

Curse of Dimensionality

- > when dimensionality increases data become sparse. (Information content is less)
- > density and distance between points, which is critical to clustering, outlier analysis becomes less meaningful.
- > The possible combinations of subspaces will grow exponentially.

Dimensionality Reduction

- > Avoid the curse of dimensionality.
- > Help eliminate irrelevant features and reduce noise.
- > Reduce time and space required for data mining.
- > Allow easier visualization.

Techniques used

PCA - Principal Component Analysis

FA - Factor Analysis

DFT - Discrete Fourier Transform

DCT - Discrete Cosine Transform

DWT - Discrete Wavelet Transform.

Time Domain - analysis of signal or data w.r.t Time.

Frequency Domain - Analysis of signal or data w.r.t. Frequency.

Example

ECG - Electro cardiogram

If a doctor maps the heartbeat with time.
say, the recording is done for 20 min.
we call it time domain signal.

on ECG a number of peaks are there.
say, in one heartbeat 4 types of peaks
or variation in amplitude occurs.

on frequency domain representation \rightarrow
How many times each peak comes is
recorded over the entire time period of
recording.

A given function or signal can be converted
between the time and frequency
domain with a pair of mathematical
operations is called a transform.

Discrete Fourier Transform (DFT)

using ~~the~~ DFT we can determine the frequency
content of a signal.

$$\xrightarrow{\text{DFT}} \quad X(K) = \sum_{n=0}^{N-1} x(n) e^{-j \frac{2\pi K n}{N}}, \quad 0 \leq K \leq N-1$$

$$\xrightarrow{\text{IDFT}} \quad x(n) = \frac{1}{N} \sum_{K=0}^{N-1} X(K) e^{j \frac{2\pi K n}{N}}, \quad 0 \leq n \leq N-1$$

Q Find the DFT of a sequence $x(n) = \{1, 1, 0, 0\}$

Solⁿ

$N = 4 = \text{length of } x(n)$

$K = 0, 1, \dots, N-1 = 0, 1, 2, 3$

Discrete Wavelet Transform (DWT) (3)

$$\underline{K=0}$$

$$X(0) = \sum_{n=0}^3 x(n) e^0 = x(0) + x(1) + x(2) + x(3) \quad (2)$$

$$= 1 + 1 + 0 + 0 = 2$$

$$\underline{K=1}$$

$$X(1) = \sum_{n=0}^3 x(n) e^{-j2\pi n/4} = x(0)e^0 + x(1)e^{-j\pi/2}$$

$$+ x(2)e^{-j\pi} + x(3)e^{-j3\pi/2}$$

$$= 1 + 1(\cos \pi/2 - j \sin \pi/2) = 1 - j$$

$$\underline{K=2}$$

$$X(2) = \sum_{n=0}^3 x(n) e^{-j2\pi \cdot 2 \cdot n/4} = \sum_{n=0}^3 x(n) e^{-j\pi n}$$

$$= x(0) + x(1)e^{-j\pi} = 1 + \cos \pi - j \sin \pi$$

$$= 1 - 1 = 0$$

$$\underline{K=3}$$

$$X(3) = \sum_{n=0}^3 x(n) e^{-j2\pi \cdot 3 \cdot n/4} = \sum_{n=0}^3 x(n) e^{-j3\pi n/2}$$

$$= x(0) + x(1)e^{-j3\pi/2} + x(2)e^{-j3\pi} + x(3)e^{-j9\pi/2}$$

$$= 1 + \cos 3\pi/2 - j \sin 3\pi/2 = 1 + j$$

$$X(K) = \{ 2, 1-j, 0, 1+j \}$$

$$\boxed{\text{IDFT}(X(K)) \rightarrow x(n)}$$

$$X(K) = \{2, 1-j, 0, 1+j\} \text{ IDFT}$$

$$e^{j\omega} = \cos\omega + j\sin\omega$$

$$e^{-j\omega} = \cos\omega - j\sin\omega$$

$$y(n) = \frac{1}{N} \sum_{K=0}^{N-1} Y(K) e^{j2\pi nK/N}, n = 0, 1, \dots, N-1$$

$$\begin{aligned} \frac{n=0}{y(0)} &= \frac{1}{4} \sum_{K=0}^3 Y(K) e^0 = \frac{1}{4} [Y(0) + Y(1) + Y(2) + Y(3)] \\ &= \frac{1}{4} [2 + 1-j + 0 + 1+j] = \frac{4}{4} = 1 \end{aligned}$$

$$\begin{aligned} \frac{n=1}{y(1)} &= \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\frac{2\pi K}{4}} = \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\frac{\pi K}{2}} \\ &= \frac{1}{4} \left[Y(0) \cdot e^0 + Y(1) e^{j\frac{\pi}{2}} + Y(2) e^{j\pi} + Y(3) e^{j\frac{3\pi}{2}} \right] \end{aligned}$$

$$= \frac{1}{4} \left[2 \cdot 1 + (1-j)(\cos\frac{\pi}{2} + j\sin\frac{\pi}{2}) + 0 + (1+j)(\cos\frac{3\pi}{2} + j\sin\frac{3\pi}{2}) \right]$$

$$= \frac{1}{4} \left[2 + (1-j)(0+j) + (1+j)(0-j) \right]$$

$$= \frac{1}{4} [2 + j - j^2 + 1 - j - j^2]$$

$$= \frac{1}{4} [2 - 2(-1)] = \frac{4}{4} = 1$$

$$\begin{aligned} \frac{n=2}{y(2)} &= \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\frac{2\pi K \cdot 2}{4}} = \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\pi K} \\ &= \frac{1}{4} \left[Y(0) e^0 + Y(1) e^{j\pi} + Y(2) e^{j2\pi} + Y(3) e^{j3\pi} \right] \end{aligned}$$

$$= \frac{1}{4} \left[2 \cdot 1 + (1-j)(\cos\pi + j\sin\pi) + 0 + (1+j)(\cos 3\pi + j\sin 3\pi) \right]$$

$$= \frac{1}{4} \left[2 + (1-j)(-1+0) + (1+j)(-1+0) \right]$$

$$= \frac{1}{4} [2 + (1-j)(-1+0) + (1+j)(-1+0)]$$

$$= \frac{1}{4} [2 - 1 + j - 1 - j] = \frac{0}{4} = 0$$

$$\begin{aligned} \frac{n=3}{y(3)} &= \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\frac{2\pi K \cdot 3}{4}} = \frac{1}{4} \sum_{K=0}^3 Y(K) e^{j\frac{3\pi K}{2}} \\ &= \frac{1}{4} \left[Y(0) e^0 + Y(1) e^{j\frac{3\pi}{2}} + Y(2) e^{j3\pi} + Y(3) e^{j\frac{9\pi}{2}} \right] \end{aligned}$$

$$= \frac{1}{4} \left[2 + (1-j)(\cos\frac{3\pi}{2} + j\sin\frac{3\pi}{2}) + 0 + (1+j)(\cos\frac{9\pi}{2} + j\sin\frac{9\pi}{2}) \right]$$

$$= \frac{1}{4} \left[2 + (1-j)(0-j) + (1+j)(0+j) \right]$$

$$= \frac{1}{4} [2 + (-j + (-1)) + (j + (-1))] = \frac{1}{4} [2 - 2] = 0$$

Discrete Wavelet Transform (DWT)

(3)

~~DWT~~ ~~is~~ ~~also~~

$$\vec{X} \rightarrow \text{DWT} \rightarrow \vec{X}'$$

(wavelet coefficient)

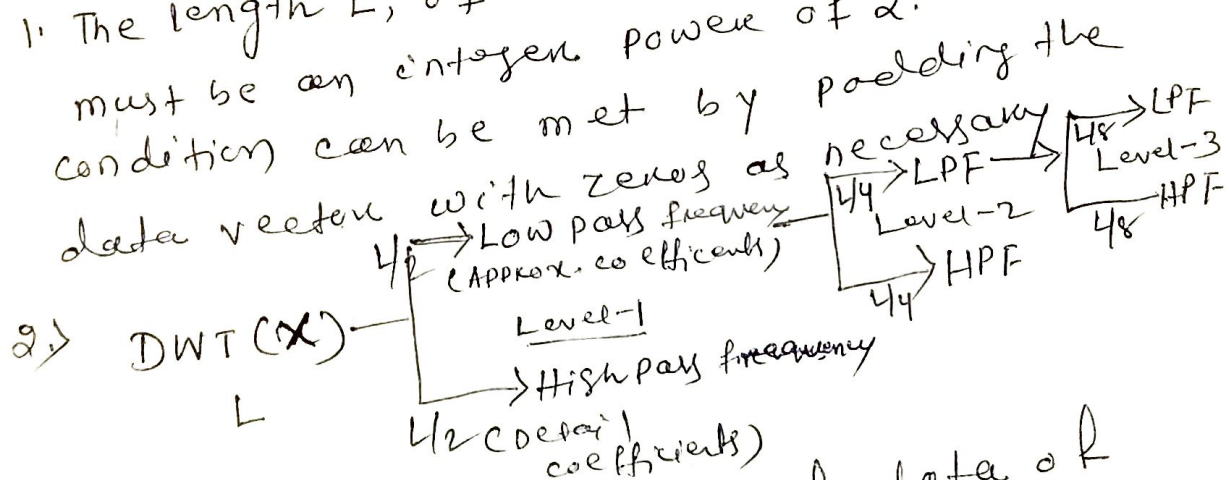
The size of $(\vec{X}) = \text{size of } (\vec{X}')$

For dimensionality reduction, we will keep the wavelet coefficients larger than some threshold value and truncate the small values or set them to zero value.

Given the wavelet coefficients, an approximation of the original data can be constructed by applying IDWT.

Method

1. The length L , of the input data vector must be an integer power of 2. This condition can be met by padding the data vector with zeros as necessary.



This results in two sets of data of length $L/2$. Generally it contains a low frequency version of the input data and high frequency content of it respectively.

3. The low frequency content again divided into two sets and the process continues until the resulting dataset obtained are of length 2.

4. selected values from the data sets obtained in above are Wavelet coefficients.

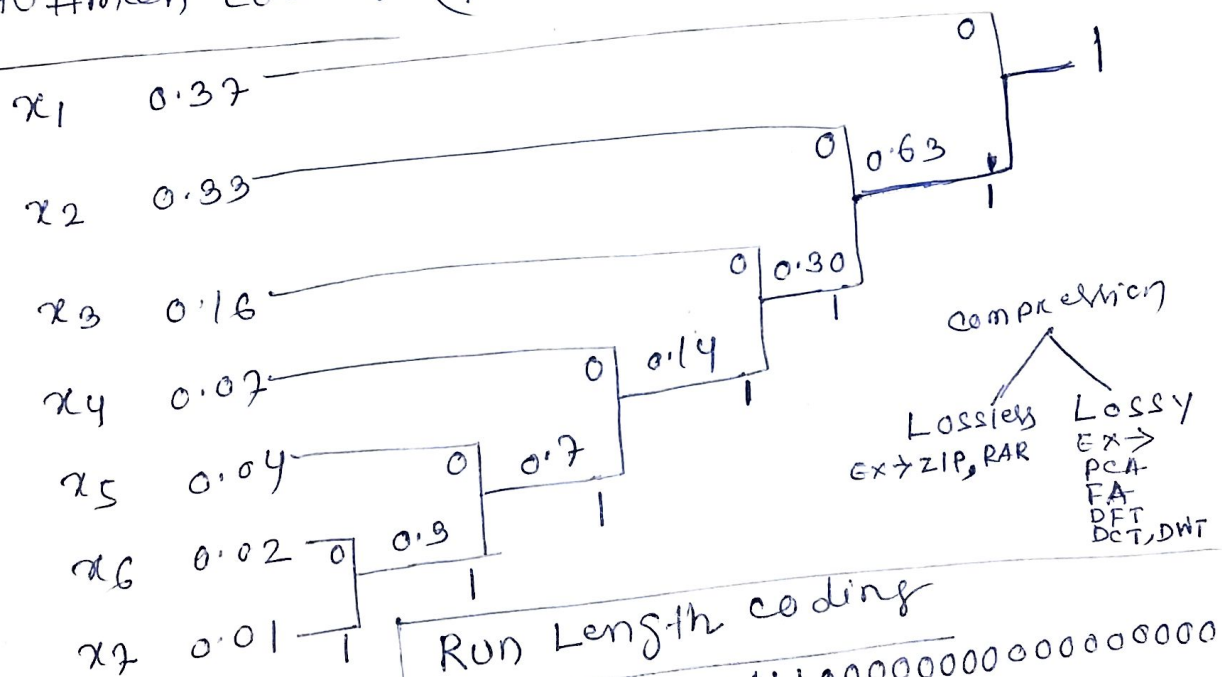
Advantages of DWT over DFT

- 1) DWT achieves better lossy compression.
- 2) Requires less space than DFT
- 3) There is only one DFT, but there are several families of DWT.
- 4) popular wavelet transforms include Haar-2, Daubechies-2, Daubechies-6.
- 5) DWT uses a hierarchical pyramid algorithm that halves the data at each iteration resulting in fast computation.

iterations resulting

Huffman coding (for Data compression)

81



Run Length coding

Run Length Code

s = 11111111111111000000000000000000
001111 (0 → 1)

Max. no. of repetition = 19, which
can be represented with 5-bits.

compression ratio = $\frac{18}{38} = 1:2.11$

codeword

$$x_1 = 0$$
$$x_2 = 10$$
$$x_3 = 110$$
$$x_4 = 1110$$
$$x_5 = 11110$$
$$x_6 = 11110$$

27 = |||| ||||