**11)** Given a dataset with repetitive patterns design a simple physical model for lossless compression. Justify how the model reduces data size and provide an example to demonstrate your approach.

**1) Physical model: Run length Encoding (RLE)**

concept

Instead of storing replaced data values individually, RLE stores:
- The values
- The number of times it repeats (run length).

model representation.

Each sequence is encoded as $(value, count)$.

**2) Compression Algorithm**

Encoding procedure

- Read the dataset size sequentially.
- count consecutive occurrence of same value.
- Replace repeated values with a pair
- continue until the end of the dataset.

Decoding procedure

- Read each pair.
- Reconstruct the original data by repeating the value count times.

### 3) Justification of size reduction.

Let:
* Original size = number of element N.
* Compressed size = number of runs R.

If data contains many repetitions.

$R \ll N$.

$$\therefore \text{compression Ratio} = N/(2R).$$

Because each run stores only:

- One value
- one count

### Key insight

- works best when data long repeated sequence
- Eliminates redundancy by grouping identical values.

### 4) Example

AAAA BB CCCCC AA.

### Step 1:
Identify runs..

- A repeated 4 times
- B repeats 2 times
- C repeat 5 times
- A repeat 2 times.

## Step 2:

encod
(A,4), (B,2), (C,5), (D,2).

## Step 3

Size compression

original length = 13

compressed length = 4 pairs = 8 values..

## 5) Code snippet

```python
def rle_encode (data):
    encode = []
    count = 1
    for i in range (1, len(data)):
        if data[i] == data[i-1]:
            count += 1
        else:
            encode.append ([data[i-1], count])
            count = 1
    encode.append ((data[-1], count))
    return encode

def rle_decode (encoded):
    decode = []
    for value, count in encoded:
        decoded.extend ([value] * count)
    return decoded
```

## 6) Advantages

- simple to implement
- Lossless
- Effective for repetitive datasets.

## 7) Limitations

- Not efficent for non-repetitive data.
- May increase size, if data has no repetition.

## 8) Conclusion.

The Run length encoding model compress data by replacing repeated sequences with (value, count) pairs. It reduces data size by eliminating redudancy, making it highly effective for datasets with repetitive patterns. It follows ~~lossless~~ lossless compression.