

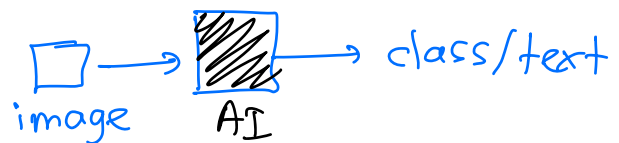
AI is vision

Intelligent is to recognize people, scenes, objects patterns

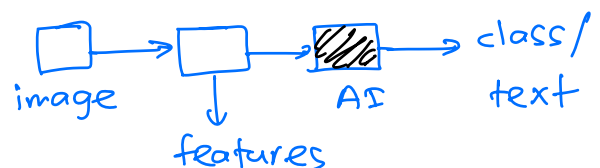
- face recognition
- Object recognition
- Auto captioning of images
- Robot navigations

Given a digital image, 2 possible ways for recognition

① AI approach



② CV/AI approach



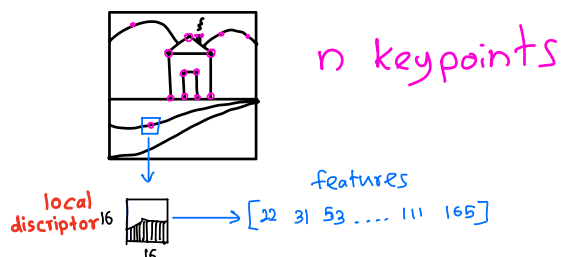
Feature extraction

① Key point oriented (SIFT, SURF,...)

→ several hundred/thousands of features

② Histogram oriented (LBP, HOG, ELP) → one histogram

Example :



for SIFT you get many feature vectors of length

$$\text{Image} = \bigcup_{i=1}^{n_{\text{keypoints}}} v_i \leftarrow \text{union sparse sampling (only sample keypoints)} \quad 128.$$

Challenges: 1) Data is too large

2) Data may not be descriptive

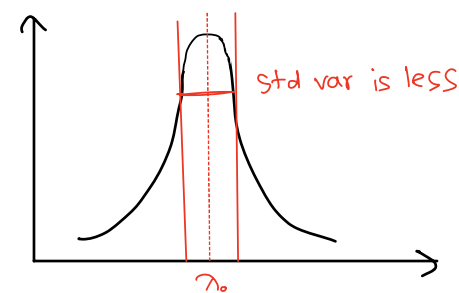
Solutions: embedding / pooling / encoding

→ Fisher vector (embedding and pooling)
→ VLAD (encoding and increasing discrimination)

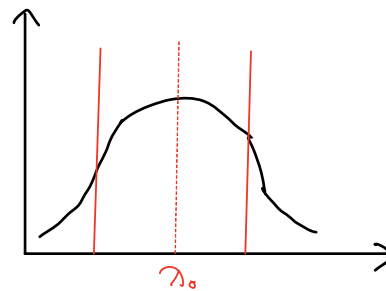
$$X = \{x_t, t=1, 2, \dots, T\}$$

u_λ = probability density function which models the generative process of elements of x .

$\lambda = [\lambda_1, \lambda_2, \dots, \lambda_m]^T \in \mathbb{R}^m$ parameters of u_λ



fits the data more accurately



Statistics : The score function is the gradient (partial derivative) wrt parameter λ of the natural

logarithm of the likelihood function.

$$\text{score function} = \underbrace{\nabla_{\lambda}}_{\text{partial derivative wrt } \lambda} \log u_{\lambda}(x)$$

a, b independent and identically distributed random variables.

$f(a, b|u) = f(a, u) * f(b, u)$ is likelihood function a, b given distribution u

$$\log(f(a, b|u)) = \log(f(a, u)) + \log(f(b, u)) \quad \begin{matrix} \log \\ \text{likelihood} \end{matrix}$$

$$\begin{aligned} \text{score function} &= \nabla_{\lambda} \log u_{\lambda}(x) \\ &= \nabla_{\lambda} \log P(x | u_{\lambda}) \end{aligned}$$

Let X be the set of D -dimensional local descriptors extracted from an image (e.g. SIFT) Don't have intelligence

$$g_{\lambda}^x = \sum_{t=1}^T L_{\lambda} \nabla_{\lambda} \log u_{\lambda}(x_t) \quad \text{fisher vector}$$

function g that has parameter λ operating on X

Fisher vector is a sum of normalized gradient statistics compute for each descriptor (feature vector)

The operation $X_t \longrightarrow f_{f_k}(x_t) = L_{\lambda} \nabla_{\lambda} \log u_{\lambda}(x_t)$

read the data and put it on another space

f uses fisher kernel for each data point x_t
is an embedding of local descriptors x_t
in a higher dimensional space which is easier for

a linear classifier.

L_λ : cholesky Decomposition

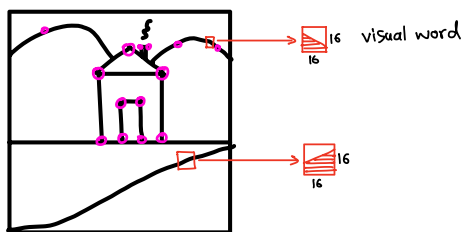
$$F_\lambda^{-1} = L_\lambda^T L_\lambda$$

$$K_{FK}(x, y) = G_\lambda^{xT} F_\lambda^{-1} G_\lambda^x$$

Fisher kernel - take 2 numbers and tell how similarity it is in higher dimensional space.

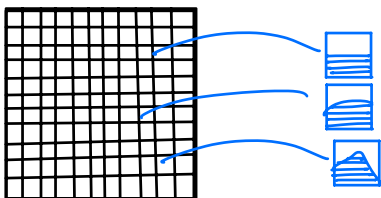
VLAD (Vector for locally aggregated descriptors)

How to recognize images? \rightarrow Bag-of-visual-words



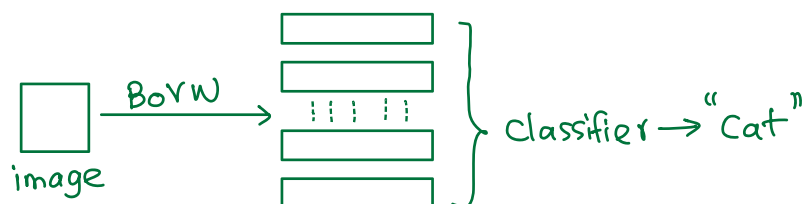
Given an image I , divided it into small cells/windows of size $n \times n$ (eg: 16×16)

dense sampling



We vectorize visual words for convenient calculations.

Theory:



Practice: - too many vectors

- Redundancy

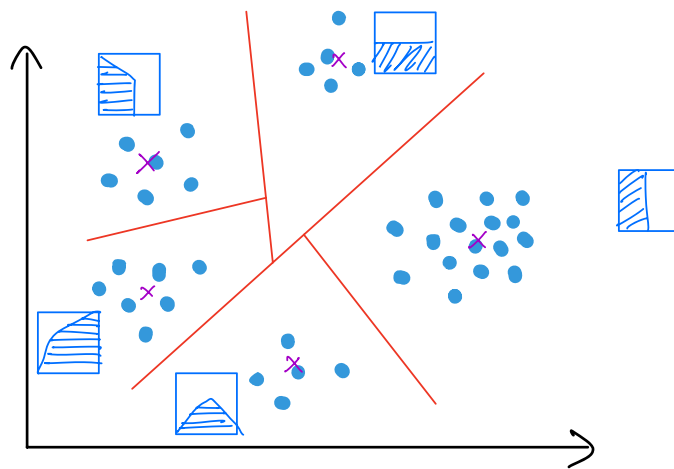
- Noise

General approach

Build a ^(dictionary) codebook $C = \{c_1, c_2, \dots, c_n\}$ from $m \gg n$ feature vectors (vectorized visual words)

idea: Use a clustering algorithm like k-means

c_i is the centre of the n classes found in the data



Core idea of VLAD

Accumulate for each visual word c_i , the difference of the vectors x assigned to c_i , $x - c_i$ ie distribution of data wrt the class centers.

$$V_{i,j} = \sum_{NN(x)=c_i} \underbrace{x_j}_{\substack{\text{jth component} \\ \text{of the descriptor } x}} - \underbrace{c_{i,j}}_{\substack{\text{cluster} \\ \text{corresponding centre}}}$$

use these in dictionary

L2 - normalization, $V: \frac{V}{\|V\|}$

Final chain

