

THE MATH BEHIND PCA

DATA

Suppose we have m data points in \mathbb{R}^n , ie. each data point has n coordinates (usually called features in machine learning) with respect to an arbitrary vector basis $\{e_1, \dots, e_n\}$:

$$pc^{(i)} = \sum_{j=1}^n X_{ij} e^{(j)} \quad i=1, \dots, n.$$

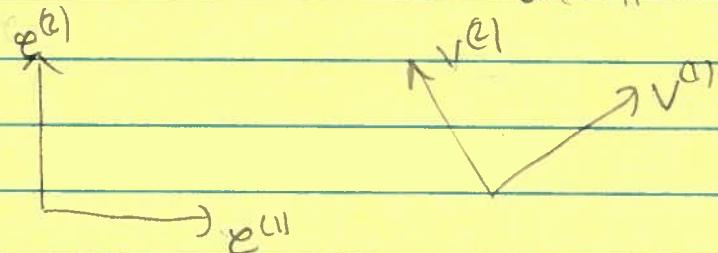
where $x^{(i)}, e^{(j)} \in \mathbb{R}^n$. Together, the data pts and the basis $e^{(i)}$ define a $m \times n$ matrix

$$X = (X_{ij})_{i=1 \dots m; j=1 \dots n}$$

Now perform SVD:

$$X = U \Sigma V^T$$

The columns of V , denoted $v^{(i)}$, define a new basis in \mathbb{R}^n :



$$V = \begin{bmatrix} 1 & 1 & 1 \\ v^{(1)} & v^{(2)} & v^{(3)} \\ 1 & 1 & 1 \end{bmatrix}$$

CHANGE
OF BASIS

In component form, the SVD says:

$$x_{ij} = \sum_k (U\Sigma)_{ik} (V^T)_{kj} \quad (*)$$

But

$$x_{ij} = \hat{x}_j^{(i)}$$

$$(V^T)_{kj} = v_{jk} = v_j^{(k)}$$

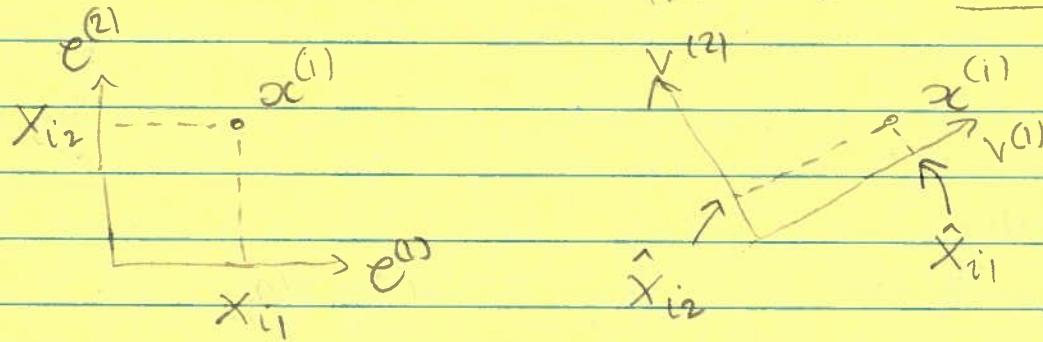
Thus $(*) \Rightarrow$

$$\hat{x}_j^{(i)} = \sum_k \hat{x}_{ik} v_j^{(k)} \Rightarrow \begin{cases} \hat{x}_i^{(i)} = \sum_k \hat{x}_{ik} v_i^{(k)} \\ i=1, \dots, n. \end{cases}$$

where:

$\hat{x}_{ik} = k^{\text{th}}$ component of $\hat{x}_i^{(i)}$ wrt basis $v^{(1)}, \dots, v^{(n)}$.

$$= (U\Sigma)_{ik}, \text{ i.e. } \boxed{\hat{x} = U\Sigma}$$



Thus V provides a new basis and $U\Sigma$ provides the coordinates w.r.t. that new basis.

COVARIANCE OF FEATURES Let X_i = random variable, realizations of which lie in the i^{th} column of X , ie. the n samples of X_i are x_{1i}, \dots, x_{ni} .

Let us now compute the covariance of X_i and X_j :

$$\text{cov}(X_i, X_j) = E[(X_i - E[X_i])(X_j - E[X_j])]$$

Now,

$$E[X_i] = \frac{x_{1i} + \dots + x_{ni}}{m} \quad (\textcircled{a})$$

Let us suppose we have "mean normalized" the data, ie:

$$x^{(i)} \leftarrow x^{(i)} - \frac{x^{(1)} + \dots + x^{(m)}}{m}$$

Then $(\textcircled{a}) = 0$ and the covariance collapses to

$$\begin{aligned} \text{cov}(X_i, X_j) &= E[X_i X_j] \\ &= \frac{1}{m} \sum_{k=1}^m x_{ki} x_{kj} \\ &= \frac{1}{m} \sum_k (x^T)_{ik} x_{kj} \end{aligned}$$

$$= \frac{1}{m} (X^T X)_{ij}$$

In general, all elements of $X^T X$ will be non-zero, ie. all features are correlated w/ one another.

Contrast that with the covariance of the new features (coordinates) defined by the SVD basis $\{e^{(i)}\}$:

$$\hat{X}^T \hat{X} = (U \Sigma)^T (U \Sigma)$$

$$= \Sigma^T U^T U \Sigma$$

$$= \Sigma^T \Sigma$$

$$= \begin{bmatrix} \sigma_1^2 & & & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & & \sigma_r^2 & 0 \\ & & \cdots & \cdots \end{bmatrix} \quad (= n \times n)$$

Thus in the SVD basis, the features are independent (uncorrelated), and their variance is:

$$\text{var}(\hat{x}_{i \cdot}) = \text{cov}(\hat{x}_{i \cdot}, \hat{x}_{i \cdot})$$

$$= (X^T X)_{ii}$$

$$= \sigma_i^2$$

In summary, the spread of the data
pts along $v^{(i)}$, as measured by $\text{var}(\hat{x}_i)$,
is σ_i^2 .

PCA / LSA (cf wiki/Latent-semantic-analysis #
Derivation)

github:

Math 105A → lab-worksheets - Lab10.

Add this to p.6 of Lab10.pdf:

The eigenvectors ~~are~~ ~~the~~ ~~are~~ are &
the ~~are~~ ~~the~~ expression of the "covariance
matrix" are $V^{(1)}$. Here's how:

$$X = \begin{bmatrix} & & & \leftarrow n \rightarrow \\ & \uparrow & & \\ & m & & \downarrow \\ & & & \end{bmatrix}$$

$$X = U \Sigma V^T$$

where $\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \sigma_{m-1}^2 & \\ & & & \sigma_m^2 \end{bmatrix}$ $\sigma_1 \geq \sigma_2 \dots \geq \sigma_m$

$$V = \begin{bmatrix} 1 & 1 \\ \sqrt{v^{(1)}} & \dots & \sqrt{v^{(n)}} \\ 1 & 1 \end{bmatrix}$$

x_i = ~~the~~ values of a R.V.
= i^{th} column of X .

$$\text{cov}(x_i, x_j) = \frac{1}{m} (x^T x)_{ij}$$

(2)

$$X^T X v^{(i)}$$

$$= (U\Sigma V^T)^T (U\Sigma V^T) v^{(i)}$$

$$= V \sum U^T \underbrace{U \sum V^T}_{\text{I}} V v^{(i)}$$

$$= V \sum \sum \underbrace{\text{I}}_{\begin{bmatrix} \sigma_1^2 & \\ & \sigma_n^2 \end{bmatrix}} \underbrace{V^T v^{(i)}}_{\begin{bmatrix} -v^{(i)} \\ \vdots \\ v^{(i)} \end{bmatrix}}$$

$$- \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix}$$

$$= V \begin{bmatrix} 0 \\ \vdots \\ \sigma_i^2 \\ \vdots \\ 0 \end{bmatrix}$$

~~$$\begin{bmatrix} 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 \end{bmatrix} = \begin{bmatrix} 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 0 \end{bmatrix}$$~~

$$= \begin{bmatrix} \sqrt{0} & \dots & \sqrt{0} & \dots & \sqrt{0} & \dots & \sqrt{0} \\ | & \ddots & | & \ddots & | & \ddots & | \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ \sigma_i^2 \\ \vdots \\ 0 \end{bmatrix}$$

$$= \sigma_i^2 \begin{bmatrix} \sqrt{0} \\ \vdots \\ \sqrt{0} \end{bmatrix} = \sigma_i^2 v^{(i)}$$

(3)

We have proved that

$$(X^T X) v^{(i)} = \sigma_i^2 v^{(i)}$$

The proof that SVD of X or eigendecomposition of $X^T X$ is equivalent to an optimization problem where one successively finds eigenvectors that maximize variance of data along them can be found in "Details" section of wikipedia article on PCA.

~~Math 1051 - Lab 10: probabilistic PCA paper by Bishop - Eqs (4) and (3) show how to generate observed data using a hidden variable s.t. MLE of parameters in the model yields the principal axes of the data, and the variance of the data along those axes!~~