

scalar is its own trace  $\alpha = \text{Tr}(\alpha)$

## Principal component analysis

collection of  $m$  points -  $\{x^{(1)} \dots x^{(m)}\}$

• Apply lossy compression to points -

These points can be encoded to represent  
a lower dimensional version of them.

For each point  $x^{(i)} \in \mathbb{R}^n$  find corresp.  
code vector  $c^{(i)} \in \mathbb{R}^d$ .

If  $d < n$  storing these code points

will occupy lesser memory than original data.

encoding function such that  $f(x) = c$

decoding function  $\pm x \approx g(f(x))$

PCA can be used as decoding func.

For a simple decoder matrix mul. can be used to map code back to  $\mathbb{R}^n$ .

$g(c) = Dc$  where  $D \in \mathbb{R}^{n \times d}$  - matrix defining decoding.

Computing the optimal code for this decoding problem so to make it simpler we constrain the columns of  $D$  to be orthogonal to each other.

The columns of  $D$  also have to be constrained to have unit norm to avoid being scaled causing to possibility of many solutions.



To generate optimal code point  $c^*$  for each input  $x$ , minimize dist between  $x$  & reconstruction  $g(c^*)$ . Use norm to measure distance (refer previous notes for norm). PCA uses  $L^2$  norm.

$$c^* = \arg \min_c \|x - g(c)\|_2$$

squared  $L^2$  norm can be used as both squared  $L^2$  norm &  $L^2$  norm are minimized by same value of  $c$  as  $L^2$  norm is non-negative & squaring operation monotonically inc. for non-negative numbers.

$$c^* = \arg \min_c \|x - g(c)\|_2^2$$

function being minimized reduces to

$$\begin{aligned} & (x - g(c))^T (x - g(c)) \\ &= x^T x - x^T g(c) - g(c)^T x + g(c)^T g(c) \\ &= x^T x - 2x^T g(c) + g(c)^T g(c) \end{aligned}$$

(because scalar  $g(c)^T x$  is equal to its transpose).

omit 1st term as it does not depend on  $c$ .

$$c^* = \arg \min_c -2x^T g(c) + g(c)^T g(c)$$

$$g(c) = Dc$$

$$\therefore c^* = \arg \min_c -2x^T Dc + c^T D^T Dc$$

$$= \arg \min_c -2x^T Dc + c^T I_1 c$$

$$= \arg \min_c -2x^T Dc + c^T c.$$

Applying vector calculus (directional derivative)

$$\nabla_c (-2x^T Dc + c^T c) = 0$$

$$-2D^T x + 2c = 0$$

$$c = D^T x.$$



From above equation  $x$  can be encoded optimally using just a matrix vector operation. To encode a

vector  $\div f(x) = D^T x$ .

PCA reconstruction  $\div r(x) = g(f(x)) = DD^T x$ .

Now to choose encoding matrix  $D$ ,

let's revisit the idea of minimizing

the  $L^2$  distance between inputs & reconstruction

Since same matrix ~~can~~<sup>will</sup> be used to decode all points, the points can't be considered in isolation. Instead the

Frobenius norm of matrix of errors computed for all points must be minimized.

$$D^* = \arg \min_D \sqrt{\sum_{i,j} (x_i^{(u)} - r(x_i^{(u)})_j)^2}$$

subject to  $D^T D = I_d$ .

To find  $D^*$ , first case  $d=1$

then  $D^*$  is single vector  $d$ .

$$\text{thus } d^* = \arg \min_d \sum_i \|x^{(i)} - d d^T x^{(i)}\|_2^2$$

subject to  $\|d\|_2 = 1$ .

A scalar is its own transpose

$$\text{so } d^* = \arg \min_d \sum_i \|x^{(i)} - x^{(i)T} d d^T\|_2^2$$

$$\|d\|_2 = 1.$$

$X \in \mathbb{R}^{m \times n}$  - matrix formed by stacking all vectors such that  $x_{a;:} = x^{(a)}$ .

$$d^* = \arg \min_d \|X - X d d^T\|_F^2, \quad d^T d = 1$$

$$= \arg \min_d \text{Tr}((X - X d d^T)^T (X - X d d^T))$$

$$= \arg \min_d \text{Tr}(X^T X - X^T X d d^T d - d d^T X X^T + d d^T X^T X d d^T)$$



$$= \arg \min_d \text{Tr}(X^T X) - \text{Tr}(X^T X d d^T) - \text{Tr}(d d^T X^T X)$$

$$+ \text{Tr}(d d^T X^T X d d^T).$$

$$= \arg \min_d -\text{Tr}(X^T X d d^T) - \text{Tr}(d d^T X^T X)$$

$$+ \text{Tr}(d d^T X^T X d d^T).$$

(as terms not changing  $d$  don't affect  $\arg \min$ )

$$= \arg \min_d -2\text{Tr}(X^T X d d^T) + \text{Tr}(d d^T X^T X d d^T)$$

(inside a trace order of matrices can be cycled).

$$= \arg \min_d -2\text{Tr}(X^T X d d^T) + \text{Tr}(X^T X d d^T d d^T).$$

remember the constraint  $d d^T = 1$

$$\Rightarrow \arg \min_d -\text{Tr}(X^T X d d^T)$$

$$= \arg \max_d \text{Tr}(X^T X d d^T)$$

$$= \arg \max_d \text{Tr}(d^T X^T X d)$$

Now the above problem can be solved using eigendecomposition (if you can't recognize the pattern see my previous notes). Optimal  $d$  is given by eigenvector  $X^T X$  across the largest eigenvalue.

Now above case is only for  $d=1$ . we can use mathematical induction to generalize.