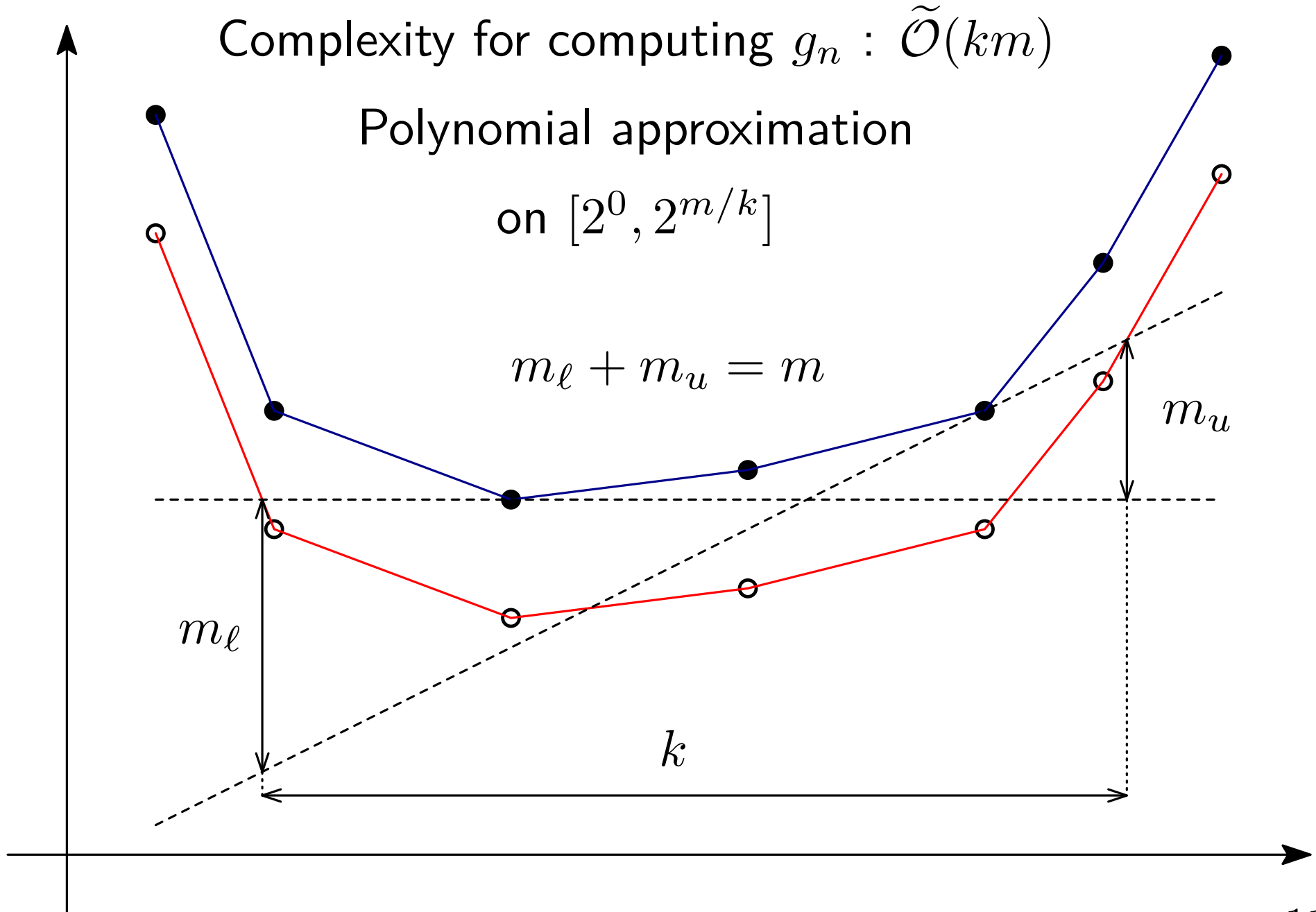
⋮

$z^{10\,000}$

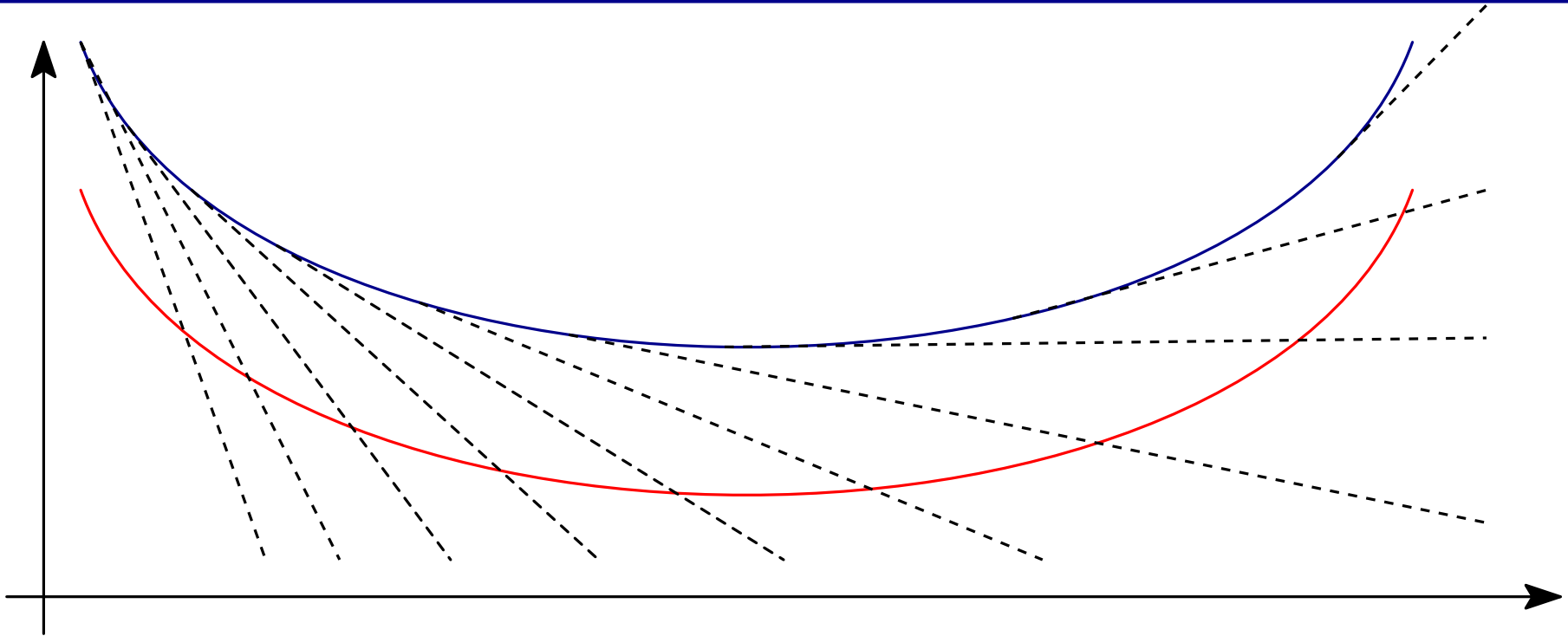
Algorithm: one interval



Algorithm: one interval



Algorithm: all intervals



Approximation algorithm

WHILE $a_j < a_{max}$

$$a_{j+1} = a_j + \frac{m}{k_j}$$

$g_{j+1} =$ Taylor approximation of degree $\mathcal{O}(m)$
on the interval $[2^{a_j}, 2^{a_{j+1}}]$

Complexity

Theorem

The piecewise approximation algorithm costs

$$\tilde{\mathcal{O}}\left(\sum k_j m\right) \in \tilde{\mathcal{O}}(dm) \text{ bit-operations}$$

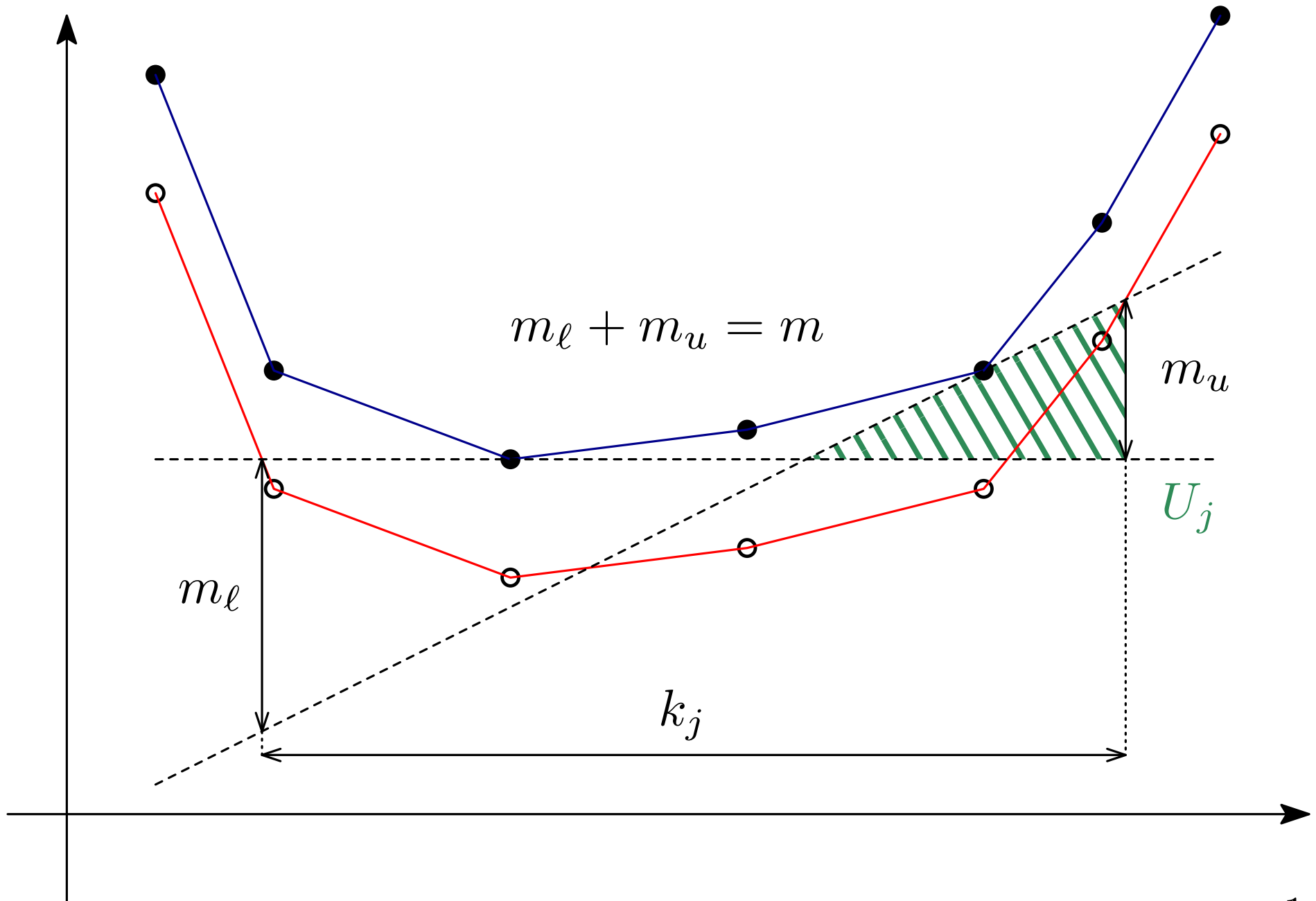
Approximation algorithm

WHILE $a_j < a_{max}$

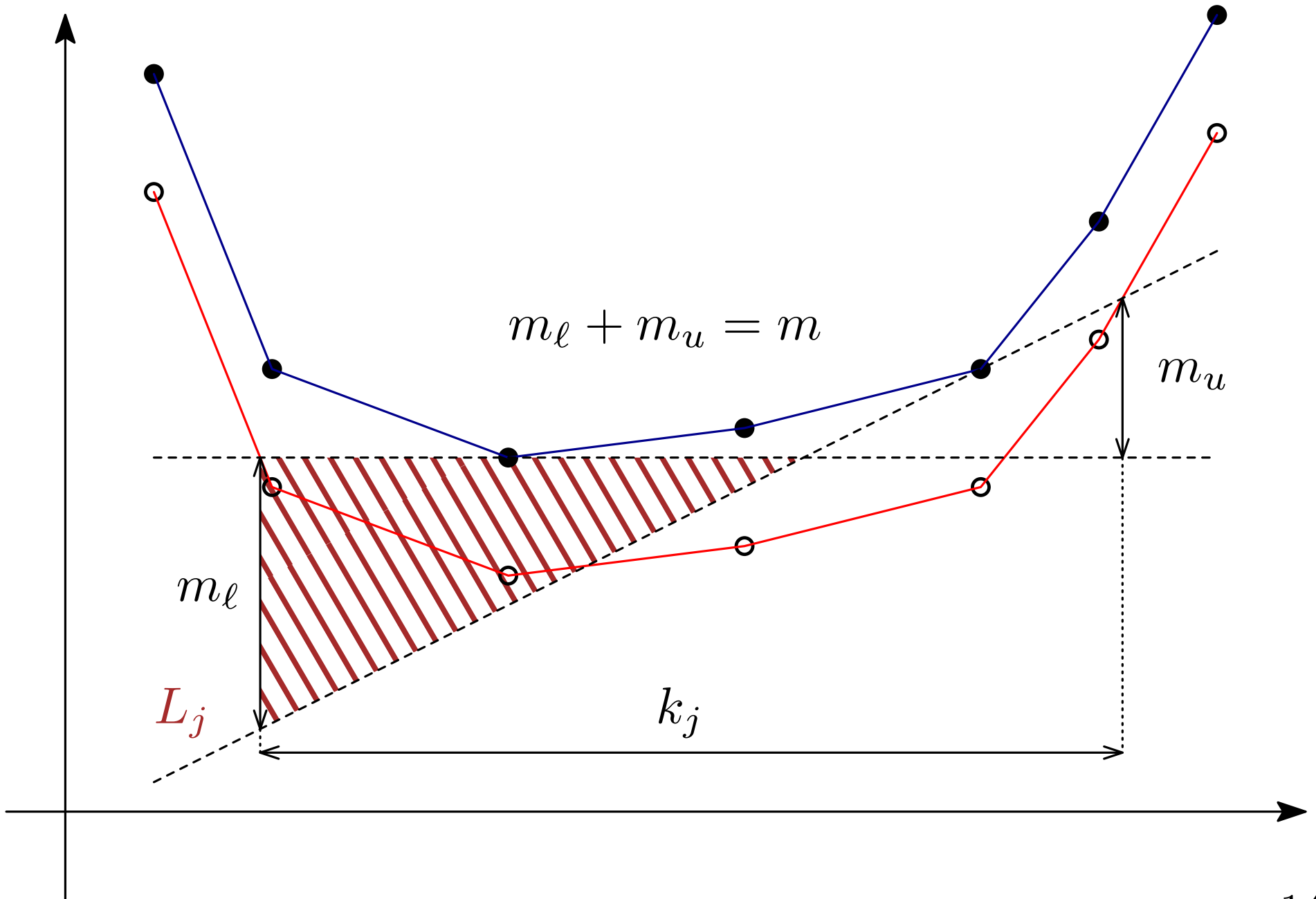
$$a_{j+1} = a_j + \frac{m}{k_j}$$

$g_{j+1} =$ Taylor approximation of degree $\mathcal{O}(m)$
on the interval $[2^{a_j}, 2^{a_{j+1}}]$

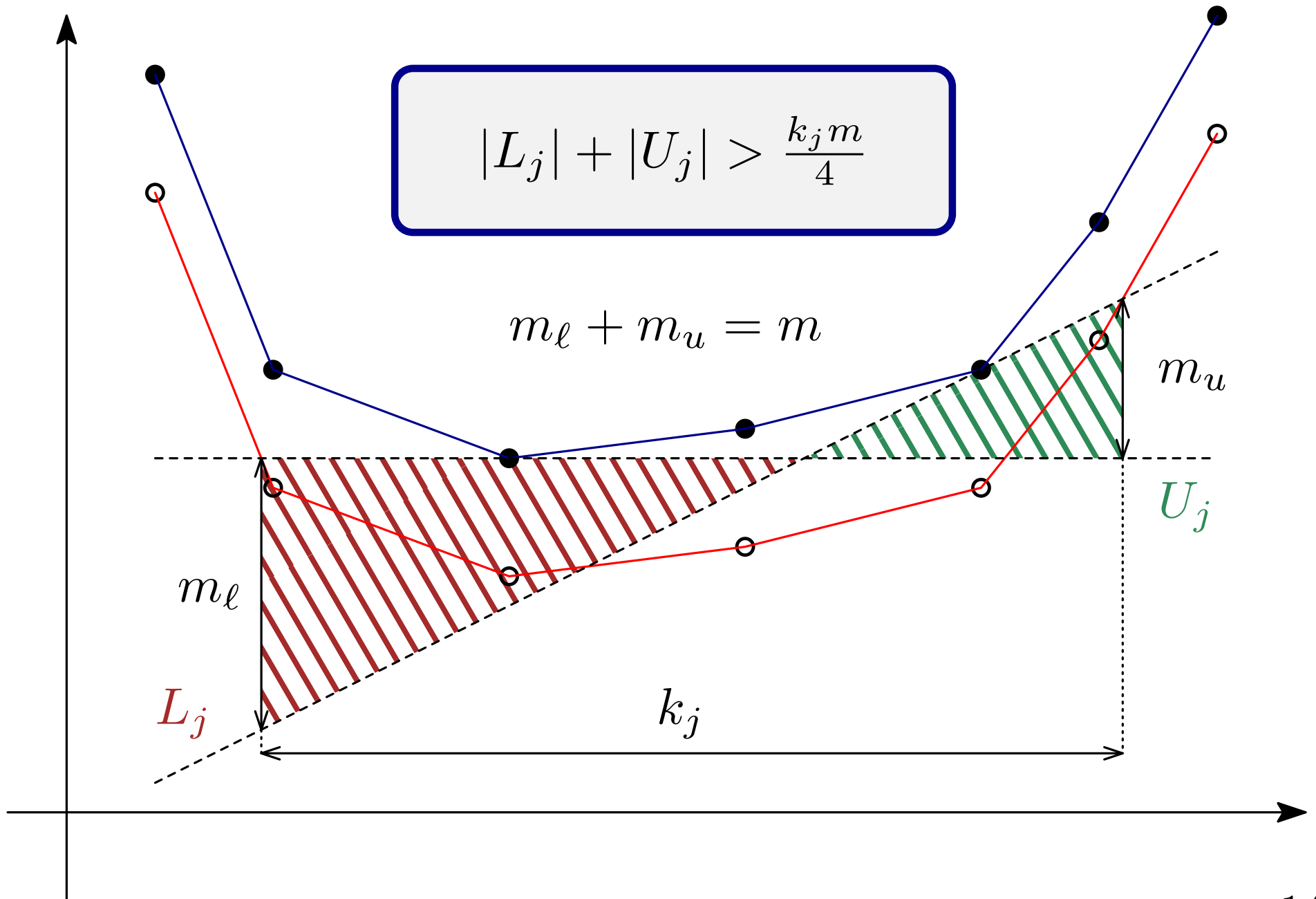
Proof



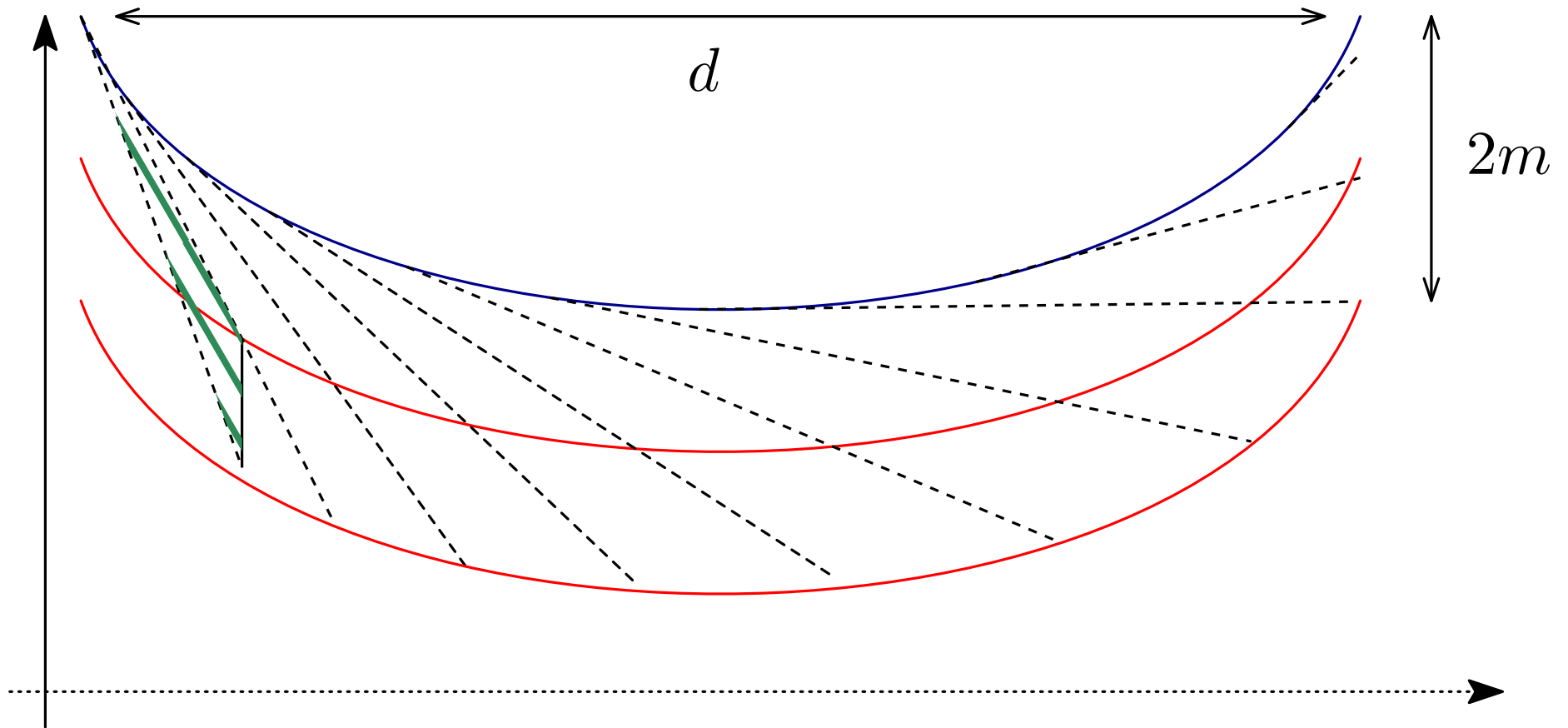
Proof



Proof



Proof

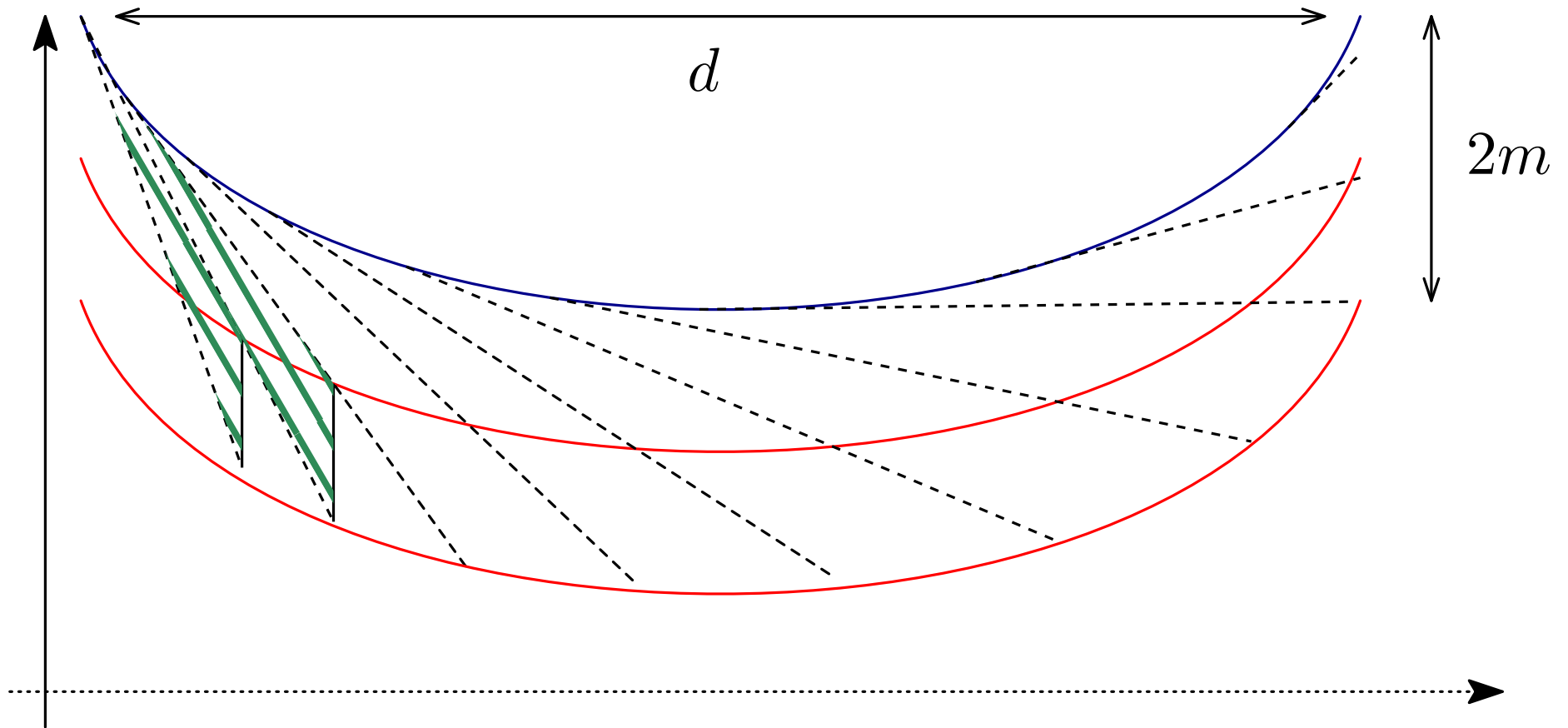


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

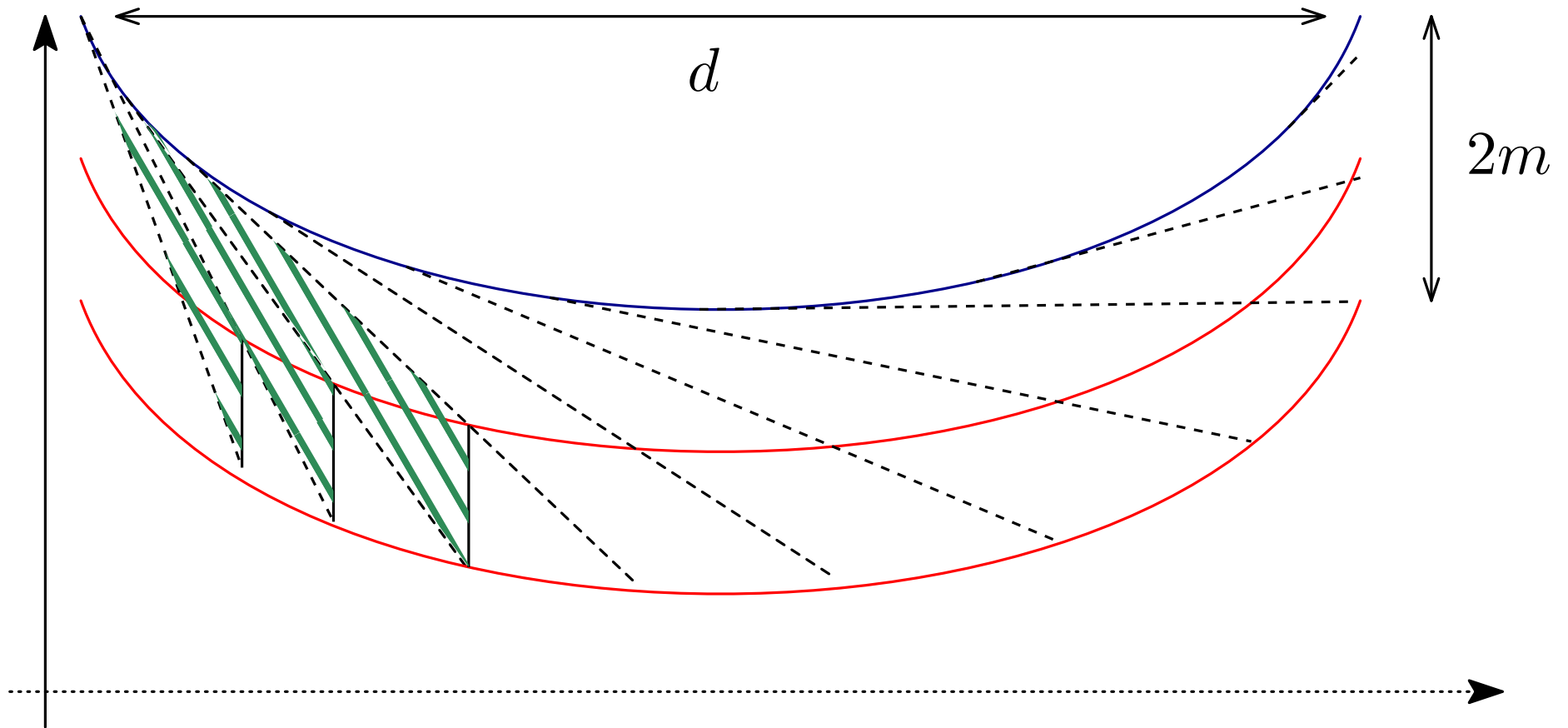


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

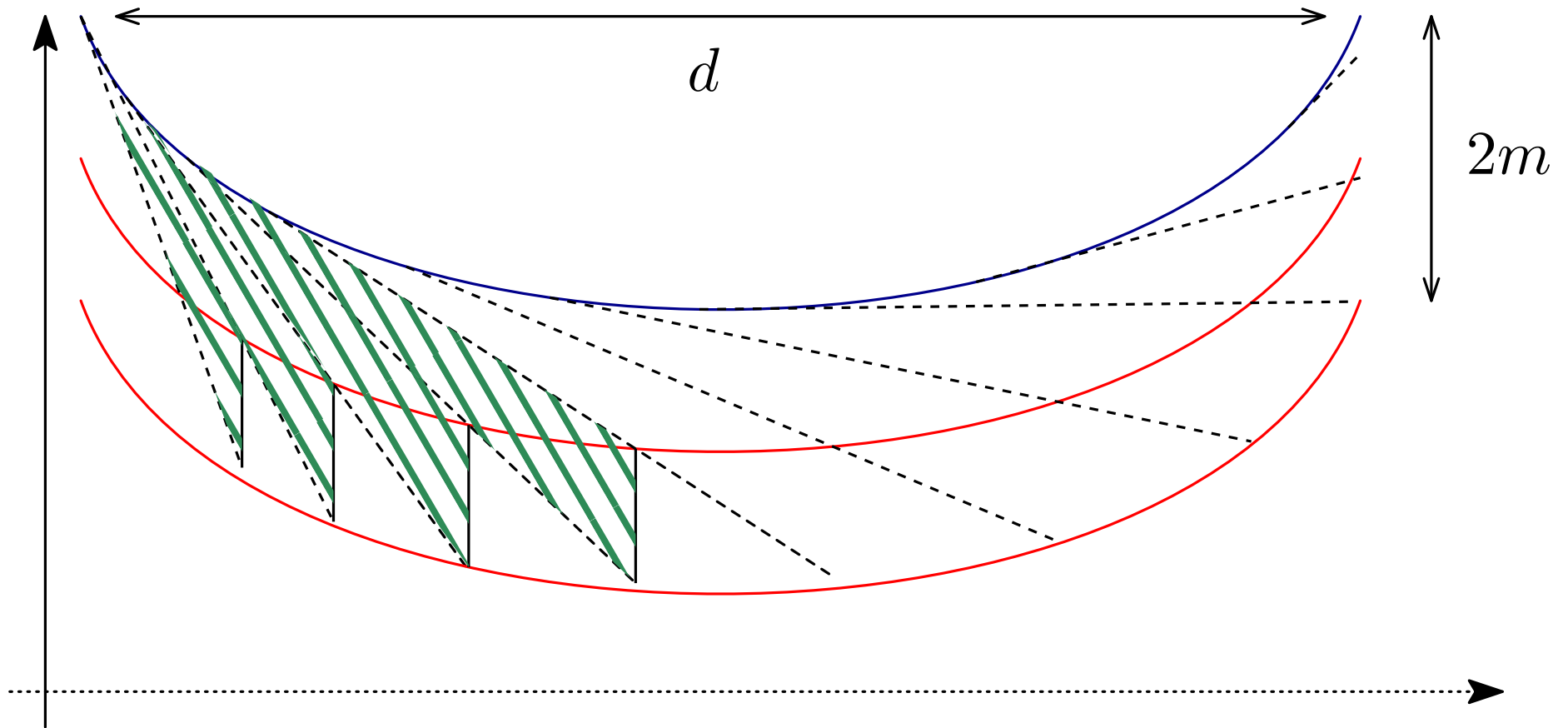


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

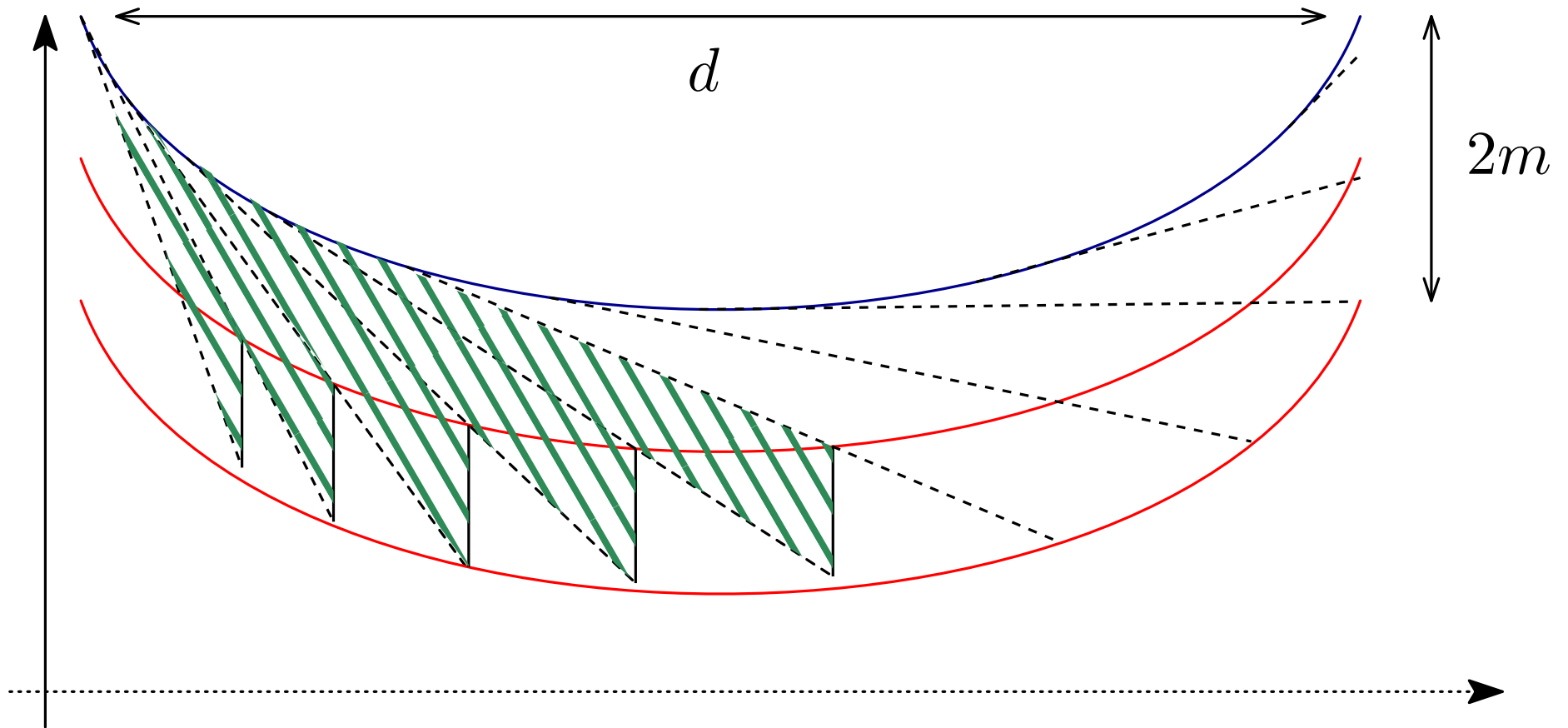


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

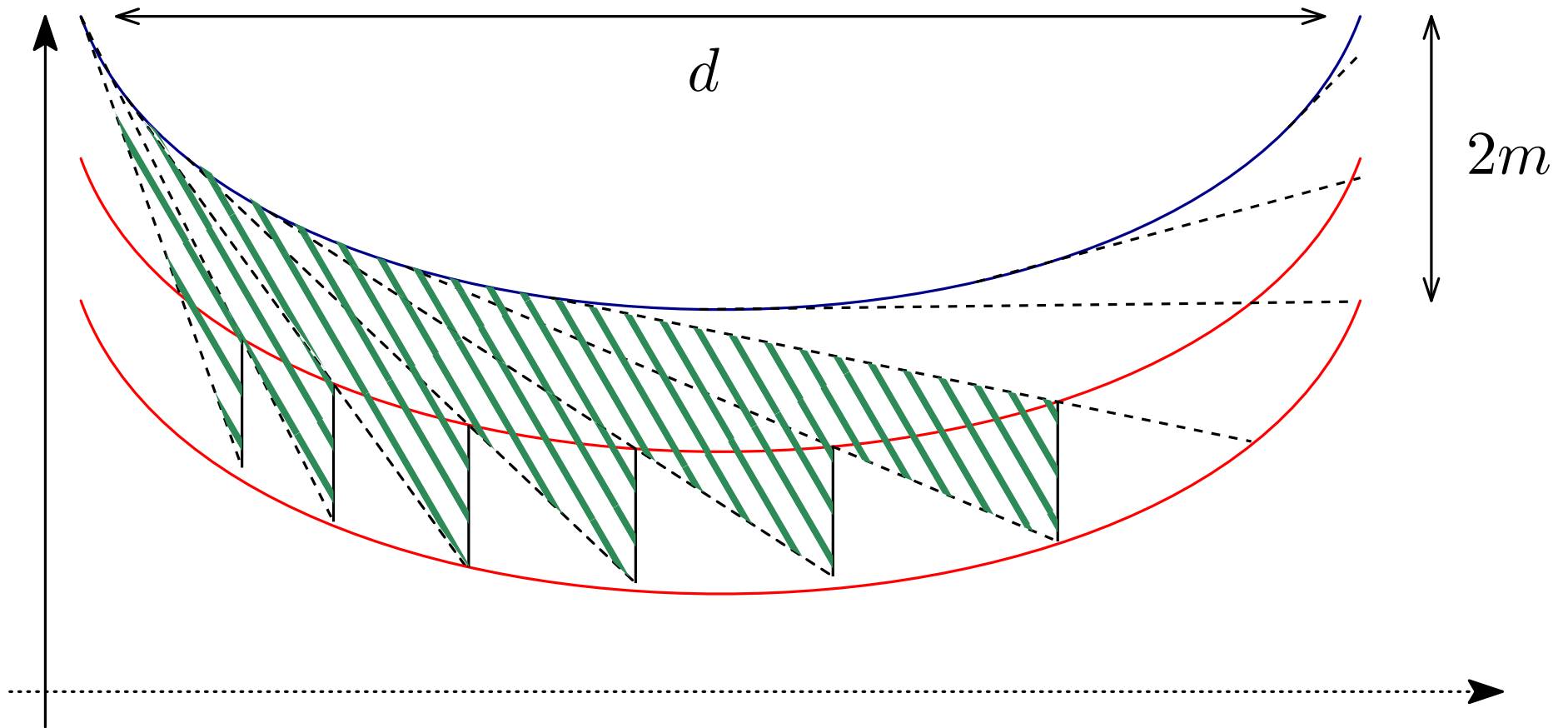


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

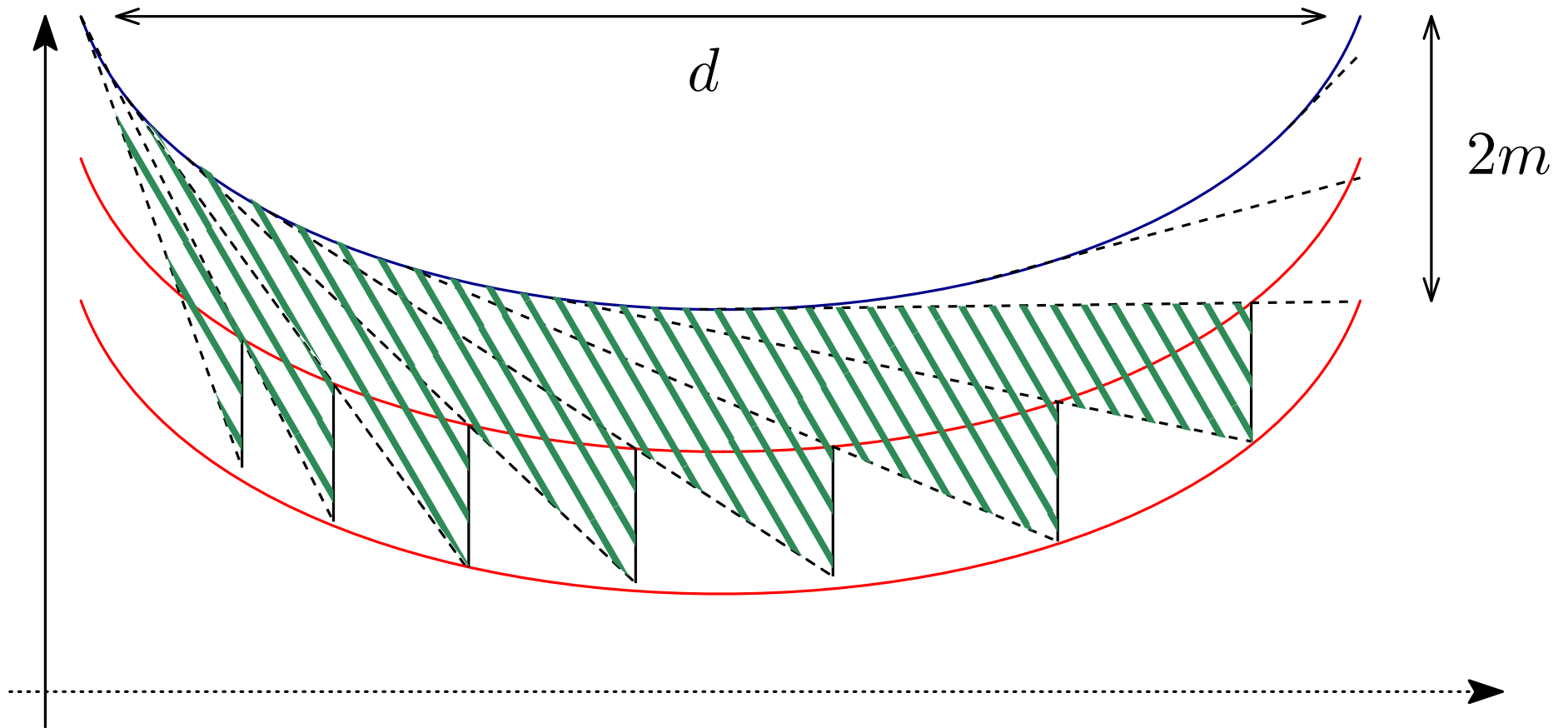


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

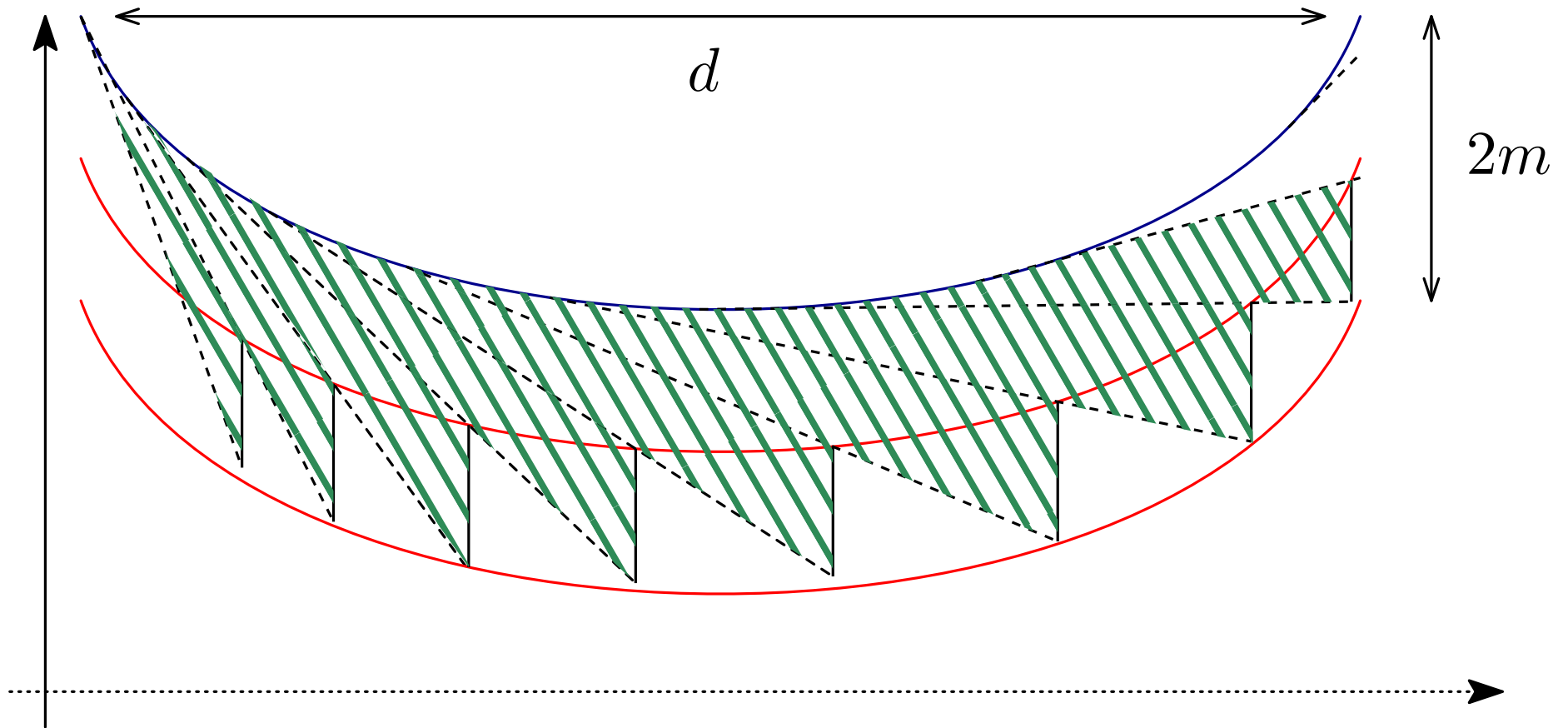


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

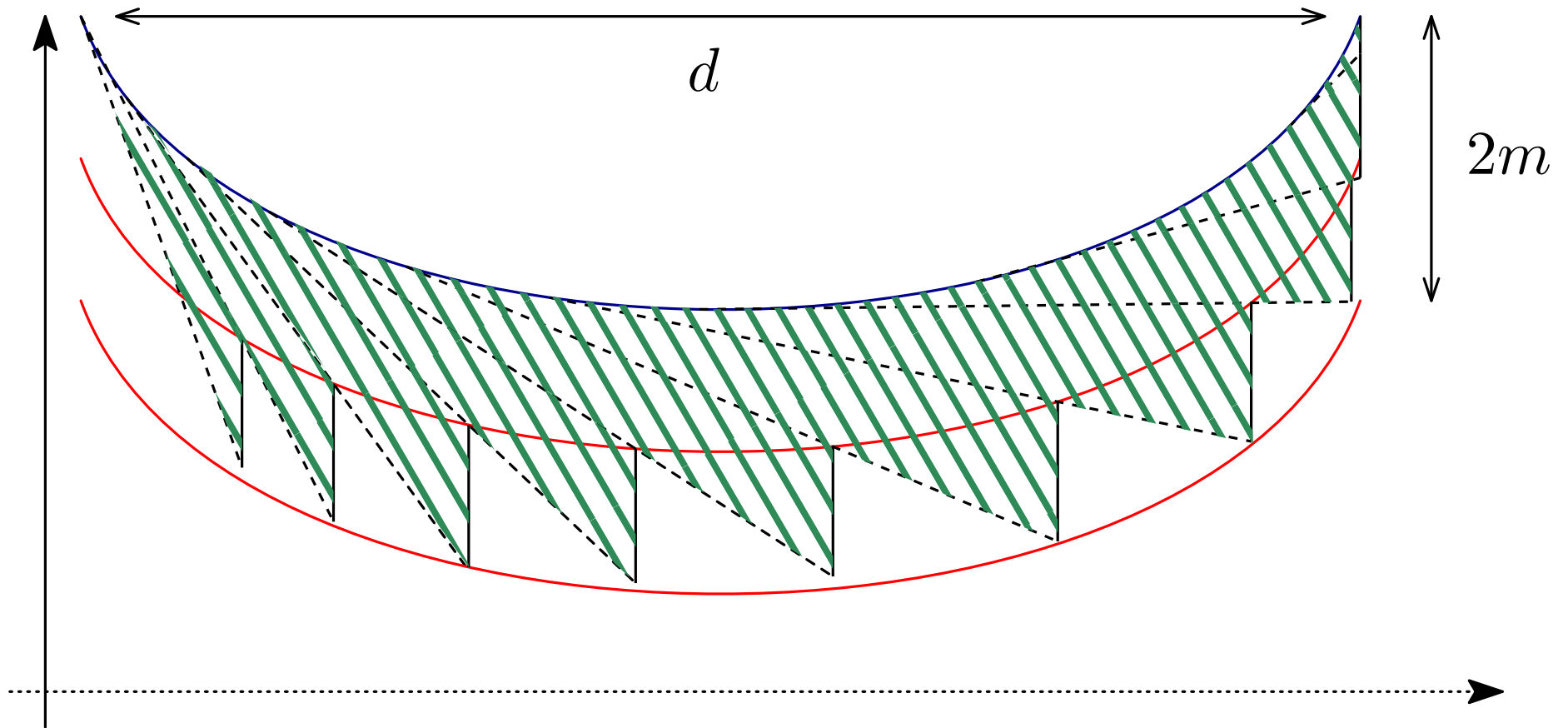


U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof



U_j piecewise disjoint

Union of U_j in B_{2m}

$$\sum |U_j| < 2dm$$

Proof

$$|L_j| + |U_j| > \frac{k_j m}{4}$$

$$\sum |L_j| < 2dm$$

$$\sum |U_j| < 2dm$$

Main lemma

$$\sum \frac{k_j m}{4} \leq \sum |L_j| + \sum |U_j| \leq 4dm$$

Benchmark

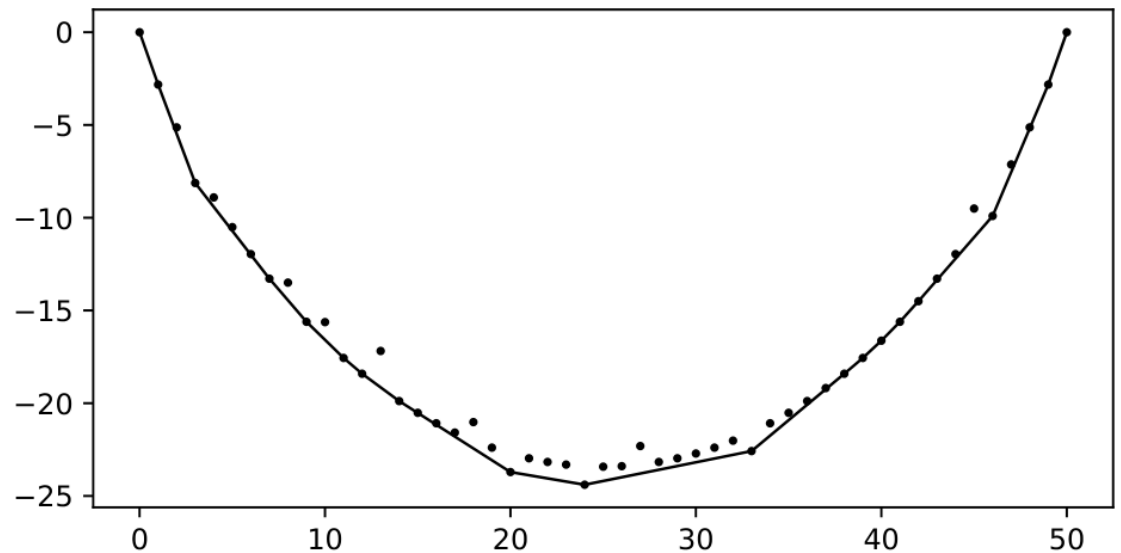
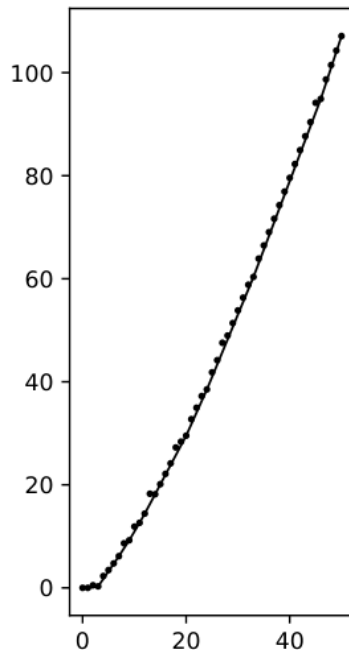
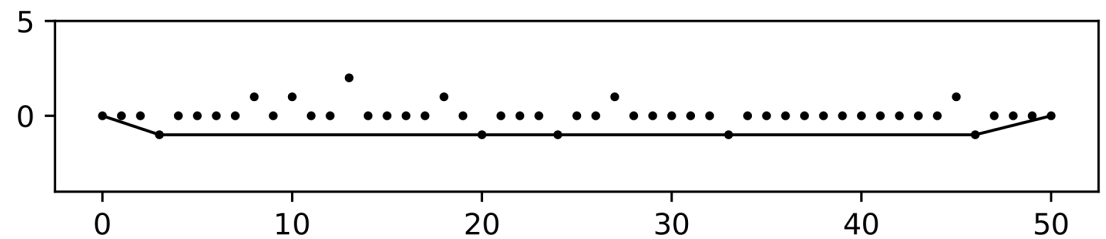
Three families of polynomials

c_j are iid random Gaussian variable with mean 0

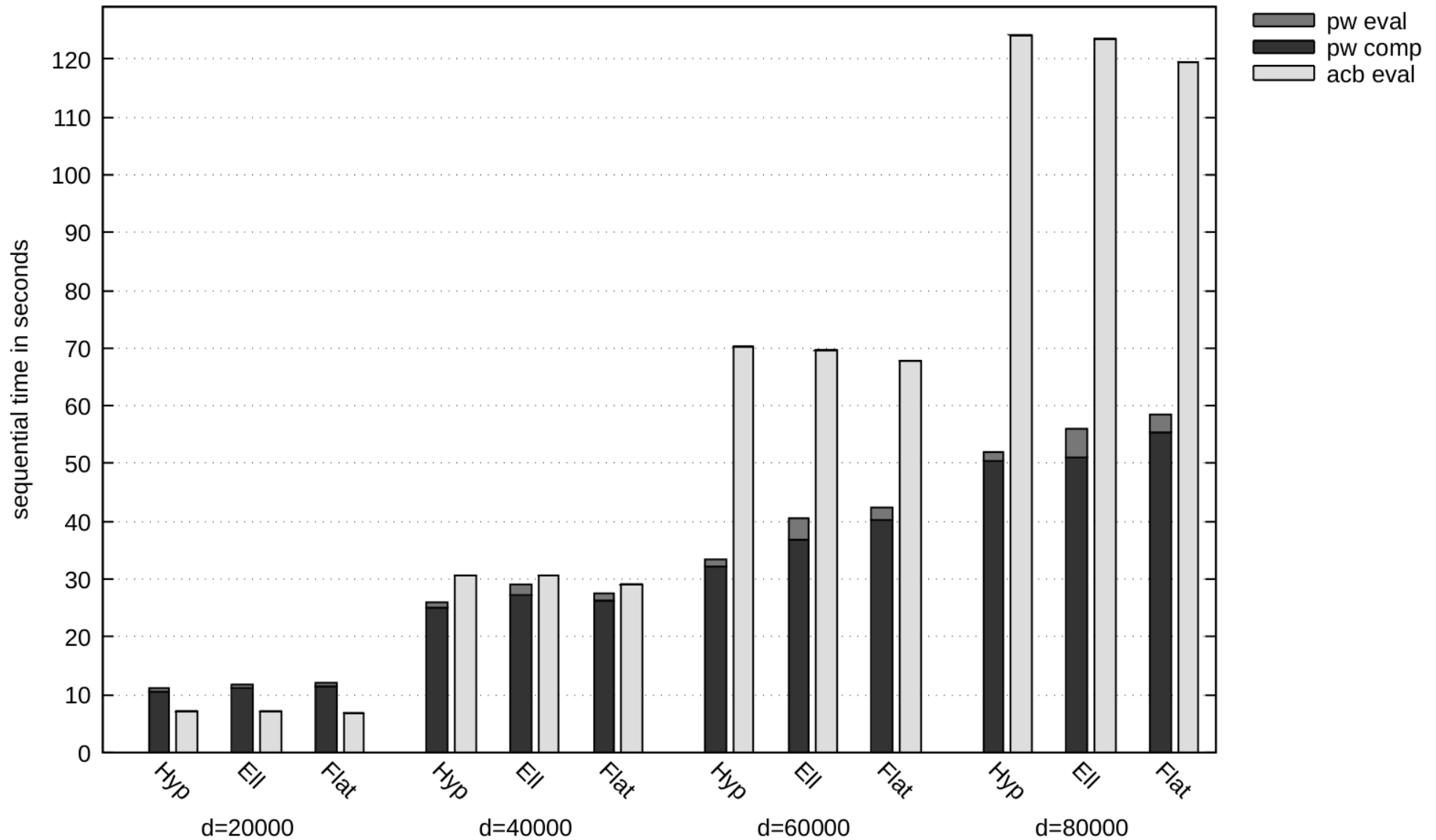
ic: $\sum c_j z^j$

ic: $\sum c_j \sqrt{\binom{d}{j}} z^j$

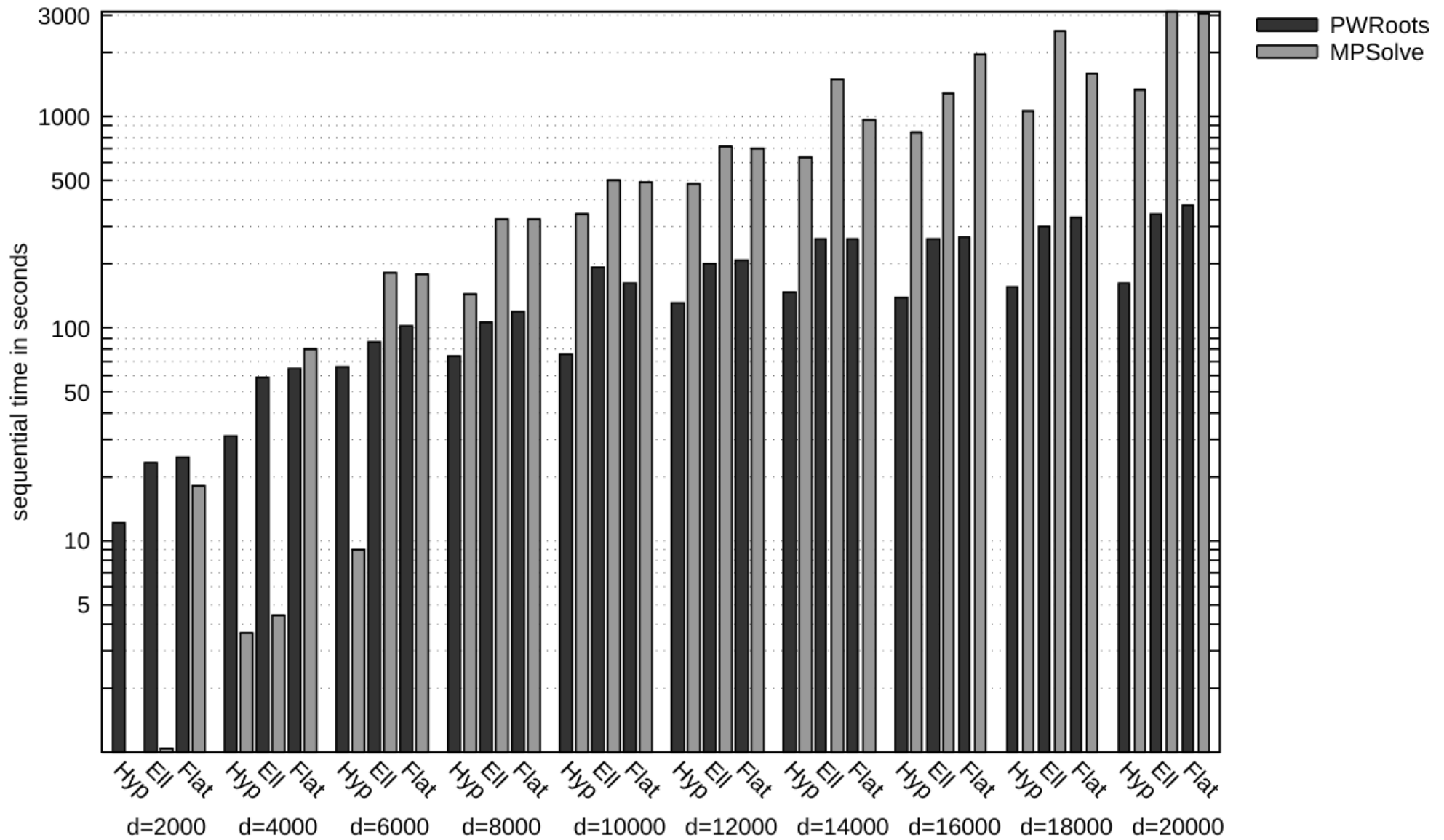
at: $\sum c_j \sqrt{\frac{1}{j!}} z^j$



Benchmark



Benchmark



Conclusion

Perspectives

- Bivariate polynomials
- Non integer exponents
- Laplace transform

Thank you!