

Terminology

data can be:

- bare numbers
- symbols
- ~~video~~ images
- audio files
- text

Preprocessing

- Filtering (noise reduction)
- Normalization
- Outlier detection
- dimensionality reduction

(we assume there is structure in the data)

Features

Attributes

- data representation

(euclidean distance fails if dimension is high, say 13,000)

hand crafted/
engineered.
eg. SIFT.

Automatically extracted
features
eg. ~~SIFT~~ CNN.

Encoding

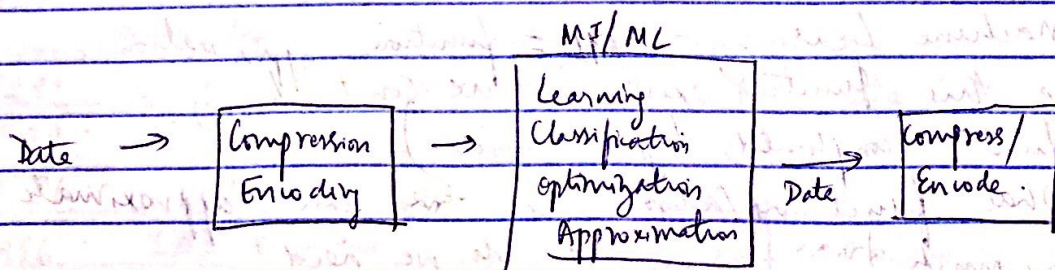
- compressing data

(purpose is not dimensionality reduction, purpose is encoding)

Conventional
Conventional
(PCA, Fisher Vector)

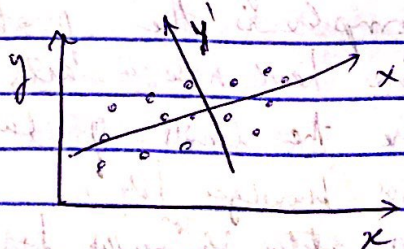
ML
(autoencoders, t-SNE)

Jan 10th 2016



- Compression PCA
- Encoding (Fisher Vector)

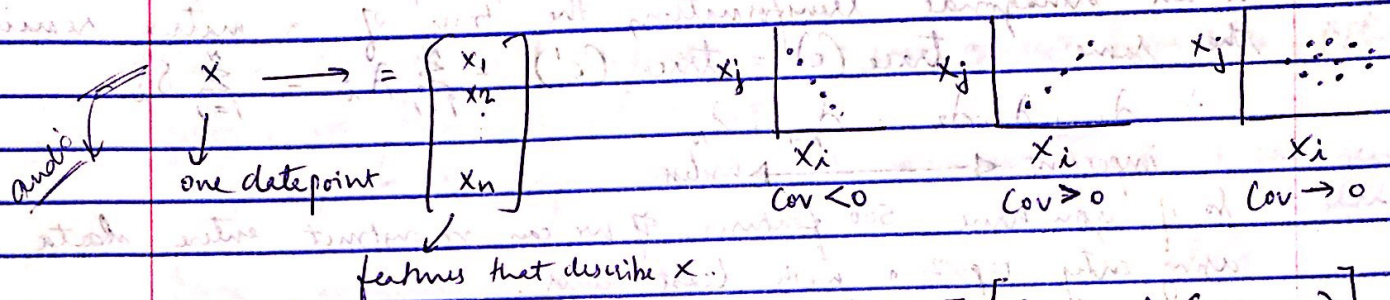
variance of data maximum when projected on x' .



Data that does not give one variance is useless! It should give change information

PCA finds X'/Y' with max. variance wrt. population mean.
PCA repeats the process until N such axes are found.

PCA finds a set of orthogonal axes that diagonalizes the covariance matrix.



$$\text{Cov. Matrix} = \sum_{i,j} = \text{Cov}(x_i, x_j) = E[(x_i - \mu_i)(x_j - \mu_j)]$$

$$\mu_i = E(x_i)$$

$$\mu_j = E(x_j)$$

Generalization of variance:

$$\Sigma = E[(X - E(X))(X - E(X))^T]$$

$$\sigma^2 = \text{var}(x) = E[(X - E(X))^2]$$

Entries of diagonal of Σ are the variances of each element of X .

Covariance matrix C :

$$C = E[(x_p - m)(x_p - m)^T]$$

Annotations: x_p = sample, m = population mean, E = Expected value.

x = scalar X = set
 x = vector X = matrix

We can estimate C as:

$$C = \frac{1}{P} \sum_{p=1}^P x_p x_p^T - m m^T$$

$$m = \frac{1}{P} \sum_{p=1}^P x_p$$

Diagonalizing C using a suitable orthogonal transformation matrix A by obtaining N orthogonal eigen vectors u_i with eigen value λ_i .

$$C u_i = \lambda_i u_i$$

$$i = 1, 2, \dots, N.$$

The vectors u_i are derived $u_i = A(X_i - m)$ $X_i = m + A u_i$

Since $A^{-1} = A^T$ (because it's orthogonal)

Hence $C' = ACA^T$

$$C' = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_N \end{bmatrix}$$

In an orthogonal transformation, the "trace" of a matrix remains the same. $\text{trace}(C) = \text{trace}(C') = \sum_{i=1}^N \lambda_i = \sum_{i=1}^N s_i^2$

$\therefore \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_{N-1}, \lambda_N$
important \longleftrightarrow useless

So if you have 500 features, if we can reconstruct entire data using only top 2, with least error.

Encoding

eg. Fisher Vector.

Applies Fisher kernel. $X = \{x_t : t=1, 2, \dots, T\}$
T observations $x_t \in T$

$x_t \in X$

μ_λ probability density function which models the generative process of elements in X

$\lambda = [\lambda_1, \lambda_2, \dots, \lambda_M] \in \mathbb{R}^M$
M # of parameters of μ_λ

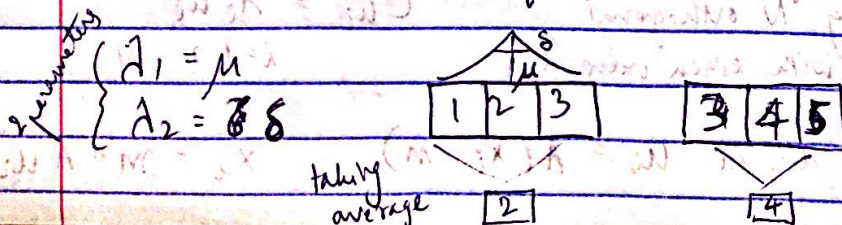
(Statistics) Score funcⁿ.

this is the gradient (partial derivatives) w.r.t the parameters λ of the logarithm (natural log) of the likelihood function (log-likelihood)

Score function: $G_\lambda^x = \nabla_\lambda \log \mu_\lambda(x)$

$= \nabla_\lambda \log P(x/\mu_\lambda)$

The gradient quantifies the contribution of individual parameters \rightarrow How to change λ (of μ_λ) to better fit data (x)?



PCA = compression.
Fisher ... = encoding

Since $G_{\lambda}^x \in \mathbb{R}^M \Rightarrow$ it does not depend on T (size of data)
where $(M \ll T)$ (depends on M , not T)

$(M) \ll T$
parameters of probability density function.
(representation) (data)

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

This is encoding, not compression. we can generate new data. After fitting, we have generative data.

Encoding because after I find M dimensions, I can generate data

From information geometry, $U = \{\mu_{\lambda} : \lambda \in \Lambda\}$

is a Riemannian manifold M_{Λ} with a local metric,

given by Fisher Information Matrix F_{λ}

Topology space with local euclidean metric.

(manifold: hyperdimensional space locally, euclidean distance (can be used, not globally))

$$F_{\lambda} = E_{x \sim \mu_{\lambda}} [G_{\lambda}^x G_{\lambda}^{xT}]$$

Using Fisher Kernel K_{FK} to measure the similarity between two samples x and y .

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

Using Cholesky decomposition: $F_{\lambda}^{-1} = L_{\lambda}^T L_{\lambda}$ lower triangular matrix

We can re-write K_{FK}

$$K_{FK}(x, y) = g_{\lambda}^{xT} g_{\lambda}^y$$

with $g_{\lambda}^x = L_{\lambda} G_{\lambda}^x$

Fisher vector

$$= L_{\lambda} \nabla_{\lambda} \log P(x/\mu_{\lambda})$$