

Distances for Sequences and Text

- **symbol distance:** $\rho(x, y) = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{otherwise} \end{cases}$

- **Hamming distance:**

$$H((x^{(1)}, \dots, x^{(p)}), (y^{(1)}, \dots, y^{(p)})) = \sum_{i=1}^p \rho(x^{(i)}, y^{(i)})$$

number of elements
p that are different

- **edit/Levenshtein distance:** append/remove/edit operations to have the 2 vectors equal

$$L((x^{(1)}, \dots, x^{(p)}), (y^{(1)}, \dots, y^{(q)})) =$$

$$\begin{cases} p & q = 0 \\ q & p = 0 \\ \min \{ \begin{array}{l} L((x^{(1)}, \dots, x^{(p-1)}), (y^{(1)}, \dots, y^{(q)})) + 1, \text{ (REMOVE 1st sequence)} \\ L((x^{(1)}, \dots, x^{(p)}), (y^{(1)}, \dots, y^{(q-1)})) + 1, \text{ (REMOVE 2nd sequence)} \\ L((x^{(1)}, \dots, x^{(p-1)}), (y^{(1)}, \dots, y^{(q-1)})) + \rho(x^{(p)}, y^{(q)}) \end{array} \} & \text{otherwise} \end{cases}$$

CHANGE symbols (edit)

Example Edit/Levenshtein Distance

		C	L	E	O	P	A	T	R	A
	0	1	2	3	4	5	6	7	8	9
C	1	0	1	2	3	4	5	6	7	8
A	2	1	1	2	3	4	4	5	6	7
E	3	2	2	1	2	3	4	5	6	7
S	4	3	3	2	2	3	4	5	6	7
A	5	4	4	3	3	3	3	4	5	6
R	6	5	5	4	4	4	4	4	4	5

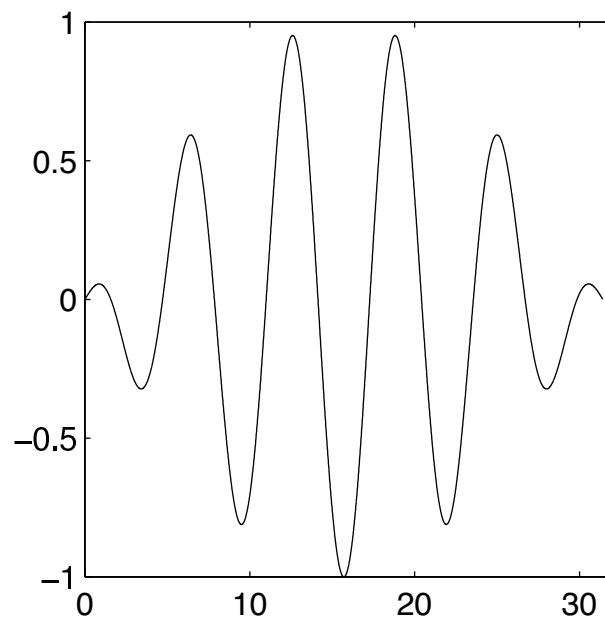
C	A	E	S		A		R	
0	1	0	1	1	0	1	0	1
C	L	E	O	P	A	T	R	A

$\Rightarrow 5$

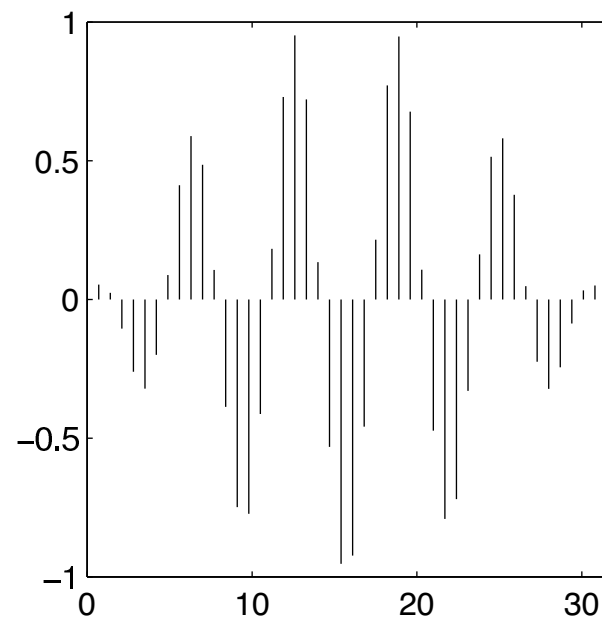
Sampling Continuous Signals

continuous functions

original signal
 $s(t)$

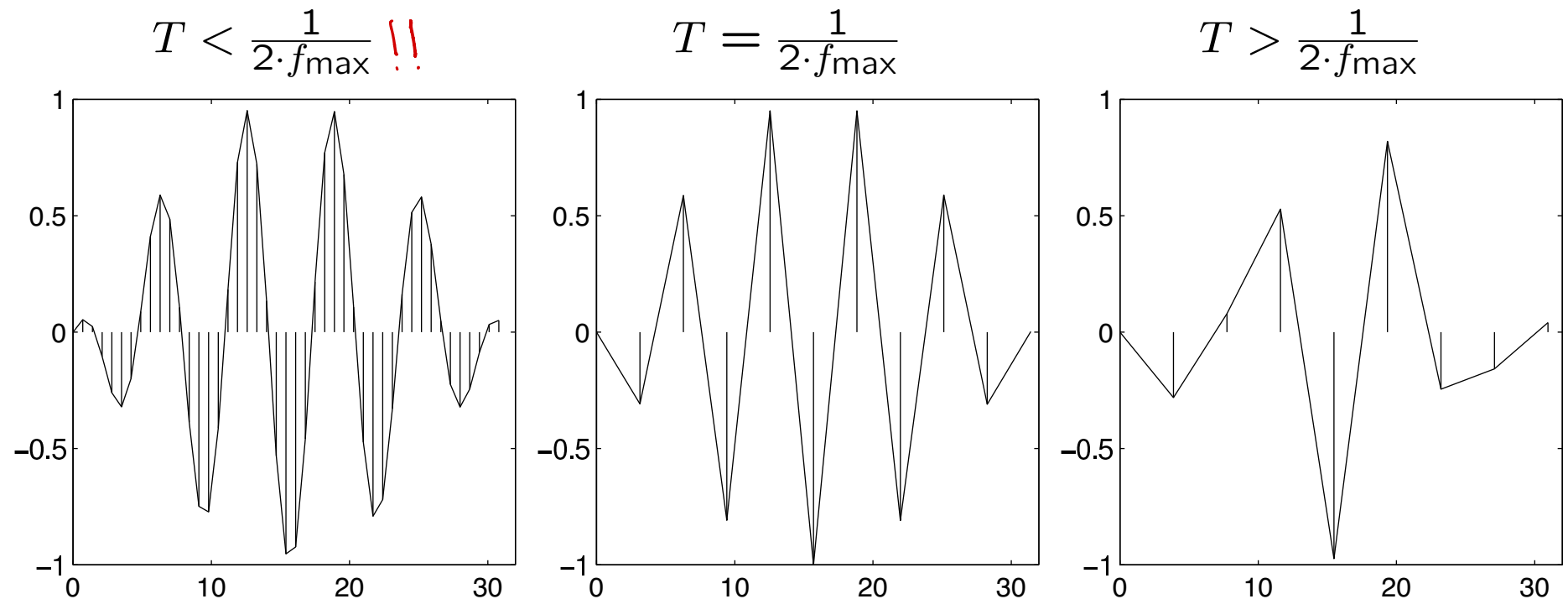


sampld signal
 $s_n = s(n \cdot T)$ T sampling time

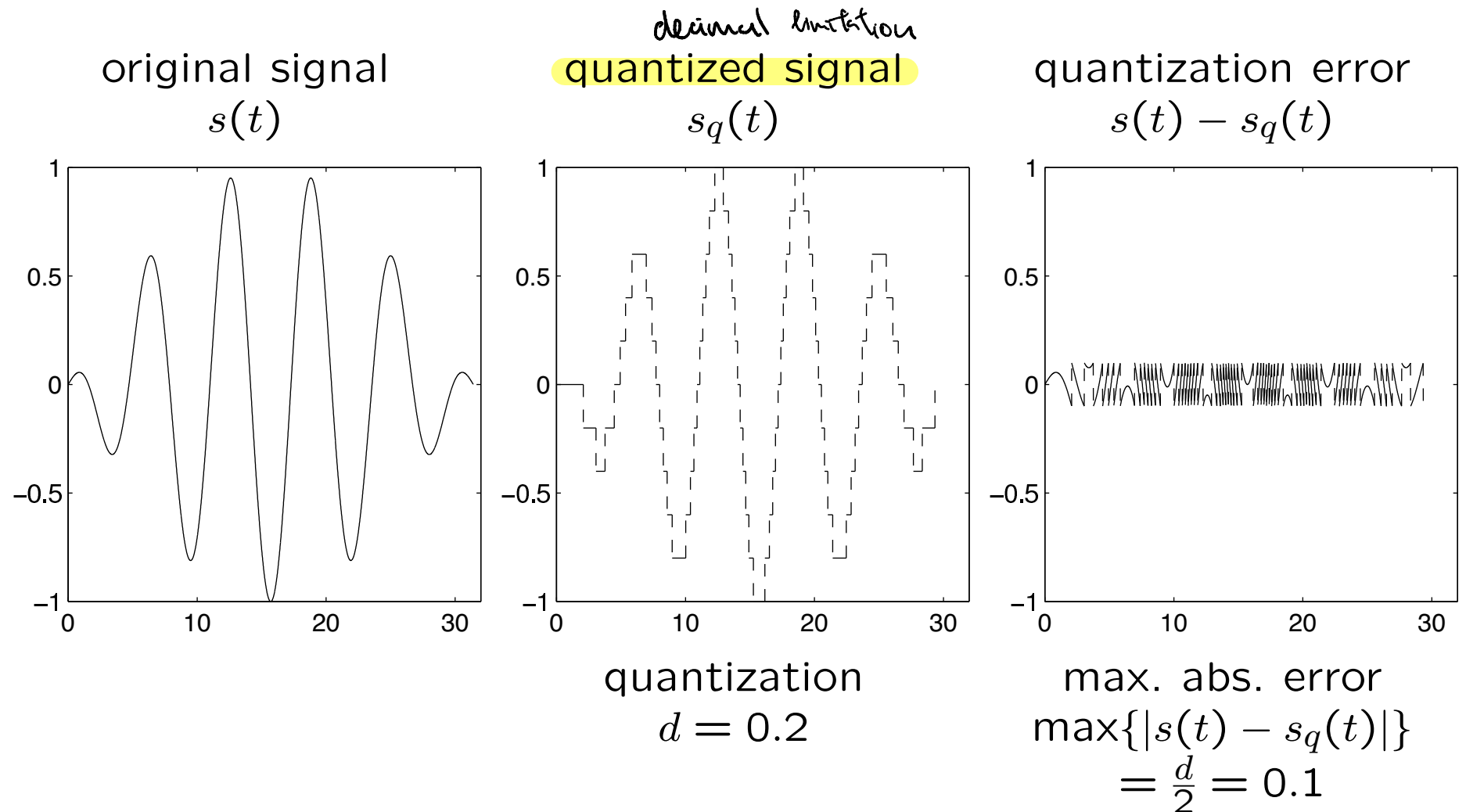


Shannon's Sampling Theorem

1. $s(t)$ band limited: Fourier spectrum $|s(j2\pi f)| = 0$ for $|f| > f_{\max}$
 2. $T_s < \frac{1}{2 \cdot f_{\max}}$ (Nyquist condition) limit
- $\Rightarrow s(t)$ can be completely reconstructed from s_n



Quantization



Chapter 3: Data Preprocessing

1. Error Types and Handling
2. Filtering
3. Standardization and Transformation
4. Data Merging