

PCA

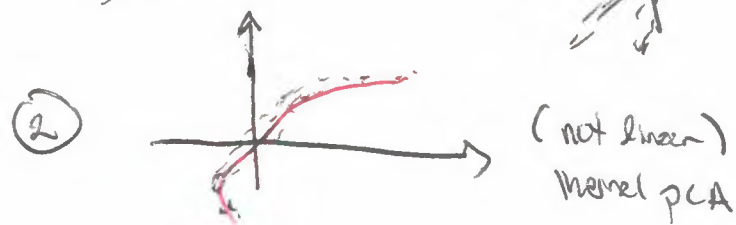
- interpretation of data can be complicated by correlation amongst features and noise in the data
- PCA can completely get rid of correlation, and eliminates the ^{most noisy} 1^{st} dims.

Assumptions

i.e. ^{an orthogonal basis} ① principle directions that correspond to actual signal are \perp to each other (necessary to diagonalize Σ)

② data behave linearly

③ directions w/ higher variance correspond to actual signal



Theory

let $X \in \mathbb{R}^{m \times n}$, and let all cols. be centered by the mean (so we can do linear algebra)

- the sample Cov-Matrix is

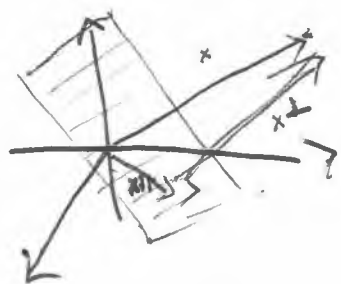
$$\Sigma = \frac{1}{m-1} X^T X; \text{ which is a real, symmetric matrix}$$

- it can be thus diagonalized in a basis of its eigenvectors: (where orthonormal by spectral thm)

$$S = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix} = \begin{bmatrix} - & v_1 & - \\ & \vdots & \\ - & v_n & - \end{bmatrix} \Sigma \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}; \text{ where } \sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2 \text{ by construction}$$

(eliminates all correlation)

- to reduce dimensionality, we project data onto-subspace defined by directions w/ highest signal ($\leq k$ principle directions), the projection can be expressed as:

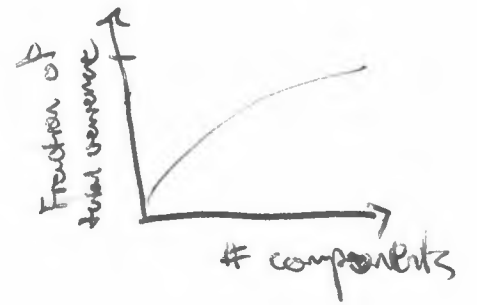


$$\begin{aligned} X &= X'' + X^\perp = \alpha_1 v_1 + \dots + \alpha_k v_k + \dots + \alpha_n v_n \\ &= \langle X, v_1 \rangle v_1 + \dots + \langle X, v_k \rangle v_k + \dots + \langle X, v_n \rangle v_n \end{aligned}$$

$$\Rightarrow \left\{ \begin{array}{l} X_{PCA} \\ m \times k \end{array} \right\} = X \begin{bmatrix} v_1 & \dots & v_k \\ \vdots & & \vdots \end{bmatrix} \begin{array}{l} n \times k \\ m \times n \quad n \times k \end{array}$$

notes

- usually use truncated SVD to obtain $\sum \sigma^2$ eigenvalues / eigenvectors of $X^T X$
since SVD is more numerically stable and ~~finding~~ computing eigenvectors of an $n \times n$ matrix might be daunting
- usually need to center (since we are looking for a subspace that goes through the center) and scale to unit variance (since 1 dimension could have high variance purely due to units).
- can get a handle on # components to use by looking at uses



- visualization in 2D
- reducing # predictors for a supervised method (or even in a clustering method)
- data compression
- can eliminate noisy dimensions w/ no meaningful full signal.

other ways of interpreting PCA

- a plane that goes through center w/ minimal sum of distances to plane.



- a low rank matrix that is as close as possible to X :

$$X_{\text{pca}} = \underset{\substack{\tilde{X} \in \mathbb{R}^{m \times n} \\ \text{rank}(\tilde{X}) = k}}{\text{argmin}} \|\tilde{X} - X\|_F^2$$

- 1st PC = direction w/ max variance
- 2nd PC = " " " orthogonal to 1st
- 3rd PC = " " " " previous
- ⋮

(Follows From SVD)