# Image processing background

## Image Processing

- Intersection of classic ML w/ SciML        Ultrasound, MRI...

    eg. medical imaging, tomography,
         astronomy, microscopy ...

- Many tasks are linked, eg.    compression
                                 denoising          which we'll
                                 inverse problems       explore

## Image Compression: JPEG (1992)    (lossy)

Rough idea:    ① split image into $8 \times 8$ patches, and from
                  now on, operate on patches separately
                       (this guarantees linear complexity and makes it
                        fast)

               ② perform a 2D ($8 \times 8$) DCT (Discrete Cosine Transform)
                   to transform to frequency space        (DCT via integer
                   (no compression/loss yet)                arithmetic!!)

               ③ set small coefficients to 0  (user defined threshold)

               ④ possibly quantize remaining coefficients    ⎤ save these
                  possibly also entropy-encode them ⎤ lossless ⎦ to file
                       "Arithmetic code"

## Why the DCT to "induce sparsity"?

to paraphrase George Box, "All (image) models are wrong,
                               some are useful"

Let's go back even further
                                                      discrete
## How to model compression:                 Set of symbols like
                                               $\{a, b, c, ..., z\}$
        Suppose we wish to compress $\{X_i\}_{i=1}^{\hat{}}$, $X_i \in \mathcal{X}$

        and we assume ① these are random,
                       ② these are iid, w/ probability $p(X)$

ex: $p(X = "a") \approx \frac{1}{26}$
    $p(X = "e") \approx \frac{1}{10}$
    $p(X = "z") \approx \frac{1}{50}$

let our alphabet $\mathcal{X}$ have $m$ symbols, so we can characterize $x_i \in \mathcal{X}$ using an integer $0, 1, ..., m-1$   (eg. $0, 1, ..., 25$)

Baseline:

encode this integer directly. Requires $\lceil \log_2(m) \rceil$ bits per symbol. To send a message $(x_1, x_2, ..., x_n)$ we need $n \cdot \lceil \log_2(m) \rceil$ bits, and the average per symbol is $\lceil \log_2(m) \rceil$

Doing better:

If $p(X = "z") = 0$, no need to encode it. If it's almost $0$, can we exploit? Yes!

Aside:

Entropy of a distribution $p$, with $X \sim p$ a r.v., is $H(X) := -\sum_{x \in \mathcal{X}} p(x) \log_2(p(x))$

Ex: if $p$ is uniform (aka uninformative) over $\mathcal{X} = \{1, 2, ..., m\}$ then $p(x) = \frac{1}{m}$ $\forall x \in \mathcal{X}$,

$H(X) = \sum p(x) \cdot \log_2\left(\frac{1}{p(x)}\right) = \sum \frac{1}{m} \log_2(m)$
$$= \log_2(m) \quad \text{Same as baseline}$$

All other distributions have lower entropy.

We can encode w/ Huffman codes or other entropy codes.

Theorem (Shannon)

It is possible to encode $(X_1, ..., X_n)$ (if iid) using $n \cdot (H(X) + \varepsilon)$ bits, for any $\varepsilon > 0$, and recover it w/ probability $1 - \varepsilon$.

**Problem**     What if $\{X_1, ..., X_n\}$ aren't iid?

Ex:    $\mathcal{X} = \{a, b, ..., z\}$

message     "Helloworld"

$X_2$ is n<u>ot</u> independent of $X_1$     <span style="color:purple">think of Wordle</span>

( it's more likely to be a vo<u>wel</u> )

So, entropy encoding is inefficient

**Strategy 1**:   larger messages,  $\mathcal{X} = \{$ all words $\}$

intractable

**Strategy 2**:  make it iid

↑ actually impossible... but we can <span style="color:red">whiten</span> it

(ie. make the <u>correlations</u> zero)

<span style="color:red">Karhunen-Loève transformation</span>

like PCA but when we treat data as random.

<span style="color:green">**Ex**</span>

Suppose   $Cov(X_i, X_j) = \sigma^2 \cdot \rho^{|i-j|}$   for some $\rho \in (-1, 1)$

then the   DCT <u>is   the   KL transform.</u>     <span style="color:purple">should be familiar to those of you who took time series</span>

So...

If $X_i$ is value of image at pixel $i$,

and    $Cov(X_i, X_j) \approx \rho^{|i-j|}$    (and extend to 2D)

then DCT wh<u>itens</u> our signal, so we can better

exploit entropy encoding.                          <span style="color:purple">⚹.</span>

# Evaluation metrics

- The OG: Mean Squared Error (MSE)

$$MSE = \frac{1}{m \cdot n} \sum_{i=1}^{m} \sum_{j=1}^{n} (\hat{X}_{ij} - X_{ij})^2$$

$MSE \geq 0$
lower better

↑ $i,j^{th}$ pixel

  - well understood
  - many nice mathematical properties

- Variant: Peak Signal-to-Noise Ratio (PSNR)

$$PSNR = 10 \cdot \log_{10} \left( \frac{max^2}{MSE} \right)$$

$PSNR \in \mathbb{R}$
higher better

max = 255 for 8-bit images

28dB + good

in decibels (dB) a relative scale
$10 \cdot \log($ power units$)$
$20 \log ($ amplitude units $)$

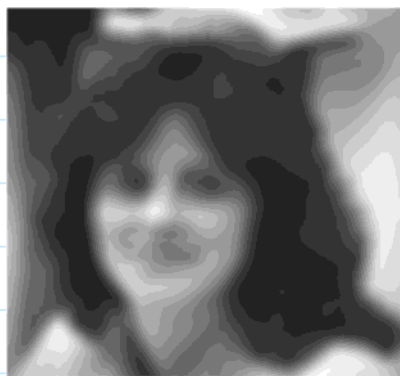- Structural Similarity Index Measure (SSIM)   '01, '04

Wang et al.
50k + citations
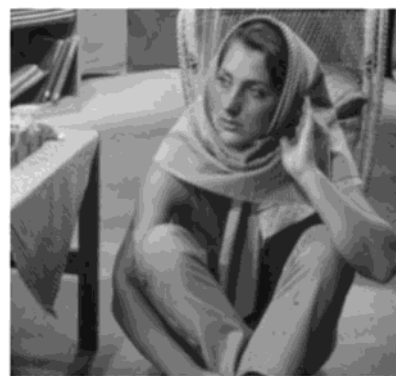
Also has an easy-to-use formula ... but

attempts to align w/ human perception

$SSIM \in (0,1]$
near 1 is better

ex: Failure of PSNR



PSNR=25.11 dB, IQA=0.0292          PSNR=25.12 dB, IQA=0.5574

Bezhadpour and Ghanbari, 2021,
https://www.researchgate.net/figure/Subjective-quality-of-two-different-contents-with-the-same-PSNR_fig2_362323114

SSIM attempts to capture luminance masking and feature masking

... like MP3's

- Learned Perceptual Image Patch Similarity (LPIPS)

lower better

Uses a trained neural net