

Eigenchangers!

Given,

We want to compute eigen-vector (u) of $X^T X$ given eigen-vector (v) of XX^T .

Note: Ignoring $\frac{1}{N}$ since eigen vectors will not change direction.

From property of eigen-vectors for a non zero eigen-value(λ) we can write:

$$XX^T v = \lambda v$$

Multiplying X^T both sides we get:

$$(X^T X)(X^T v) = \lambda(X^T v)$$

Again, from eigen-vector property, we can say that $X^T v$ is an eigen-vector of $X^T X$ for any non-zero eigen-value λ .

For zero eigen-value we will have:

$$X^T X v = 0 \implies v^T X^T X v = 0 \implies X v = 0$$

So eigen-vectors for zero-eigen values can be obtained from null-space of X matrix in similar way to that of non-zero eigen value given we know v from null space.

Observation : Eigen-values of $X^T X$ and XX^T will be same and are square of singular values we obtain in diagonal matrix of SVD decomposition.

Advantage of doing PCA this way :

For doing PCA we need to do $O(D^3)$ computation, but if $(D > N)$ then it is more computationally expensive as compared to $O(N^3 + ND)$ which is required in this method.

This concludes the answer to above question.

A General Activation Function

Given,

$$h(x) = \frac{x}{1 + \exp(-\beta x)}$$

1. Linear Function $\{\beta = 0\}$

When $\beta = 0$,

$$h(x) = \frac{x}{1 + \exp(0)} = \frac{x}{2}$$

So it becomes a linear function.

2. ReLU Function $\{\beta \rightarrow \infty\}$

When $\beta \rightarrow \infty$ then,

$$h(x) = \frac{x}{1 + \exp(-\infty x)} = \begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$$

So it becomes ReLU activation function.

NOTE : There are different forms of this activation function depending on what we choose β to be. It can become reflections of ReLU along Y-axis as well and many other non-linear forms. That's why with introduction of one variable we have obtained a more general and non-linear activation function.

This concludes the answer to this question.

Student Name: Abhishek Kumar
 Roll Number: 18111002
 Date: November 8, 2018

Mixtures meets Neural Nets!

Given,

$$\begin{aligned} \{x_n, y_n\}_{n=1}^N, \quad x_n \in \mathbb{R}^D, \quad y_n \in \{0, 1\} \\ z_n \sim \text{multinoulli}(\{\pi_k\}_{k=1}^K) \\ y_n \sim \text{Bernoulli}(\sigma(w_{z_n}^T x_n)) \end{aligned}$$

Now we can write out marginal probability as:

$$\begin{aligned} p(y_n = 1|x_n) &= \sum_{k=1}^K p(y_n = 1, z_n = k|x_n) \\ &= \sum_{k=1}^K p(y_n = 1|z_n = k, x_n) p(z_n = k|x_n) \\ &= \sum_{k=1}^K \pi_k \sigma(w_k^T x_n) \\ &= \sum_{k=1}^K \pi_k \sigma\left(\sum_{d=1}^D w_{kd} x_{nd}\right) = \sum_{k=1}^K \pi_k h_k \end{aligned}$$

Now this can be represented as a Neural Network with a single hidden layer, W and π as weight matrices for first and second layer respectively, σ as activation function for hidden layer and x_n as input and $p(y_n = 1|x_n)$ as output respectively. Output of hidden layer node is $h_k = \sigma(\sum_{d=1}^D w_{kd} x_{nd})$

See figure below:

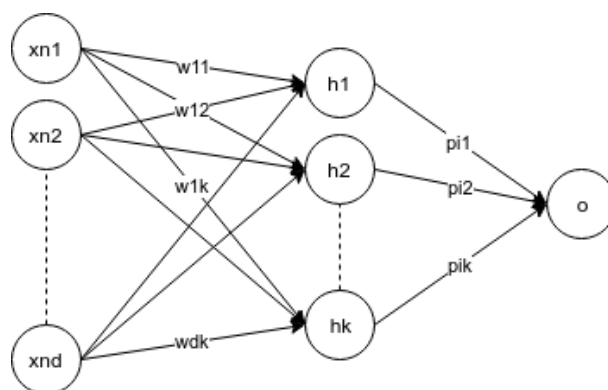


Figure 1: Equivalent Neural Network

This Concludes the answer to this problem.

Student Name: Abhishek Kumar
 Roll Number: 18111002
 Date: November 8, 2018

Probabilistic Formulation of Matrix Factorization with Side Information

Given,

We have a matrix factorization problem over X which represents users and items as rows and columns respectively.

Model Assumptions

$$\begin{aligned} X_{nm} &\sim \mathcal{N}(X_{nm} | \theta_n + \phi_m + u_n^T v_m, \lambda_x^{-1}) \\ u_n &\sim \mathcal{N}(u_n | W_u a_n, \lambda_u^{-1} I_k) \\ v_m &\sim \mathcal{N}(v_m | W_v b_m, \lambda_v^{-1} I_k) \end{aligned}$$

where, $\{X_{nm}, \theta_n, \phi_m\} \in \mathbb{R}$, $\{u_n, v_m\} \in \mathbb{R}^k$, $W_u \in \mathbb{R}^{k \times D_u}$, $W_v \in \mathbb{R}^{k \times D_v}$, $a_n \in \mathbb{R}^{D_u}$, $b_m \in \mathbb{R}^{D_v}$

We want to estimate $\Theta = \{\{u_n, \theta_n\}_{n=1}^N, \{v_m, \phi_m\}_{m=1}^M, W_u, W_v\}$ given X where Ω is set of known indexes.

Loss Function

We want to do ALT-OPT optimization over this modelling so we can write is as maximization problem over know observation by deducing MAP as follows:

$$\begin{aligned} \hat{\Theta} &= \arg \max_{\Theta} p(X_{\Omega}, U, V) \\ &= \arg \max_{\Theta} p(X_{\Omega} | U, V) p(U, V) \\ &= \arg \max_{\Theta} p(X_{\Omega} | U, V) p(U) p(V) \\ &= \arg \max_{\Theta} [\log(p(X_{\Omega} | U, V)) + \log(p(U)) + \log(p(V))] \\ &= \arg \min_{\Theta} L(\Theta) \end{aligned}$$

Now we can write based on our assumptions .

$$\begin{aligned} \log(p(X_{\Omega} | U, V)) &= C_0 - \sum_{\{n,m\} \in \Omega} \frac{\lambda_x}{2} (X_{nm} - \theta_n - \phi_m - u_n^T v_m)^2 \\ \log(p(U)) &= C_1 - \sum_{n=1}^N \frac{\lambda_u}{2} (u_n - W_u a_n)^T (u_n - W_u a_n) \\ \log(p(V)) &= C_2 - \sum_{m=1}^M \frac{\lambda_v}{2} (v_m - W_v b_m)^T (v_m - W_v b_m) \end{aligned}$$

Now we can write our loss function or cost function ignoring constants as follows:

$$L(\Theta) = \sum_{\{n,m\} \in \Omega} \frac{\lambda_x}{2} (X_{nm} - \theta_n - \phi_m - u_n^T v_m)^2 + \sum_{n=1}^N \frac{\lambda_u}{2} (u_n - W_u a_n)^T (u_n - W_u a_n) \\ + \sum_{m=1}^M \frac{\lambda_v}{2} (v_m - W_v b_m)^T (v_m - W_v b_m)$$

ALT-OPT updates

Now we can take derivative equate to 0 to get update equations as follows:

$$\frac{\partial L}{\partial u_n} = \left(\sum_{m \in \Omega_{r_n}} \lambda_x (X_{nm} - \theta_n - \phi_m - u_n^T v_m) (-v_m) \right) + \lambda_u (u_n - W_u a_n) = 0 \\ \Rightarrow \boxed{u_n = \left[\sum_{m \in \Omega_{r_n}} \lambda_x v_m v_m^T + \lambda_u I_k \right]^{-1} \left[\sum_{m \in \Omega_{r_n}} \lambda_x (X_{nm} - \theta_n - \phi_m) v_m + \lambda_u W_u a_n \right]}$$

$$\frac{\partial L}{\partial \theta_n} = \sum_{m \in \Omega_{r_n}} \lambda_x (X_{nm} - \theta_n - \phi_m - u_n^T v_m) (-1) = 0 \\ \Rightarrow \boxed{\theta_n = \frac{1}{|\Omega_{r_n}|} \left[\sum_{m \in \Omega_{r_n}} (X_{nm} - \phi_m - u_n^T v_m) \right]}$$

$$\frac{\partial L}{\partial W_u} = \sum_{n=1}^N \lambda_u (u_n - W_u a_n) (-a_n^T) = 0 \\ \Rightarrow \boxed{W_u = \left[\sum_{n=1}^N u_n a_n^T \right] \left[\sum_{n=1}^N a_n a_n^T \right]^{-1}}$$

Similarly and symmetrically we get :

$$\boxed{\begin{aligned} v_m &= \left[\sum_{n \in \Omega_{c_m}} \lambda_x u_n u_n^T + \lambda_v I_k \right]^{-1} \left[\sum_{n \in \Omega_{c_m}} \lambda_x (X_{nm} - \theta_n - \phi_m) u_n + \lambda_v W_v b_m \right] \\ \phi_m &= \frac{1}{|\Omega_{c_m}|} \left[\sum_{n \in \Omega_{c_m}} (X_{nm} - \theta_n - u_n^T v_m) \right] \\ W_v &= \left[\sum_{m=1}^M v_m b_m^T \right] \left[\sum_{m=1}^M b_m b_m^T \right]^{-1} \end{aligned}}$$

This concludes the answer to given problem.

Programming Problem 1

Given to apply PPCA with constant sigma using ALT-OPT method. which is nothing but PCA model itself without eigen-decomposition and orthogonality constraint over basis vectors.

ALT-OPT PPCA const. Sigma

Alternating optimization over PPCA with constant sigma basically boils down to two steps:

1. Update Z keeping W constant

$$\hat{Z} = XW[W^TW]^{-1}$$

2. Update W keeping Z constant

$$\hat{W} = XZ[Z^TZ]^{-1}$$

Doing this repeatedly for 200 iterations we will come to following solution:

Original Images were:

Original Images

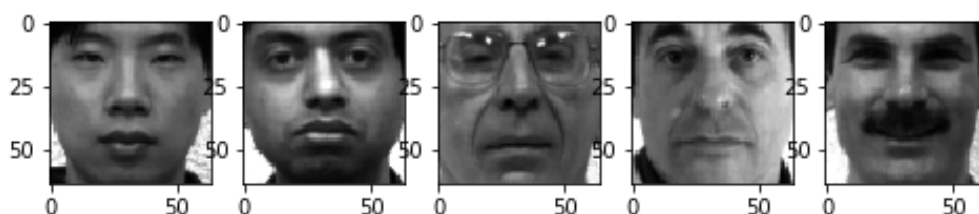


Figure 2: Original Images

Reconstructed Images



Figure 3: Reconstructed Images

Basis Images (Columns of W)

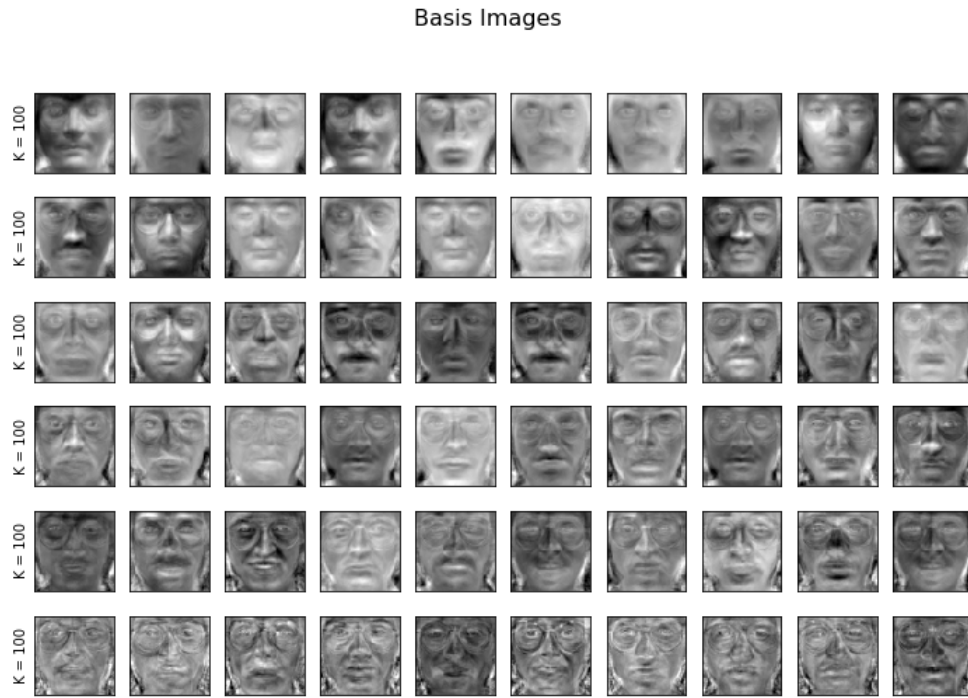


Figure 4: 10 Basis Images

Observation

The Reconstructed Images are worse than the original ones for low value of K but for large values they seem fine. Also face reconstructions are better with large K since it has more face parameters to take combinations.

Observation

The more i increase the number of basis size (K) the features that are learned are more like edges of faces as compared to low K size where it's mostly faces themselves which decreases the fluidity of models with less K and hence their accuracy.

Programming Problem 2

Given to implement PCA and TSNE with Kmeans Clustering over mnist dataset.

PCA Clustering

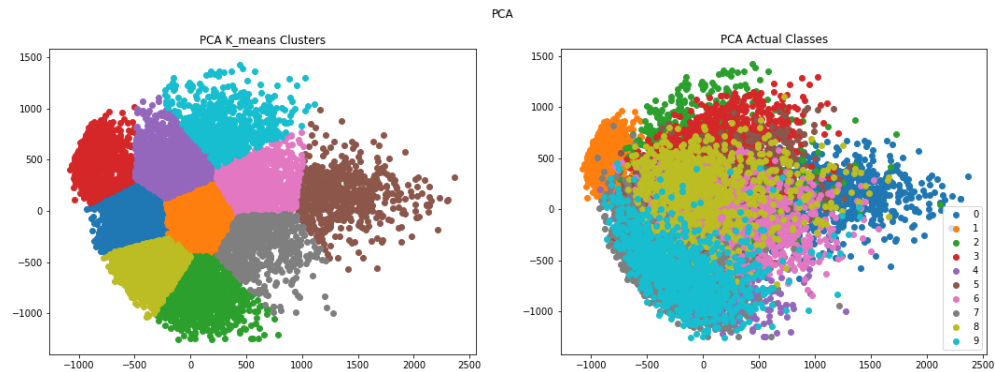


Figure 5: PCA Clustering vs Actual Classes

TSNE Clustering

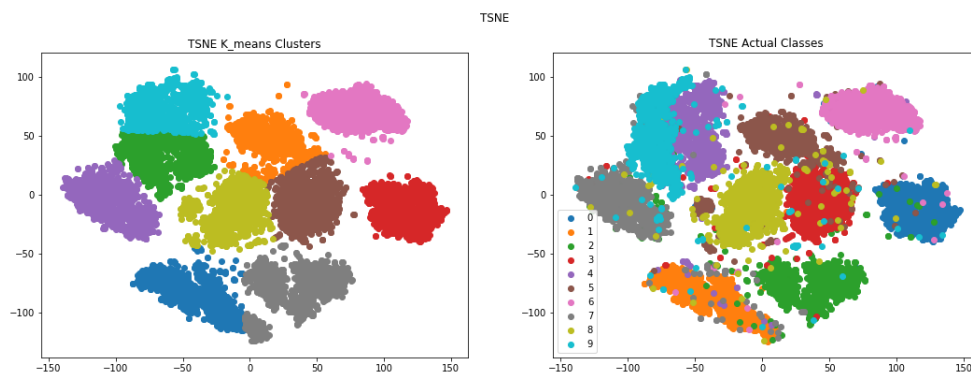


Figure 6: TSNE Clustering vs Actual Classes

Observation :

TSNE gives a much better clustering as compared to PCA for classification purpose since PCA simply looks at variance whereas TSNE looks at neighbourhood and maintains the relationship.

Some random runs over kmeans after PCA and TSNE reduction.



Figure 7: TSNE vs PCA Clusterings

We can observe how PCA clusters are merged and TSNE clusters are somewhat far apart and so we can conclude that TSNE gives better results visually as compared to PCA. This concludes the answer to this question.