

Data Reduction Fundamentals

1

Principal Component Analysis

- * Based on the Singular Value Decomposition
- * Non-parametric approach towards denoising
- * Unlikely leads to shrinkage in one dimension of the data

$$\text{Variance } \sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Q Why is statistical variance divided by " $n-1$ " and not " n "?

That is, why not $S^2 = \frac{\sum_i (x_i - \bar{x})^2}{n}$? Observe that

$$S^2 = \frac{1}{n} \left[\sum_i (x_i^2 + \bar{x}^2 - 2\bar{x}x_i) \right] = \frac{1}{n} \sum_i x_i^2 - \bar{x}^2$$

Now, if we compute S^2 over many mini-batches of size n ,

$$\mathbb{E}[S^2] = \frac{1}{n} \sum_i \mathbb{E}[x_i^2] - \mathbb{E}[\bar{x}^2]$$

If $x_i \sim \mathcal{N}(\mu, \sigma^2)$ then $\mathbb{E}[x_i^2] = \sigma^2 + \mu^2$

$$\text{and } \mathbb{E}[\bar{x}^2] = \mathbb{E}\left[\left(\frac{1}{n} \sum_i x_i\right)^2\right] = \frac{1}{n^2} \mathbb{E}\left[\sum_i x_i^2 + 2 \sum_i x_i \bar{x}\right]$$

$$= \frac{1}{n^2} \left[(\sigma^2 + \mu^2) + 2 \sum_i \mathbb{E}[x_i \bar{x}] \right] = \cancel{\frac{n-1}{n} \sigma^2}$$

$n(n-1)\sigma^2$ (Assumes x_i, \bar{x} are iid)

$$= \bar{x}^2 + \frac{\sigma^2}{n}$$

$$\Rightarrow \mathbb{E}[S^2] = \sigma^2 + \bar{x}^2 - \left(\bar{x}^2 + \frac{\sigma^2}{n} \right) = \sigma^2 \left(\frac{n-1}{n} \right)$$

Thus, to match ensemble average of i.i.d variances with the theoretical ("generator") average, we

define $S^2 = \frac{\sum_i (x_i - \bar{x})^2}{n-1}$ instead.

→ This argument is due to Bessel but there are a lot of non technical reasons as well to prefer $(n-1)$ in the denominator

consequently, define covariance as $\text{cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$

If $\text{cov}(X, Y) > 0$ we say variables are positively correlated, $\text{cov}(X, Y) < 0$ neg. correlated.

$\text{cov}(X, Y) = 0$ variables are independent.

(This is wrong! Independence $\Rightarrow \text{cov}(X, Y) = 0$)

Part dependent X, Y can be made to have $E(XY) = E(X)E(Y)$!

e.g. Take RV X such that $E(X) = 0$ and $E(X^3) = 0$

i.e., $X \sim N(0, \sigma^2)$ then $\text{cov}(X, Y) = E(X^2) = \sigma^2$

$$f(Y = X^2)$$

but Y, X are

clearly
dependent

Linear Independence:

If vectors x_1, x_2, \dots, x_k are LI then if

$$\sum c_j x_j = 0 \Rightarrow c_1 = c_2 = \dots = c_k = 0$$

be written

That is no vector can be written as a combination of the rest. Basically linear dependence is a measure of redundancy among the vectors.

Span: The span of a set of vectors is the set of all vectors generated by linearly combining them.

Basis: If the spanning vectors are linearly independent and # of such vectors = dimension of vector space then the span is the basis for that vector space

Orthogonal / Orthonormal:

If the basis vectors form a set have mutually \neq dot product and if they have each norm 1

Using the Gram-Schmidt ritual you can take any linearly independent set of vectors and produce an orthogonal set of vectors. (2)

Start with v_1, v_2, \dots, v_n . Define $\text{proj}_v u$

$$\text{Then } u_k = v_k - \sum_{j=1}^{k-1} \frac{\langle u, v_j \rangle}{\langle v_j, v_j \rangle} v_j.$$

The idea behind PCA:

Noise information cannot have high variance or else it would be something interesting is, a signal.

Consequently, if we can find the basis of representation for a dataset in which the axes are aligned towards directions of high variance (signal axes) we can read off qualitative information more easily.

$$P \underset{\text{matrix}}{x} = \underset{\text{matrix}}{y} \quad P_i \in \mathbb{R}^m \quad x_i, y_i$$

$$\begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_m \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} = \begin{bmatrix} y_1 & y_2 & \dots & y_n \end{bmatrix}$$

$$\underline{x}^T = (\langle p_1, x \rangle, \langle p_2, x \rangle, \dots, \langle p_m, x \rangle)$$

Note that this is equivalent to scaling + stretching of x over m axes

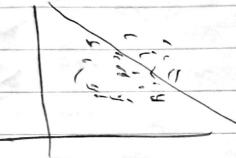
That is, we are changing the basis of x from \mathbb{R}^m to \mathbb{P} and the new representation is $y \in \mathbb{R}^m$ spanned by P .

What is key here is that change of basis does not modify the data itself, it just provides a new viewpoint \rightarrow It is the same as rotation/translation of axes in 2D, 3D.

As described earlier, we want to represent our data in a basis where the noise variance is minimized. That is, if we call

$$LNR = \frac{\sigma_{\text{signal}}^2}{\sigma_{\text{noise}}^2}$$

we want to
re-express X such
that
 LNR is maximized.

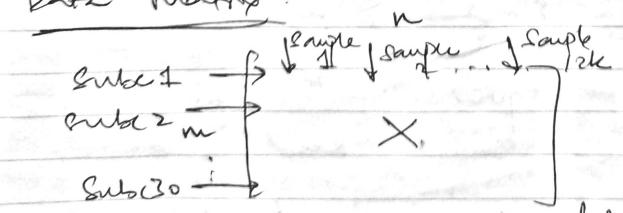


here there is
no basis where
SNR is high because
noise variance stays
high (cannot distinguish
from signal)

there ~~exists~~ exists a
rotation of axes where
the highly "correlated"
samples lie on

This is our first
principal component

Data Matrix:



Basically row
dimension is object
dimension, column
dimension is
time dimension

We subtract the mean across row dimension
(over time) to center our data.

(3)

$$\text{Thus, we have } \hat{\Sigma}_x = \frac{1}{n-1} \mathbf{X} \mathbf{X}^T$$

(Note that we
are now
calling
 \mathbf{X} as row vectors i.e., $\mathbf{X} \in \mathbb{R}^{n \times n}$).

- Covariance matrix is symmetric square
- The diagonal is packed with variance of individual subchannels
- The off diagonal is filled with cross variance or correlation b/w different subchannels.

Signal Energy is on diagonal, noise energy is off diagonal, can we find a change of basis transformation such that $\hat{\Sigma}_x \rightarrow \hat{\Sigma}_y$ is diagonalized.

We want change of basis / PCA matrix to be orthonormal i.e., $\langle p_i, p_j \rangle = \delta_{ij}$ so that we can find directions of descending variance by iteratively eliminating previously exposed orthogonal directions.

$$\begin{aligned}\hat{\Sigma}_y &= \frac{1}{n-1} \mathbf{Y} \mathbf{Y}^T = \frac{1}{n-1} (\mathbf{P} \mathbf{X}) (\mathbf{P} \mathbf{X})^T \\ &= \frac{1}{n-1} \mathbf{P} (\mathbf{X} \mathbf{X}^T) \mathbf{P}^T \\ &\Rightarrow \mathbf{P} \hat{\Sigma}_x \mathbf{P}^T.\end{aligned}$$

So, we need to find \mathbf{P} s.t for any $\hat{\Sigma}_x$, $\hat{\Sigma}_y$ is diagonalized.

Recall, eigenvalue diagonalization \rightarrow any square matrix $\Sigma_x = \Lambda D \Lambda^T$ where D is a diag mat \times composed of eigenvalues and Λ has corresponding eigenvectors.

④ The Karhunen-Loeve expansion is a generalization of the e.v basis of cov matrices for arbitrary random process

If we pick $P = \Lambda^T$ then $\Sigma_y = D$ as
 $P^{-1} = P^T = \Lambda^{-1}$.

In summary, the change of basis is the eigenbasis of the covariance matrix of the signal data
 Another ic, the new representation in the Karhunen-Loeve basis yields a diagonal cov.

Therefore, if we sort the eigenvalues in a descending order, we have the most to least significant principal components

If you drop the last k eigenvectors, you have a lower dimensional representation of your data.

$$\text{Proportion of Variance} \rightarrow \frac{\lambda_1 + \lambda_2 + \dots + \lambda_r}{\lambda_1 + \lambda_2 + \dots + \lambda_m}$$

where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$.

when you swap eigenvalue decoupling with its parent the SVD you obtain LSA / Linear Dimensional

via the Eckart-Young-Murphy Analysis
 Method (Next Lecture)