## Prob 1

let's assume we have a eigenvector $\underline{v}_i$ of

Matrix $S = \frac{1}{N} \underline{X} \underline{X}^T$

$\therefore$ It will satisfy $\quad S\underline{v} = \lambda \underline{v}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \overset{\downarrow}{\text{eigenvalue}}$

$\Rightarrow \quad \frac{1}{N}\left(\underline{X}\underline{X}^T\right) \underline{v} = \lambda \underline{v}$

$\Rightarrow \quad \frac{1}{N}\left(\underline{X}^T\underline{X}\right)\left(\underline{X}^T\underline{v}\right) = \lambda \left(\underline{X}^T\underline{v}\right) \quad \left\{\begin{array}{l}\text{premultiply } \underline{X}^T \\ \text{both side}\end{array}\right.$

Substitute $\quad \underline{u} = \underline{X}^T\underline{v}$

$\Rightarrow \quad \frac{1}{N}\left(\underline{X}^T\underline{X}\right)\underline{u} = \lambda \underline{u}$

$\therefore \quad \underline{u} = \underline{X}^T\underline{v}$ is nothing else but an

$\quad\quad\quad$ eigenvector for $\frac{1}{N}\left(\underline{X}^T\underline{X}\right)$

In normal case to compute $K$ eigenvectors for $\frac{1}{N}(X^TX)$
complexity will be $O\left(KD^2\right)$
but in this case order will be $O\left(KN^2\right) + O\left(KND\right)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \overset{\downarrow}{\underset{\text{of } \frac{1}{N}(XX^T)}{\text{for decomp.}}} \quad \overset{\downarrow}{\underset{\text{matrix}}{\text{for multi}}}$

$\therefore$ overall complexity $O\left(KND\right)$
$\quad\quad\quad\quad\quad\quad\quad\quad \downarrow$ which is less than $O\left(KD^2\right)$

$\quad\quad\quad\quad$ as $\boxed{N < D}$

Scanned by CamScanner

$$h(x) = x \sigma(\beta x)$$

\# if we choose $\beta = 0$

$$\Rightarrow \sigma(0) = 1 \Rightarrow \quad \text{a linear}$$
activation
function.

$$\Rightarrow \quad h(x) = \frac{x}{2} \Rightarrow$$

\# if we choose $\beta \to \infty$

$$\Rightarrow \text{for } x < 0 \Rightarrow \sigma(\beta x) = \sigma(-\infty) = 0$$

$$\& \text{ for } x > 0 \Rightarrow \sigma(\beta x) = \sigma(\infty) = 1$$
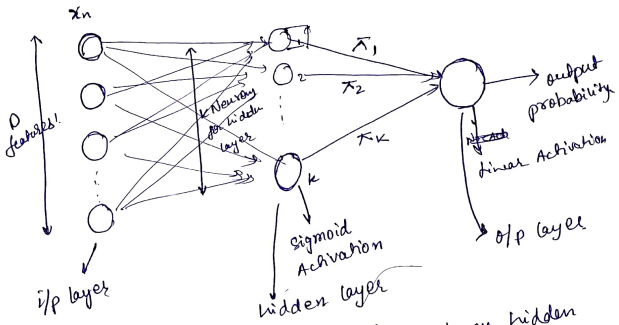
$$\Rightarrow \quad h(x) = \begin{cases} x & ; \quad x > 0 \\ \\ 0 & ; \quad x < 0 \end{cases} \quad \to \text{RELU}$$
$$\downarrow$$
Activation
f^n

$$P(y_n = 1 \mid x_n) = \prod_{k=1}^{K} P(y_n = 1 \mid z_{n\emptyset} = k, x_n)\, P(z_n = k \mid x_n)$$

∴ Above equation can be modelled as a multiplication of two probabilistic f(n)s.

⇒ 2 layer NN with one hidden layer.



$x_n$

O features!

i/p layer

K Neurons for hidden layer

$\overline{\Lambda}_1$

$\pi_2$

$\pi_K$

sigmoid Activation

hidden layer

output probability

Linear Activation

o/p layer

∴ basically $\sigma(w_{zn}^T x_n) \rightarrow$ acts as an hidden neurons.

Activation is sigmoids.

∴ Weights for first layer will be $(w_{z1}, w_{z2} \cdots w_{zk})$ for $k^{th}$ neurons.

∴ weight for outputs will be $[\pi_1, \overline{\pi}_2 \cdots \overline{\pi}_K]^T$

multiplying the probability of $P(z_n \mid x_n)$

Prob. 4

$$P\left(u_n, v_m, \theta_n, \phi_m \middle| \overset{\frown}{\smile} x\right)$$

$$= \prod_{n=1, m=1}^{M, N} P\left(x_{nm} \middle| u_n, v_m, \theta_n, \phi_m\right) P(u_n) P(v_m)$$

$\Rightarrow 0$

$$= \prod_{n, m} \sqrt{\frac{\lambda_x}{2\pi}} \, \exp\left\{-\frac{\lambda_x}{2}\left(x_{nm} - \theta_n - \phi_m - u_n^T v_m\right)^2\right\}$$

$$\times \sqrt{\frac{\lambda_u}{2\pi}} \left\{\exp\left(-\frac{\lambda_u}{2}\left\|u_n - W_u a_n\right\|^2\right)\right.$$

$$\times \sqrt{\frac{\lambda_v}{2\pi}} \, \exp\left(-\frac{\lambda_v}{2}\left\|v_m - W_v a_m\right\|^2\right)$$

$$\Rightarrow NLL = \sum_{n, m}\left\{\frac{\lambda_x}{2}\left(x_{nm} - \theta_n - \phi_m - u_n^T v_m\right)^2\right\}$$

$$+ \sum_n\left\{\frac{\lambda_u}{2}\left\|u_n - W_u a_n\right\|^2\right\}$$

$$+ \sum_m\left\{\frac{\lambda_v}{2}\left\|v_m - W_v a_m\right\|^2\right\}$$

NLL → Cost function

⚡ Now, for $W_u$, $W_v$ → minimization will be similar to linear regression model.

which will give :—

$$W_u^T = (A^T A)^{-1} A^T U$$

where; $A$ → $N \times q$ matrix of all $a_n$'s

$U$ → $N \times k$ matrix of all $u_n$'s

$$\Rightarrow \quad W_u = U^T A (A^T A)^{-1}$$

Similarly, $\quad W_v = V^T B (B^T B)^{-1}$

# for $\theta_n$ & $\theta_m$ | taking the derivatives yield.

$$\sum_m -\lambda_r \left[ x_{nm} - \theta_n - \phi_m - u_n^T v_m \right] = 0$$

$$\Rightarrow \quad \theta_n = \frac{1}{M} \sum_m \left( x_{nm} - \phi_m - u_n^T v_m \right)$$

Similarly $\phi_m = \frac{1}{N} \sum_n \left( x_{nm} - \theta_n - u_n^T v_m \right)$

\# Now taking the derivative w.r.t. $v_m$ weget.

$$\sum_n -\lambda_x \left( x_{nm} - \theta_n - \phi_m - u_n^T v_m \right) u_n$$

$$+ \lambda_v \left( v_m - W_e\, b_m \right) = 0$$

$$\Rightarrow \left( \sum_n \lambda_x u_n u_n^T + \lambda_v I_{\theta k} \right) v_m = \lambda_v W_e\, b_m$$
$$+ \sum_n \lambda_x \left( x_{nm} - \theta_n - \phi_m \right) u_n$$

we can write $\sum_n u_n u_n^T = U^T U$

$$\Rightarrow \boxed{ v_m = \left( \lambda_x U^T U + \lambda_v I_k \right)^{-1} \left[ \sum_n \lambda_n \left( x_{nm} - \theta_n - \phi_m \right) u_n + \lambda_v W_e\, b_m \right] }$$

Similarly, for $U_m$

$$\boxed{ U_m = \left( \lambda_x V^T V + \lambda_u I_k \right)^{-1} \left[ \sum_m \lambda_x \left( x_{nm} - \theta_n - \phi_m \right) v_m + \lambda_u W_u a_n \right] }$$

$\therefore$ Using the equation shew above. ALT-OPT algo
will be :—

• Initial



PTO

## ALT - OPT

→① Initialise latent variables $U_n = \hat{U}_n$ ; $V_m = \hat{V}_m$

→② Solve for $\hat{W}_u$, $\hat{W}_v$, $\hat{\theta}_n$, $\hat{\theta}_m$ using the $eq^n$.

→③ Solve $\hat{U}_n$, & $\hat{V}_m$ using the $eq^n$

→④ Go to step ② if not converge.

*Student Name:* Shashi Kant Gupta
*Roll Number:* 160645
*Date:* November 17, 2018

**Programming Part 1:**
On increasing $k$ there is improvement in the reconstruction of images! Morevover the images appear to be more clearer. Explanation is simple higher the value of $k$ more the no of features retained in $z$.
Please find the plots from next page.

For K = 10

For K = 20

For K = 30

# For K = 40

For K = 50

For K = 100

For K = 10

For K = 20

For K = 30

For K = 40

For K = 50

For K = 100

**Programming Part 2:**
Visually tSNE appears to be better, clusters are more seperated as well as on running 10 different initiallisations errors in tSNE seems to be less than that of PCA!
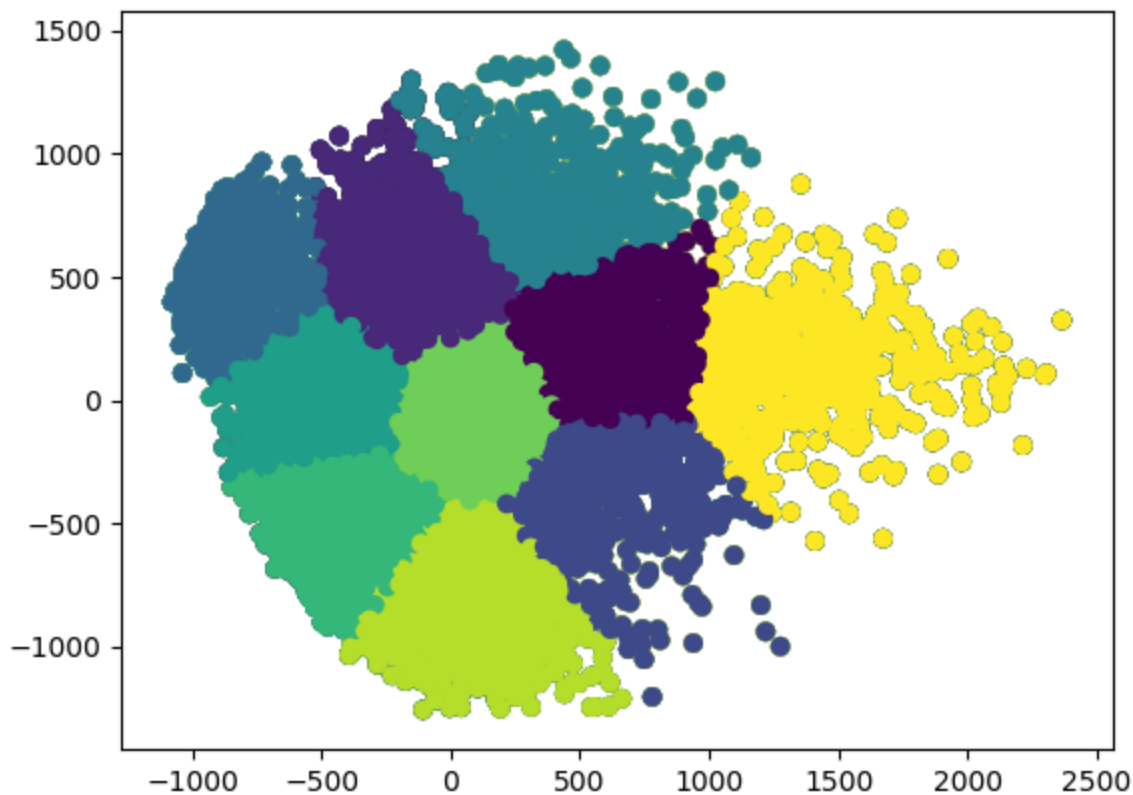Please find the plots from next page.

PCA Based Init: 0
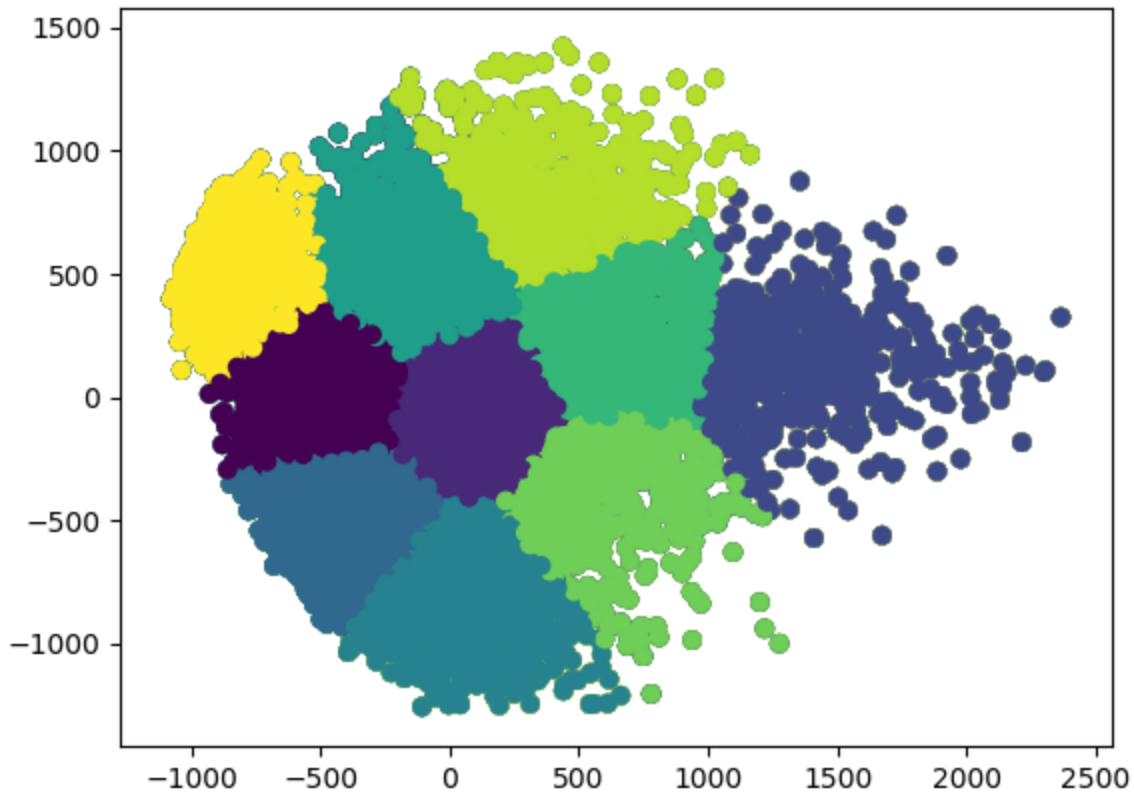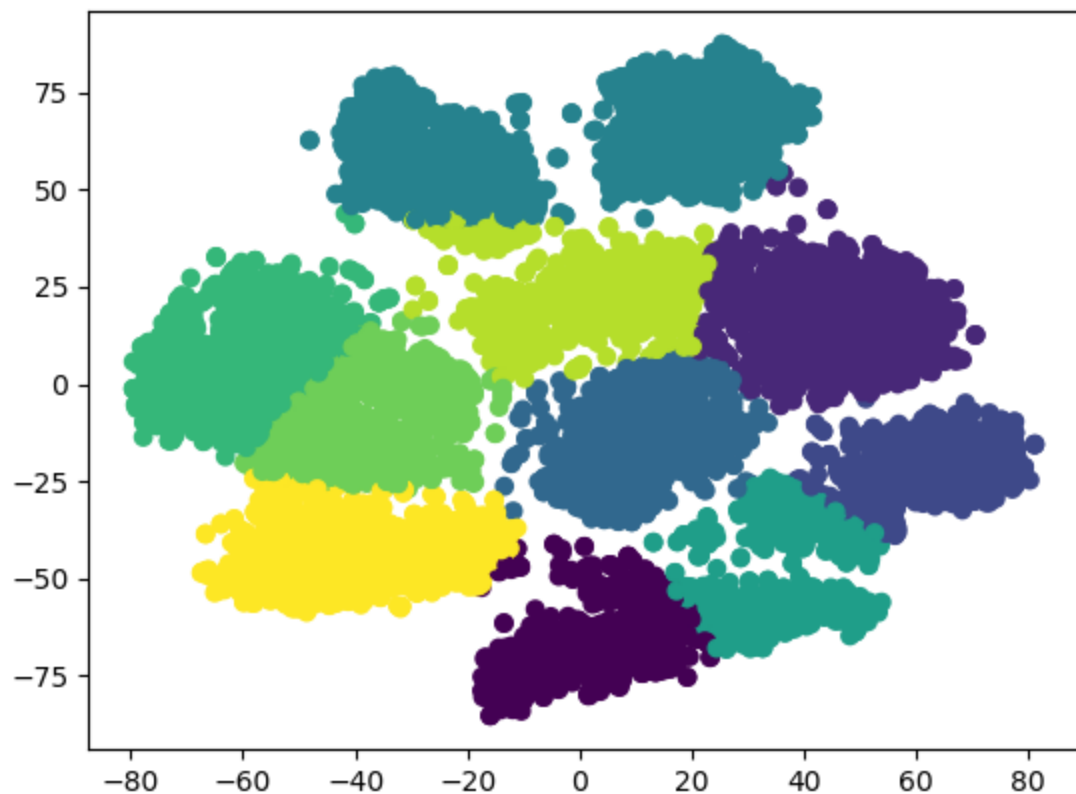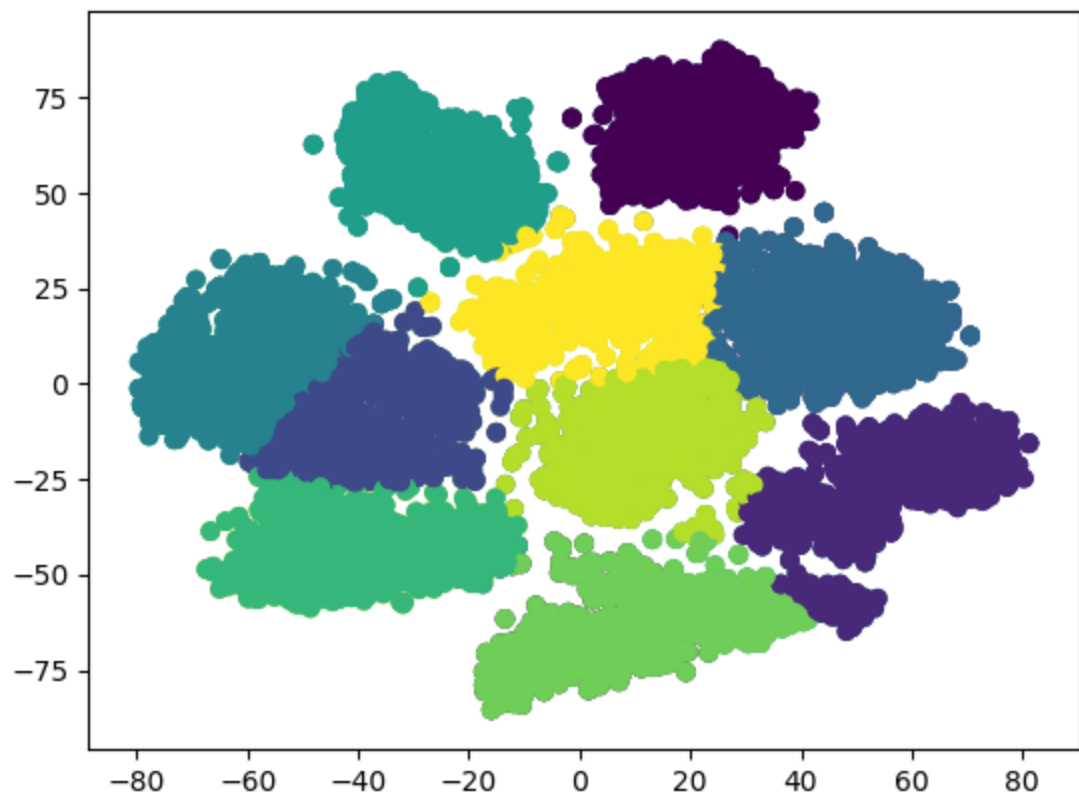
PCA Based Init: 1

PCA Based Init: 2

PCA Based Init: 3
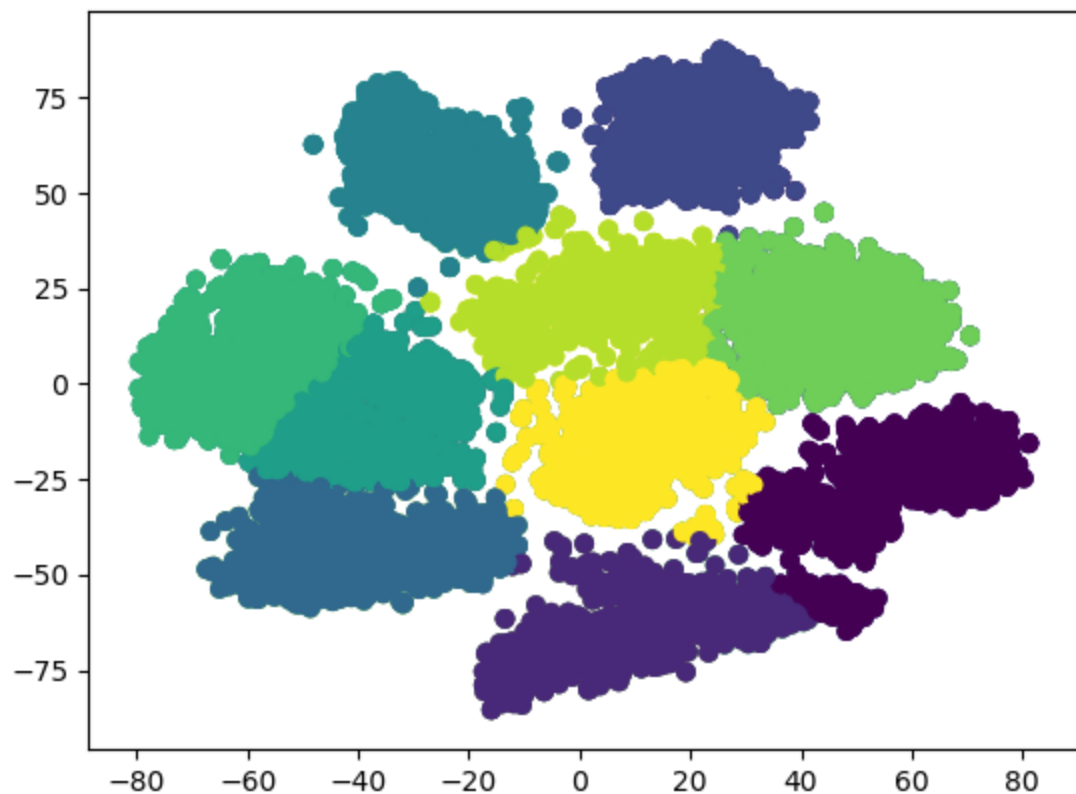
PCA Based Init: 4

PCA Based Init: 5

PCA Based Init: 6

PCA Based Init: 7

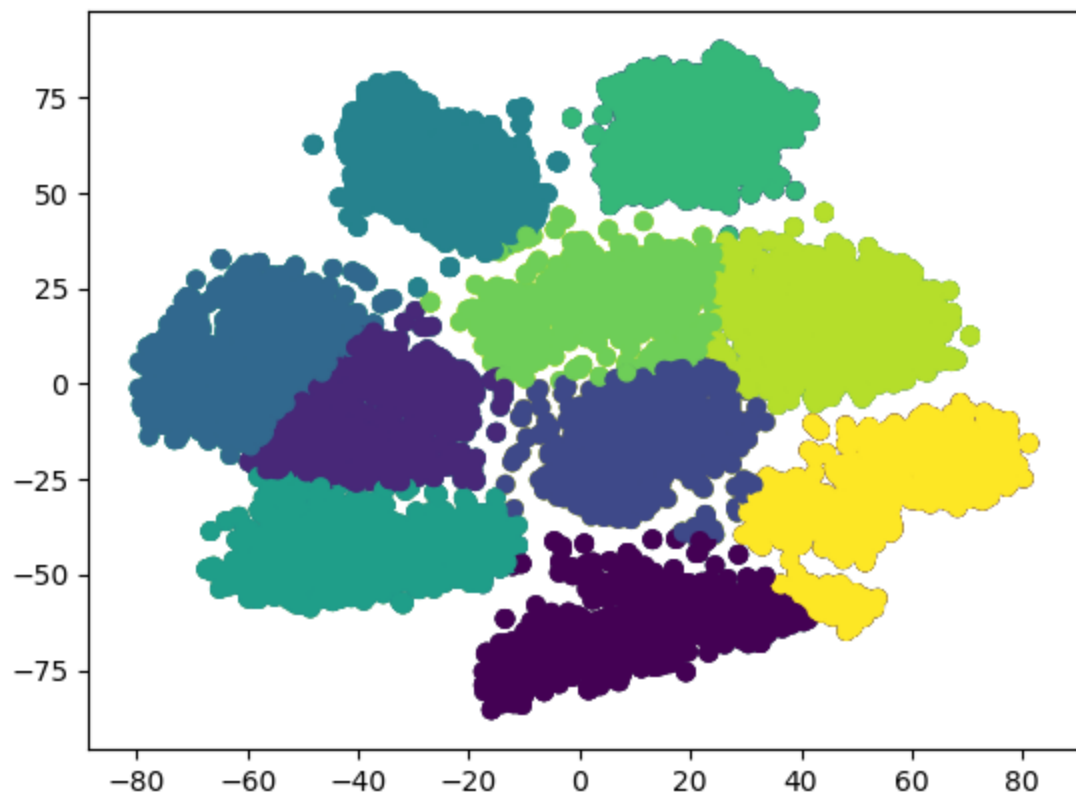PCA Based Init: 8

PCA Based Init: 9
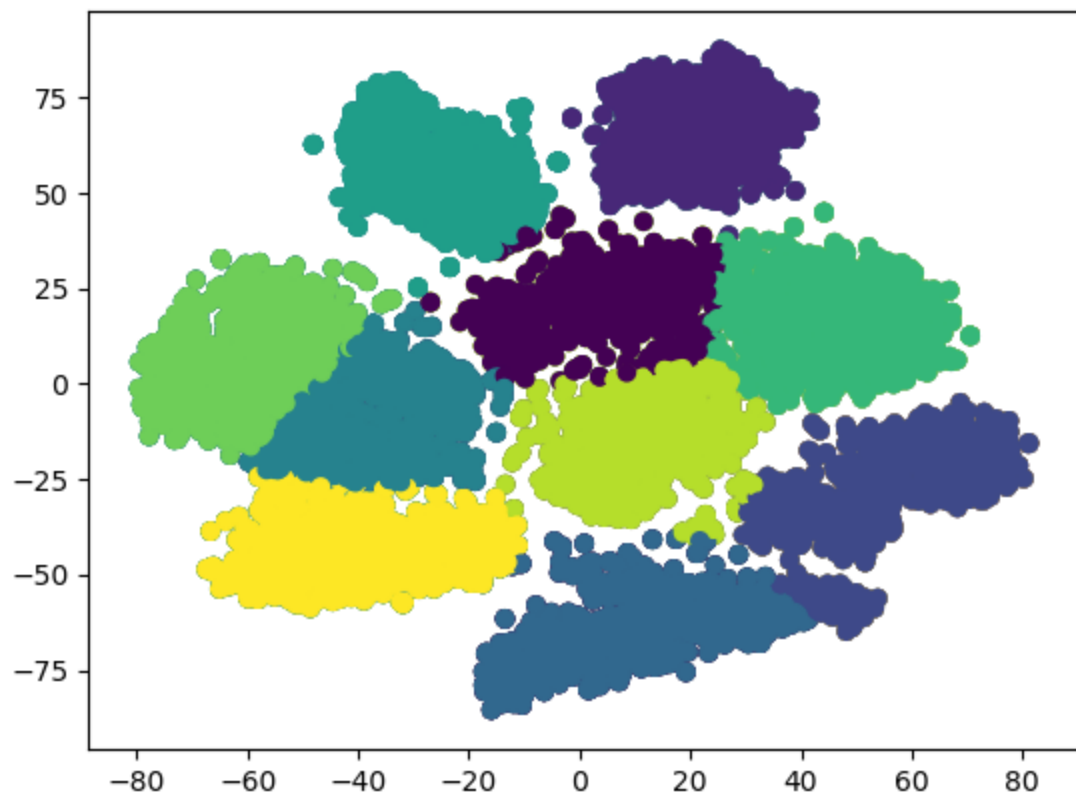
tSNE Based Init: 0

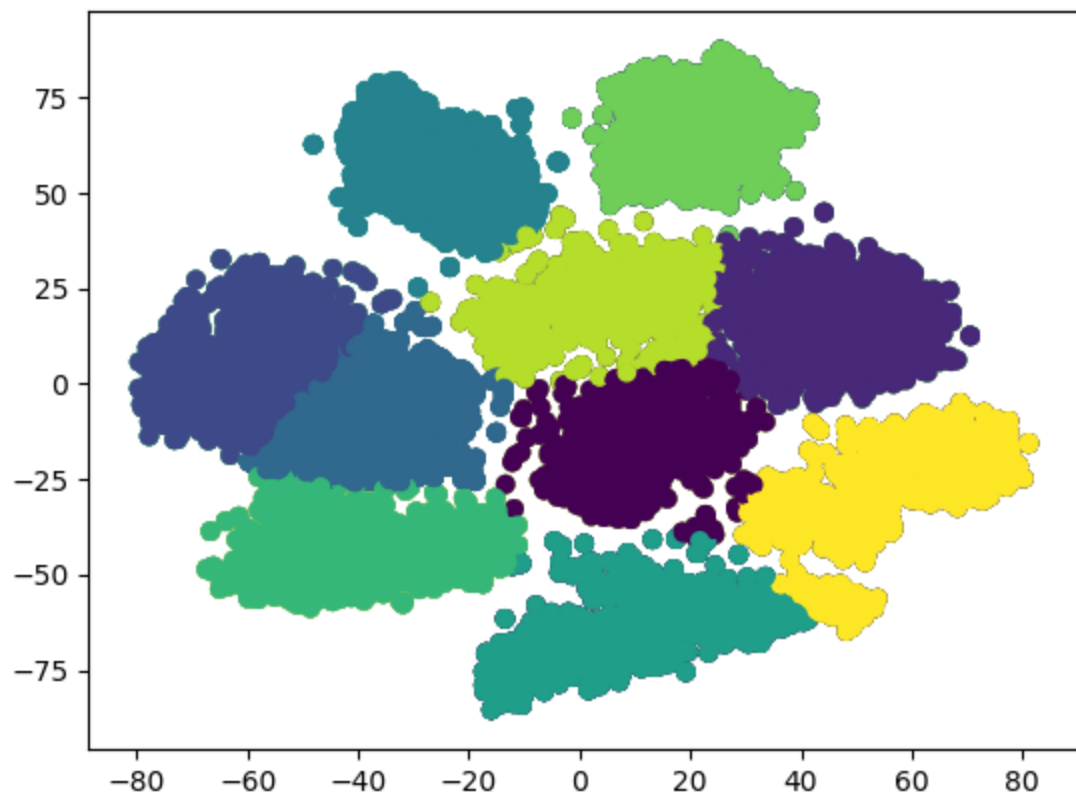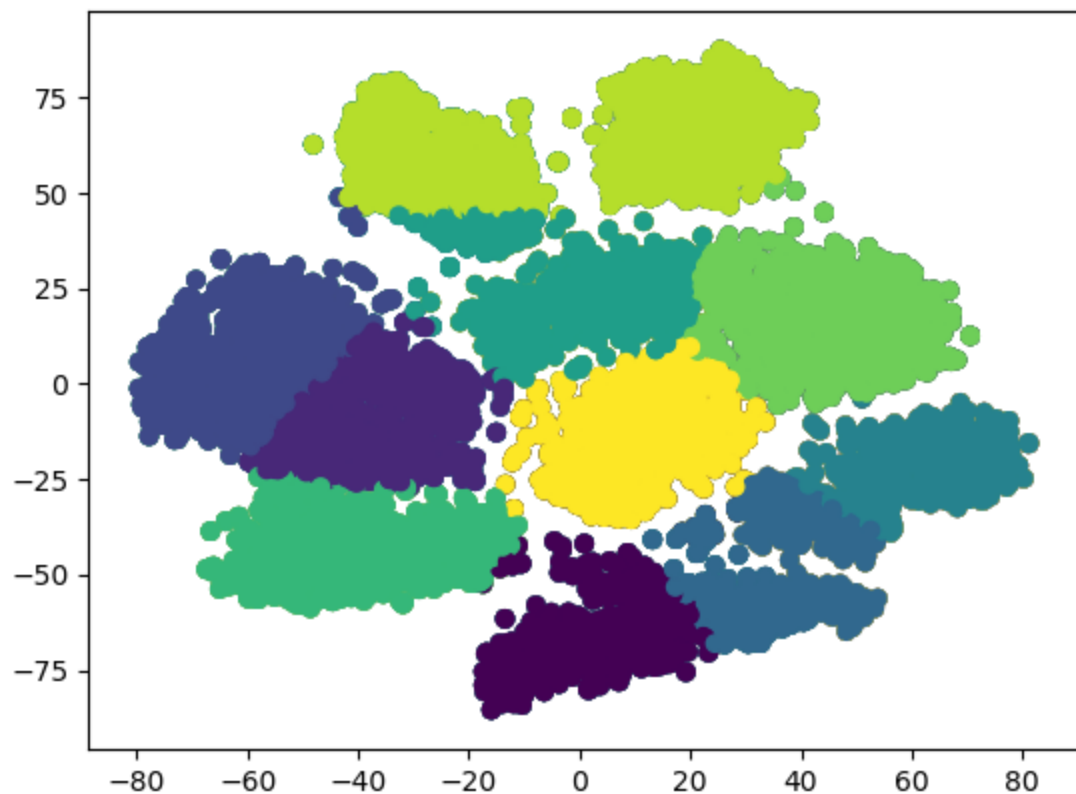tSNE Based Init: 1

tSNE Based Init: 2

tSNE Based Init: 3
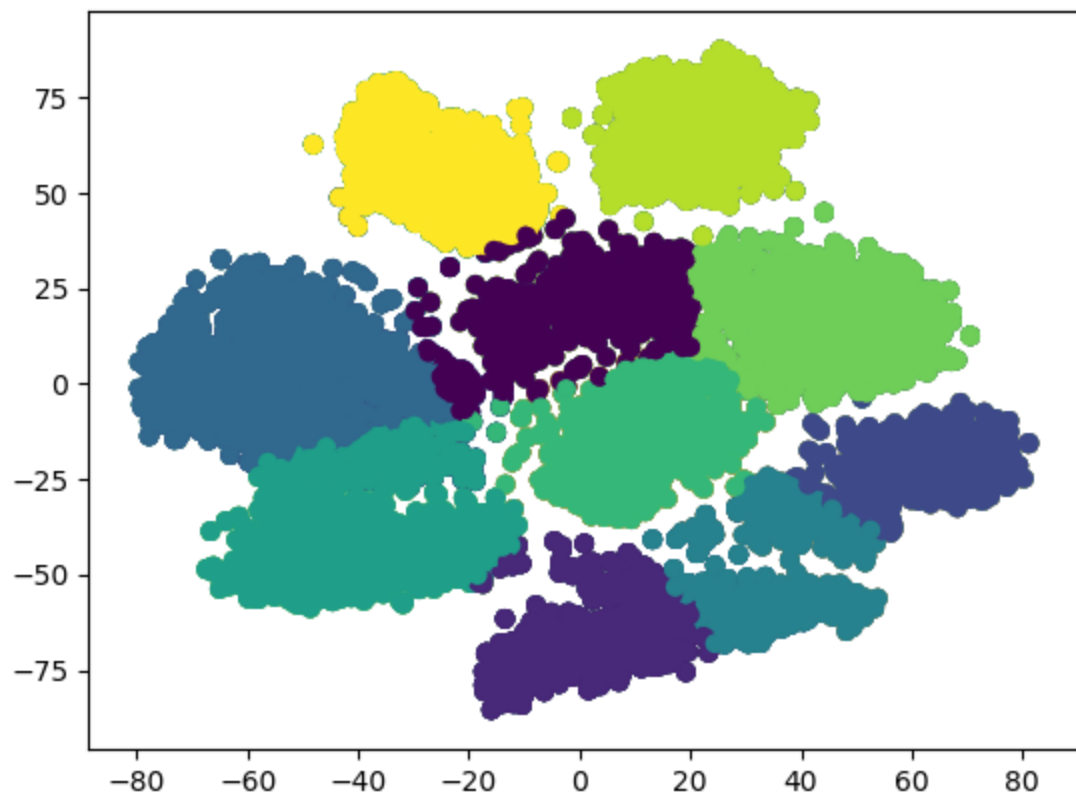
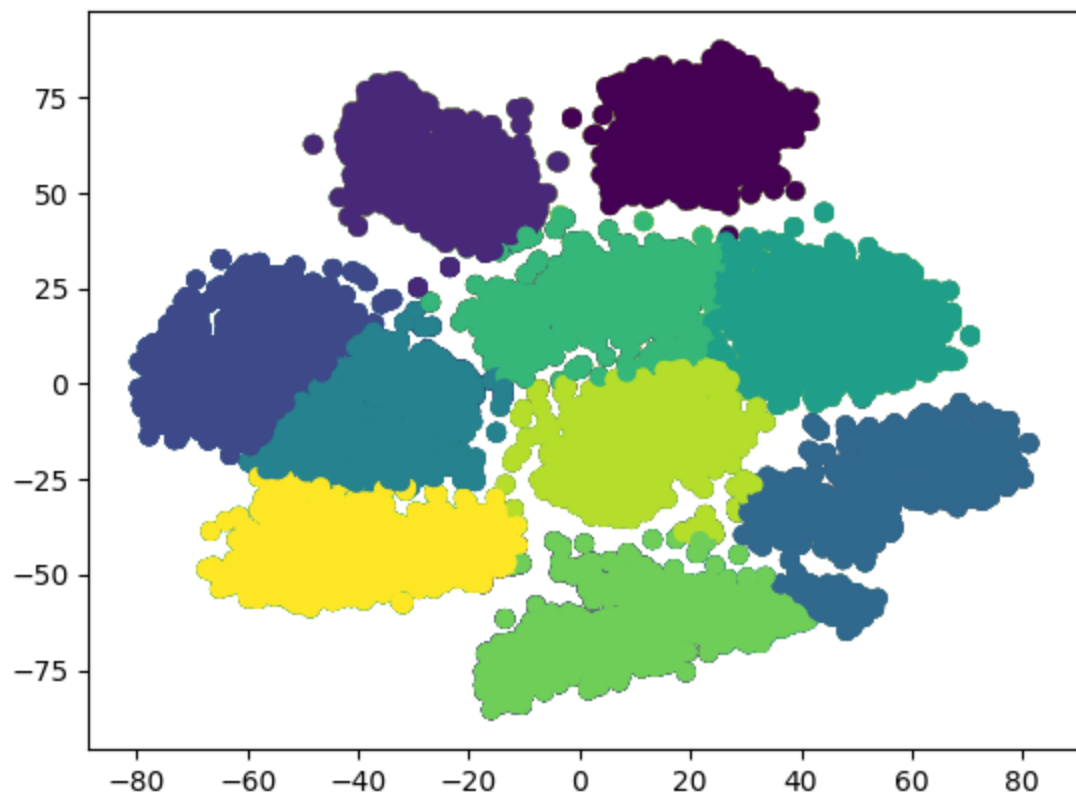tSNE Based Init: 4

tSNE Based Init: 5

tSNE Based Init: 6

tSNE Based Init: 7

tSNE Based Init: 8

tSNE Based Init: 9