

Theory Questions

① (a) we would like to find  $\vec{w}$  s.t.

$$X\vec{w} = \vec{y} \quad \text{rank}(X) = n$$

$$X \in \mathbb{R}^{n \times d} \quad n < d$$

$$\vec{y} \in \mathbb{R}^n$$

$$\vec{w} \in \mathbb{R}^d$$

denote  $u_i$ : the  $i$ -th column vector in  $X$ . Then

$$X\vec{w} = \sum_{i=1}^d w_i u_i$$

but within the set  $\{u_i\}_{i=1}^d$ , there exist  $n$  orthogonal vectors that can span a space of dimension  $n$ . Therefore, we are able to compose any vector  $\vec{y} \in \mathbb{R}^n$  with this linear combination and in particular  $\vec{y} = X\vec{w}$ .

(b) we would like to solve

$$\min_w \|\vec{w}\|^2 \rightarrow \min_w \frac{1}{2} \|\vec{w}\|^2$$

$$\text{s.t. } X\vec{w} = \vec{y} \rightarrow X\vec{w} - \vec{y} = 0$$

we define the Lagrangian:

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum_i \lambda_i (X\vec{w} - \vec{y})_i = \frac{1}{2} \|\vec{w}\|^2 - \lambda^T (X\vec{w} - \vec{y})$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w} = w - (\lambda^T x)^T \stackrel{!}{=} 0 \Rightarrow w = X^T \lambda$$

$\Rightarrow$  if we denote  $x_i$  the  $i$ -th column vector  $\Rightarrow w^* = \sum_{i=1}^n \lambda_i x_i$

$\vec{\lambda} \in \mathbb{R}^n$   
 $\vec{x}_i \in \mathbb{R}^d$

and here we have  $n$  vectors that represent all the data points

(c)  $v^T w^* = v^T X^T \lambda =$

$v \in \mathbb{R}^d$   
 $x_i \in \mathbb{R}^d$   
 $\lambda \in \mathbb{R}^n$

$$(v_1, v_2, \dots, v_d) \begin{pmatrix} | & & | \\ \vec{x}_1 & \cdots & \vec{x}_n \\ | & \cdots & | \end{pmatrix} \vec{\lambda}$$

$$= (v^T x_1, v^T x_2, \dots, v^T x_n) \vec{\lambda}$$

$$= (K(v, x_1), K(v, x_2), \dots, K(v, x_n)) \vec{\lambda}$$

$$= \sum_{i=1}^n K(v, x_i) \lambda_i$$

$$(2) \text{ (a)} \quad V_i = \phi(x_i) - \frac{1}{m} \sum_{t=1}^m \phi(x_t)$$

$$K_{ij} = \langle V_i, V_j \rangle =$$

$$(\phi(x_i) - \frac{1}{m} \sum \phi(x_t))^T (\phi(x_j) - \frac{1}{m} \sum \phi(x_t))$$

$$= \phi(x_i)^T \phi(x_j) - \frac{1}{m} \sum_t (\phi(x_i)^T \phi(x_j) + \phi(x_i)^T \phi(x_t))$$

$$+ \frac{1}{m^2} \sum_{t \neq i} \phi(x_t)^T \phi(x_t) =$$

$$= \bar{K}_{ij} - \frac{1}{m} \sum_t (\bar{K}_{tj} + \bar{K}_{it}) + \frac{1}{m^2} \sum_{t,t'=1}^m \bar{K}_{tt'}$$

(b) the covariance matrix in the dimension of  $\phi(x)$ :

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \phi(x_i) \phi(x_i)^T$$

The set  $\{u_i\}$  are the eigenvectors of the  $\Sigma$ , (principal components) thus:

$$\Sigma u_i = \sigma_i u_i \quad (\sigma_i \text{ is the eigenvalue of } u_i)$$

$$\left( \frac{1}{m} \sum_{i=1}^m \phi(x_i) \phi(x_i)^T \right) u_i = \sigma_i u_i$$

$$\text{Define } \alpha_{ij} = \frac{\phi(x_i)^T u_j}{m \sigma_j}$$

$$\Rightarrow u_j = \sum_{i=1}^m \alpha_{ji} \phi(x_i)$$

We have  $u_j$  as a linear combination of  $\phi(x_i)$  where  $\vec{\alpha}_j$  is the coefficients vector. In order to find  $\alpha_j$  for a particular  $j$ , we have to solve the equation

$$\sum u_j = \tau_j u_j \quad (\tau_j \neq 0 \text{ for } \rho C) \Rightarrow$$

$$\left( \frac{1}{m} \sum_i \phi(x_i) \phi(x_i)^T \right) \left( \sum_k \alpha_{jk} \phi(x_k) \right) = \tau_j \sum_i \alpha_{ji} \phi(x_i)$$

$$\frac{1}{m} \sum_i \phi(x_i) \sum_k \alpha_{jk} \phi(x_k)^T \phi(x_k) = \tau_j \sum_i \alpha_{ji} \phi(x_i)$$

$$\frac{1}{m} \sum_i \phi(x_i) \sum_k \alpha_{jk} K(x_i, x_k) = \tau_j \sum_i \alpha_{ji} \phi(x_i)$$

Multiplying from left with some vector  $\phi(x_l)^T$

$$\frac{1}{m} \sum_i \phi(x_l)^T \phi(x_i) \sum_k \alpha_{jk} K(x_i, x_k) = \tau_j \sum_i \alpha_{ji} \phi(x_l)^T \phi(x_i)$$

$$\frac{1}{m} \sum_i K(x_l, x_i) \sum_k \alpha_{jk} K(x_i, x_k) = \tau_j \sum_i \alpha_{ji} K(x_l, x_i)$$

$$\left( \frac{1}{m} \bar{K} \bar{K} \bar{\alpha}_j \right)_l = \tau_j (\bar{K} \bar{\alpha}_j)_l \quad \forall l = 1 \dots m$$

$$\bar{K} \bar{\alpha}_j = \bar{K} m \tau_j \bar{\alpha}_j \Rightarrow \bar{K} \bar{\alpha}_j = m \tau_j \bar{\alpha}_j$$

$\Rightarrow$  we have  $\bar{\alpha}_j$  if we diagonalize  $\bar{K}$

$$(c) \quad \langle u_j, \phi(x') \rangle = u_j^\top \phi(x') = (\text{from } b)$$

$$= \sum_{i=1}^m \alpha_{ji} \phi(x_i)^\top \phi(x') = \sum_{i=1}^m \alpha_{ji} K(x_i, x')$$

So if we can calculate kernel we can project any  $\phi(x')$  on the subspace spanned by  $\{u_j\}_{j=1}^k$

If the kernels are calculated with complexity  $O(d) \Rightarrow$   
 $\langle u_j, \phi(x') \rangle$  is calculated with  $O(nd)$   
 if we would like to project on the subspace of  $\{u_j\}_{j=1}^k$  than the complexity becomes  $O(knd)$

$$(3) \quad X \sim \lambda e^{-\lambda x} \quad L(x) = \prod_{i=1}^n p(x_i)$$

$$\log(L(x)) = \sum_{i=1}^n \log p(x_i) = \sum_{i=1}^n \log \lambda - \lambda x_i \\ = n \log \lambda - \sum_{i=1}^n \lambda x_i$$

$$(a) \quad 0 \stackrel{!}{=} \frac{\partial \log(L\lambda)}{\partial \lambda} = \frac{n}{\lambda} - \sum_i x_i$$

$$\Rightarrow \hat{\lambda} = \frac{n}{\sum_{i=1}^n x_i} = \frac{1}{\bar{x}}$$

(b) If we sample  $n$  (iid) values from exp. distribution with  $\lambda = 1$  the posterior would then be

$$p(\lambda = \lambda | X_1 = x_1, \dots, X_n = x_n) = \frac{P(\lambda = \lambda) \prod_{i=1}^n P(X_i = x_i | \lambda = \lambda)}{P(X_1 = x_1, \dots, X_n = x_n)}$$

$$\propto P(\lambda = \lambda) \prod_{i=1}^n P(X_i = x_i | \lambda = \lambda) = \\ e^{-\lambda} \prod_{i=1}^n \lambda e^{-\lambda x_i} \quad P(\lambda = \lambda) \propto \exp(1)$$

We can try maximizing the log of the posterior and then we are maximizing

$$\log(e^{-\lambda} \prod_{i=1}^n \lambda e^{-\lambda x_i}) = -\lambda + \sum_{i=1}^n (\log(\lambda) + \log(e^{-\lambda x_i})) \\ = -\lambda + n \log(\lambda) - \lambda \sum_i x_i$$

$$\Rightarrow 0 = \frac{\partial}{\partial \lambda} \left( n \log \lambda - \lambda \left( \sum_{i=1}^n x_i - 1 \right) \right) =$$

$$= \frac{n}{\lambda} - \sum_{i=1}^n x_i - 1$$

$$\Rightarrow \lambda = \frac{n}{\sum_{i=1}^n x_i + 1}$$

Easy to see that this is indeed

a maximum  $\frac{\partial^2}{\partial \lambda^2} P(\lambda = \lambda | X_1, \dots, X_n = x_n) < 0$

c) For  $n$  samples and  $K$  possible values of  $\lambda$ , we have  $10$ -parameters of distributions  $Z_i$  - the dist. type of  $i$ -th sample

$$Q(\theta, \theta^t) = \sum_{i=1}^n \sum_{k=1}^K P(Z_i = k | X_i = x_i, \theta^t) \log(P(X_i = x_i, Z_i = k, \theta))$$

Now, if  $Z_i \sim \cup \{1, \dots, K\}$  we have

$$P(X_i = x_i, Z_i = z_i, \theta) = P(X_i = x_i | Z_i = k, \theta) P(Z_i = k, \theta)$$

$$= \lambda_k e^{-\lambda_k x_i} \cdot \frac{1}{K}$$

$$Q(\theta, \theta^t) = \sum_{i=1}^n \sum_{k=1}^K P(Z_i = k | X_i = x_i, \theta^t) \log\left(\frac{1}{K} \lambda_k e^{-\lambda_k x_i}\right)$$

keeping only terms dependent on  $\lambda_k$  ~

$$\sum_{i=1}^n \underbrace{\sum_{k=1}^K P(Z_i = k | X_i = x_i, \theta^t) (\log \lambda_k - \lambda_k x_i)}_{= C_{i, k}^t}$$

$$\sum_i \sum_k C_{i, k}^t (\log \lambda_k - \lambda_k x_i)$$

$\Rightarrow$  deriving with respect to  $\lambda_k$

$$\frac{\partial Q(\theta, \theta^t)}{\partial \lambda_k} = \sum_i C_{i,k}^T \left( \frac{1}{\lambda_k} - x_i \right) \stackrel{!}{=} 0$$

$$\lambda_k = \sum_{i=1}^n C_{i,k}^T \cdot \frac{1}{\sum_{i=1}^n C_{i,k}^T x_i}$$

This is indeed a maximum:

$$\frac{\partial^2 Q}{\partial \lambda^2} = -\frac{1}{\lambda^2} \sum_i C_{i,k}^T < 0$$

(ii) a) Log likelihood:

$$\log \left( \prod_{i=1}^n P(X_i = x_i) \right) = \sum_{i=1}^n \log(P(X_i = x_i)) =$$

$$\sum_i \log \left( \sum_{r=0}^{31} P(X_i = x_i | Z_r = r) P(Z_r = r) \right) =$$

$$\sum_i \log \left( \sum_r \Theta_r \underbrace{P(X_i = x_i | Z_r = r)}_{\checkmark} \right) = *$$

$X$  is a 5-bit representation of  $Z_r$ ,  
with 2 (out of 5) randomly selected  
digits replaced with "2".

So given  $r_i$ , each possible  $x_i$   
has a chance of  $\binom{5}{2}^{-1} = \frac{1}{10}$  to be selected

Define  $I_{Z_r, X_i} = \begin{cases} 1 & Z_r \rightarrow x_i \text{ is a possible} \\ & \text{transformation} \\ 0 & \text{else.} \end{cases}$

Indicates if  $x_i - I_{Z_r, X_i}$   
could be built with  $Z_r$   
and 2-digits replacements.

$$* = \sum_i \log \left( \sum_r \Theta_r I_{Z_r, X_i} \cdot \frac{1}{10} \right)$$

(b) We can think of a mixture of models with parameters  $\{\theta_r\}_{r=0}^{|S|}$  and each sample  $x_i$  is generated from one of these models. For each EM step then, we have ( $\theta \equiv \{\theta_r\}$ )

$$\begin{aligned} Q(\theta, \theta^t) &= \sum_{i=1}^n \sum_{r=0}^{|S|} P(z_i = z_r | x_i = x_i, \theta^t) \log(P(x_i = x_i | z_i = z_r, \theta)) \\ &\equiv \sum_i \sum_r C_{ir}^t \log(P(x_i = x_i | z_i = z_r, \theta)) \\ &= \sum_i \sum_r C_{ir}^t \log(\theta_r I_{z_r, x_i} \cdot \frac{1}{10}) \end{aligned}$$

Now, we need to maximize  $Q(\theta, \theta^t)$  given the constraint:

$$\sum_r \theta_r = 1. \Rightarrow \text{The Lagrangian:}$$

$$L(\theta, \lambda) = \sum_i \sum_r C_{ir}^t \log\left(\frac{\theta_r}{10} I_{z_r, x_i}\right) - \lambda \left(\sum_{r=0}^{|S|} \theta_r - 1\right)$$

$$\frac{\partial L}{\partial \theta_r} = \sum_i C_{ir}^t \frac{1}{\frac{\theta_r}{10} I_{z_r, x_i}} \cdot \frac{I_{z_r, x_i}}{10} - \lambda$$

$$= \frac{1}{\theta_r} \sum_i C_{ir}^t - \lambda = 0 \Rightarrow$$

$$\theta_r = \frac{1}{\lambda} \sum_i C_{ir}^t$$

indeed a maximum

$$\frac{\partial^2 L}{\partial \theta_r^2} = -\frac{1}{\theta_r} \sum_i C_{ir}^t < 0$$

$\lambda$  could be found using the constraint:

$$1 = \sum_{r=0}^{31} \theta_r = \frac{1}{\lambda} \sum_{r=0}^{31} \sum_{i=1}^n C_{i,r}^T =$$

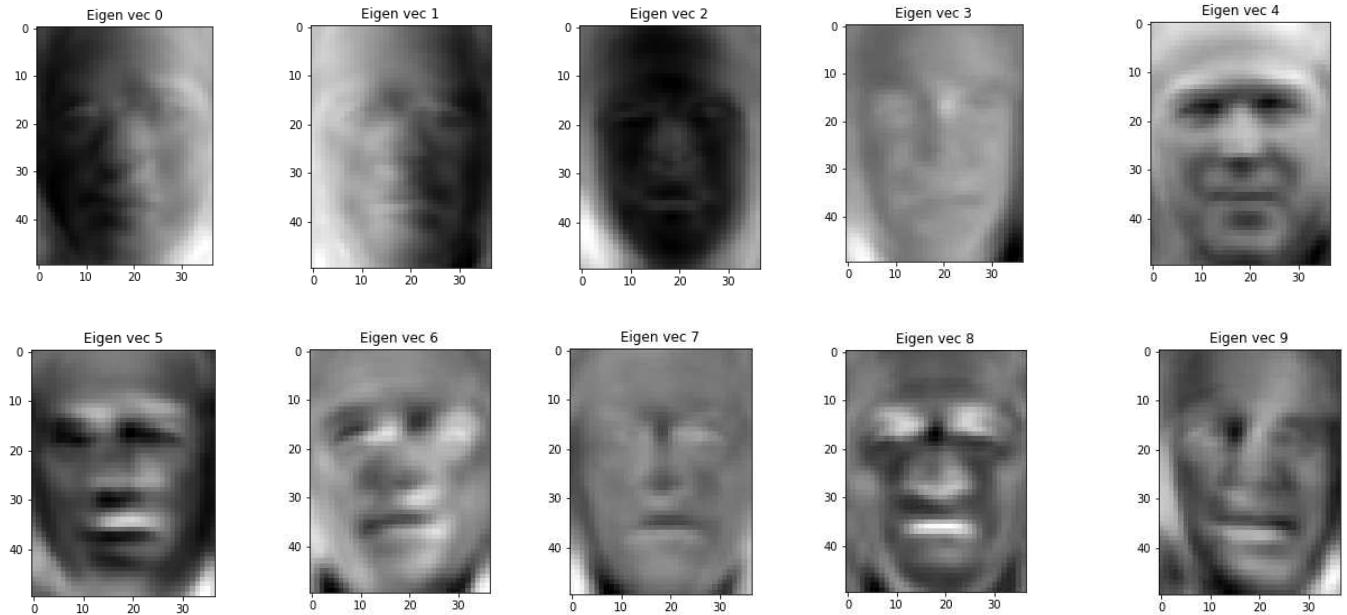
$$\frac{1}{\lambda} \sum_{i=1}^n \underbrace{\sum_{r=0}^{31} C_{i,r}^T}_1 = \frac{n}{\lambda} \Rightarrow \lambda = n //$$

$$\Rightarrow \theta_r = \frac{1}{n} \sum_i C_{i,r}^T //$$

## Programming assignment – PCA

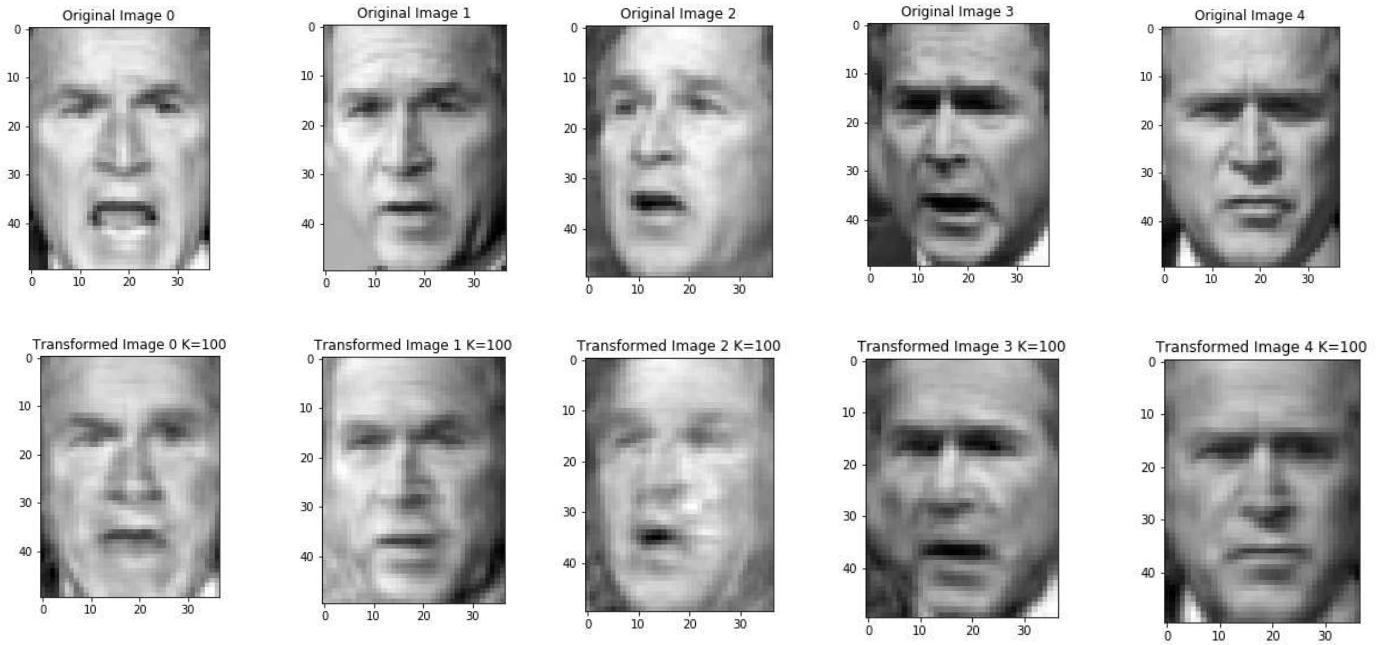
a. See scripts

b. 10 eigenvectors for George W Bush:

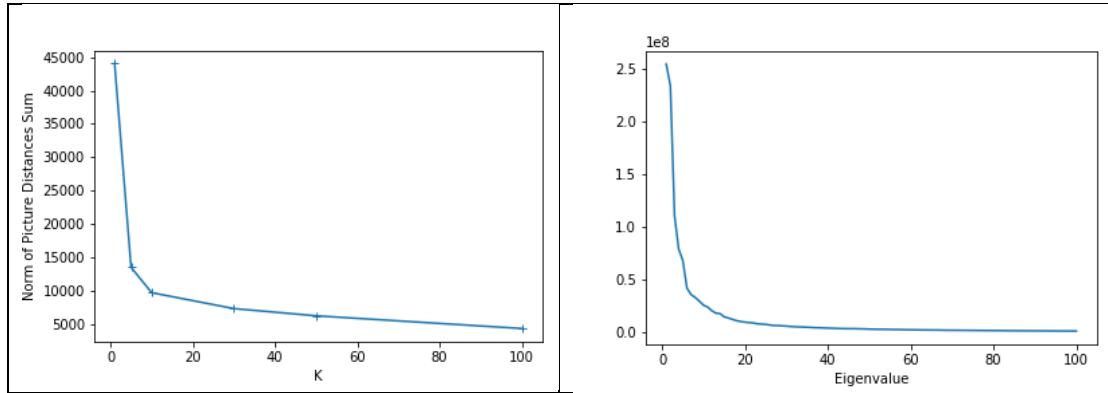


The Eigenvectors are linear combination of features (in our case, pixels) of all the samples. Those linear combinations can be viewed as new features which has the maximal variance across all samples. We can see that some of the eigenvectors give greater weights to pixels in the nose area and some to pixels in the eyes area in the pictures. Some also emphasize the shading of the pictures. each such “trait” represented by the eigenvectors is considered as a trait with high variance within the data (the original George W Bush pictures).

c. 5 Original Pictures (upper Pictures) vs Transformed Pictures for k = 100(lower pictures)

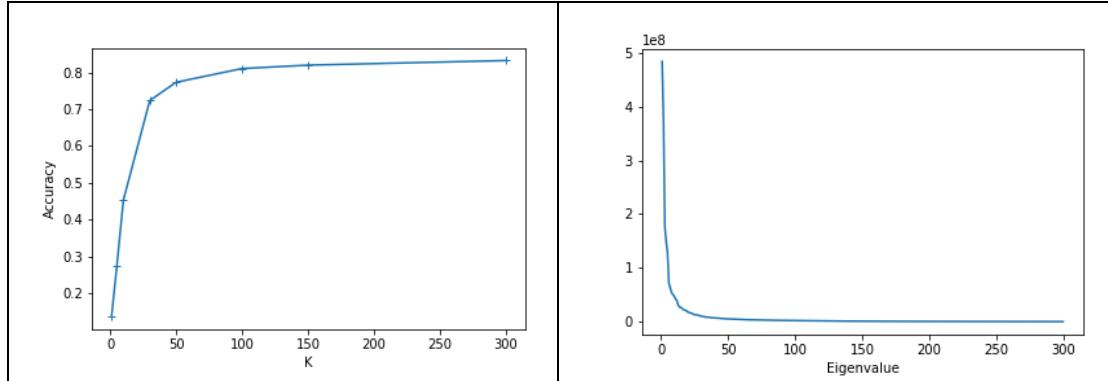


We can see that projecting the Pictures on a subspace spanned by 100 eigenvectors, we are still able to reconstruct the pictures so that they are recognizable by the human eye. For  $K < 50$  we did not had such remarkable results (the reconstructed pictures were of too poor quality). We can examine how closer the transformed pictures are to the original ones by looking at the sum of distances wrt k



We can see that the more eigenvectors we use, the closer the pictures are. But, we do reach a point from which adding more eigenvectors does not improve the distances as it did for smaller values of k. This is due to the fact that the eigenvalues themselves decrease with k and reach a value that is significantly smaller than the first few eigenvalues (see eigenvalues wrt k plot as well), meaning the data has smaller variance in those additional eigenvectors.

d. Test accuracy wrt k while using 75% of the whole data as training data and 25% as test data:



Here as well we can see that the accuracy increases with  $k$  as expected, but reaches saturation around  $k=50$ . In this value of  $k$  we can also see that the eigenvectors reach saturation at a value which is much smaller than the first few eigenvalues. This means that adding more eigenvectors does not improve our ability to distinguish between pictures as it did in smaller  $k$ 's because those eigenvectors represent smaller variance.