

PCA Two Ways

Nick Landolfi
Stanford University

Outline

- ▶ background: affine sets, projections, extremal trace
- ▶ minimum-residual affine set
- ▶ maximum-variance affine set
- ▶ examples with protein data

Affine sets

- ▶ a set $M \subset \mathbf{R}^n$ is *affine* if it contains the lines through any two of its points
 - ▶ i.e., $(1 - \lambda)x + \lambda y \in M$ for all $x, y \in M, \lambda \in \mathbf{R}$
 - ▶ other terminology: affine subspace, linear variety, affine variety, flat set
- ▶ the affine sets are the solution sets of linear equations
 - ▶ given conforming A and b the set $\{x \in \mathbf{R}^n \mid Ax = b\}$ is affine, and vice versa
- ▶ the affine sets are translated subspaces
 - ▶ if M is affine, there exists a unique $a \in \mathbf{R}^n$ and subspace $S \subset \mathbf{R}^n$ so that $M = a + S$
 - ▶ notation $a + S$ means $\{a + x \mid x \in S\}$; dimension of M is dimension of S
- ▶ concrete representation for $M = a + S$ is $a + \text{range}(U)$ where $\text{range}(U) = S$ and $U^\top U = I$

Projection onto affine set

- ▶ given $a \in \mathbf{R}^n$ and $U \in \mathbf{R}^{n \times k}$ with $U^\top U = I$
- ▶ question: what is the projection of $x \in \mathbf{R}^n$ onto $a + \text{range}(U)$
- ▶ find $z \in \mathbf{R}^k$ to minimize

$$\|a + Uz - x\| = \|Uz - (x - a)\|$$

- ▶ solution is $z^* = U^\top(x - a)$
- ▶ projection is $Uz^* + a = UU^\top x + (I - UU^\top)a$

Extremal trace problem

▶ problem: given $A = A^\top$, find $U \in \mathbf{R}^{n \times k}$ to

$$\text{maximize} \quad \text{trace}(U^\top A U)$$

$$\text{subject to} \quad U^\top U = I$$

▶ solution: pick first k (orthonormal) eigenvectors

▶ let $A = Q \Lambda Q^\top$ be an eigendecomposition with $\lambda_1 \geq \dots \geq \lambda_n$

▶ then $U^* = \begin{bmatrix} q_1 & \dots & q_k \end{bmatrix}$ is a solution

▶ “a solution”, since any permutation obtains same objective value

Extremal trace diagonalized problem

- ▶ $A = Q\Lambda Q^\top$ with $\lambda_1 \geq \dots \geq \lambda_n \geq 0$
- ▶ parameterize columns of U by basis Q ; i.e., $U = QZ$ where $Z \in \mathbf{R}^{n \times k}$
 - ▶ columns of Z give coordinates of U in basis Q
 - ▶ there is one-to-one correspondence between U and Z
- ▶ in new coordinates, we find Z to

$$\begin{aligned} & \text{maximize} && \text{trace}(Z^\top \Lambda Z) \\ & \text{subject to} && Z^\top Z = I \end{aligned}$$

- ▶ since $U^\top U = 1$ if and only if $Z^\top Z = 1$ and $U^\top AU = Z^\top \Lambda Z$
- ▶ we have *diagonalized* the problem; changed coordinates to Q

Extremal trace diagonalized objective

- ▶ we have

$$\text{trace}(Z^T \Lambda Z) = \sum_{j=1}^k \lambda_j \tilde{z}_j^T \tilde{z}_j = \sum_{j=1}^k \lambda_j \|\tilde{z}_j\|^2 \leq \sum_{i=1}^k \lambda_i$$

since

- ▶ $\|\tilde{z}_i\|^2 \leq \|Z\|^2 = 1$; i.e., the rows of an orthonormal matrix have norm bounded by 1
- ▶ $\sum_{i=1}^n \|\tilde{z}_i\|^2 = \|Z\|_F^2 = k$; i.e., the sum of squares elements of an orthonormal matrix is bounded by k
- ▶ we can achieve this upper bound by selecting $Z^* = [e_1 \ \cdots \ e_k]$
- ▶ this choice corresponds to $U^* = QZ^* = [q_1 \ \cdots \ q_k]$

Minimum-residual affine set

- ▶ given dataset $x_1, x_2, \dots, x_m \in \mathbf{R}^n$
 - ▶ define $X = \begin{bmatrix} x_1 & \cdots & x_m \end{bmatrix}$, $\bar{x} = (1/m)X1$, and $\bar{X} = X - (1/m)X11^\top = (I - (1/m)11^\top)X$
- ▶ we want to find the k -dimensional affine set “closest to” data
- ▶ problem: find $a \in \mathbf{R}^n$ and $U \in \mathbf{R}^{n \times k}$ (giving affine set $M_{a,U}$) to minimize

$$\sum_{i=1}^m \|x_i - \text{proj}_{M_{a,U}}(x_i)\|^2$$

- ▶ solution: pick $a^* = \bar{x}$ and U to have columns first k eigenvectors of $\bar{X}\bar{X}^\top$
 - ▶ eigenvectors of $\bar{X}\bar{X}^\top$ are first k left singular vectors of \bar{X} (right singular vectors of \bar{X}^\top)

Minimum-residual affine set, offset

▶ fix $U \in \mathbf{R}^{n \times k}$, $U^\top U = I$

▶ find $a \in \mathbf{R}^n$ to minimize

$$\sum_{i=1}^m \|x_i - UU^\top x_i - (I - UU^\top)a\|^2 = \left\| \begin{bmatrix} I - UU^\top \\ \vdots \\ I - UU^\top \end{bmatrix} a - \begin{bmatrix} I - UU^\top & \cdots & 0 \\ & \ddots & \\ 0 & \cdots & I - UU^\top \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \right\|^2$$

▶ solution a^* must satisfy normal equations, helps to notice $(I - UU^\top)^2 = (I - UU^\top)$

▶ normal equations are $(I - UU^\top)a^* = (I - UU^\top)\bar{x}$

▶ $a^* = \bar{x}$ works, and *does not depend on U*

Minimum-residual affine set, subspace

► assume $M = \bar{x} + \text{range}(U)$

► find $U \in \mathbf{R}^{n \times k}$ to

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m \|(I - UU^\top)(x_i - \bar{x})\|^2 \\ & \text{subject to} && U^\top U = I \end{aligned}$$

► $\|(I - UU^\top)(x_i - \bar{x})\|^2 = \|x_i - \bar{x}\|^2 - \|U^\top(x_i - \bar{x})\|^2$, first term constant w.r.t. U

► $\sum_{i=1}^m \|U^\top(x_i - \bar{x})\|^2 = \|U^\top \bar{X}\|_F^2 = \text{trace}(\bar{X}^\top UU^\top \bar{X}) = \text{trace}(U^\top \bar{X} \bar{X}^\top U)$

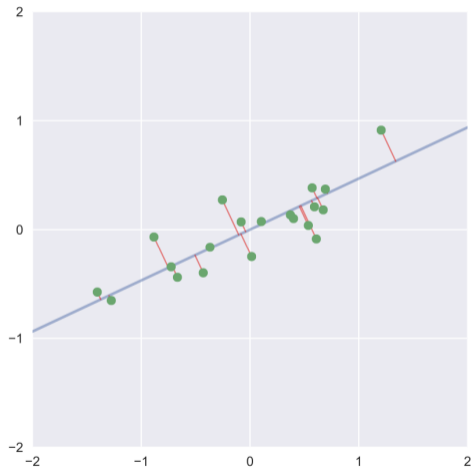
► so we want to find $U \in \mathbf{R}^{n \times k}$ to

$$\begin{aligned} & \text{maximize} && \text{trace}(U^\top \bar{X} \bar{X}^\top U) \\ & \text{subject to} && U^\top U = I \end{aligned}$$

► an extremal trace problem, solution is first k eigenvectors of $\bar{X} \bar{X}^\top$

Total least squares

- ▶ measure distances orthogonal to line



Maximum-variance affine set

- ▶ given data set $x_1, \dots, x_m \in \mathbf{R}^n$
- ▶ we want to find the k -dimensional affine subspace in which our data has “maximal variance”
- ▶ for affine subspace M , define the *projected mean* and *projected variance* by

$$\bar{x}(M) = \frac{1}{m} \sum_{i=1}^m \text{proj}_M(x_i) \quad \text{and} \quad \nu(M) = \frac{1}{m} \sum_{i=1}^m \|\text{proj}_M(x) - \bar{x}(M)\|^2$$

- ▶ problem: find $a \in \mathbf{R}^n$ and $U \in \mathbf{R}^{n \times k}$

$$\begin{aligned} & \text{maximize} && \nu(a + \text{range } U) \\ & \text{subject to} && U^\top U = I \end{aligned}$$

- ▶ solution pick columns of U to be the first k eigenvectors $\bar{X}\bar{X}^\top$, any $a \in \mathbf{R}^n$ works

Maximum-variance affine set solution

- ▶ express $\bar{x}(a + \text{range } U)$ as

$$\frac{1}{m} \sum_{i=1}^m UU^\top x_i + (I - UU^\top)a = UU^\top \bar{x} + (I - UU^\top)a$$

- ▶ drop the constant $1/m$ and write the objective

$$\sum_{i=1}^m \|\text{proj}_M(x) - \bar{x}(S)\|^2 = \sum_{i=1}^m \|UU^\top x_i - UU^\top \bar{x}\|^2 = \sum_{i=1}^m \|UU^\top(x_i - \bar{x})\|^2$$

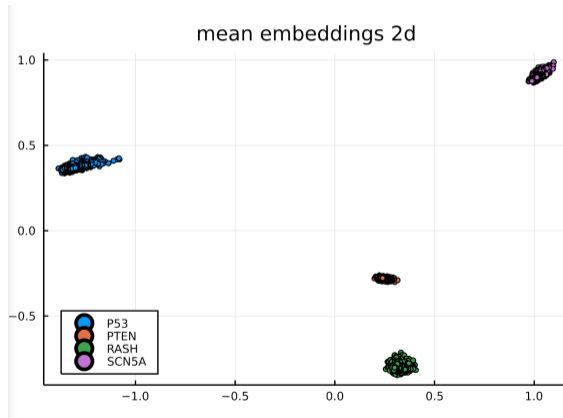
- ▶ since U is orthonormal, $\|UU^\top(x_i - \bar{x})\| = \|U^\top(x_i - \bar{x})\|$, a familiar expression
- ▶ the variance of the projected points does not depend on a
- ▶ so we want to find $U \in \mathbf{R}^{n \times k}$ to

$$\begin{aligned} & \text{maximize} && \text{trace}(U^\top \bar{X} \bar{X}^\top U) \\ & \text{subject to} && U^\top U = I \end{aligned}$$

- ▶ an extremal trace problem, pick the first k eigenvectors of $\bar{X} \bar{X}^\top$

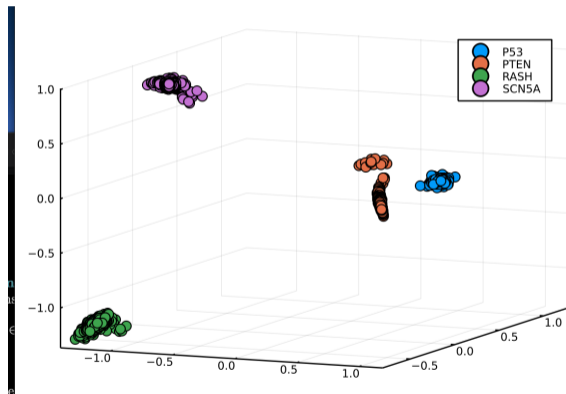
Protein embeddings 2d

- ▶ train a big neural network which maps proteins to vectors in \mathbf{R}^{1024}



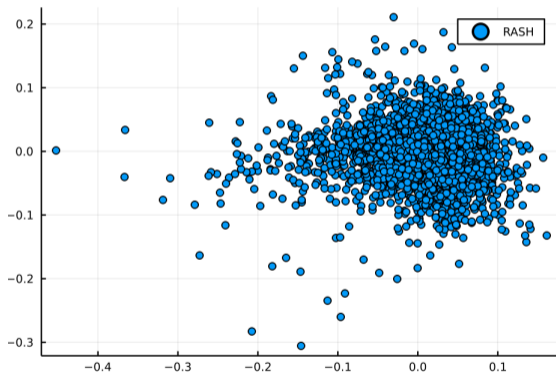
Protein embeddings 3D

- ▶ train a big neural network which maps proteins to vectors in \mathbf{R}^{1024}



Rash embeddings 3D

► RASH protein family



Rash embeddings 3D

► RASH protein family

