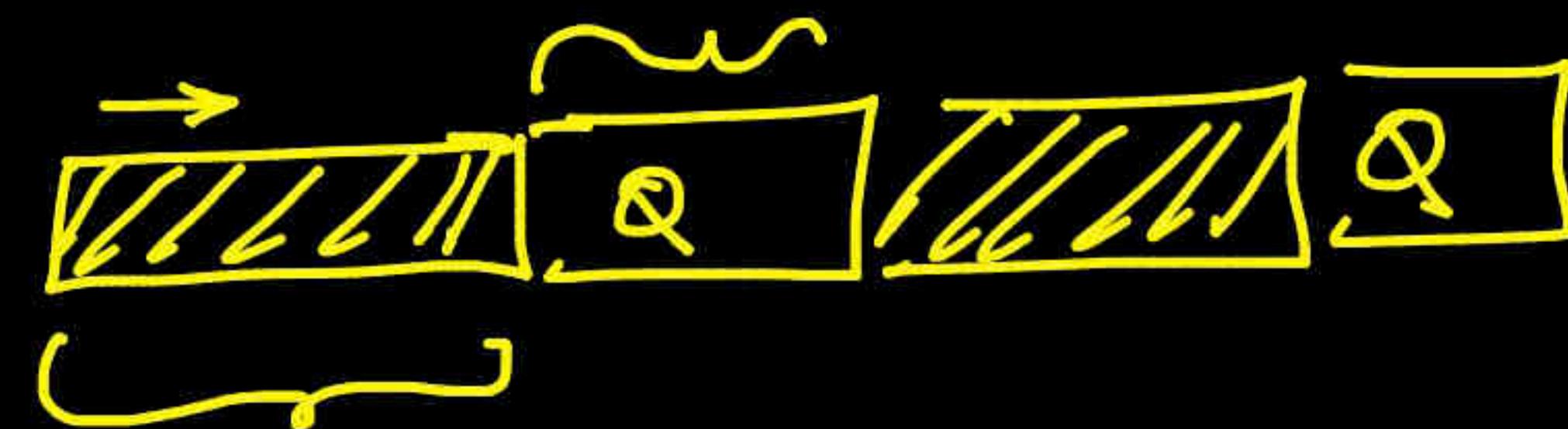


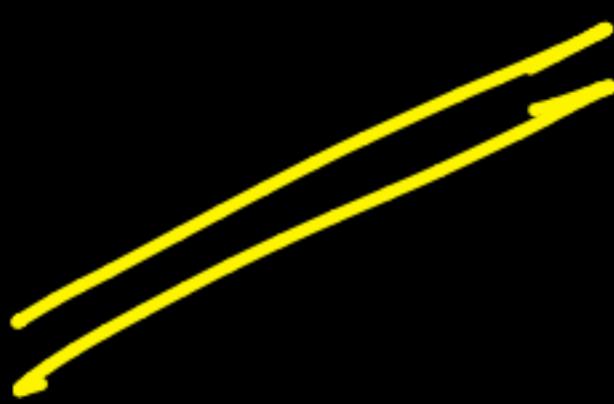
~~OPS:~~

## Question - Tabs



"END: · · · · "

# Vector & Matrix derivatives



①

$$y_{m \times 1} = \underbrace{A_{m \times n}}_{\text{post-mul}} \underbrace{x_{n \times 1}}_{(\text{trick})}$$

$$\frac{\partial y}{\partial x} = A \Rightarrow \frac{d(Ax)}{dx} = A$$

$\alpha$

$\alpha = \underbrace{y^T A}_{l \times m} \underbrace{x}_{m \times n} \xrightarrow{\text{post-mul}} \underbrace{(y_{m \times 1})}_{l \times 1}$

Scalar

$$\frac{\partial \alpha}{\partial x} = y^T A$$

(3)

$$\text{pre} \rightarrow \tilde{\alpha} = \underbrace{y^T A}_{l \times M} \underbrace{x}_{M \times n} \underbrace{z}_{n \times 1}$$

$$\underline{(A \cdot B)^T = B^T A^T}$$

$$\frac{\partial \tilde{\alpha}}{\partial y} = \tilde{(Ax)}^T = x^T A^T$$

(earlier)

$\frac{\partial \text{Loss}}{\partial w}$  scalar  
 $\hookrightarrow \text{vec } \bar{w}$

$$\nabla_{\omega} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \end{bmatrix}$$

{ Typical proof:

To prove:  $\frac{\partial \underline{\alpha}}{\partial \underline{y}} = \underline{x^T A^T}$  ;  $\underline{\alpha} = \underbrace{\underline{y^T A x}}_{\underline{w^T}}$

$$(A \cdot B)^T = B^T A^T$$

Let  $\underline{w} = \underline{y^T A}$   $\Rightarrow (\underline{w^T})^T = (\underline{y^T A})^T$   
 $\Rightarrow \underline{\alpha} = \underline{w^T x}$        $\downarrow$   
 $\Rightarrow \underline{\alpha^T} = (\underline{w^T \cdot x})^T$   
 $\Rightarrow \underline{\alpha^T} = \underline{x^T w}$

⊕ {  $\alpha_{1 \times 1}$  : scalar }  
 $\checkmark \alpha^T = \alpha$

$$\alpha = \alpha^T = x^T w$$

$$\Rightarrow \alpha = x^T A^T y \xrightarrow{m \leftarrow \text{post}}$$

$$\frac{\partial \alpha}{\partial y} = x^T A^T$$

4 ✓

$$d = \underbrace{\alpha^T A \alpha}_{\text{pre}} \underbrace{\alpha^T}_{\text{post}} \quad (\text{PCA - earlier})$$

$\alpha \in \mathbb{R}^n$

$$\frac{\partial d}{\partial \alpha} = \alpha^T A + (A\alpha)^T$$

$$= \alpha^T A + \alpha^T A^T$$

$$= \alpha^T (A + A^T)$$

$A$ : Symm. Matrix

$$A_{ij} = A_{ji}$$

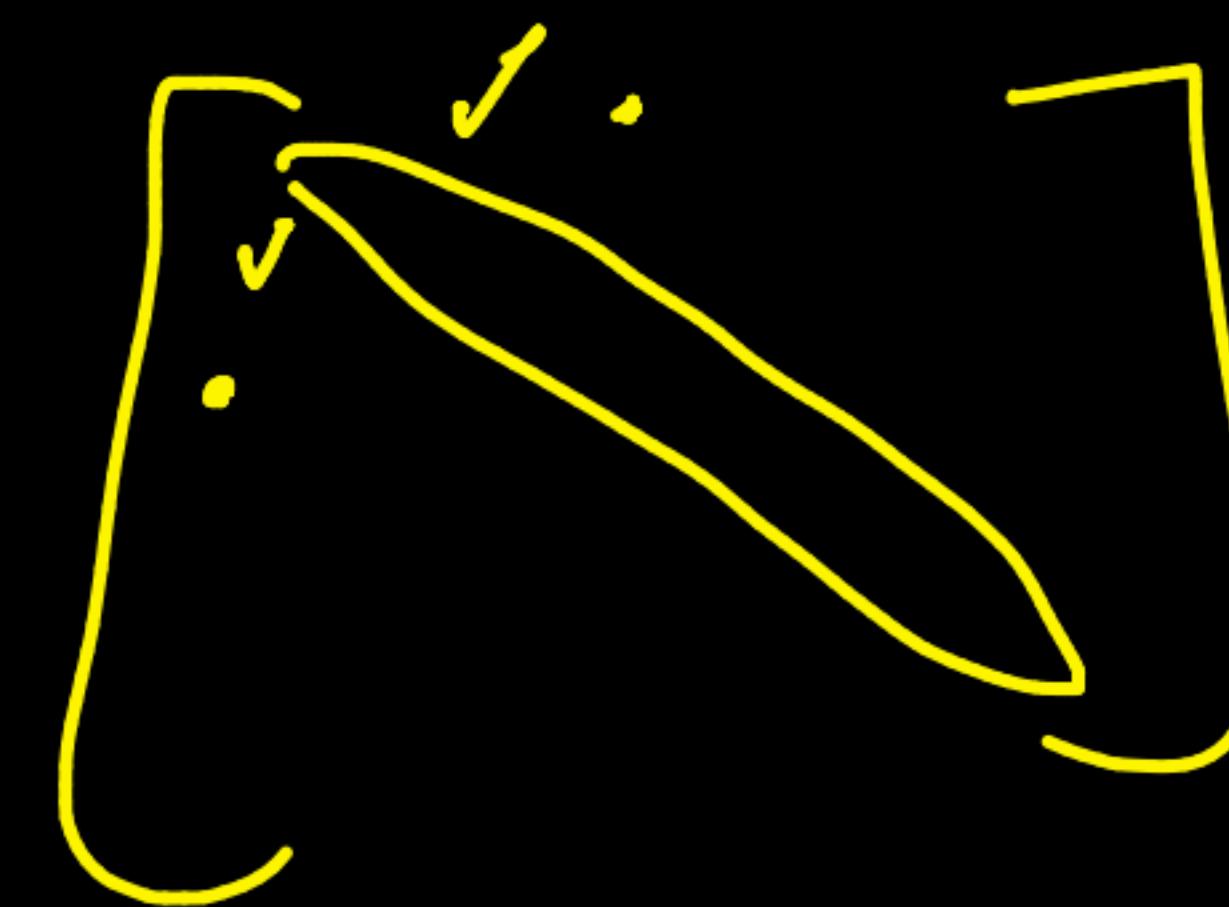
$$A^T = A$$

$$\Rightarrow \underline{\underline{2\alpha^T A}} \text{ or } \underline{\underline{2\alpha^T A^T}}$$

$$\underline{\underline{2\alpha^T A}}$$

|

Cov-Matrix

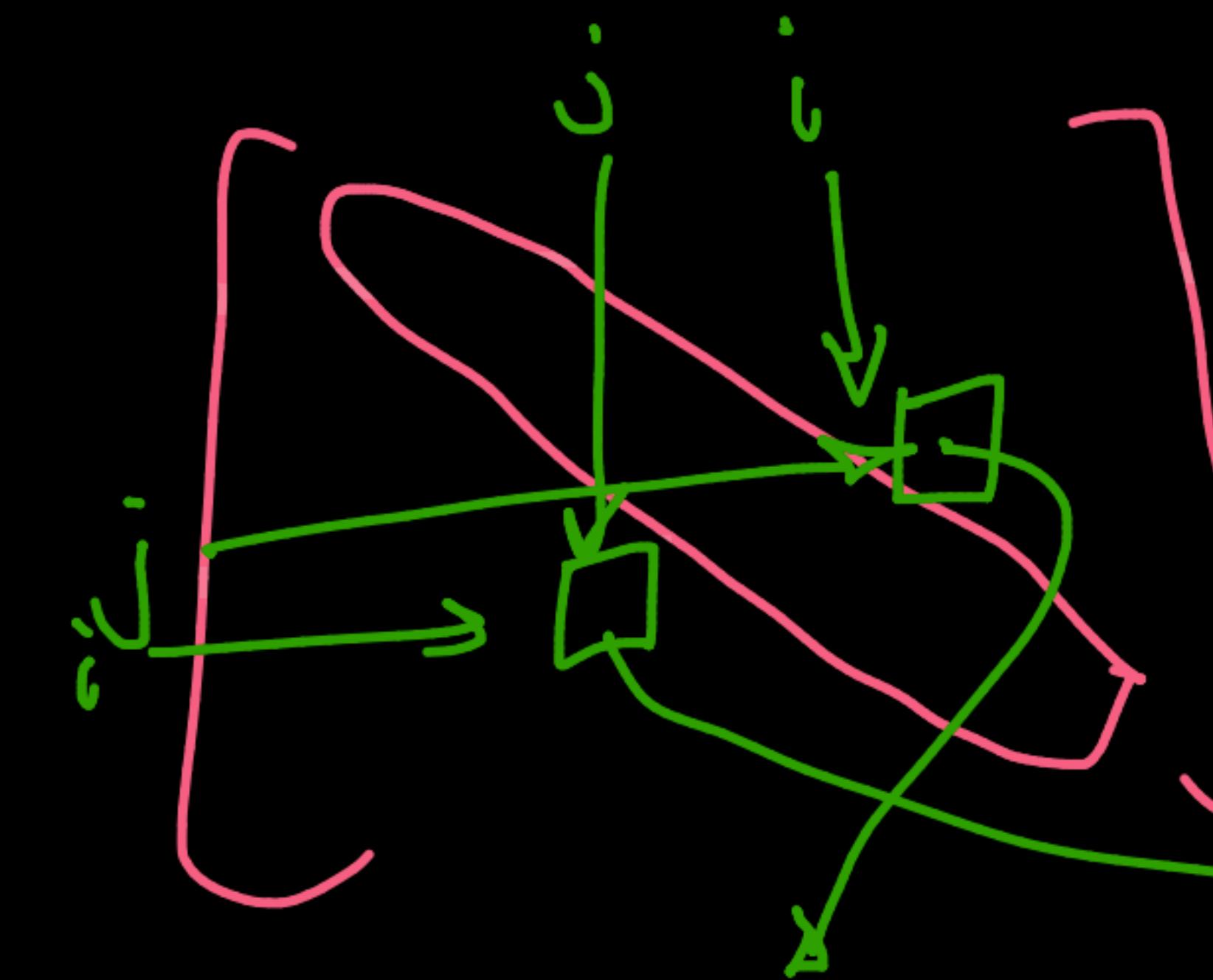


$$A_{ij}^T = A_{ji}$$

$$A = A^T$$

} Symm. Matrix

Cov. Matrix



$$\left\{ \begin{array}{l} \text{Cov}(\tilde{f}_i, \tilde{f}_i) = \text{Cov}(f_i, f_i) \\ \text{=====} \end{array} \right.$$

5

$$\alpha = \underbrace{x^T}_{1 \times 1} \underbrace{x}_{1 \times n} \xrightarrow{\text{post}} \underbrace{n \times 1}_{\text{pre}}$$

$$\|w\|_2^2 = 1$$

$$\Downarrow$$

$$w^T w = 1$$

$$\begin{aligned}\frac{\partial \alpha}{\partial x} &= x^T + x^T \\ &= \underline{\underline{2x^T}}\end{aligned}$$

$$\min_w \text{Loss} + \lambda(w^T w)$$

Most-widely used vector & Matrix der



$$\mathcal{L} = \frac{\omega}{2} \mathbf{x}^T \mathbf{y}$$

$\equiv$

$\mathbf{x} \in \mathbb{R}^{n \times 1}$        $\mathbf{y} \in \mathbb{R}^{m \times 1}$

doesn't occur much in ML

$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}^T} = \mathbf{A} \mathbf{x}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{y}} = \mathbf{x}^T \mathbf{A}^T$$

$\equiv$

$$\nabla_{\mathbf{y}} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_1} \\ \frac{\partial \mathcal{L}}{\partial y_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial y_d} \end{bmatrix}$$

$$\nabla_{\mathbf{y}^T} \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_1} & \frac{\partial \mathcal{L}}{\partial y_2} & \dots \end{bmatrix}$$

$$(\nabla_{\mathbf{y}} \mathcal{L})^T = \nabla_{\mathbf{y}^T} \mathcal{L}$$



## Statement of theorem [edit]

Bayes' theorem is stated mathematically as the following equation:<sup>[4]</sup>

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

posterior

where  $A$  and  $B$  are events and  $P(B) \neq 0$ .

$P(A | B)$  is a conditional probability: the probability of event  $A$  occurring given that  $B$  is true. It is also called the posterior probability of  $A$  given  $B$ .

- $P(B | A)$  is also a conditional probability: the probability of event  $B$  occurring given that  $A$  is true. It can also be interpreted as the likelihood of  $A$  given a fixed  $B$  because  $P(B | A) = L(A | B)$ .
- $P(A)$  and  $P(B)$  are the probabilities of observing  $A$  and  $B$  respectively without any given conditions; they are known as the marginal probability or prior probability.

likelihood

## Proof [edit]

### For events [edit]

Bayes' theorem may be derived from the definition of conditional probability:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \text{ if } P(B) \neq 0,$$

where  $P(A \cap B)$  is the probability of both  $A$  and  $B$  being true. Similarly,

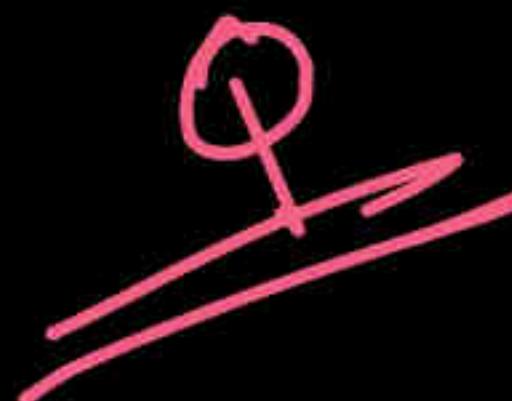
$$P(B | A) = \frac{P(A \cap B)}{P(A)}, \text{ if } P(A) \neq 0,$$

Solving for  $P(A \cap B)$  and substituting into the above expression for  $P(A | B)$  yields Bayes' theorem:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}, \text{ if } P(B) \neq 0.$$

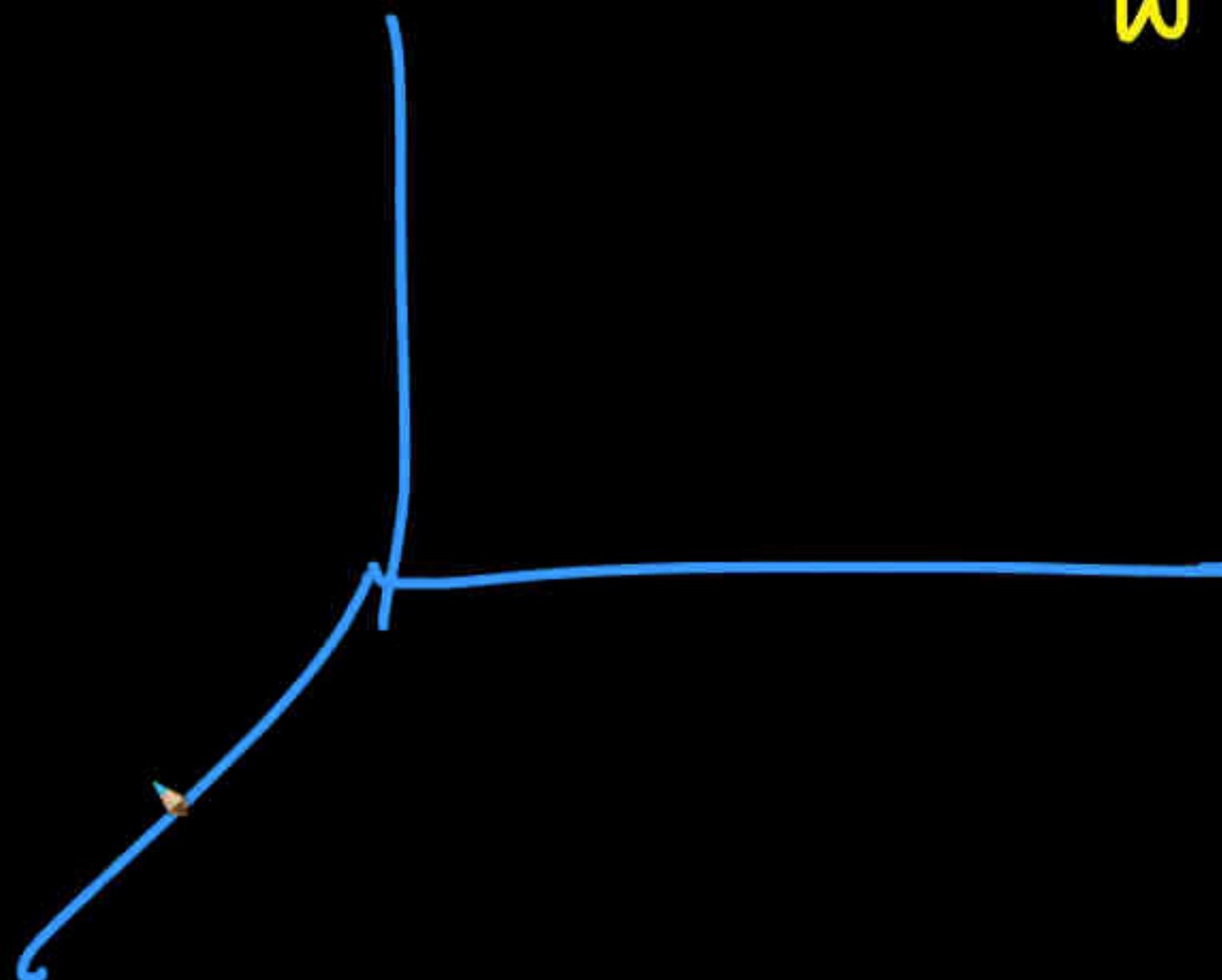
### For continuous random vari...





ML: vedn̄ -dev

$$\nabla_{\vec{w}} \text{Loss} = \left[ \begin{array}{c} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{array} \right]$$



MatrixCalc... x plot(x^2+y... x PCA\_tSNE... x JMLR\_200... x Student's t... x Kullback-L... x How to Use... x sklearn.ma... x t-SNE: The... x 1802.0342... x UMAP Dim... x Basic UMAP... x New Tab x +

google.com/search?q=plot(x%5E2%2By%5E2&rlz=1C5CHFA\_enIN958IN958&oq=plot(x%5E2%2By%5E2&aqs=chrome..69i57j0i30l9.3644j0j4&sourceid=chrome&ie=UTF-8

Google plot(x^2+y^2) X |

All Images News Maps Videos More Tools

About 73,60,00,000 results (0.66 seconds)

Graph for  $x^2+y^2$

From: To:

x	-10.0000	10.0000
y	-10.0000	10.0000
z	-63.1513	189.482

17 / 17

google.com/search?q=plot(x^2+y^2&amp;rlz=1C5CHFA\_enIN958IN958&amp;oq=plot(x^2+y^2&amp;qs=chrome..69i57j0i30i9.3644j0j4&amp;sourceid=chrome&amp;ie=UTF-8



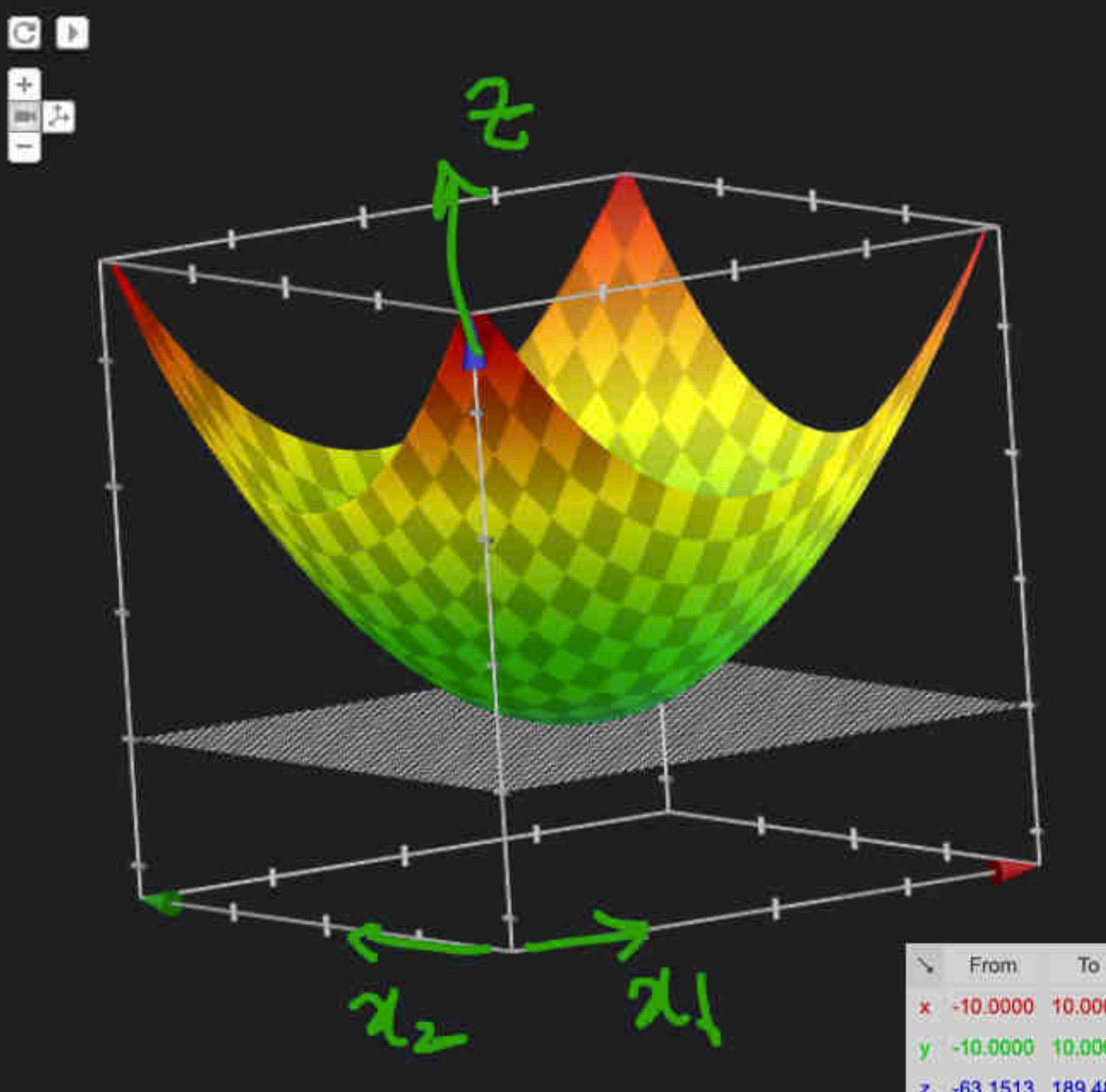
Google

plot(x^2+y^2



All Images News Maps Videos More Tools

About 73,60,00,000 results (0.66 seconds)

Graph for  $x^2+y^2$ 

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

scalar

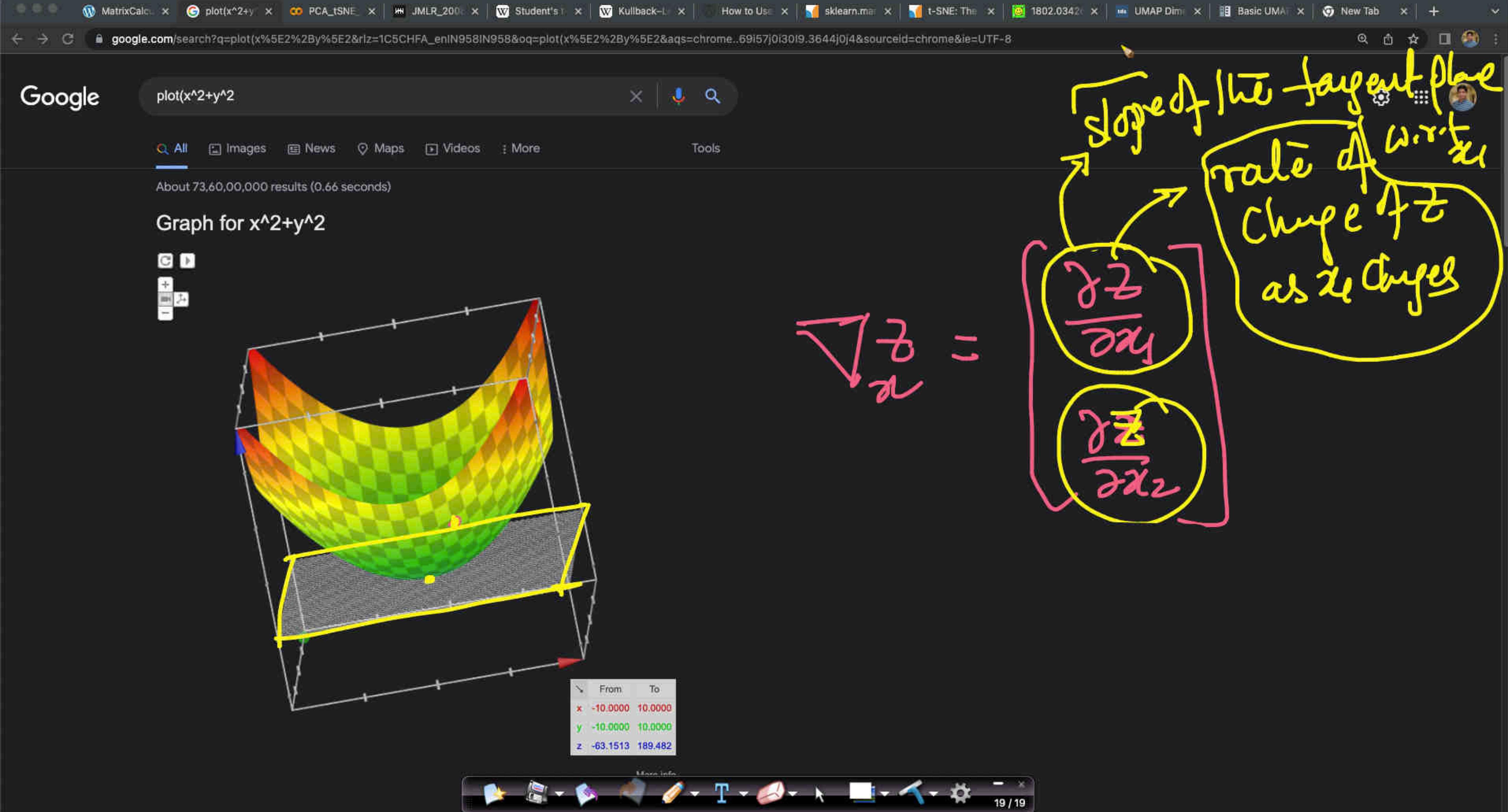
$$z = x^2 + y^2$$

$$z = x_1^2 + x_2^2$$

$$z = \mathbf{x}^\top \mathbf{x}$$

$$\nabla_z z$$

$$= \begin{bmatrix} \frac{\partial z}{\partial x_1} \\ \frac{\partial z}{\partial x_2} \end{bmatrix}$$



## PCA-recap

*data points*

$$X = \begin{bmatrix} \vdots & & & & \end{bmatrix}^T \quad \begin{matrix} 1 & 2 & \dots & d \\ \leftarrow x_i \rightarrow \end{matrix}$$

$n \times d$

standardized

$$S : \text{cov}(X) = X^T X$$

$d \times d$

$\hookrightarrow$  symmetric

optimization

$$\min_u \sum_{i=1}^n (\bar{u}^\top x_i)^2$$

$$\text{s.t. } \bar{u}^\top \bar{u} = 1$$

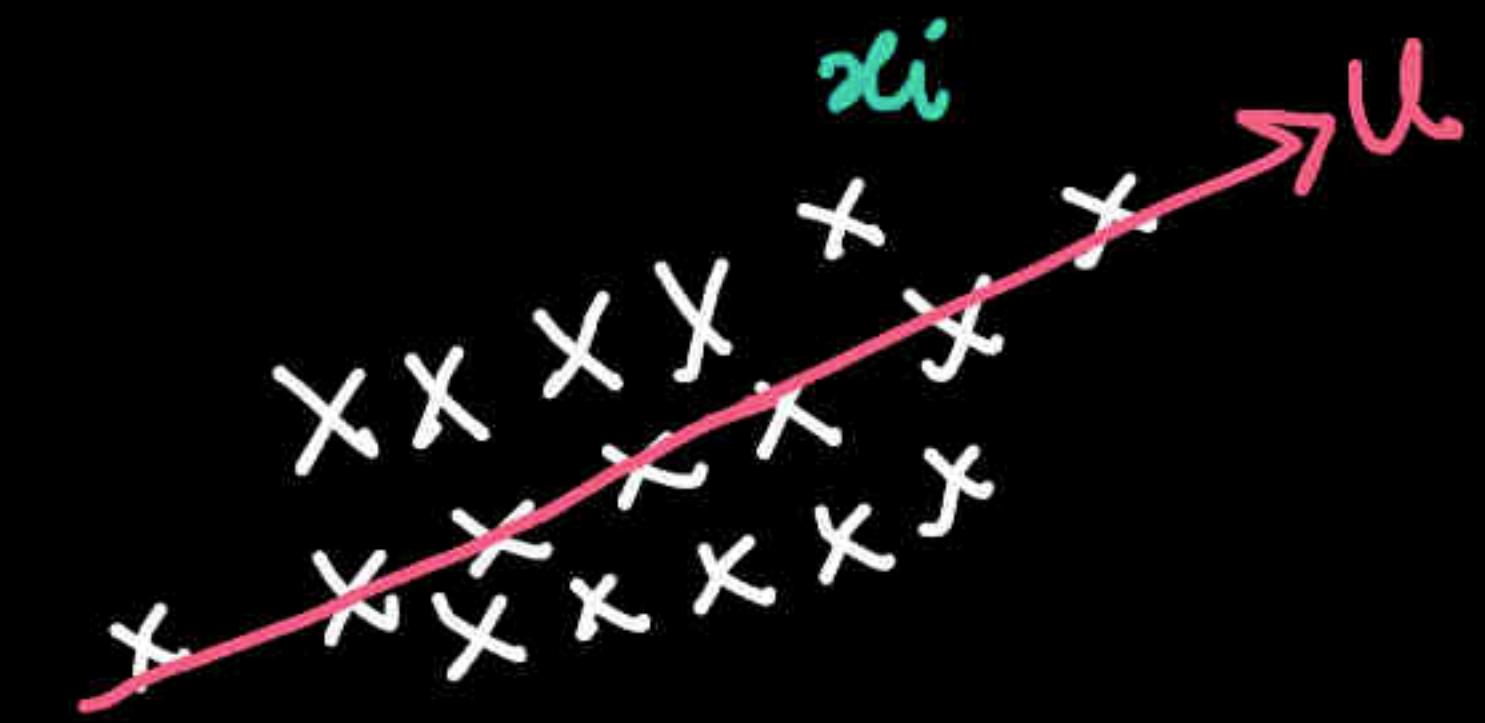
||

$$\min_u \bar{u}^\top S \bar{u}$$

$S \in \mathbb{R}^{d \times d}$

$$\text{s.t. } \bar{u}^\top \bar{u} = 1$$

geom



$$\bar{x} = \text{mean}(x_i) \rightarrow \vec{0}$$

$$\bar{u}^\top \bar{x} \rightarrow 0$$

$$\min_u \overbrace{u^T S u + \lambda u^T u}^{\mathcal{L}} \rightarrow \text{One } \underline{\text{minima}}$$

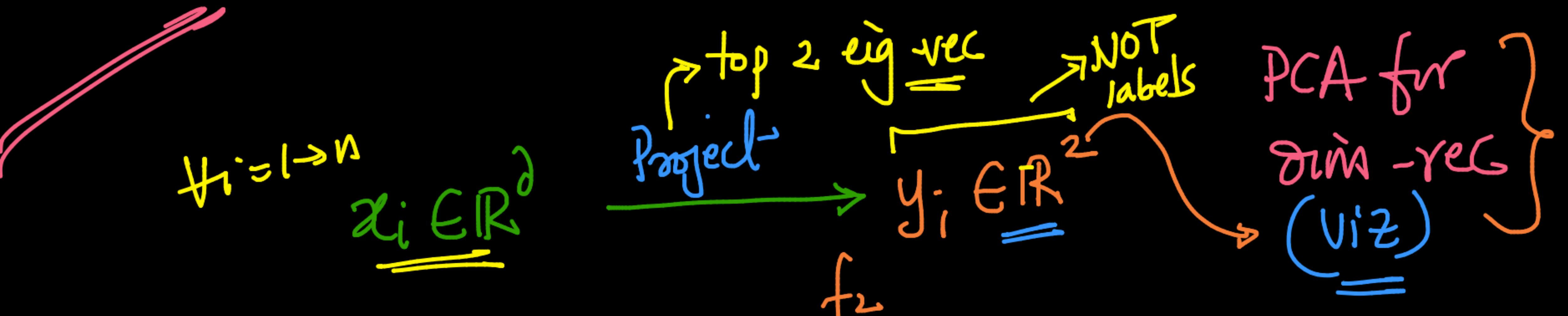
$$@ u^* \quad \frac{\partial \mathcal{L}}{\partial u} = 0$$

$$= \cancel{\lambda u^T S} + \cancel{\lambda u^T} = 0$$

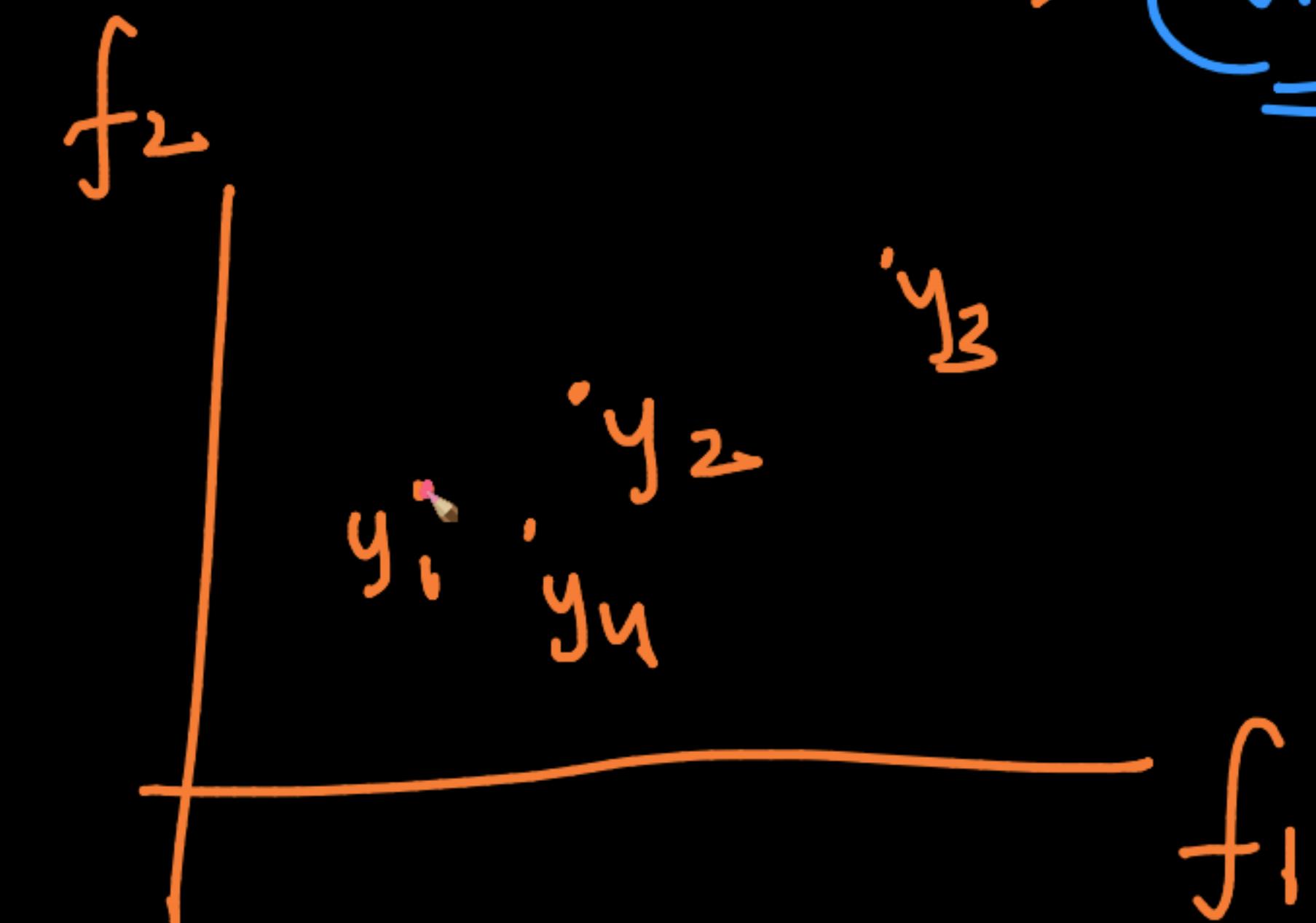
$$u^T S = -\cancel{\lambda} u^T$$

$S^T u = u \cancel{\lambda}$   
 $\cancel{S} u = \cancel{\lambda} u$   
 eigenval & eige vec

best  $u \rightarrow \cancel{\text{eigvec}}$  corresponding to largest eig-val



$\underline{\underline{u_1 \text{ & } u_2}}$  : 2 largest eig-vec



$$\begin{bmatrix} u_1^T x_i & u_2^T x_i \end{bmatrix}^T = y_i$$

$\downarrow$

$y_{i,1}$

$\downarrow$

$y_{i,2}$

$\dots$

$$f_2 = u_2$$

$y_{i,2}$

$\dots$

$y_{i,n}$

$f_1$

$= u_1$

I

PCA: linear mapping

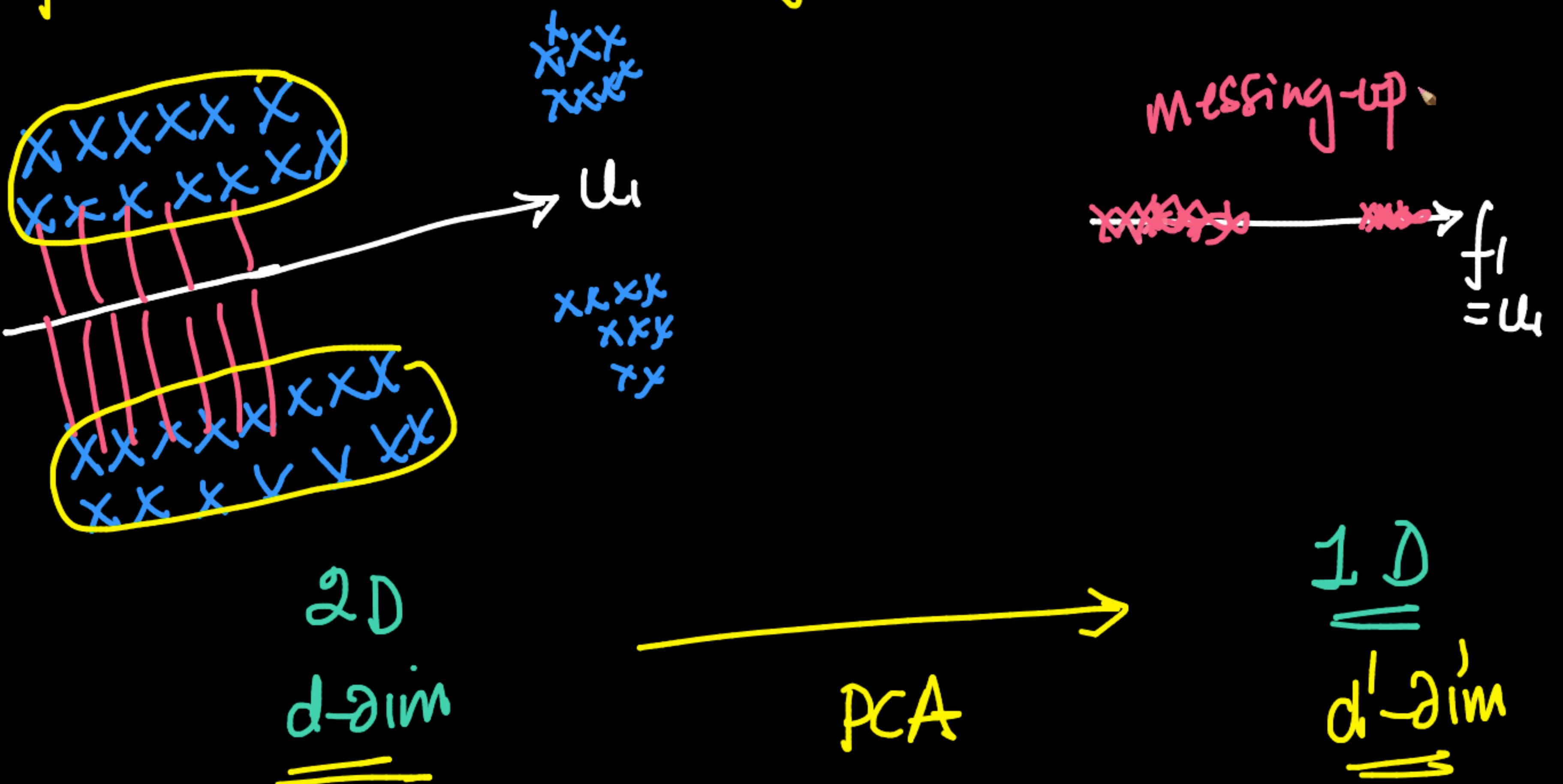
$$y_i = \begin{bmatrix} \xleftarrow{\quad U_1^T \quad} \\ \xleftarrow{\quad U_2^T \quad} \end{bmatrix} x_i \underset{2 \times 1}{\underset{d \times 1}{\text{---}}} \underset{2 \times d}{\text{---}}$$

$$\tilde{y}_i = U \cdot \tilde{x}_i \xrightarrow{\quad \sim \quad} \underset{d \text{-dim}}{\text{---}} \quad \text{Linear-mapping}$$

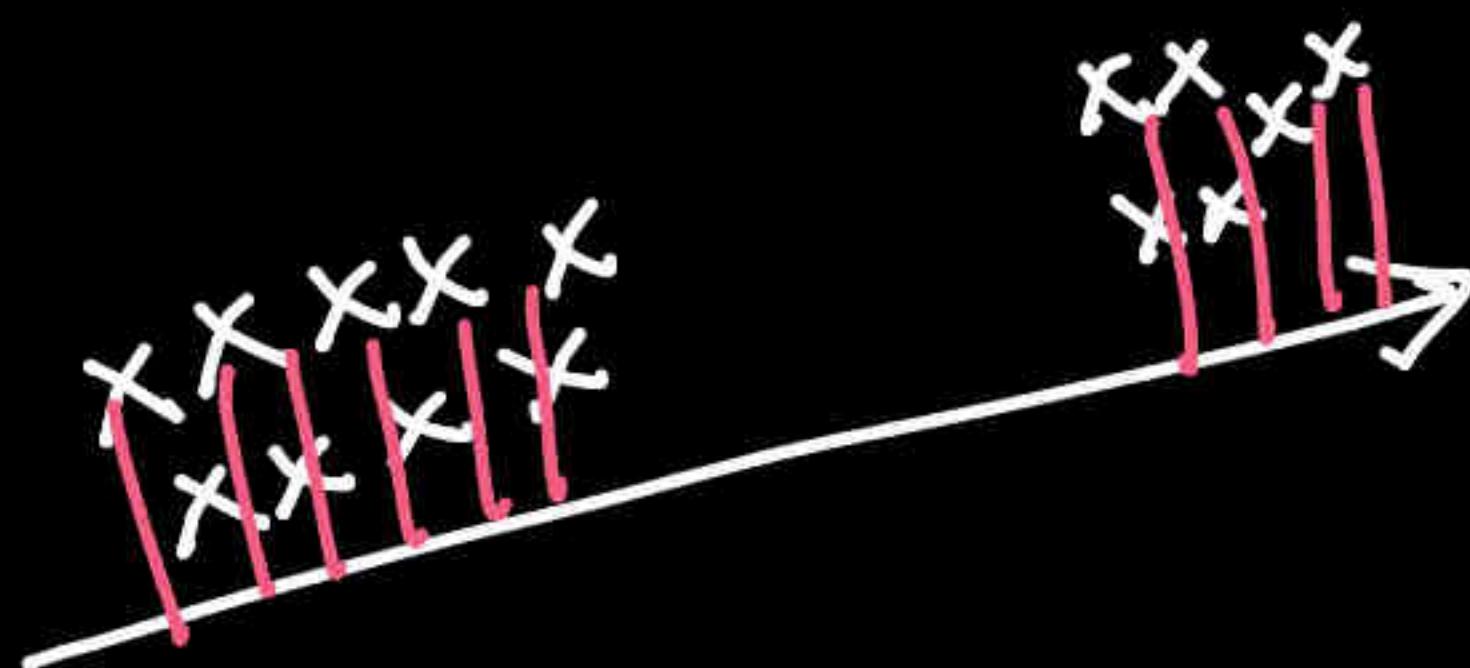
2-dim rep. of  $x_i$

2

does Not ~~often~~ preserve the local/neighborhood info



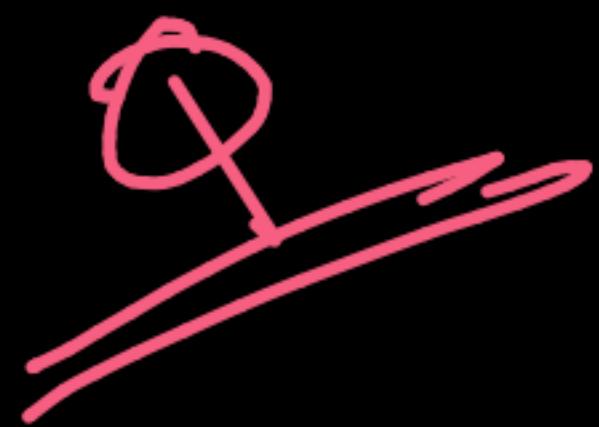
preserving neighbors



2D

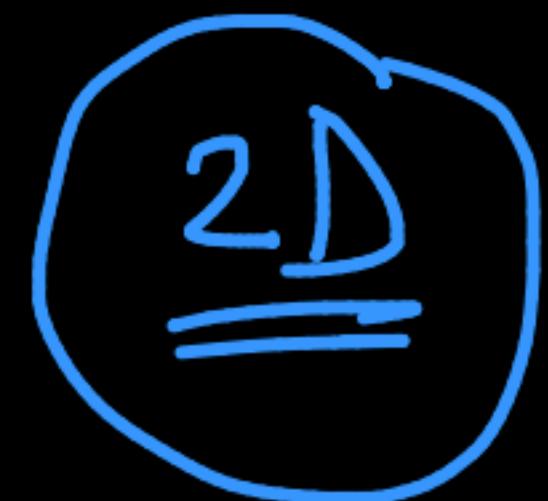


1D

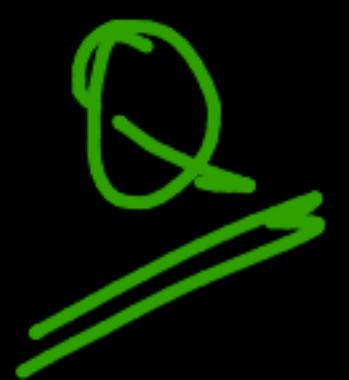


high dim viz using clustering

$$\underline{x_i \in \mathbb{R}^d}$$



- ✓ {
  - polar plots
  - parallel coordinate plots
  - hierarchical clustering  
(tree)



PCA for regression?

$x_i \in \mathbb{R}^d$



$\hat{\beta}: V \cdot V \cdot \text{large}$   
 $\downarrow$  PCA → loss of info

$d'$  → build model

Heavy (large)  $\lambda$   
Regularization

→ can be avoided

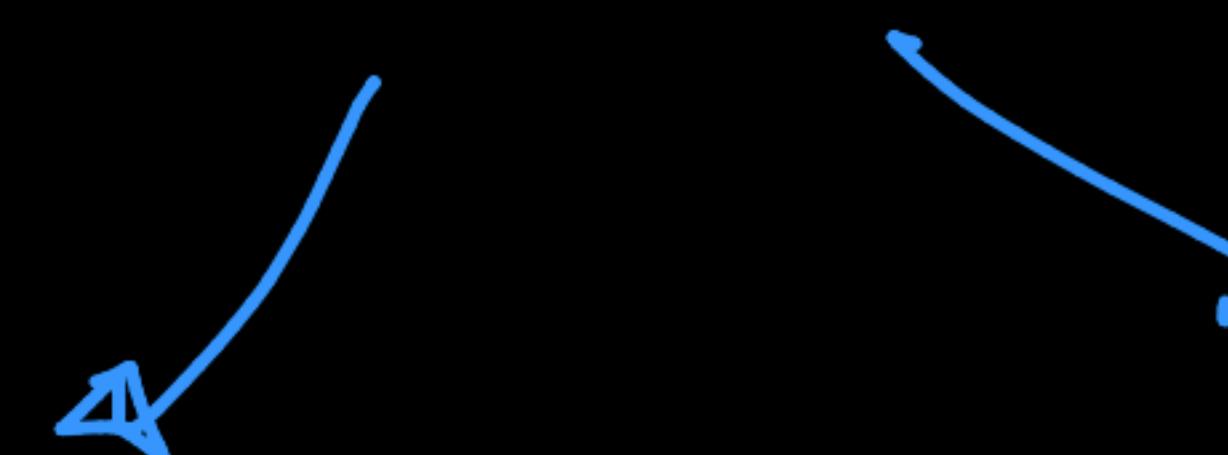
can be done





(3)

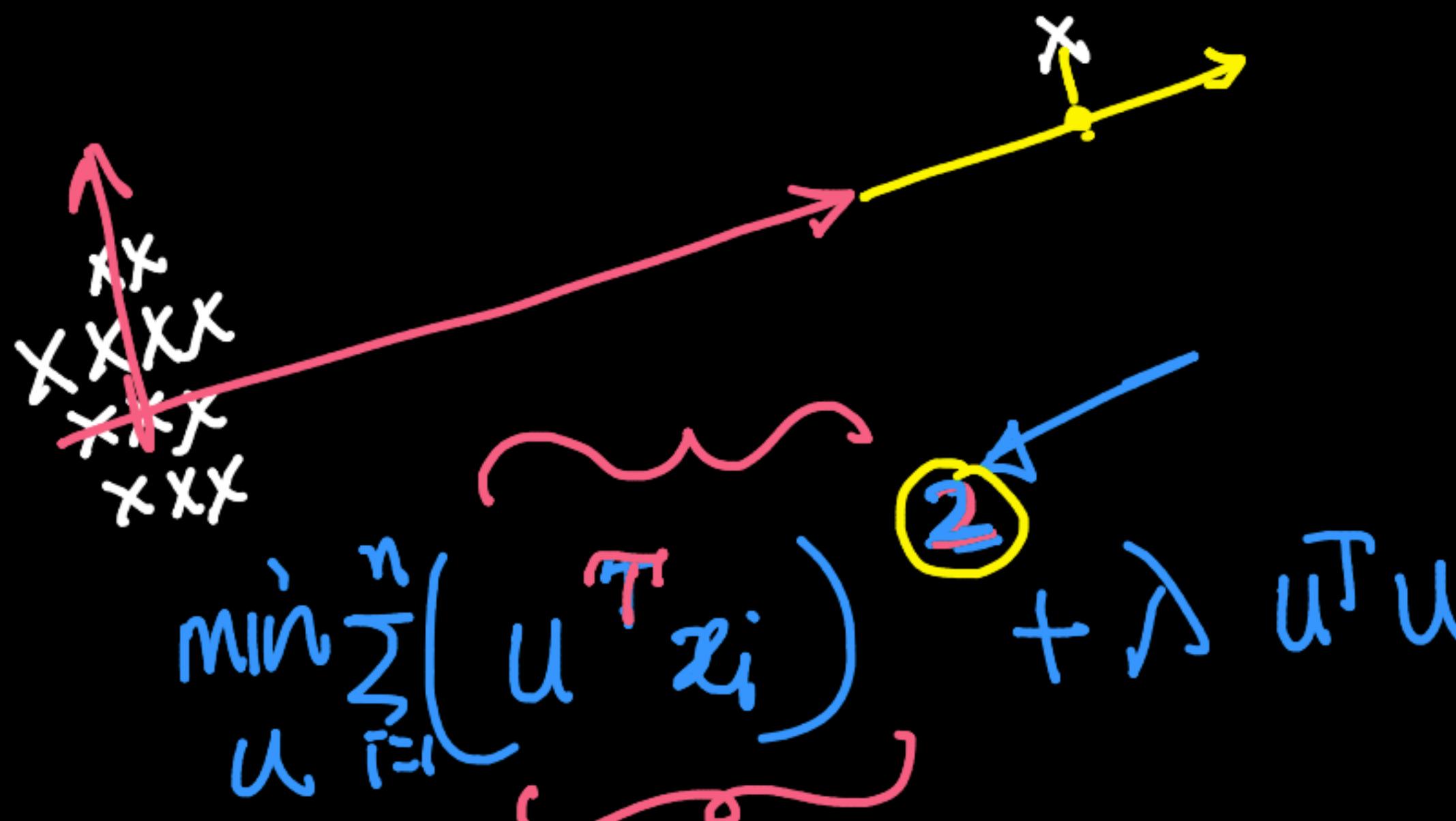
PCA: v. low noise tolerance



How to solve it?

Why?

geom:



Math:

→ remove outliers }  
using stats method

→ RANSAC-like  
Li-PCA → abs-value

$$\min_u \sum_{i=1}^n |u^T x_i| + \lambda u^T u$$

L1

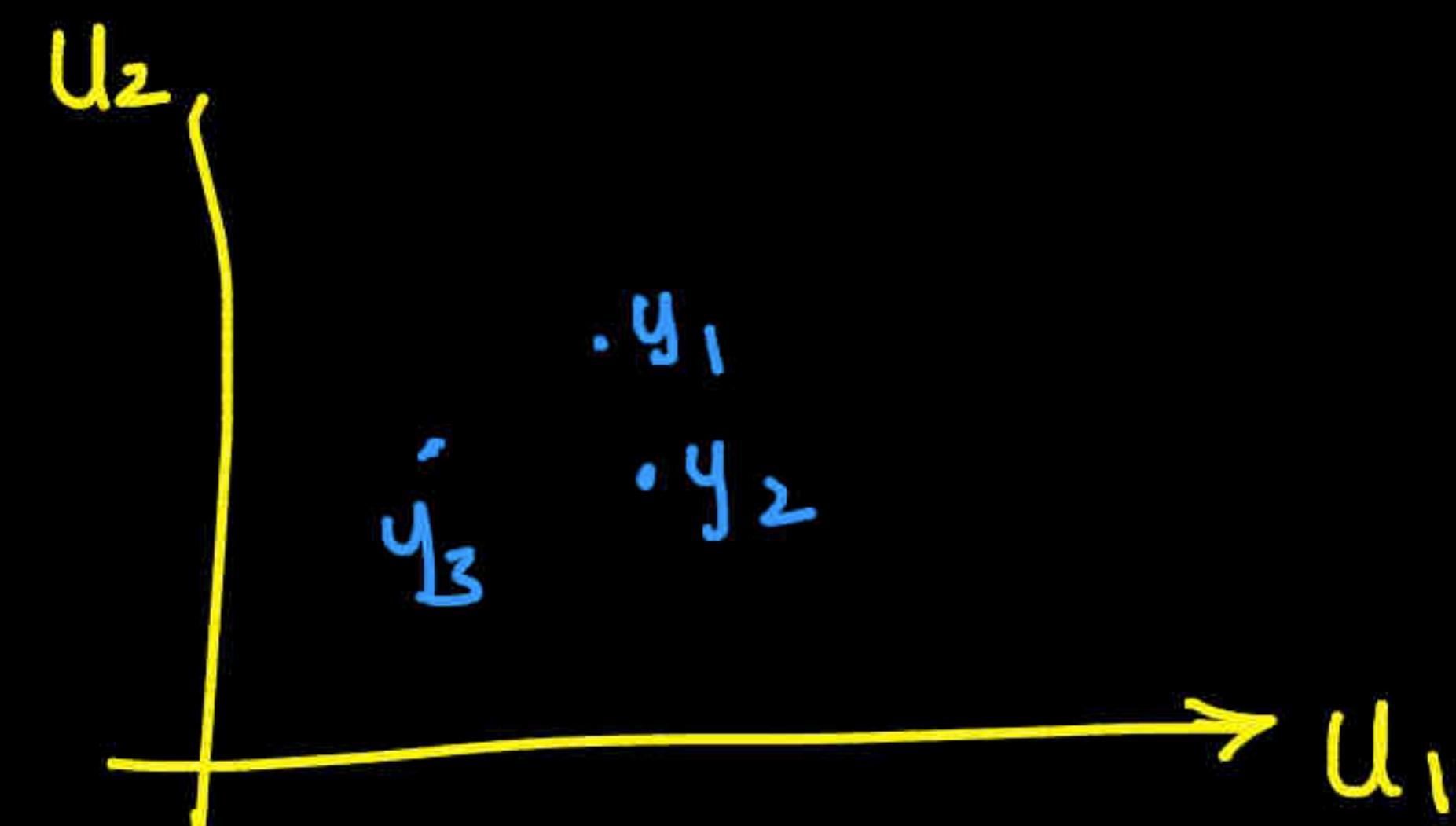
Q  
||

$$x_i \in \mathbb{R}^d$$

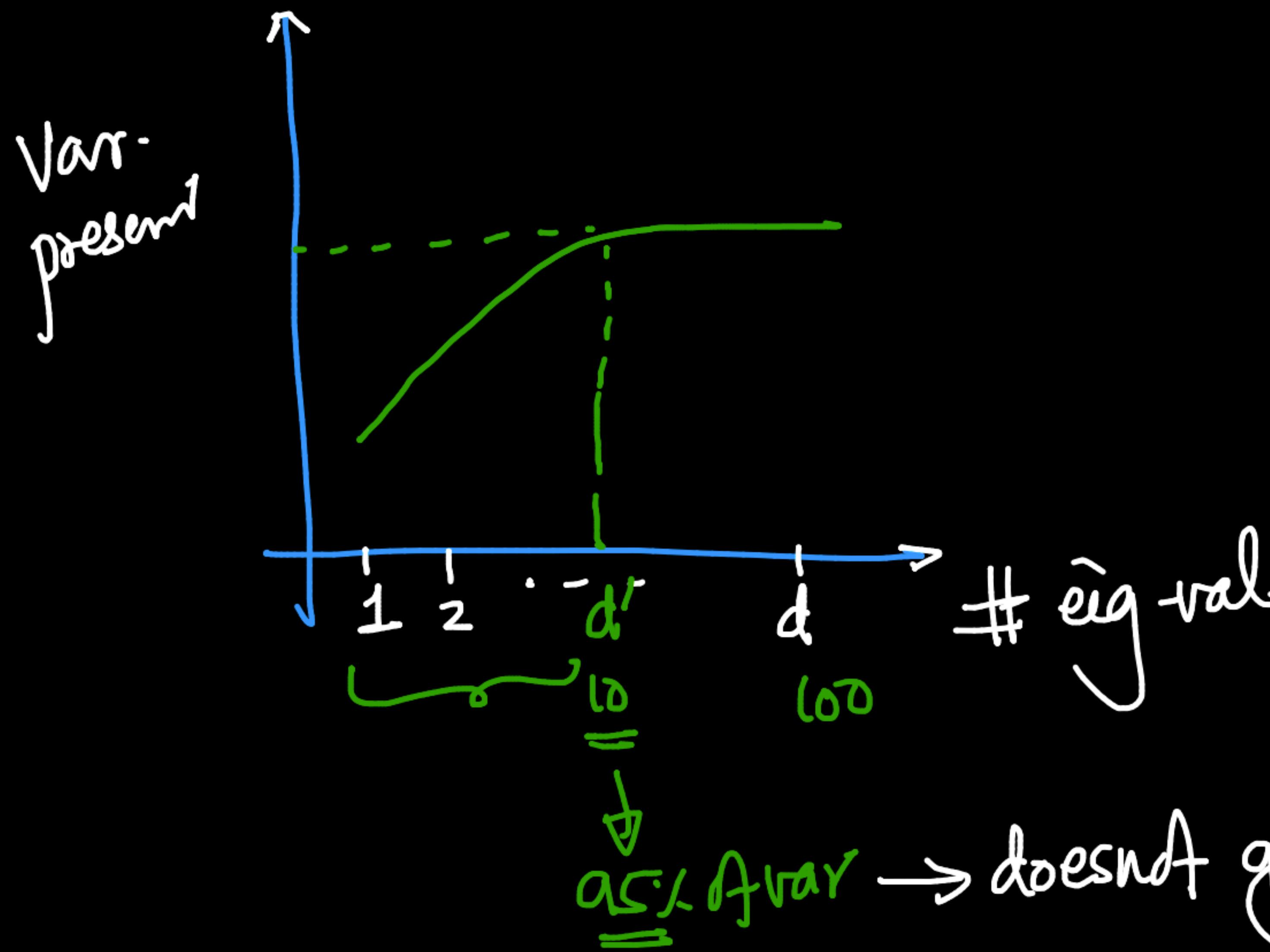
PCA

$$y_i \in \mathbb{R}^2 \rightarrow 2 \text{ or } 3$$

Vit  
||



2D  
||

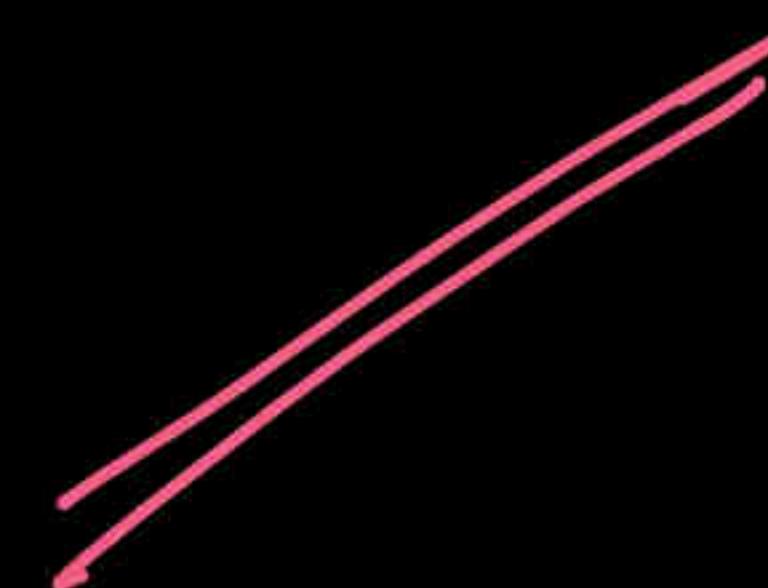


dim-red

$d \rightarrow d'$

$d' < d$

doesn't guarantee Neighbr  
preservation



+ Code + Text

Reconnect



```
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import cm
import seaborn as sns

from openTSNE import TSNE
from umap import UMAP
from sklearn import decomposition
```

PCA: Matrix decompo  
= SVD; NNMF ...  
(later)

```
[ ] from sklearn.datasets import load_digits
digits = load_digits()
```

```
[ ] digits.keys()
```

```
dict_keys(['data', 'target', 'frame', 'feature_names', 'target_names', 'images', 'DESCR'])
```

```
[ ] X = digits.data
Y = digits.target
```

```
[ ] X.shape
```



MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=yUxdlwz\_AmA

+ Code + Text Reconnect

from sklearn import decomposition

[ ] from sklearn.datasets import load\_digits  
digits = load\_digits()

[ ] digits.keys()  
dict\_keys(['data', 'target', 'frame', 'feature\_names', 'target\_names', 'images', 'DESCR'])

[ ] X = digits.data  
Y = digits.target

[ ] X.shape  
(1797, 64)

[ ] digits.images.shape  
(1797, 8, 8)

[ ] plt.gray()  
plt.imshow(digits.images[10])



MatrixCalculus x PCA\_tSNE\_ x JMLR\_2008. x Student's t- x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=yUxdlwz\_AmA

+ Code + Text Reconnect

from sklearn import decomposition

[ ] from sklearn.datasets import load\_digits  
digits = load\_digits()

[ ] digits.keys()

[ ] dict\_keys(['data', 'target', 'frame', 'feature\_names', 'target\_names', 'images', 'DESCR'])

[ ] X = digits.data  
Y = digits.target

[ ] X.shape  
(1797, 64)

[ ] digits.images.shape  
(1797, 8, 8)

[ ] plt.gray()  
plt.imshow(digits.images[10])

Reconnect

Up Down Left Right Home Stop Refresh Back Forward

38 / 38

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13o1Y9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=yUxdlwz\_AmA

+ Code + Text Reconnect  

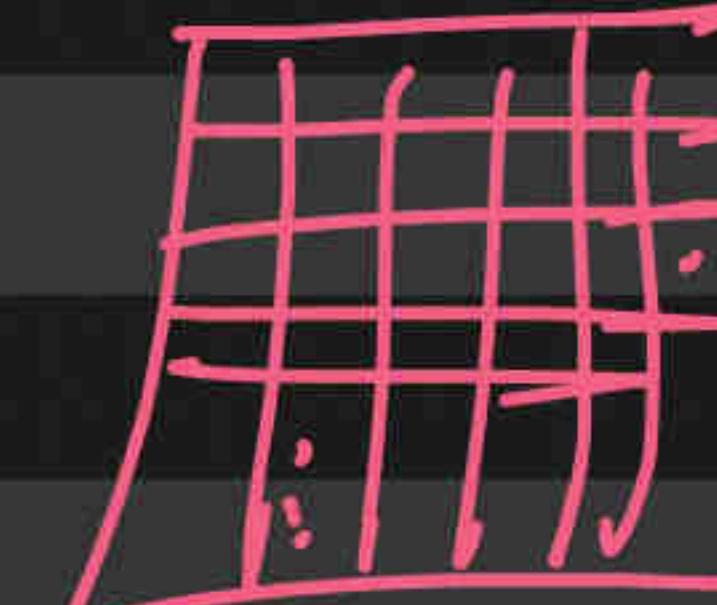
```
dict_keys(['data', 'target', 'frame', 'feature_names', 'target_names', 'images', 'DESCR'])
```

[ ] X = digits.data  
Y = digits.target

[ ] X.shape  
(1797, 64)

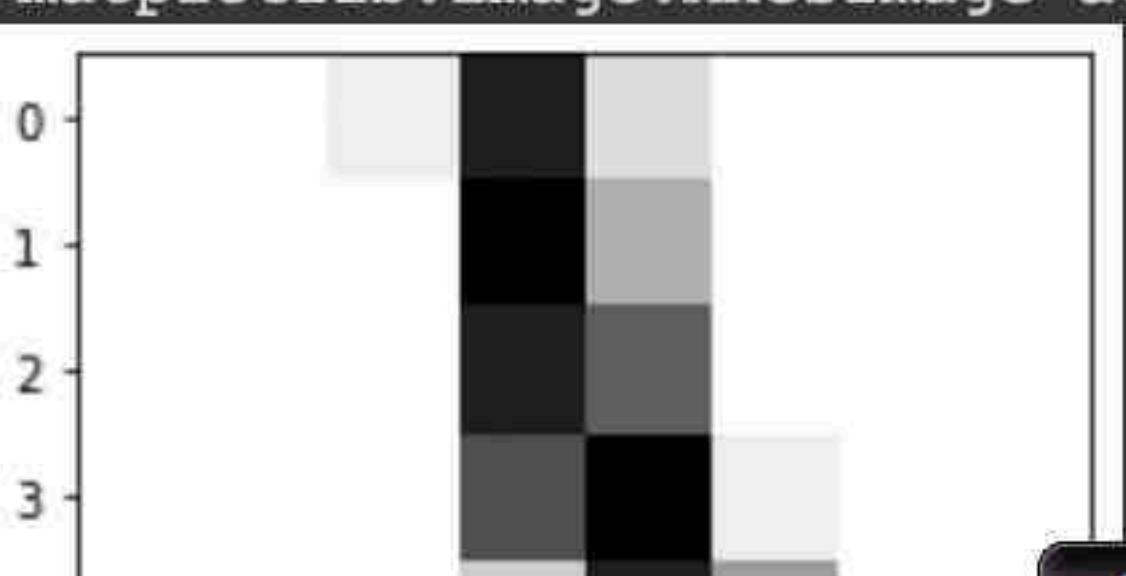
[ ] digits.images.shape  
(1797, 8, 8)

plt.gray()  
plt.imshow(digits.images[1000], cmap=plt.cm.gray\_r)



0 — 255  
B W

<matplotlib.image.AxesImage at 0x7f5b003e5110>



39 / 39

+ Code + Text

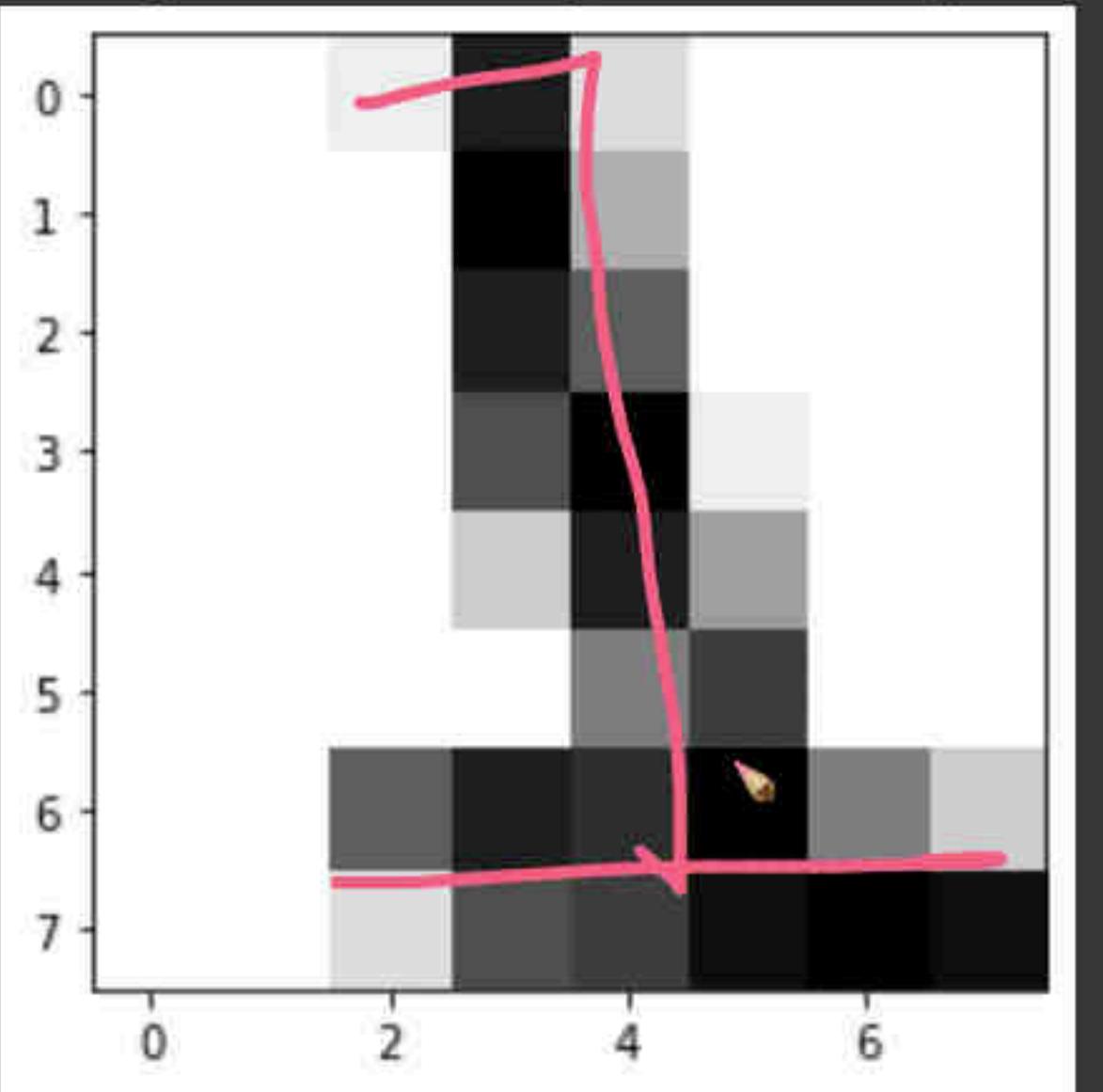
Reconnect



```
[ ] plt.gray()  
plt.imshow(digits.images[1000], cmap=plt.cm.gray_r)
```

{x}  

```
<matplotlib.image.AxesImage at 0x7f5b003e5110>
```

<>  

```
[ ] #PCA  
%%time  
pca_2D = decomposition.PCA(n_components=2)  
pca_2D.fit(X)  
Z1 = pca_2D.transform(X)
```

+ Code + Text

Reconnect



[ ] Y = digits.target

[ ] X.shape

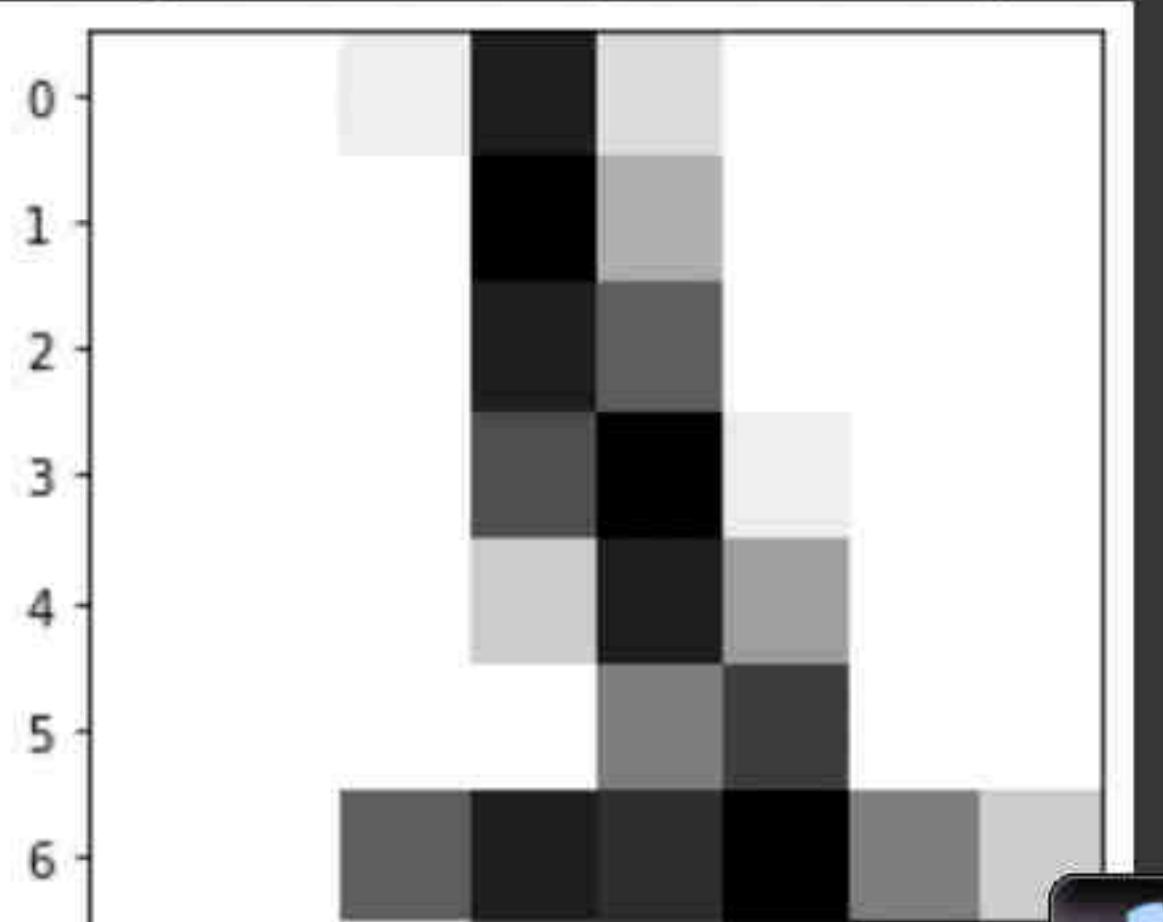
{x} (1797, 64)

[ ] digits.images.shape

[ ] (1797, 8, 8)

[ ] plt.gray()  
plt.imshow(digits.images[1000], cmap=plt.cm.gray\_r)

&lt;matplotlib.image.AxesImage at 0x7f5b003e5110&gt;



64x64

[920]

[080]

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=yUxdlwz\_AmA

+ Code + Text Reconnect  

```
dict_keys(['data', 'target', 'frame', 'feature_names', 'target_names', 'images', 'DESCR'])
```

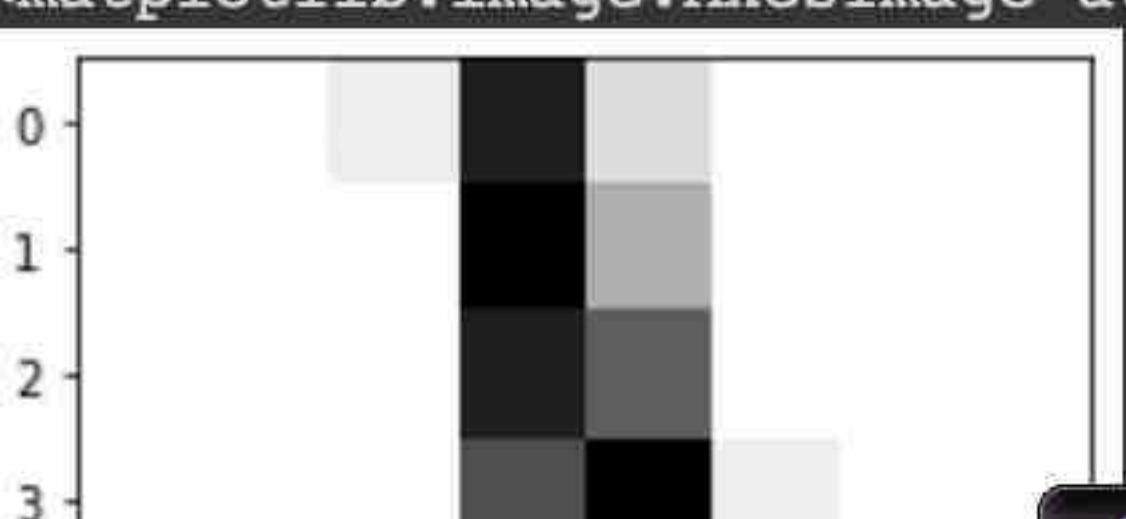
{x} [ ] X = digits.data  
Y = digits.target

[ ] X.shape  
**(1797, 64)**

[ ] digits.images.shape  
**(1797, 8, 8)**

▶ plt.gray()  
plt.imshow(digits.images[1000], cmap=plt.cm.gray\_r)

[ ] <matplotlib.image.AxesImage at 0x7f5b003e5110>



42 / 42

 + Code  + Text

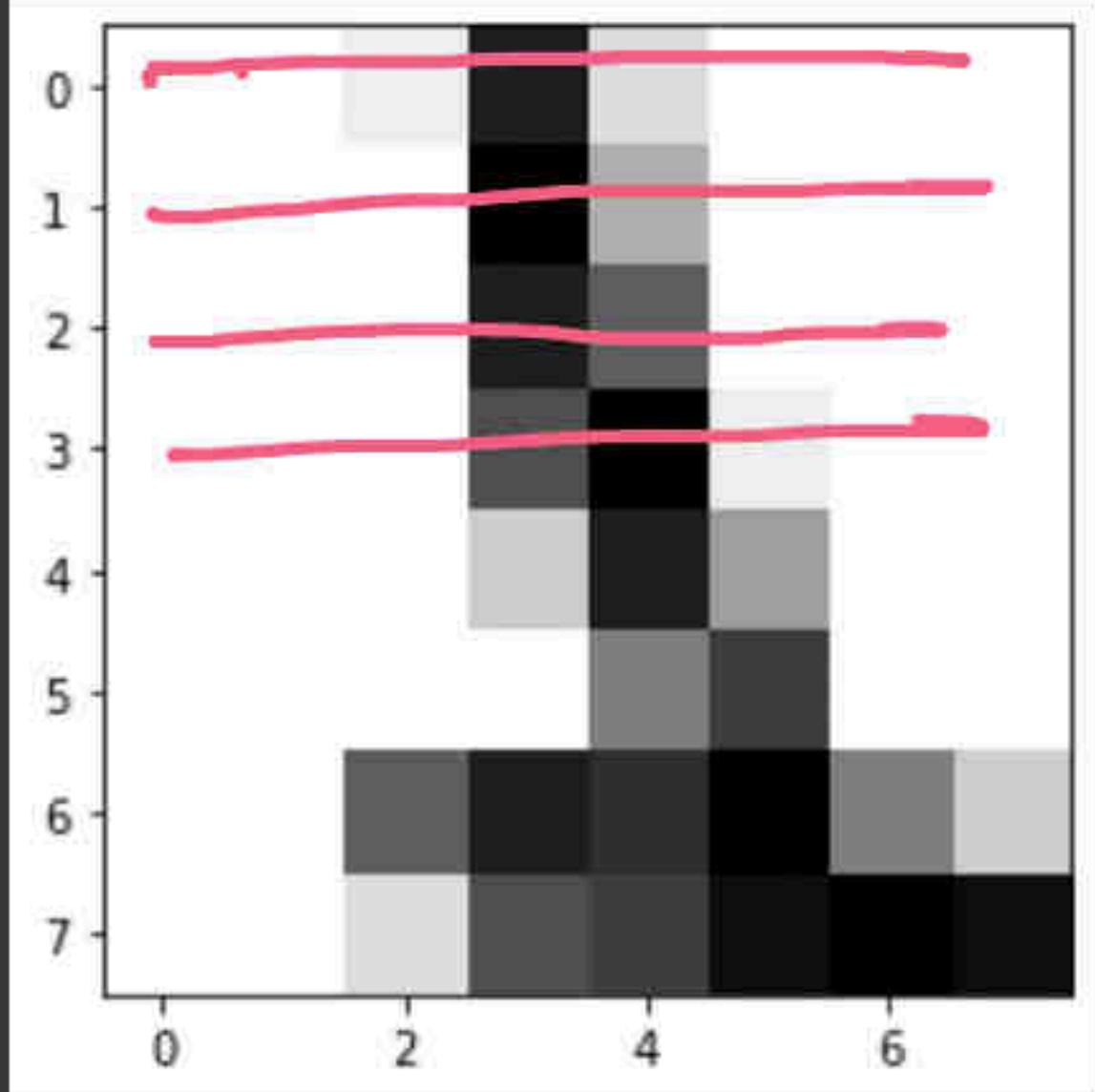
[ ] digits.images.shape

(1797, 8, 8)

{x}

[ ] plt.gray()  
plt.imshow(digits.images[1000], cmap=plt.cm.gray\_r)

&lt;matplotlib.image.AxesImage at 0x7f5b003e5110&gt;



255 255  
8x8 → 64

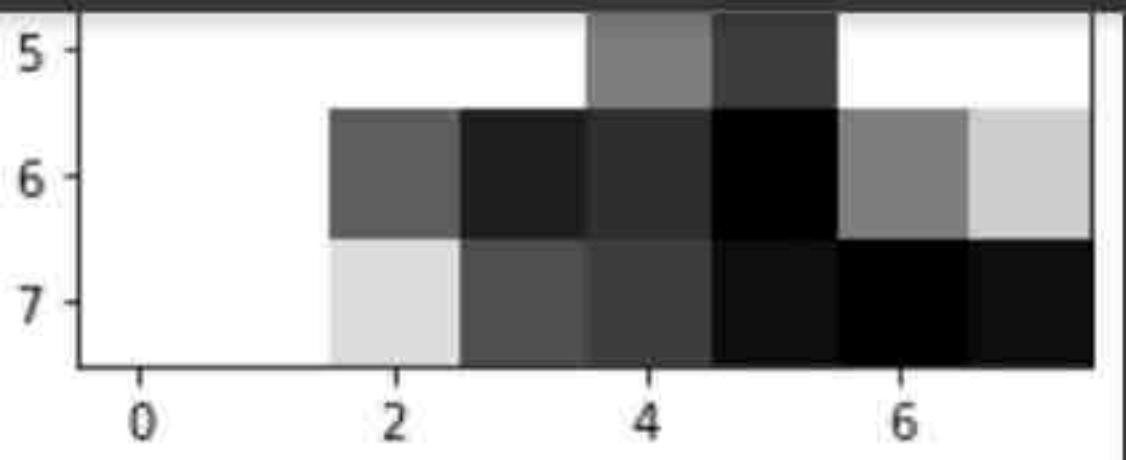
<> [ ] #PCA  
%%time  
pca\_2D = decomposition.PCA(

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oIY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=zpy9Pd8nAjTQ

+ Code + Text Reconnect

[ ]



{x}

[ ] #PCA  
%%time  
{  
pca\_2D = decomposition.PCA(n\_components=2)  
pca\_2D.fit(X)  
Z1 = pca\_2D.transform(X)

64.2ms → 2.3ms

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms  
Wall time: 18.4 ms

[ ] from matplotlib.colors import ListedColormap  
cmap = ListedColormap(sns.husl\_palette(len(np.unique(Y))))

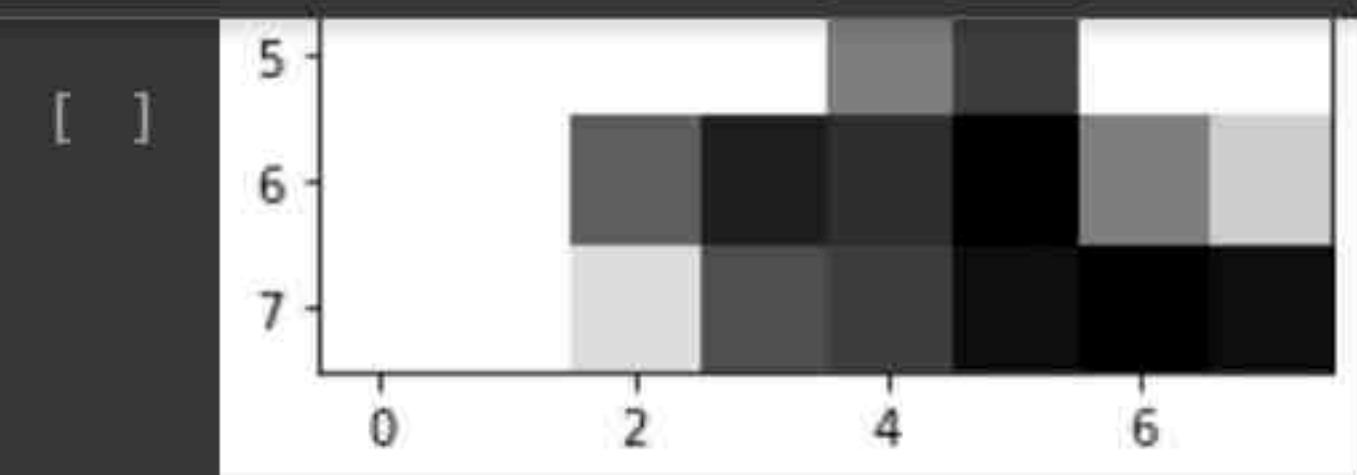
▶ fig, ax = plt.subplots()  
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')

30 9

44 / 44

+ Code + Text

Reconnect



X : 1797 x 64

```
[ ] #PCA  
%%time  
pca_2D = decomposition.PCA(n_components=2)  
pca_2D.fit(X)  
Z1 = pca_2D.transform(X)
```

64 dim → 2 dim

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms

Wall time: 18.4 ms

```
[ ] from matplotlib.colors import ListedColormap  
cmap = ListedColormap(sns.husl_palette(len(np.unique(Y))))
```

```
▶ fig, ax = plt.subplots()  
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')
```

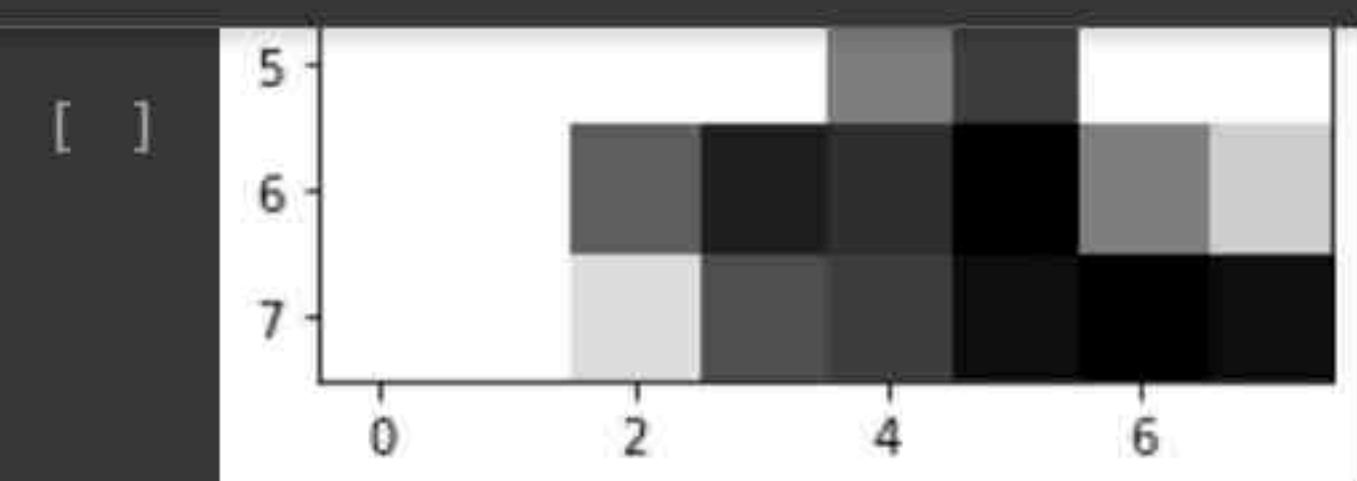
30

9



+ Code + Text

Reconnect



[ ] #PCA  
%%time  
{  
pca\_2D = decomposition.PCA(n\_components=2)  
pca\_2D.fit(X)  
z1 = pca\_2D.transform(X) → *میں جسے*

$$z_i = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} x_i$$

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms

Wall time: 18.4 ms

[ ] from matplotlib.colors import ListedColormap  
cmap = ListedColormap(sns.husl\_palette(len(np.unique(Y))))

▶ fig, ax = plt.subplots()  
im = ax.scatter(z1[:,0], z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')

30

9

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=zpy9Pd8nAjTQ

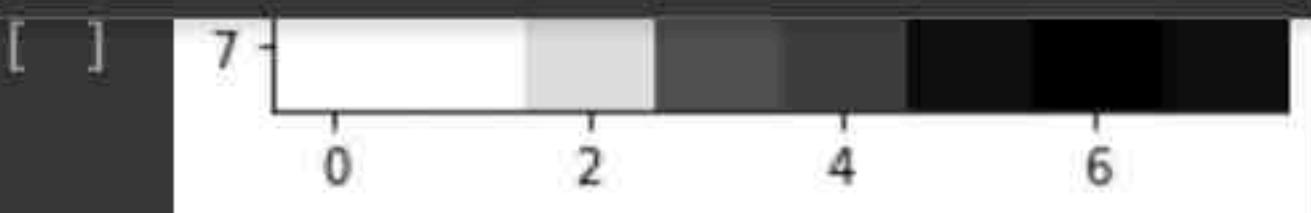
+ Code + Text Reconnect

```
[ ] plt.gray()
plt.imshow(digits.images[1000], cmap=plt.cm.gray_r)
```

<matplotlib.image.AxesImage at 0x7f5b003e5110>

[ ] #PCA
%%time
pca\_2D = decomposition.PCA(n\_components=2)
pca\_2D.fit(X)
Z1 = pca\_2D.transform(X)

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms

 + Code  + Text

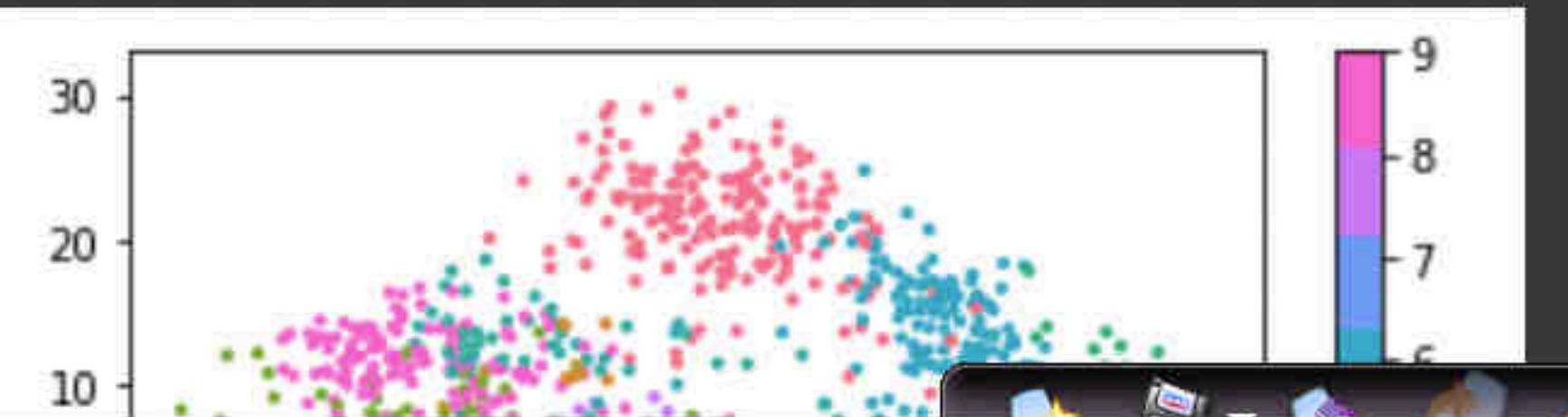
{x}

```
[ ] #PCA
{ %%time
  pca_2D = decomposition.PCA(n_components=2)
  pca_2D.fit(X)
  Z1 = pca_2D.transform(X)
```

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms  
Wall time: 18.4 ms

```
[ ] from matplotlib.colors import ListedColormap
cmap = ListedColormap(sns.husl_palette(len(np.unique(Y))))
```

```
▶ fig, ax = plt.subplots()
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')
cbar = fig.colorbar(im, ax=ax, label='Digit')
```



MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=zpy9Pd8nAjTQ

+ Code + Text Reconnect

[ ] 7

{x} [ ] #PCA  
%%time  
pca\_2D = decomposition.PCA(n\_components=2)  
pca\_2D.fit(X)  
Z1 = pca\_2D.transform(X)

64.21ms → 281ms  
24

CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms  
Wall time: 18.4 ms

[ ] from matplotlib.colors import ListedColormap  
cmap = ListedColormap(sns.husl\_palette(len(np.unique(Y))))

▶ fig, ax = plt.subplots()  
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')

30  
20  
10

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=zpy9Pd8nAJTQ

+ Code + Text Reconnect

```
[ ] z1 = pca_2D.transform(X)
CPU times: user 18.2 ms, sys: 8.02 ms, total: 26.2 ms
Wall time: 18.4 ms
```

{x}

```
[ ] from matplotlib.colors import ListedColormap
cmap = ListedColormap(sns.husl_palette(len(np.unique(Y))))
```

□

```
▶ fig, ax = plt.subplots()
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')
cbar = fig.colorbar(im, ax=ax, label='Digit')
```

Digit

50 / 50

MatrixCalculus x PCA\_tSNE\_U x JMLR\_2008 x Student's t-distribution x Kullback-Leibler x How to Use t-SNE x sklearn.man x t-SNE: The e x 1802.03426 x UMAP Dimensionality Reduction x Basic UMAP x New Tab x +

colab.research.google.com/drive/13oY9wIU97WrwgEwKa9XijJr0RF\_CbjH#scrollTo=zpy9Pd8nAjTQ

+ Code + Text Reconnect

Wall time: 18.4 ms

```
[ ] from matplotlib.colors import ListedColormap  
cmap = ListedColormap(sns.husl_palette(len(np.unique(Y))))
```

```
[ ] fig, ax = plt.subplots()  
im = ax.scatter(Z1[:,0], Z1[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')
```

$x_i$

$Z_1$  in 2D

64dim

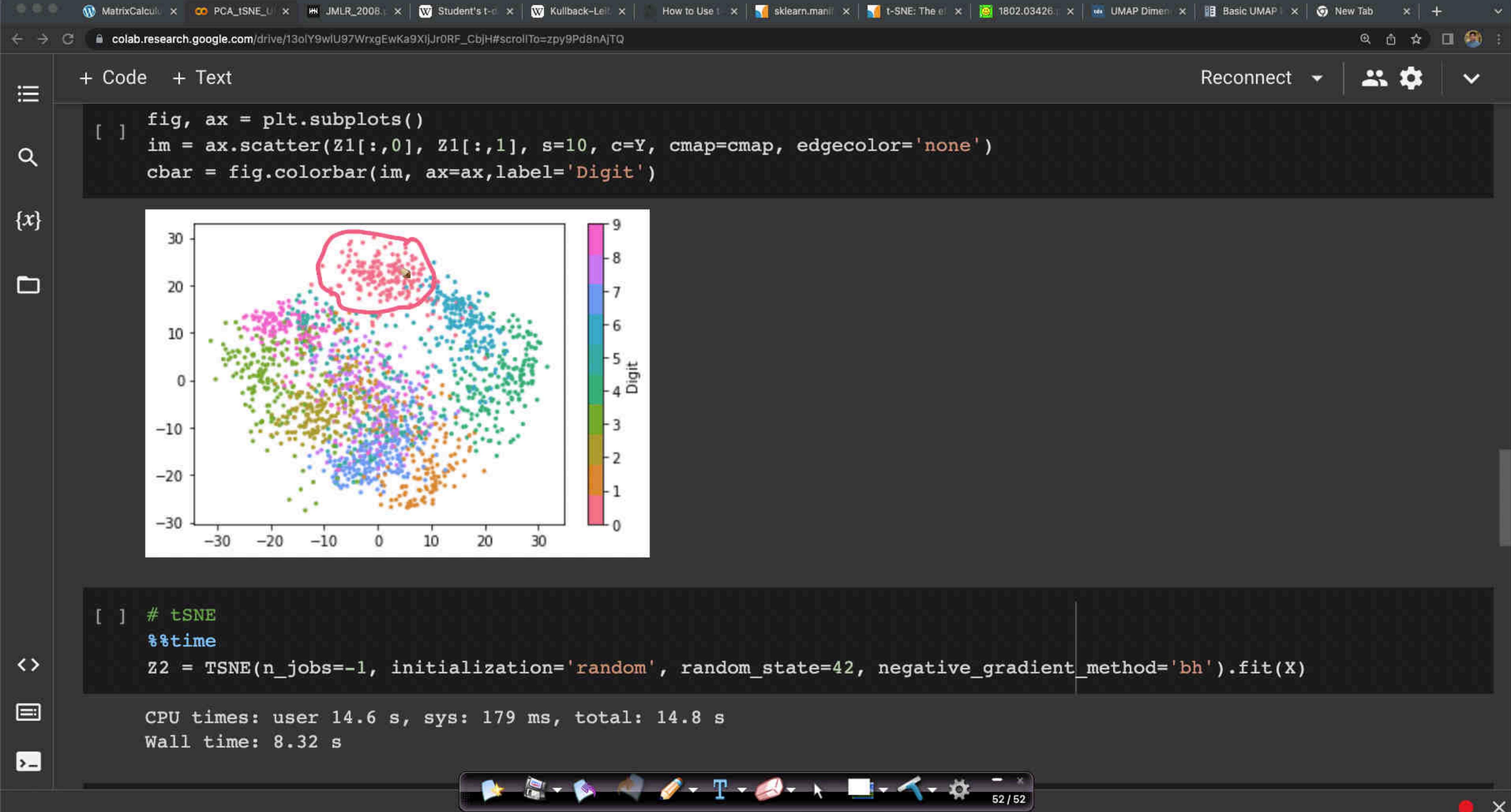
30  
20  
10  
0  
-10  
-20  
-30

-30 -20 -10 0 10 20 30

9  
8  
7  
6  
5  
4  
3  
2  
1  
0

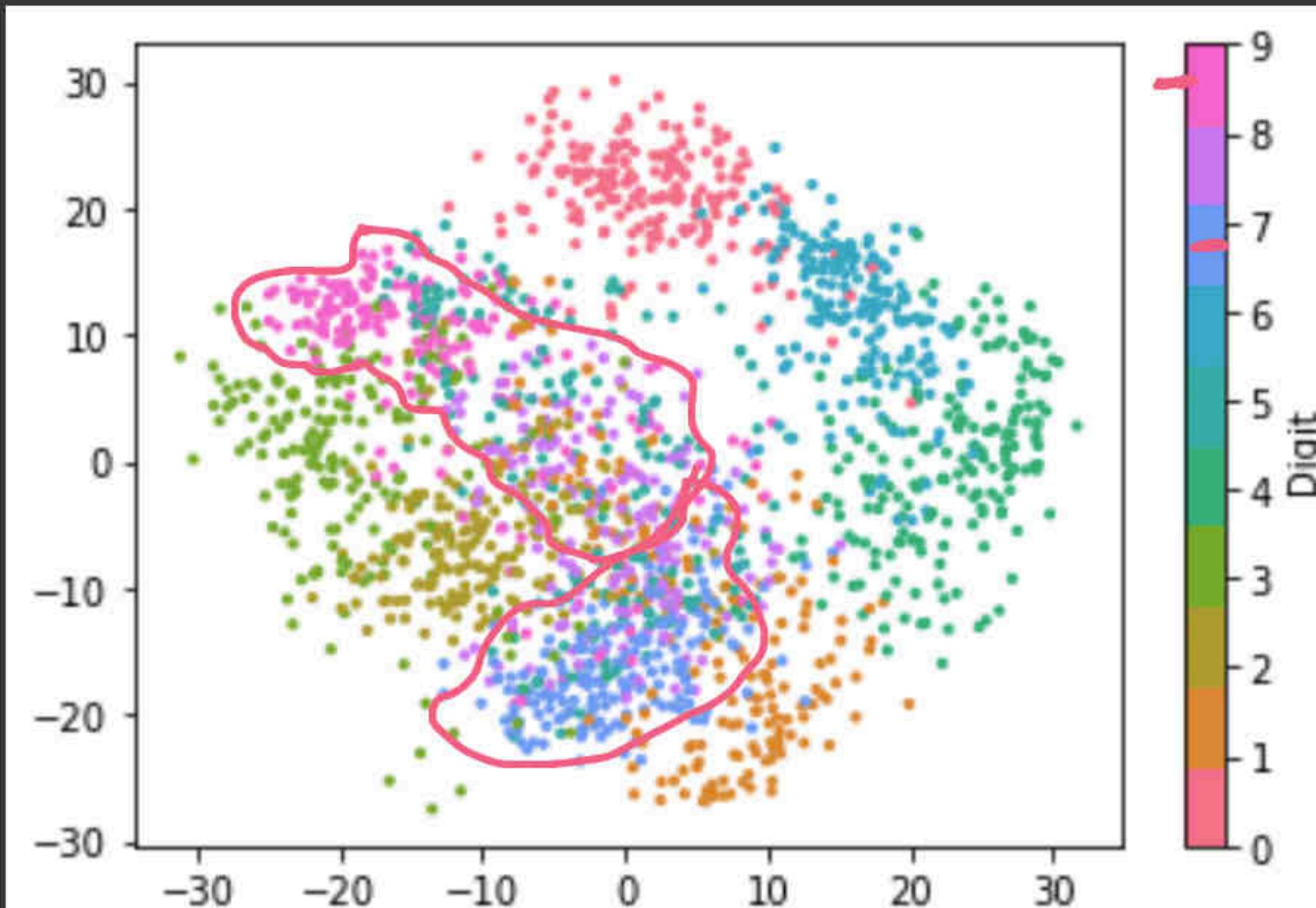
Digit

51 / 51





Reconnect



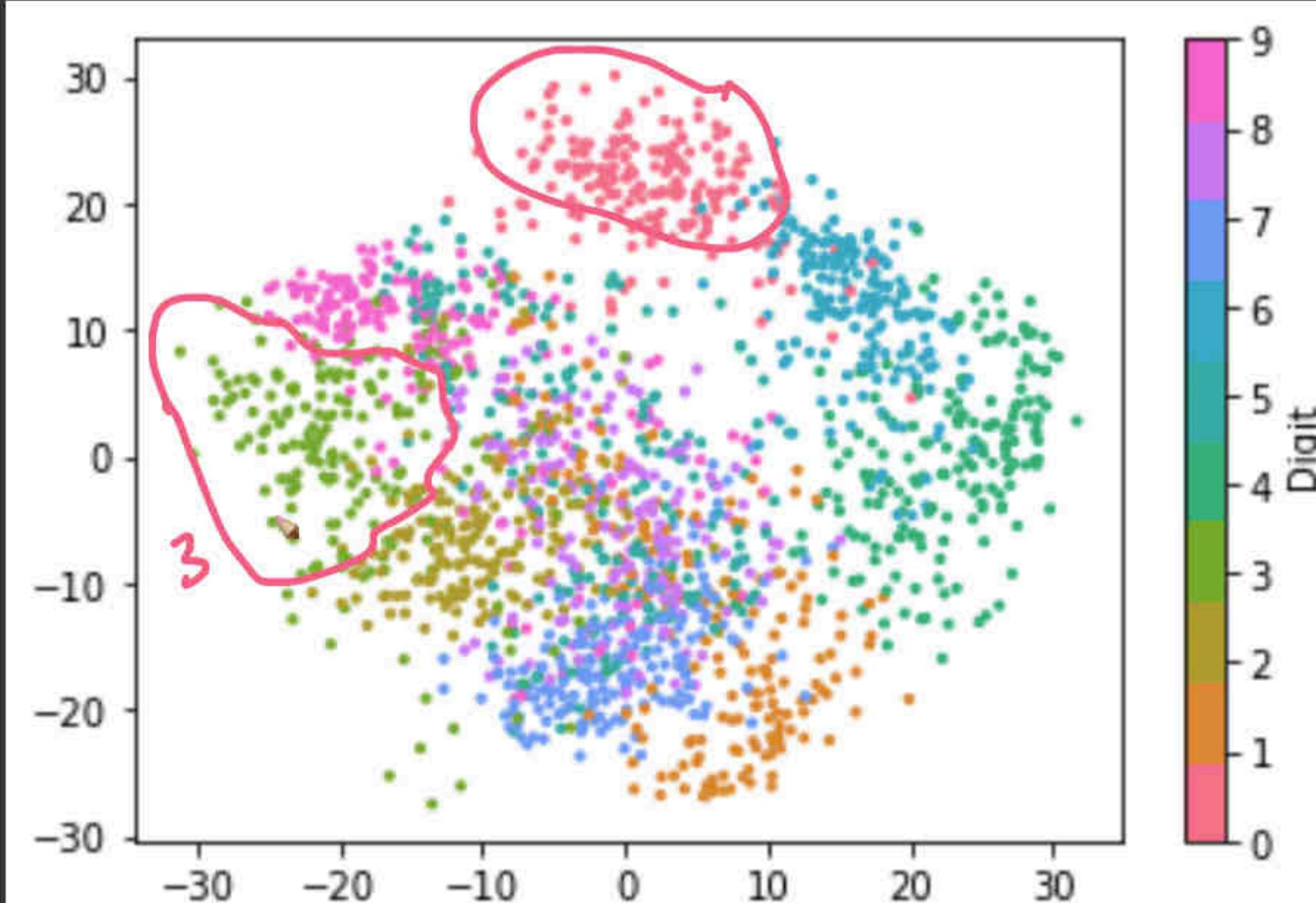
7 9

# tSNE

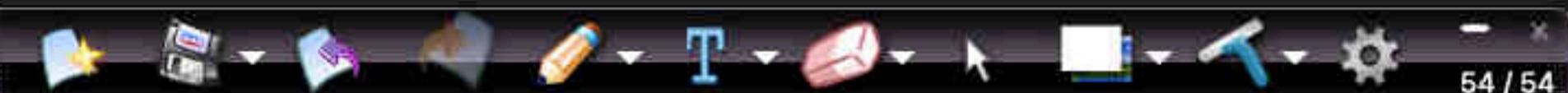




Reconnect

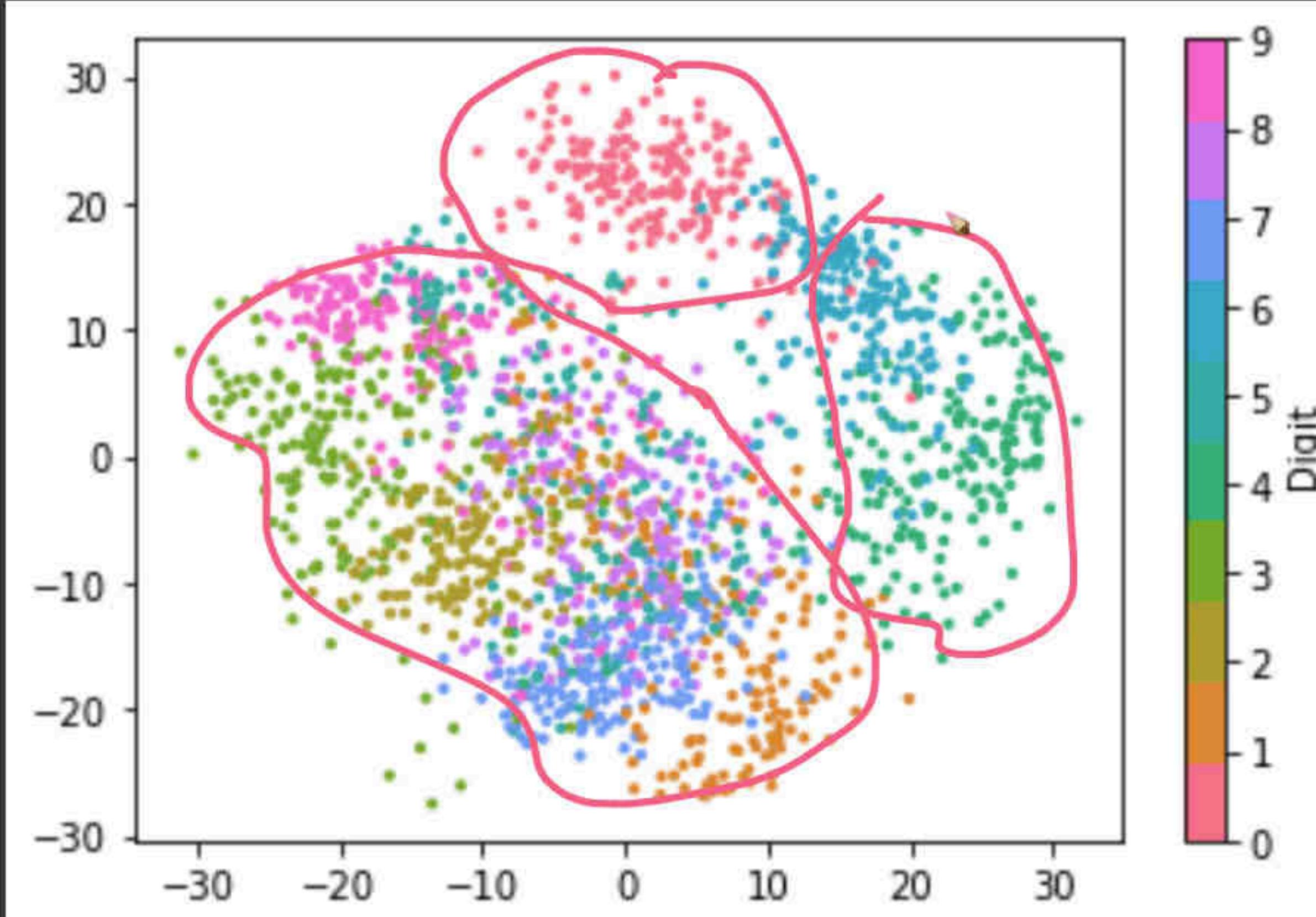


# tSNE



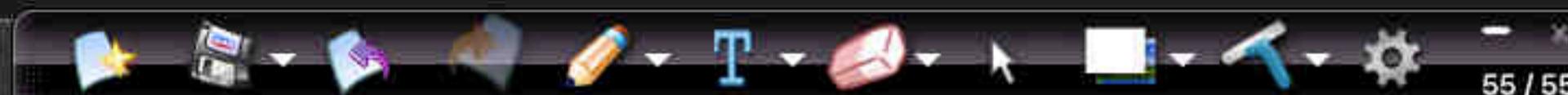


Reconnect



$x_i$ 's 64 dim  
↓  
2-dim

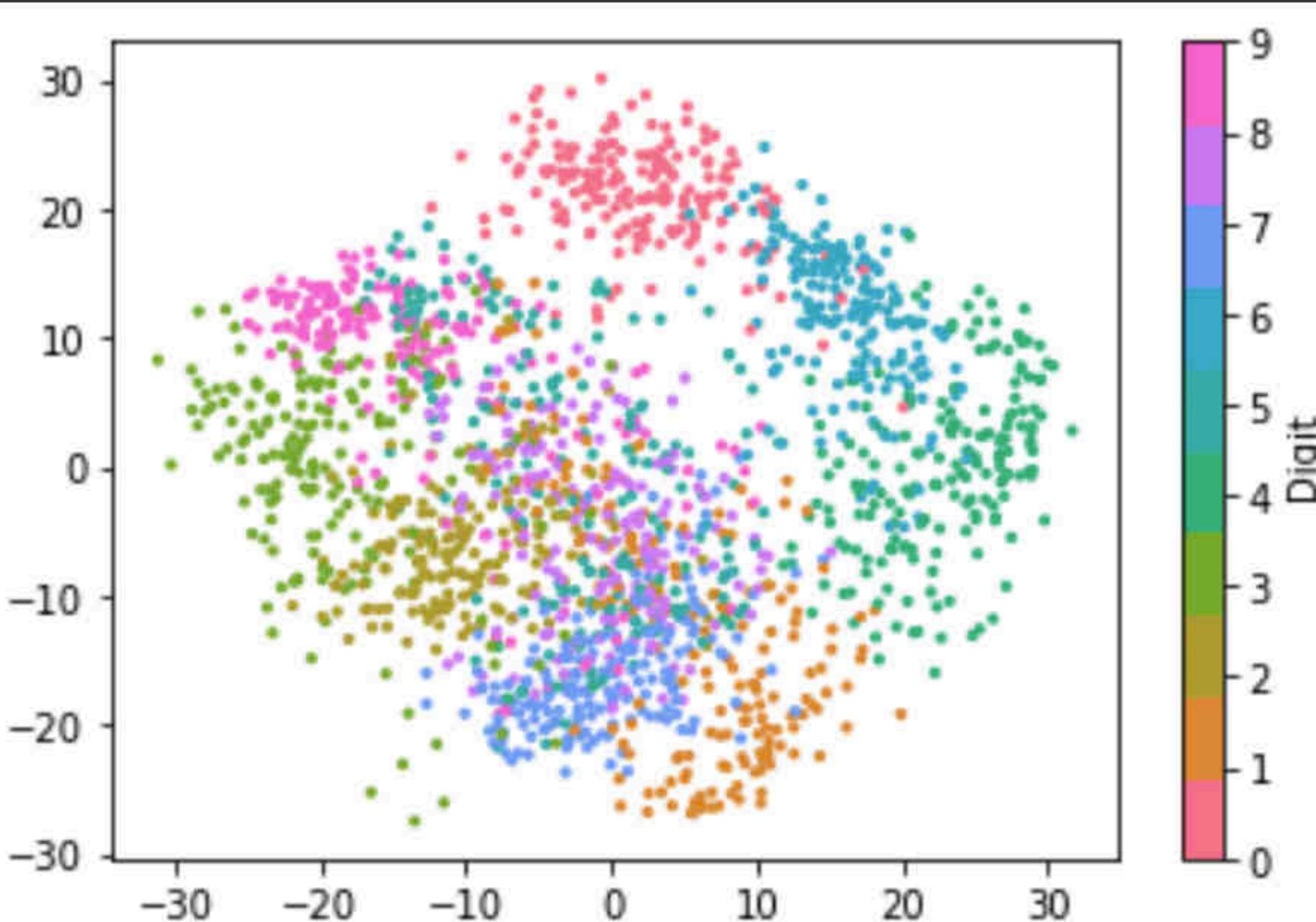
# tSNE



Reconnect



+ Code



PCA  
64 → 2dIM  
no Neigh Probs

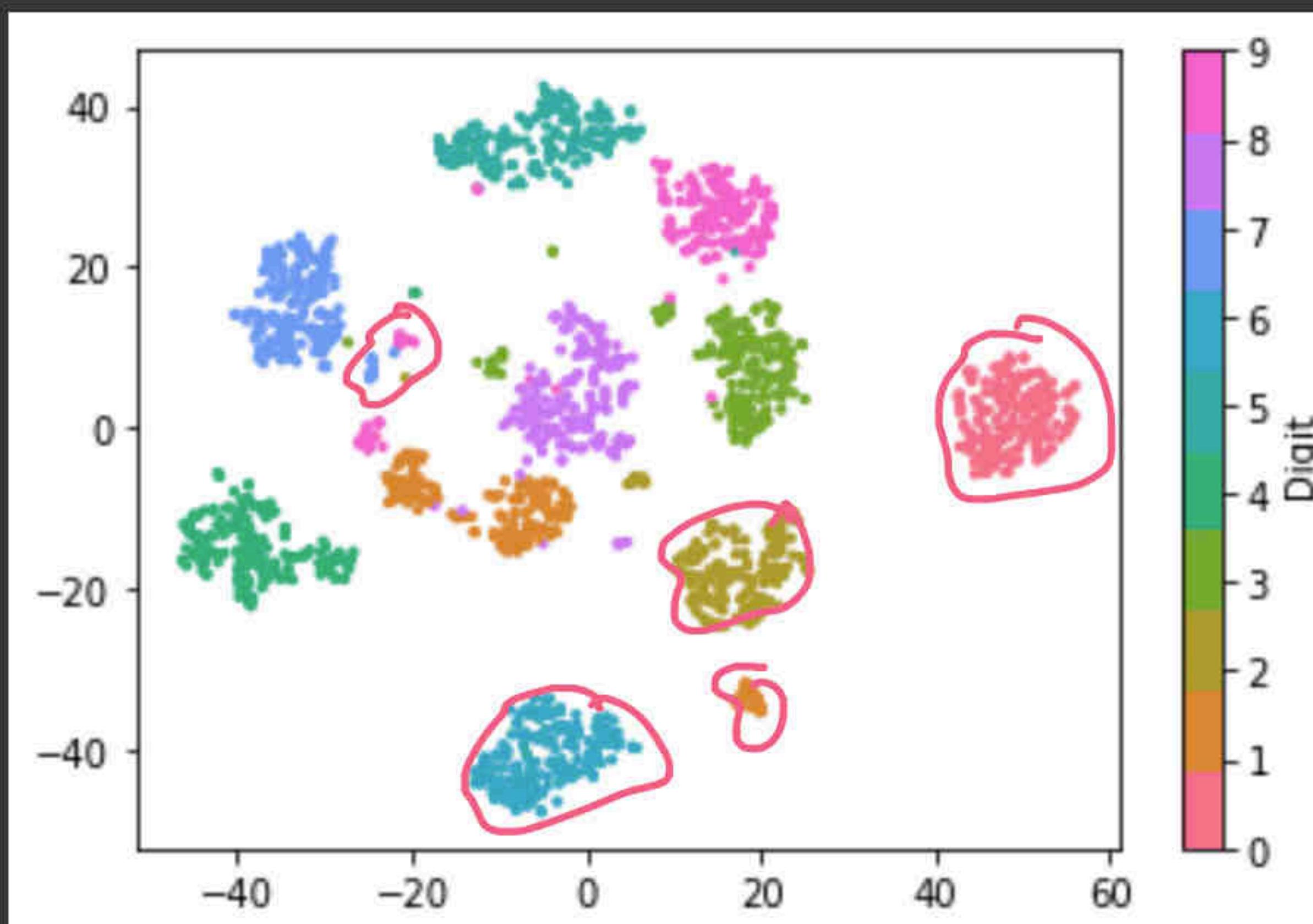
[ ] # tSNE



+ Code + Text

Reconnect

cmap = fig.colorbar(m, ax=ax, label='Digit')



tsNE  
tsNE

# Visualizing Data using t-SNE

**Laurens van der Maaten**

*TiCC*

*Tilburg University*

*P.O. Box 90153, 5000 LE Tilburg, The Netherlands*

LVDMAATEN@GMAIL.COM

**Geoffrey Hinton**

*Department of Computer Science*

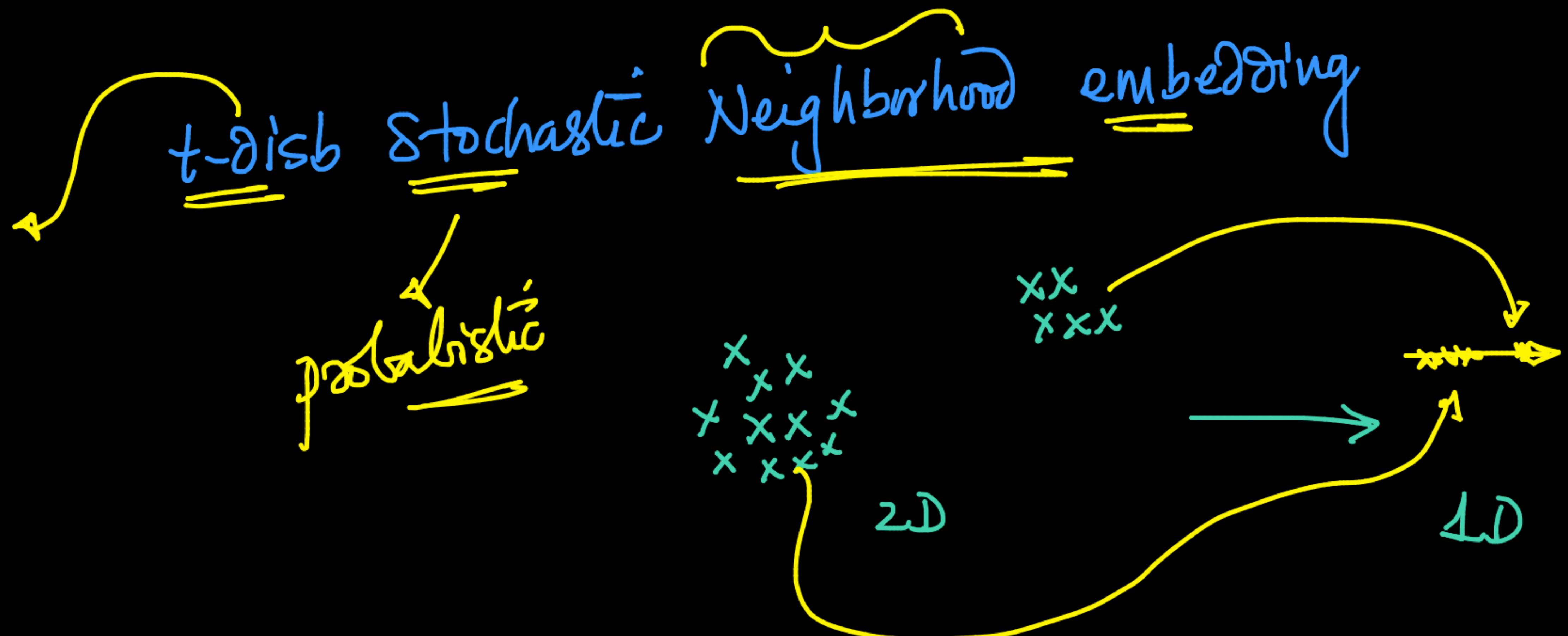
*University of Toronto*

*6 King's College Road, M5S 3G4 Toronto, ON, Canada*

HINTON@CS.TORONTO.EDU

**Editor:** Yoshua Bengio

t-SNE:



Cool-idea!

Preserve all pairwise distances in a Neighbourhood as best as possible

$$x_i \xrightarrow{x_j} \\ x_j \in N(x_i)$$

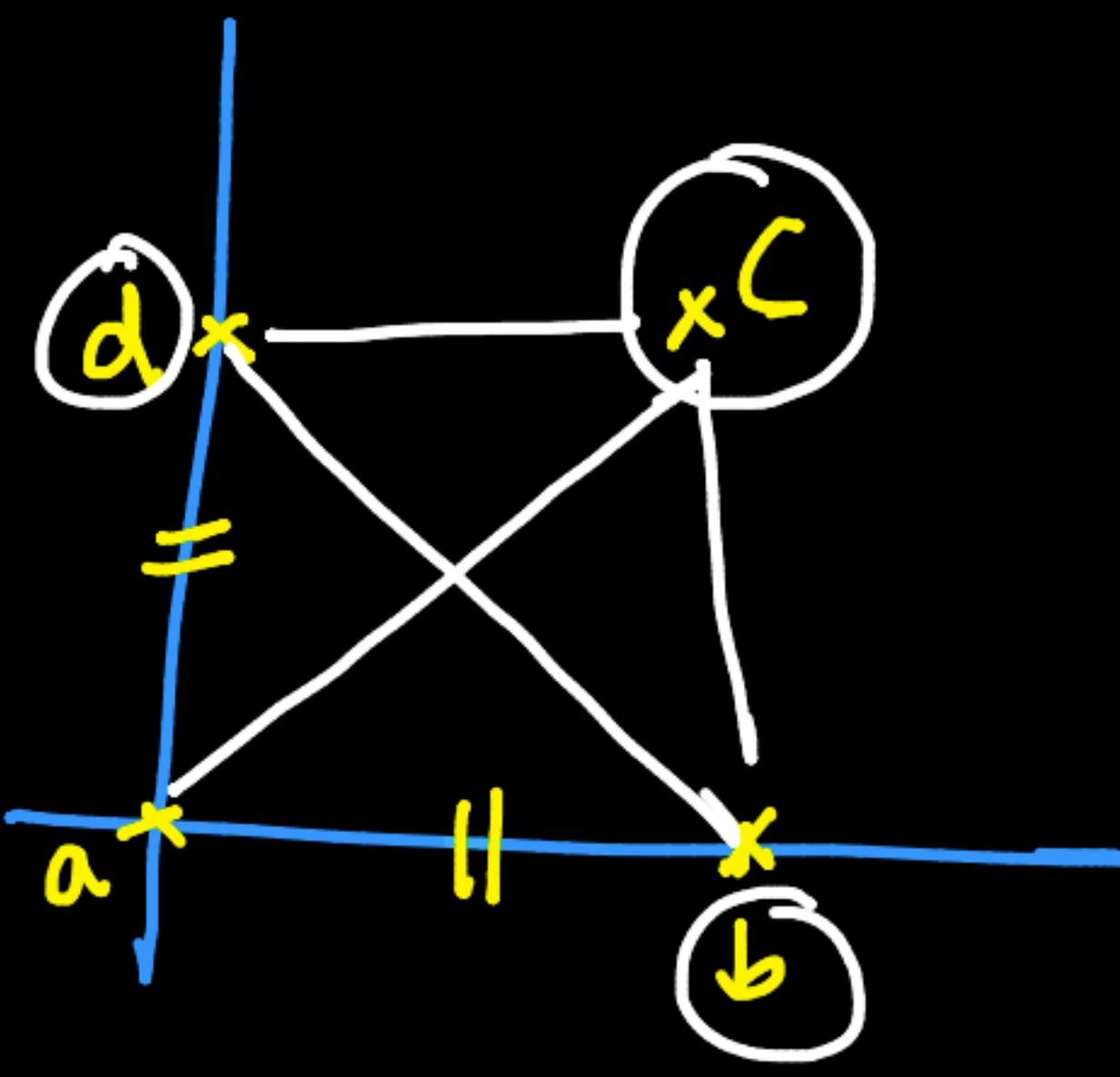
$$d \xrightarrow{x_i \in \mathbb{R}^d}$$

$$d' = 2 \text{ (typically)} \\ y_j \in \mathbb{R}^{d'}$$

$$\hat{x}_i \xrightarrow{\hat{x}_j}$$

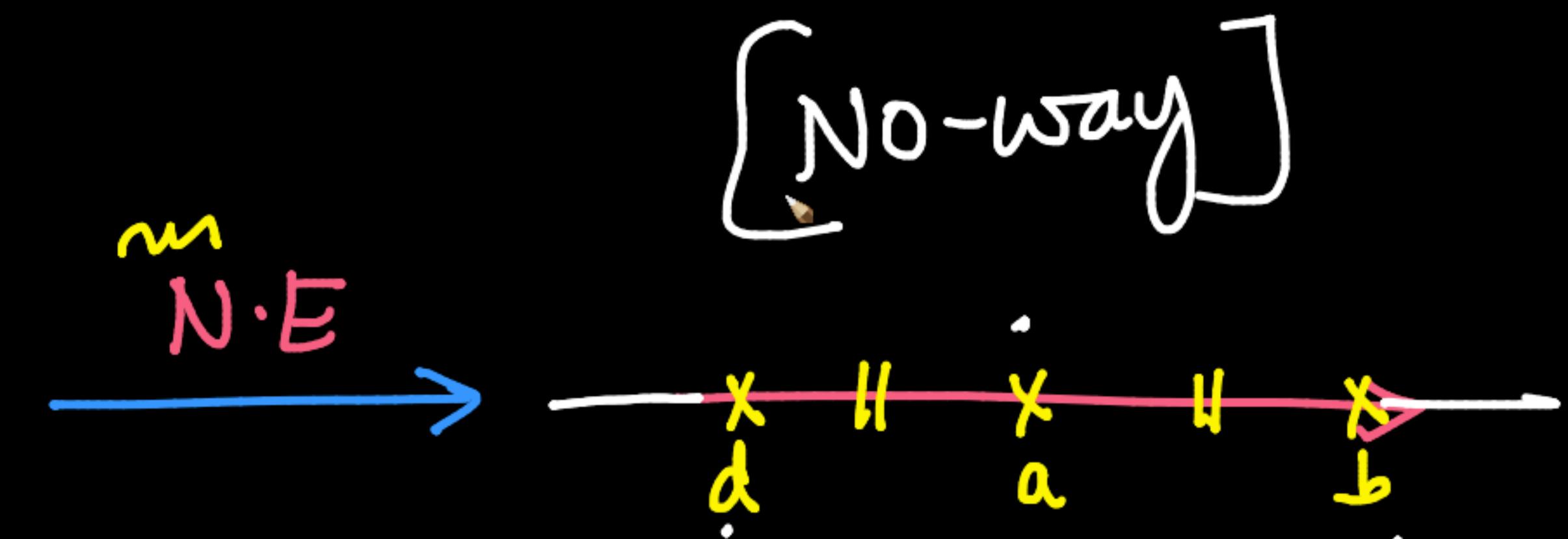
Counting

$$ac > ad = ab$$



2D

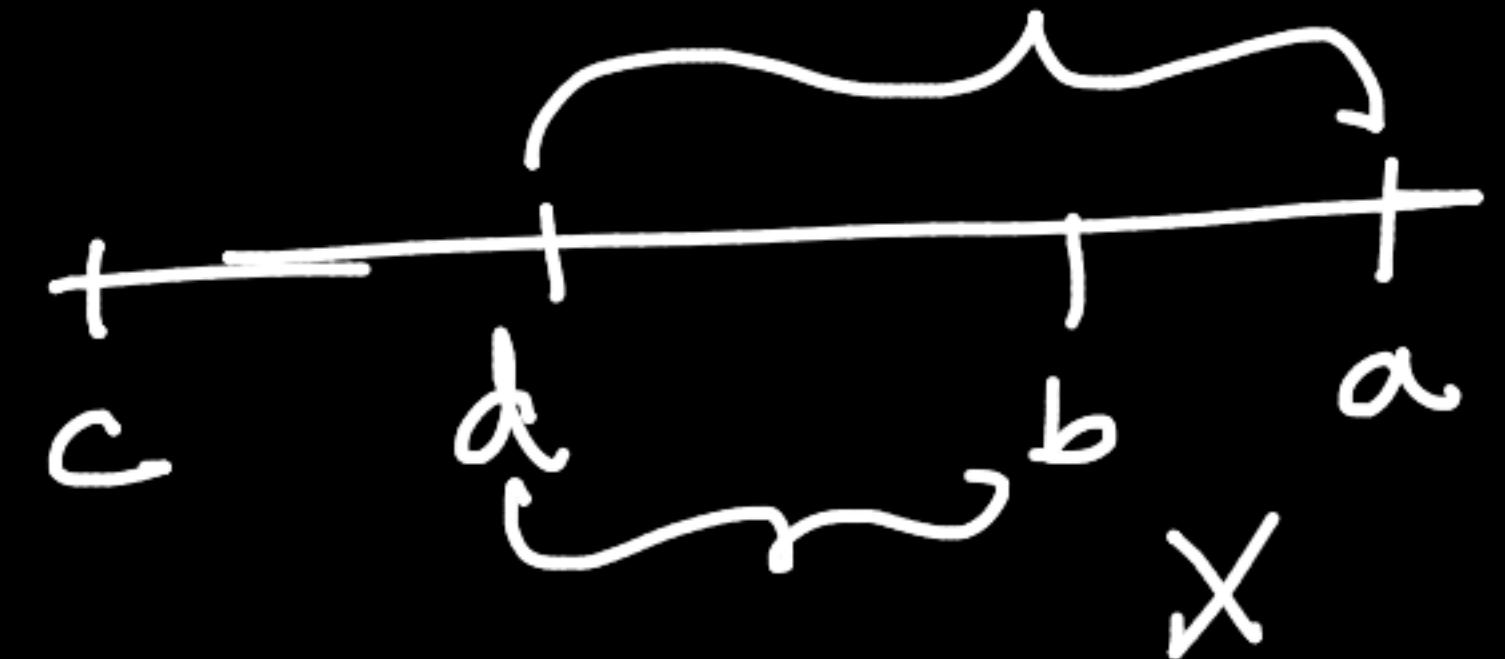
N·E



1D

Let

$$\mathcal{N}(A) = \{b, c, d\}$$





Math

$$x_i \in \mathbb{R}^d \xrightarrow{\text{tSNE}} y_i \in \mathbb{R}^{d'} \xrightarrow{\text{NOT-labels}}$$

i  
SNE

the low-dimensional map  $q_{ij}$  are given by

$d'$ -dim

$$\sum_i \sum_j p_{ij} = 1$$

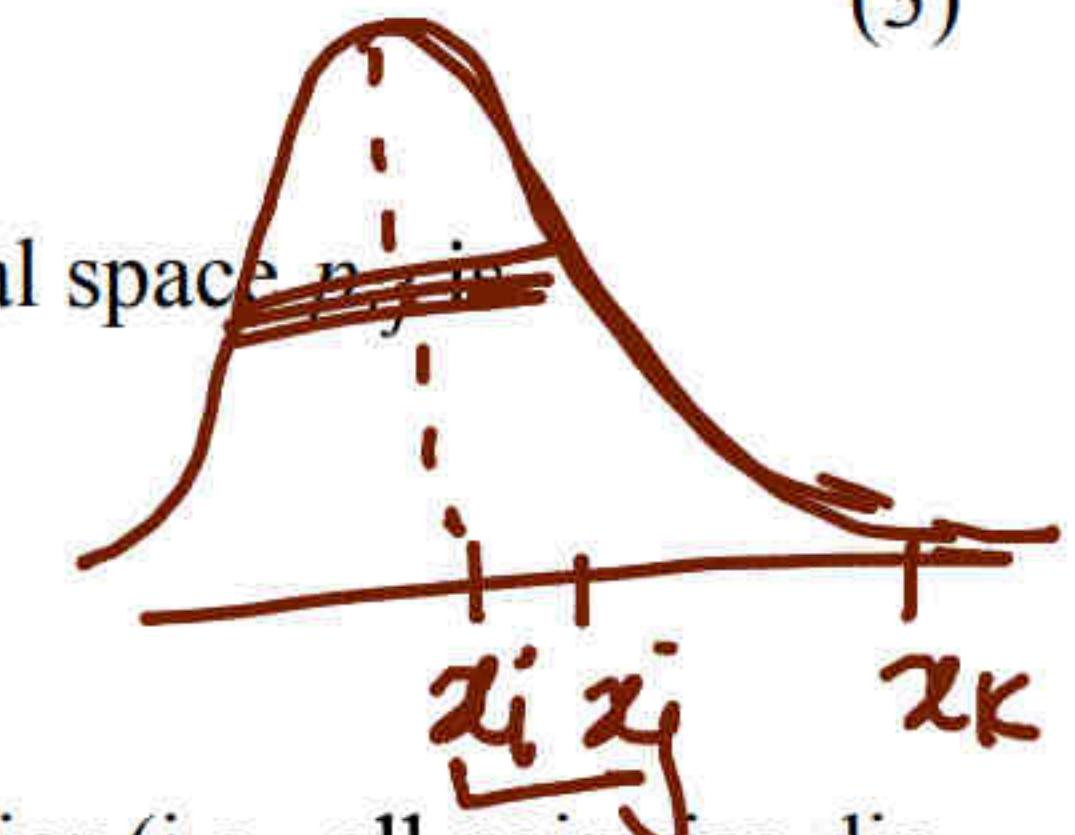
$$q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)}, \quad (3)$$

The obvious way to define the pairwise similarities in the high-dimensional space

$d$ -dim

[post that  $x_i$  &  $x_j$  are near]

$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)}$$



but this causes problems when a high-dimensional datapoint  $x_i$  is an outlier (i.e., all pairwise distances  $\|x_i - x_j\|^2$  are large for  $x_i$ ). For such an outlier, the values of  $p_{ij}$  are extremely small for all  $j$ , so the location of its low-dimensional map point  $y_i$  has very little effect on the cost function. As a result, the position of the map point is not well determined by the positions of the other map points. We circumvent this problem by defining the joint probabilities  $p_{ij}$  in the high-dimensional space to be the symmetrized conditional probabilities, that is, we set  $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$ . This ensures that  $\sum_j p_{ij} > \frac{1}{2n}$  for all  $x_i$  makes a significant contri-

MatrixCalculus x | PCA\_tSNE x | JMLR\_2008 x | Student's t-distribution x | Kullback-Leibler x | How to Use t-SNE x | sklearn.manifold x | t-SNE: The e x | 1802.03426 x | UMAP Dimensionality Reduction x | Basic UMAP x | New Tab x | +

lvdmaaten.github.io/publications/papers/JMLR\_2008.pdf

6 / 27 | - 250% + | ☰

the low-dimensional map  $q_{ij}$  are given by

$$q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)}, \quad (3)$$

The obvious way to define the pairwise similarities in the high-dimensional space  $p_{ij}$  is

$\checkmark$   $p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)}$

1 prob  $\sum_j p_{ij} = 1$   
2  $x_i$  &  $x_j$  belong to  $N$

but this causes problems when a high-dimensional datapoint  $x_i$  is an outlier (i.e., all pairwise distances  $\|x_i - x_j\|^2$  are large for  $x_i$ ). For such an outlier, the values of  $p_{ij}$  are extremely small for all  $j$ , so the location of its low-dimensional map point  $y_i$  has very little effect on the cost function. As a result, the position of the map point is not well determined by the positions of the other map points. We circumvent this problem by defining the joint probabilities  $p_{ij}$  in the high-dimensional space to be the symmetrized conditional probabilities, that is, we set  $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$ . This ensures that  $\sum_j p_{ij} > \frac{1}{2n}$  for all  $x_i$  makes a significant contribution.

$x_i$   $\xrightarrow{\text{tSNE}}$   $y_i$   $\cdot 2\text{dim}$

$d$

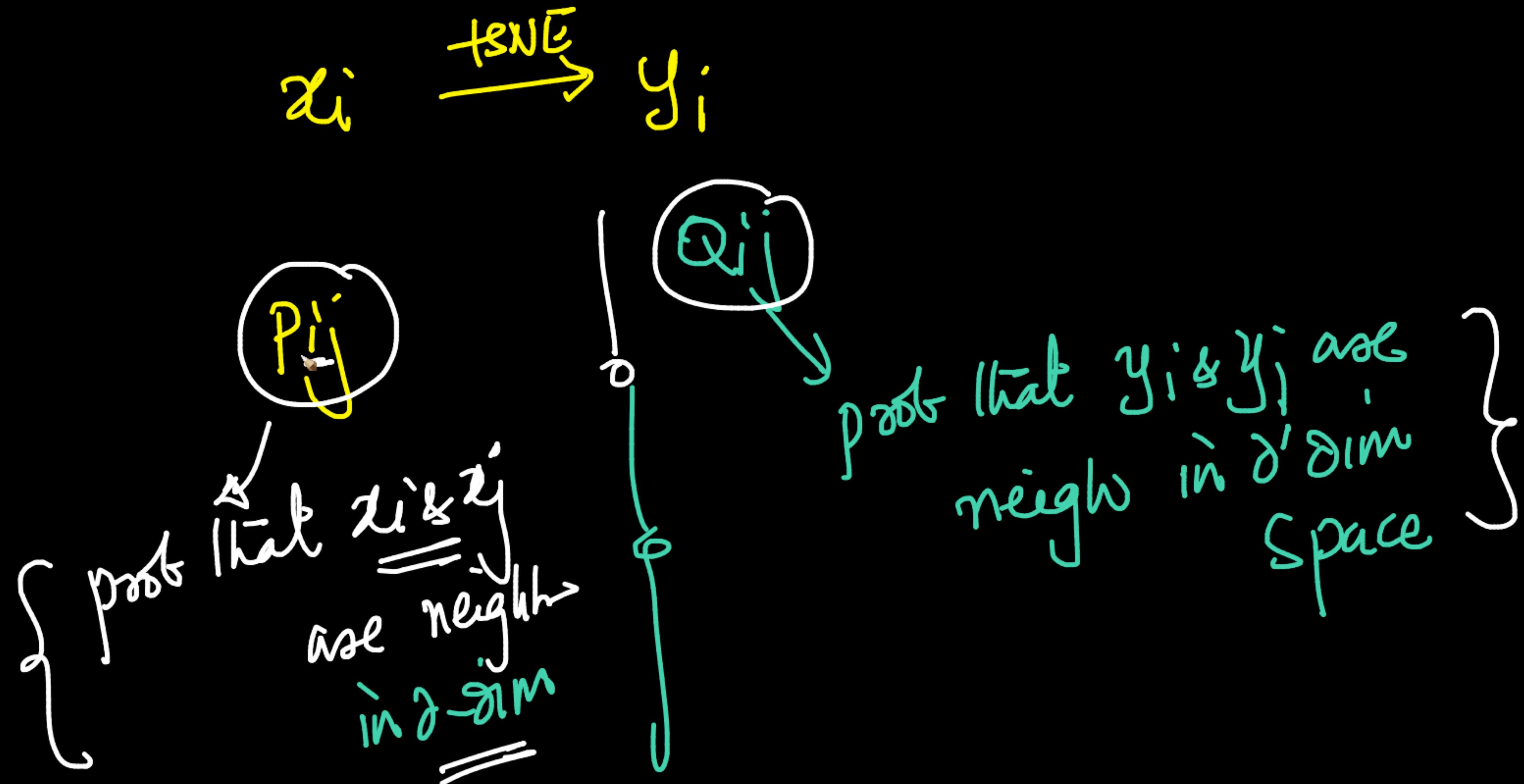
the low-dimensional map  $q_{ij}$  are given by

$$\uparrow q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)}, \quad (3)$$

The obvious way to define the pairwise similarities in the high-dimensional space  $p_{ij}$  is

$$\uparrow p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)},$$

but this causes problems when a high-dimensional datapoint  $x_i$  is an outlier (i.e., all pairwise distances  $\|x_i - x_j\|^2$  are large for  $x_i$ ). For such an outlier, the values of  $p_{ij}$  are extremely small for all  $j$ , so the location of its low-dimensional map point  $y_i$  has very little effect on the cost function. As a result, the position of the map point is not well determined by the positions of the other map points. We circumvent this problem by defining the joint probabilities  $p_{ij}$  in the high-dimensional space to be the symmetrized conditional probabilities, that is, we set  $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$ . This ensures that  $\sum_j p_{ij} > \frac{1}{2n}$  for all  $x_i$  makes a significant contri-



$H_{ij}$ 

$$P_{ij} \approx Q_{ij}$$

 $P, Q$ 

KL-div

→ differentiable

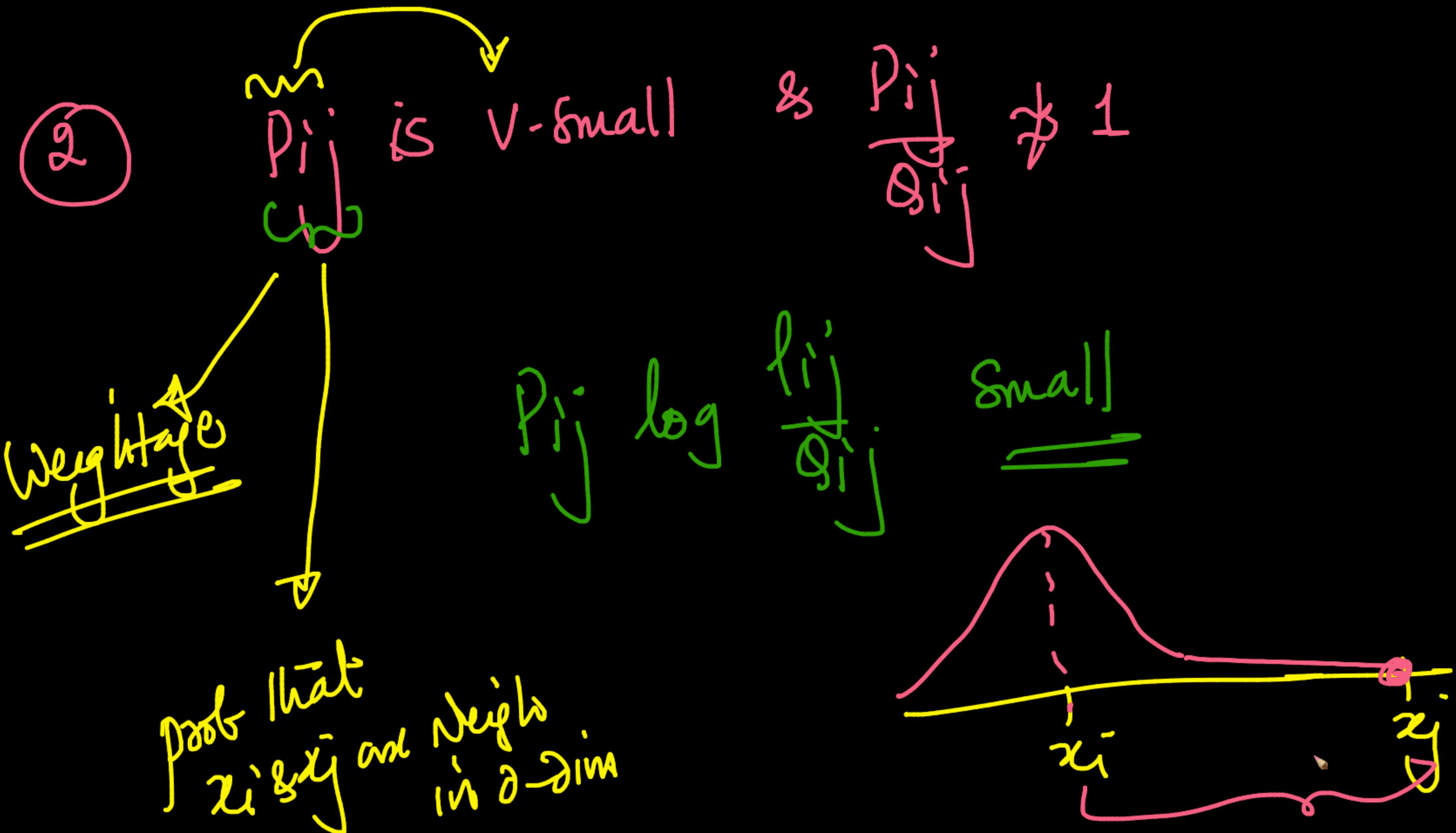
→ measure similarity  
b/w dist

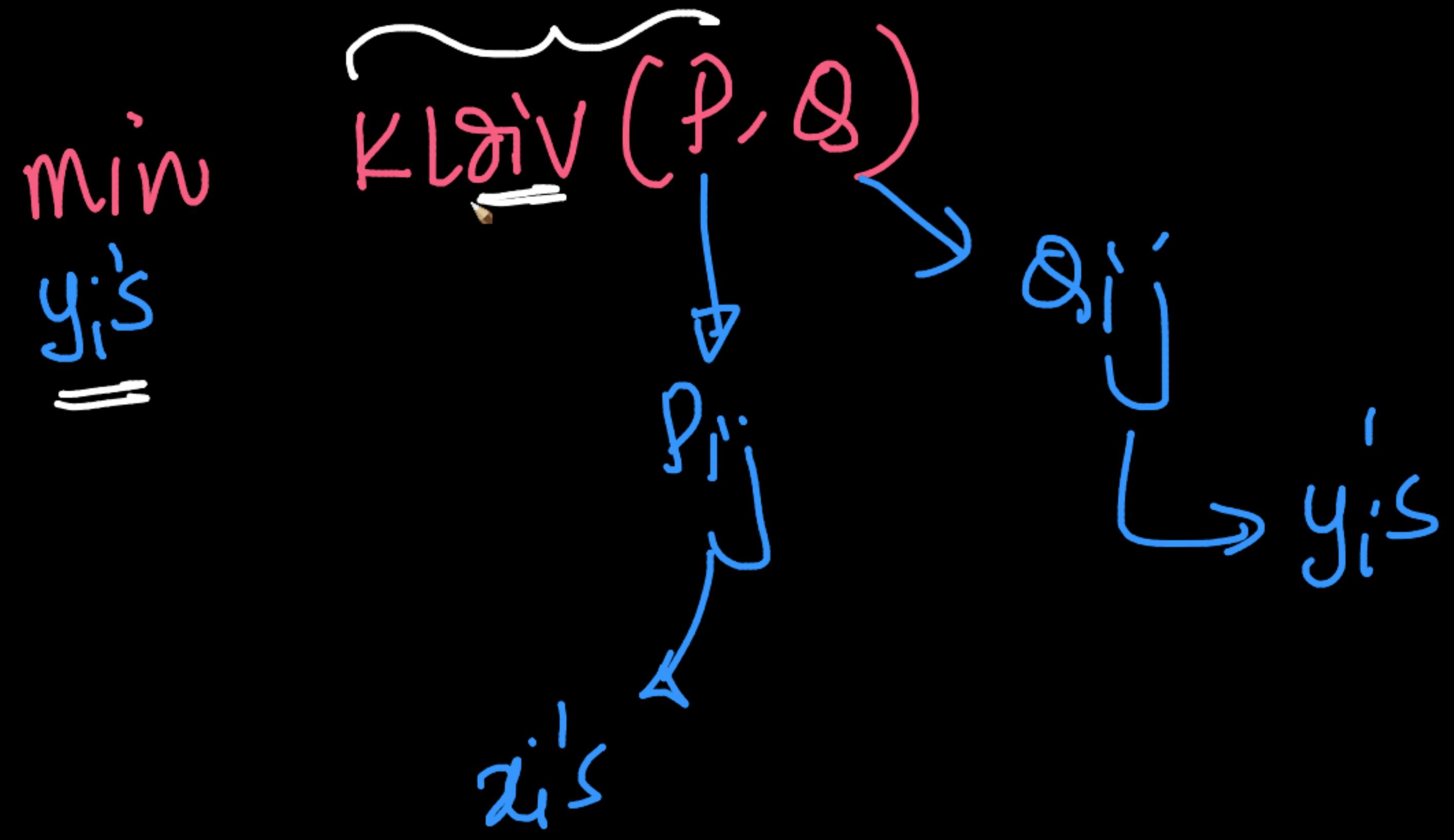
$$KL\text{-div}(P, Q) \xrightarrow{\substack{P_{ij} \\ Q_{ij}}} \text{relative Entropy} \xrightarrow{\cancel{\text{relative Entropy}}} \text{COMM. theory}$$

$$= \sum_i \sum_j P_{ij} \log \left( \frac{P_{ij}}{Q_{ij}} \right)$$

✓ if  $P_{ij} \approx Q_{ij}$   $\frac{P_{ij}}{Q_{ij}} \approx 1 \Rightarrow P_{ij} \log \frac{P_{ij}}{Q_{ij}} \rightarrow 0$   
 (i)

$P \approx Q$  then  $KL\text{-div}(P, Q) \rightarrow 0$

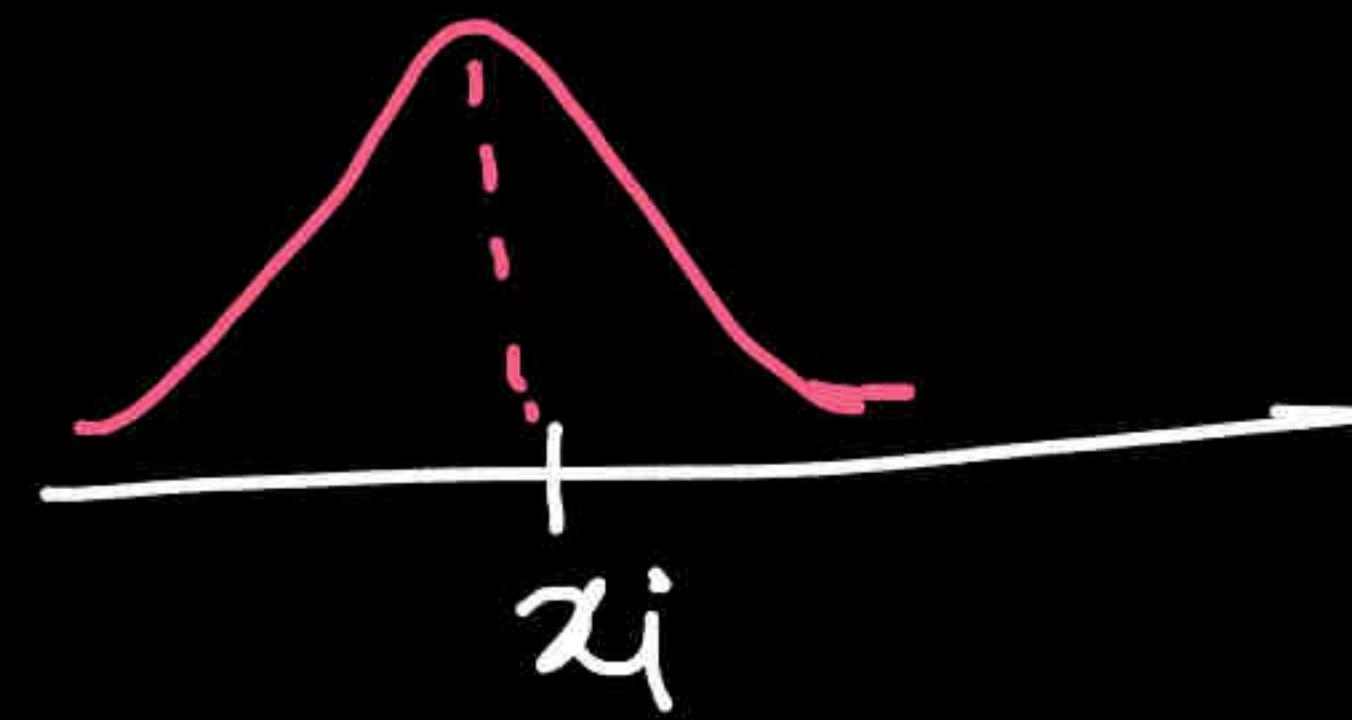




# t-SNE

$\hat{P}_{ij}$  : Gaussian - function

$\hat{P}_{ij} \neq$  t-SNE with  
 $\delta\alpha_j = 1$



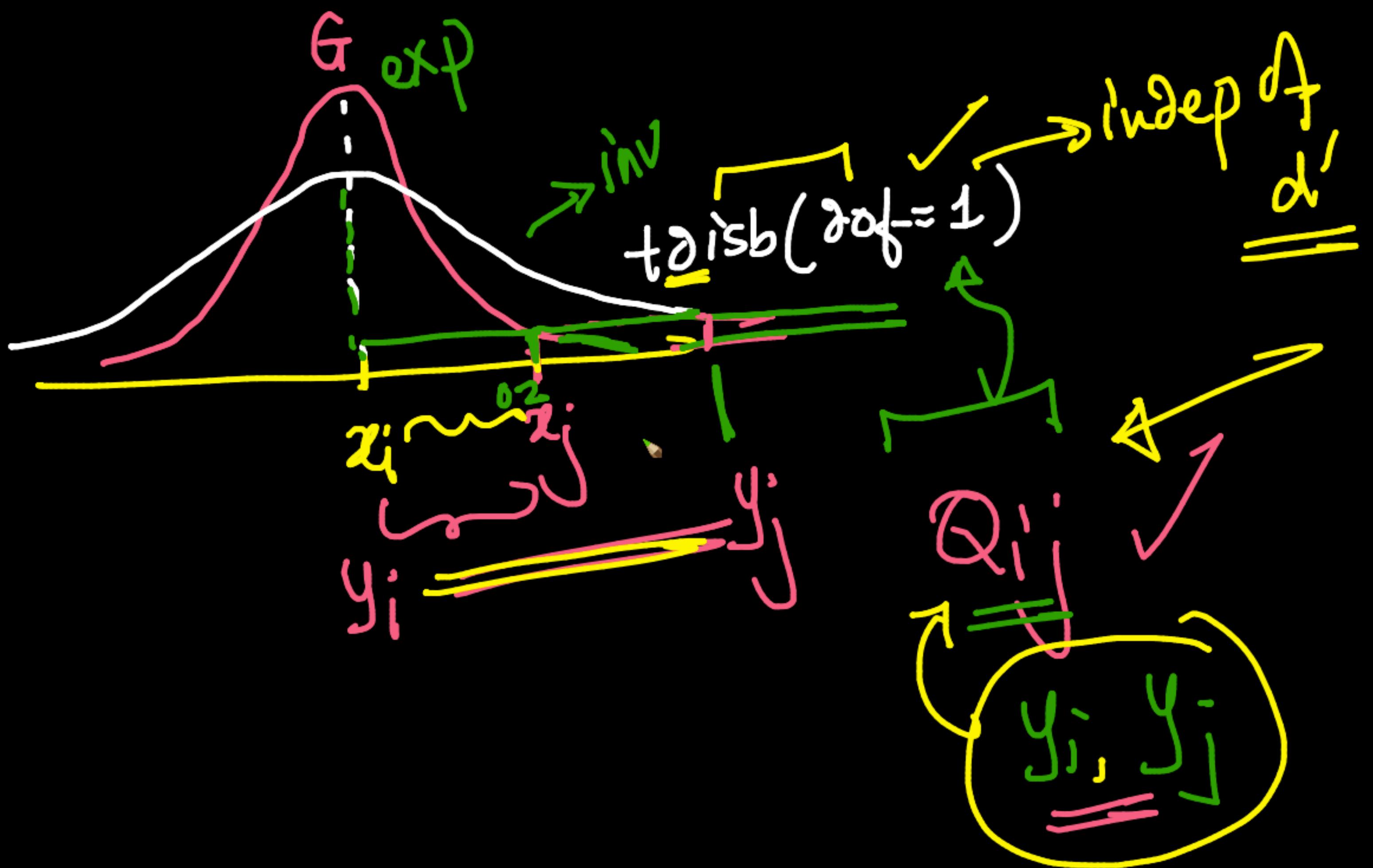
$x_i, x'_j$

Crowding

$\hat{P}_{ij}$

$x_i \quad x'_j$

Gaussian



convert distances into probabilities using a Gaussian distribution. In the low-dimensional map, we can use a probability distribution that has much heavier tails than a Gaussian to convert distances into probabilities. This allows a moderate distance in the high-dimensional space to be faithfully modeled by a much larger distance in the map and, as a result, it eliminates the unwanted attractive forces between map points that represent moderately dissimilar datapoints.

In t-SNE, we employ a Student t-distribution with one degree of freedom (which is the same as a Cauchy distribution) as the heavy-tailed distribution in the low-dimensional map. Using this distribution, the joint probabilities  $q_{ij}$  are defined as

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}. \quad (4)$$

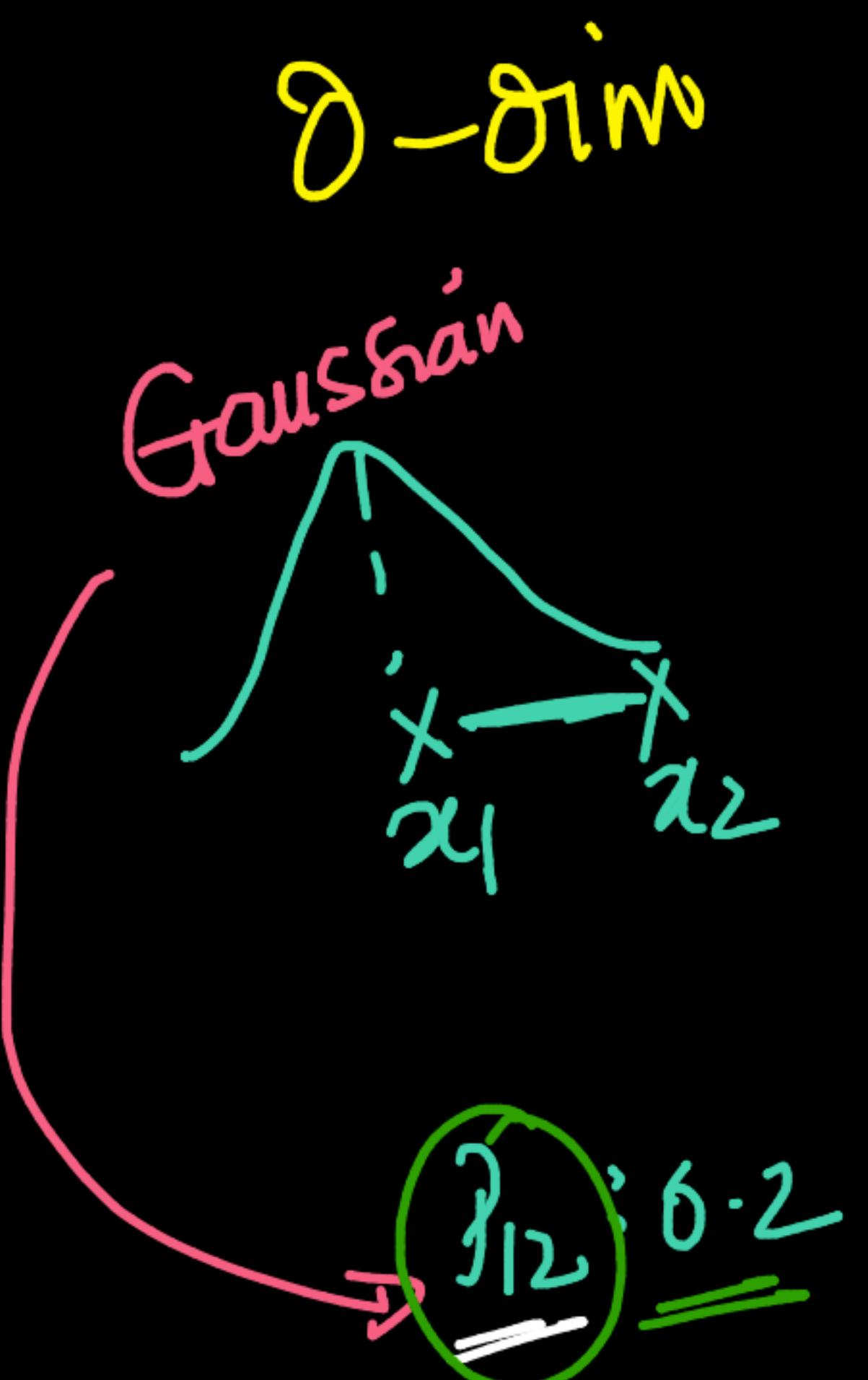
We use a Student t-distribution with a single degree of freedom, because it has the particularly nice property that  $(1 + \|y_i - y_j\|^2)^{-1}$  approaches an inverse square law for large pairwise distances  $\|y_i - y_j\|$  in the low-dimensional map. This makes the map's representation of joint probabilities (almost) invariant to changes in the scale of the map for map points that are far apart. It also means that large clusters of points that are far apart interact in just the same way as individual points, so the optimization operates in the same way at all but the finest scales. A theoretical justification for our

convert distances into probabilities using a Gaussian distribution. In the low-dimensional map, we can use a probability distribution that has much heavier tails than a Gaussian to convert distances into probabilities. This allows a moderate distance in the high-dimensional space to be faithfully modeled by a much larger distance in the map and, as a result, it eliminates the unwanted attractive forces between map points that represent moderately dissimilar datapoints.

In t-SNE, we employ a Student t-distribution with one degree of freedom (which is the same as a Cauchy distribution) as the heavy-tailed distribution in the low-dimensional map. Using this distribution, the joint probabilities  $q_{ij}$  are defined as

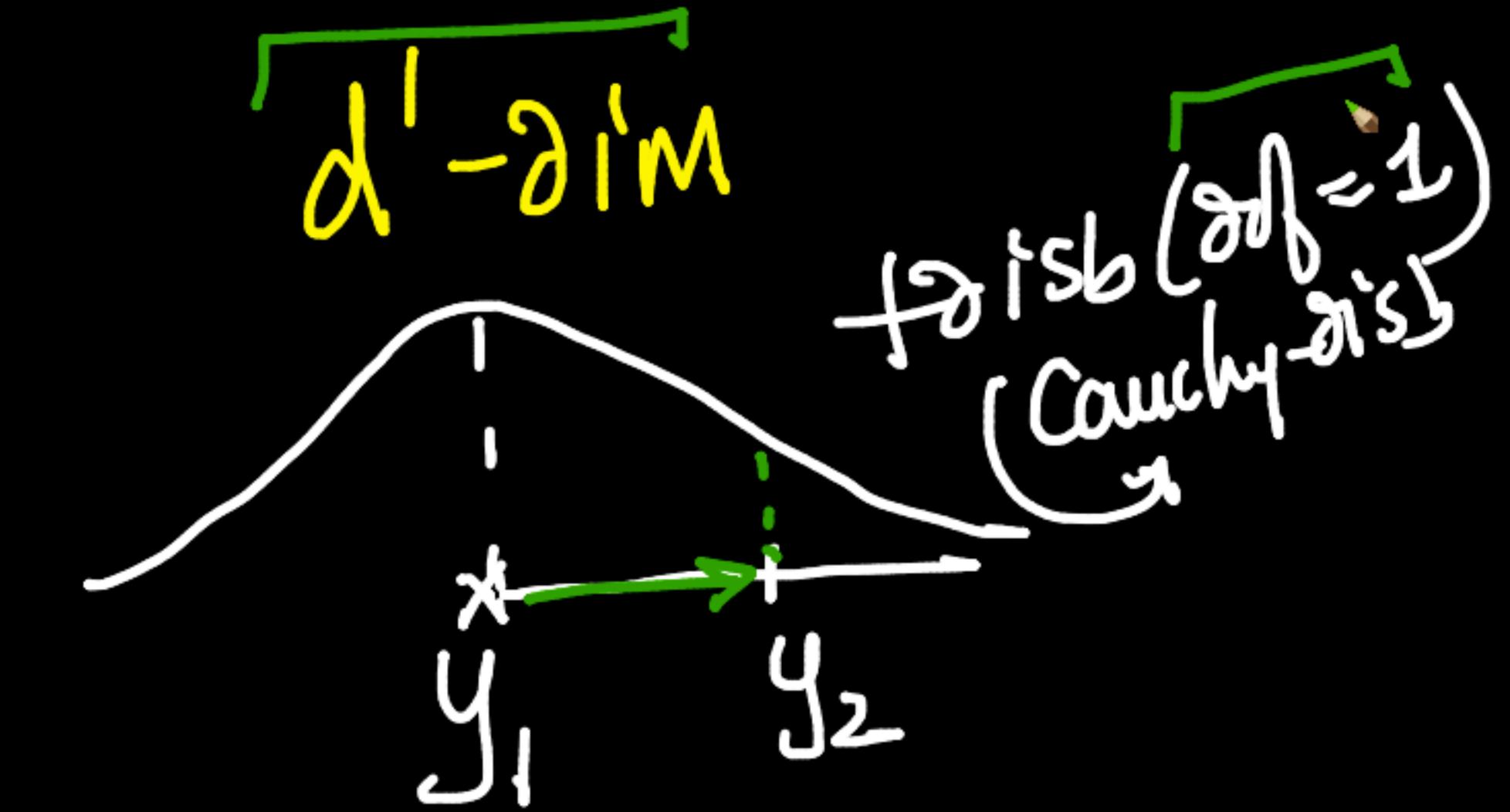
$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}. \quad (4)$$

We use a Student t-distribution with a single degree of freedom, because it has the particularly nice property that  $(1 + \|y_i - y_j\|^2)^{-1}$  approaches an inverse square law for large pairwise distances  $\|y_i - y_j\|$  in the low-dimensional map. This makes the map's representation of joint probabilities (almost) invariant to changes in the scale of the map for map points that are far apart. It also means that large clusters of points that are far apart interact in just the same way as individual points, so the optimization operates in the same way at all but the finest scales. A theoretical justification for our



mitigates  
Crossing

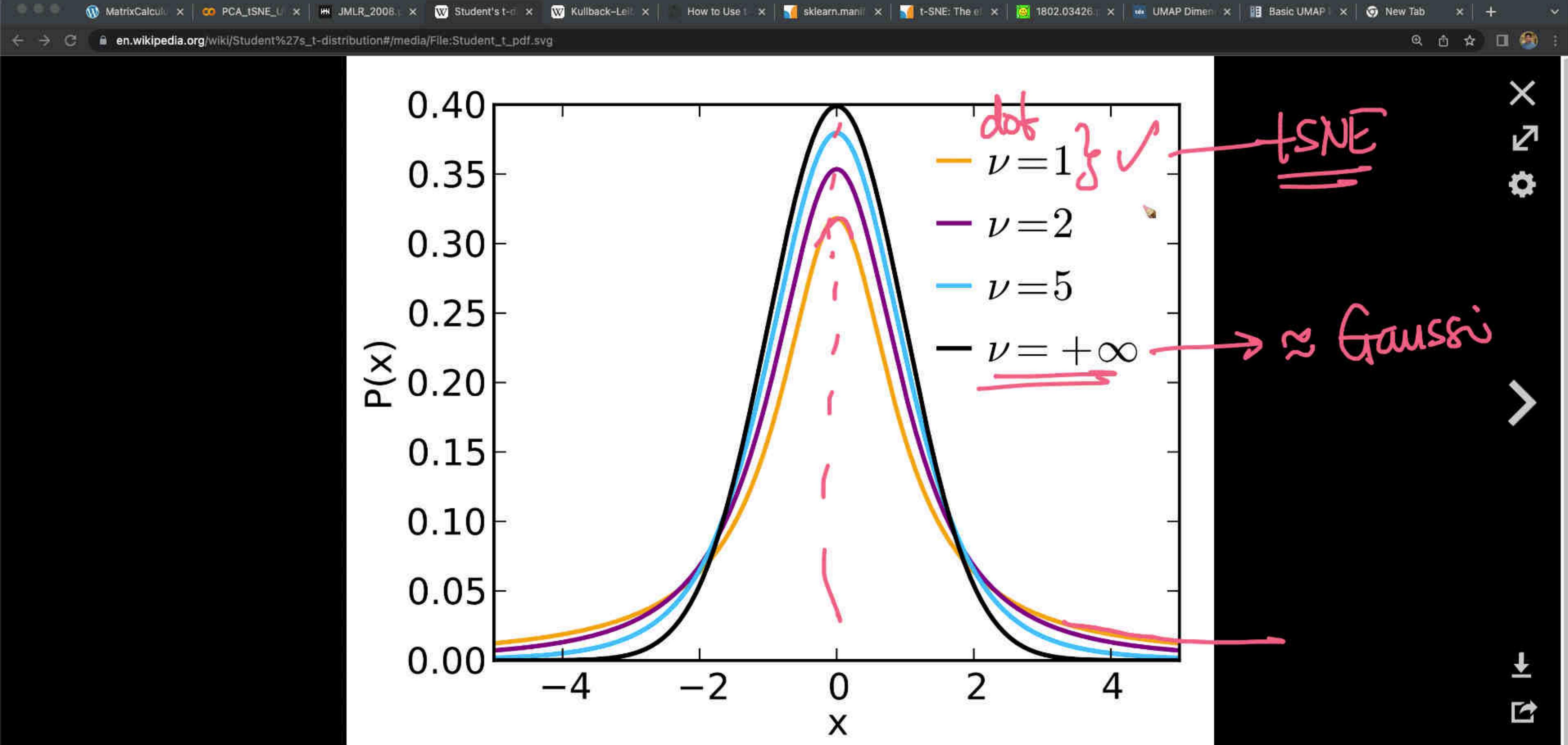
A small diagram consisting of a horizontal line with two points labeled  $y_1$  and  $y_2$ . A green arrow points from  $y_1$  towards  $y_2$ .



$Q_{12} \approx 0.2$

help me somewhat,  
overcome Crossing  
Some "wiggle" ~~zoom~~

A large circular callout contains handwritten text. It starts with  $Q_{12} \approx 0.2$ , followed by the sentence "help me somewhat, overcome Crossing", and ends with "Some 'wiggle' ~~zoom~~".



convert distances into probabilities using a Gaussian distribution. In the low-dimensional map, we can use a probability distribution that has much heavier tails than a Gaussian to convert distances into probabilities. This allows a moderate distance in the high-dimensional space to be faithfully modeled by a much larger distance in the map and, as a result, it eliminates the unwanted attractive forces between map points that represent moderately dissimilar datapoints.

In t-SNE, we employ a Student t-distribution with one degree of freedom (which is the same as a Cauchy distribution) as the heavy-tailed distribution in the low-dimensional map. Using this distribution, the joint probabilities  $q_{ij}$  are defined as

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}. \quad (4)$$

We use a Student t-distribution with a single degree of freedom, because it has the particularly nice property that  $(1 + \|y_i - y_j\|^2)^{-1}$  approaches an inverse square law for large pairwise distances  $\|y_i - y_j\|$  in the low-dimensional map. This makes the map's representation of joint probabilities (almost) invariant to changes in the scale of the map for map points that are far apart. It also means that large clusters of points that are far apart interact in just the same way as individual points, so the optimization operates in the same way at all but the finest scales. A theoretical justification for our

✓ Perplexity  $\approx$  (intuitively) effective # neighbors whose distances you are preserving

typically:  $\underline{\underline{5-50}}$

Optimzn- problem  $\rightarrow$  single optima  
↳ very time-taking  
(burnes - but approx)

slow

$$\mathcal{D} = \{x_i\}_{i=1}^n$$

t-SNE

$$\mathcal{D}' = \{y_i\}_{i=1}^n$$





Prev Up Next

## scikit-learn 1.1.1

[Other versions](#)

Please [cite us](#) if you use the software.

[sklearn.manifold.TSNE](#)

Examples using  
[sklearn.manifold.TSNE](#)

## sklearn.manifold.TSNE

```
class sklearn.manifold.TSNE(n_components=2, *, perplexity=30.0, early_exaggeration=12.0,  
learning_rate='warn', n_iter=1000, n_iter_without_progress=300, min_grad_norm=1e-07,  
metric='euclidean', metric_params=None, init='warn', verbose=0, random_state=None,  
method='barnes_hut', angle=0.5, n_jobs=None, square_distances='deprecated')
```

[\[source\]](#)

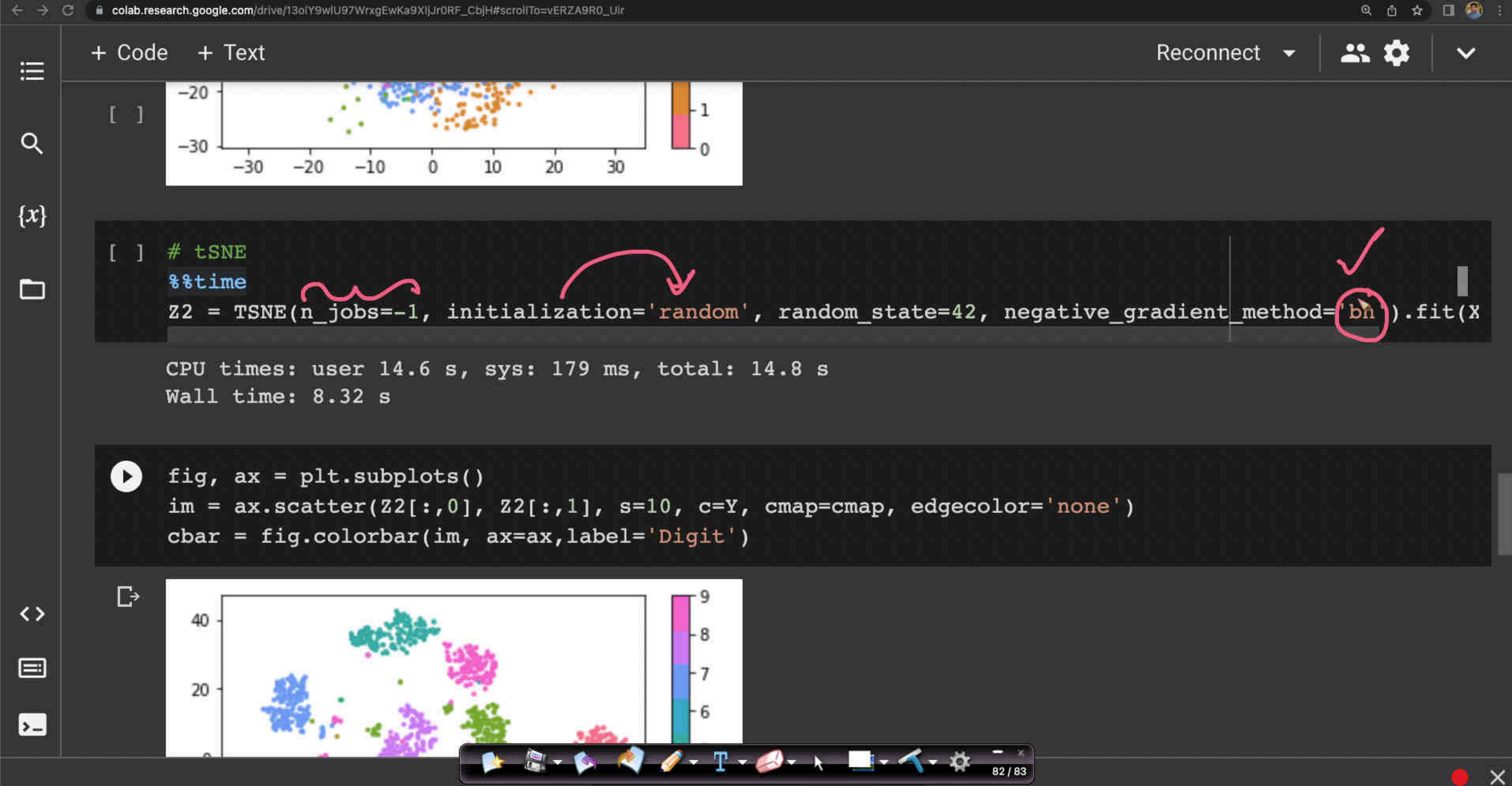
T-distributed Stochastic Neighbor Embedding.

t-SNE [1] is a tool to visualize high-dimensional data. It converts similarities between data points to joint probabilities and tries to minimize the Kullback-Leibler divergence between the joint probabilities of the low-dimensional embedding and the high-dimensional data. t-SNE has a cost function that is not convex, i.e. with different initializations we can get different results.

It is highly recommended to use another dimensionality reduction method (e.g. PCA for dense data or TruncatedSVD for sparse data) to reduce the number of dimensions to a reasonable amount (e.g. 50) if the number of features is very high. This will suppress some noise and speed up the computation of pairwise distances between samples. For more tips see Laurens van der Maaten's FAQ [2].

Read more in the [User Guide](#)

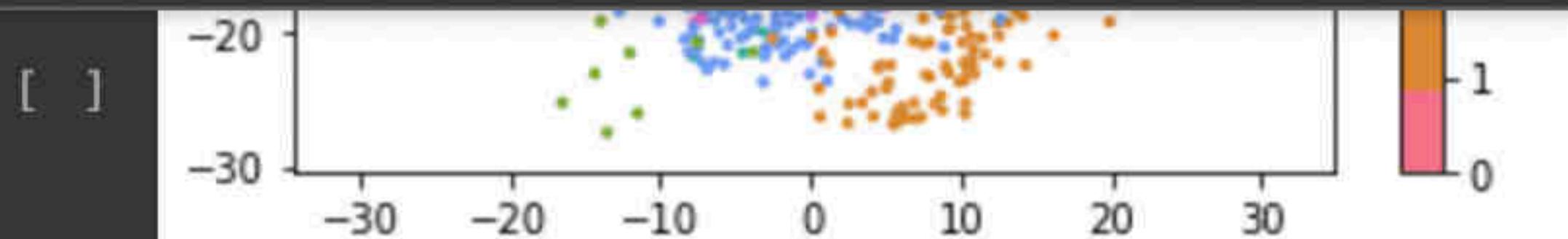




+ Code

+ Text

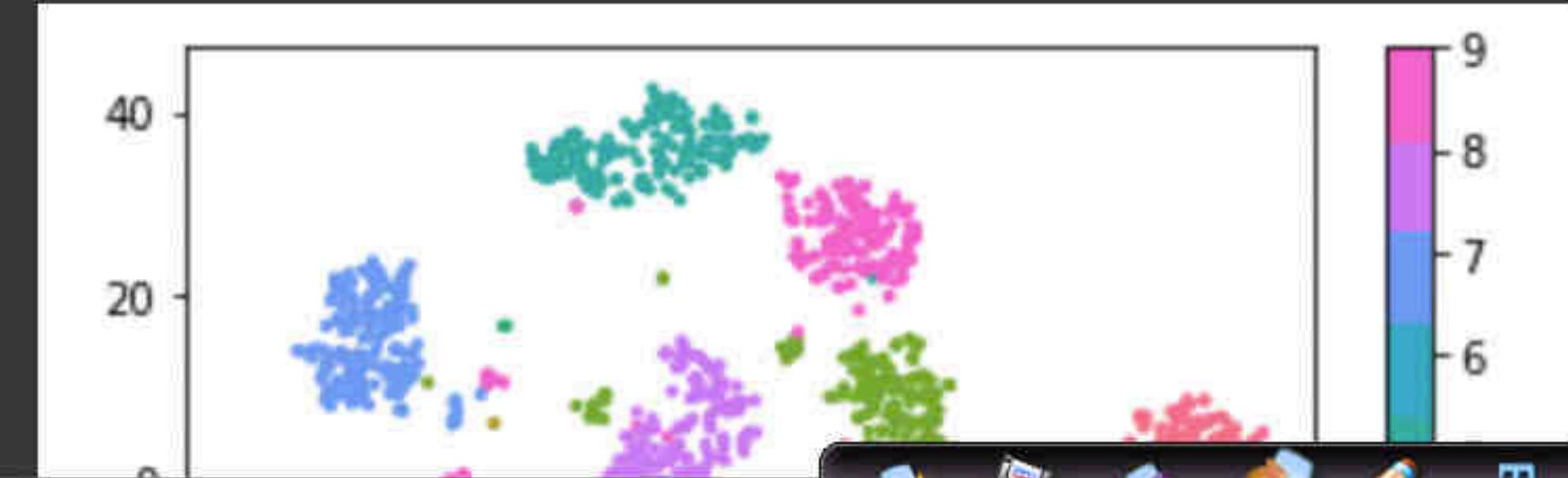
Reconnect



```
# tSNE  
%%time  
Z2 = TSNE(n_jobs=-1, initialization='random', random_state=42, negative_gradient_method='bh').fit(X)
```

CPU times: user 14.6 s, sys: 179 ms, total: 14.8 s  
Wall time: 8.32 s

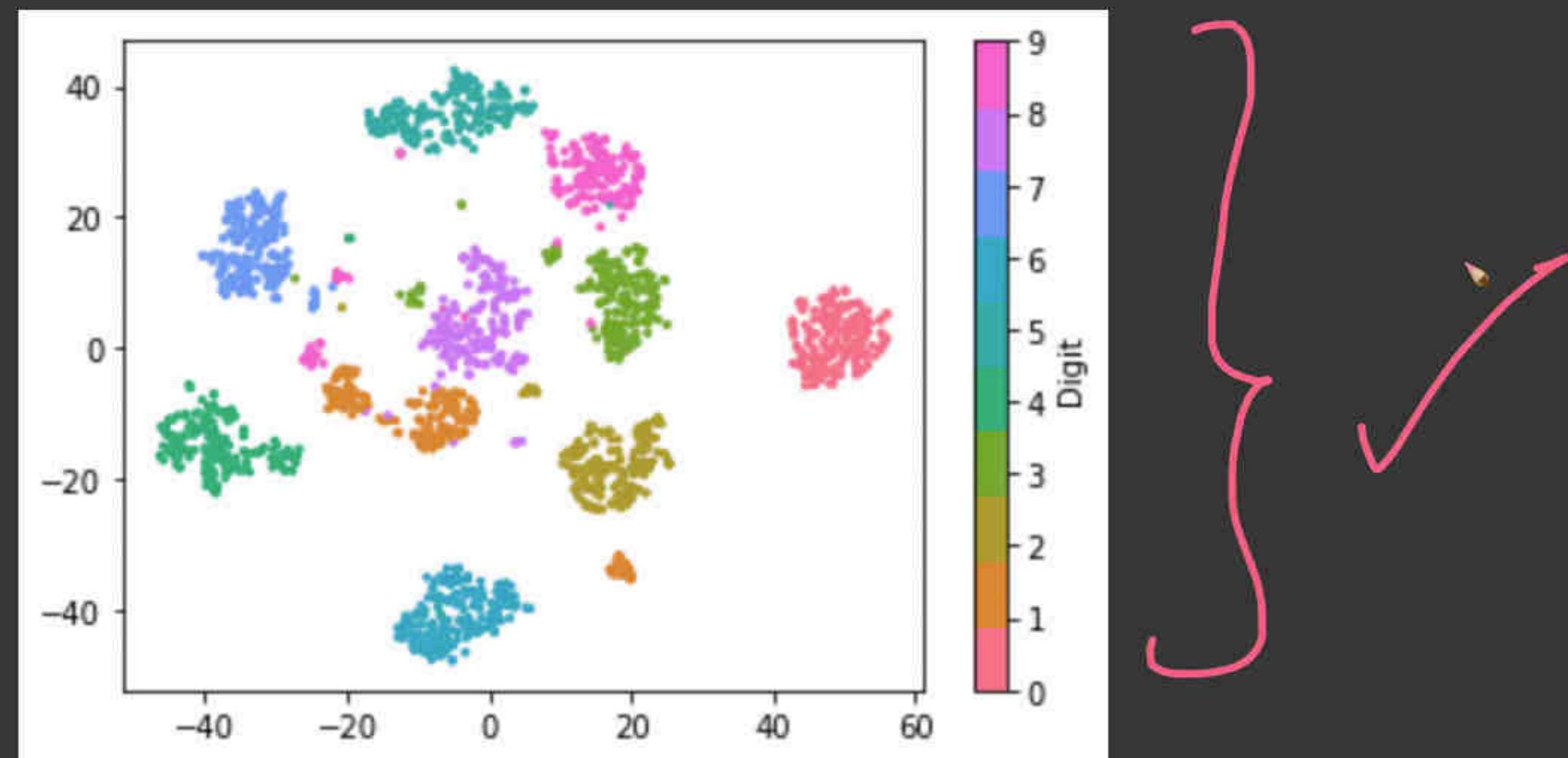
```
fig, ax = plt.subplots()  
im = ax.scatter(Z2[:,0], Z2[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')  
cbar = fig.colorbar(im, ax=ax, label='Digit')
```





+ Code + Text

```
[ ] im = ax.scatter(Z2[:,0], Z2[:,1], s=10, c=Y, cmap=cmap, edgecolor='none')
cbar = fig.colorbar(im, ax=ax, label='Digit')
```

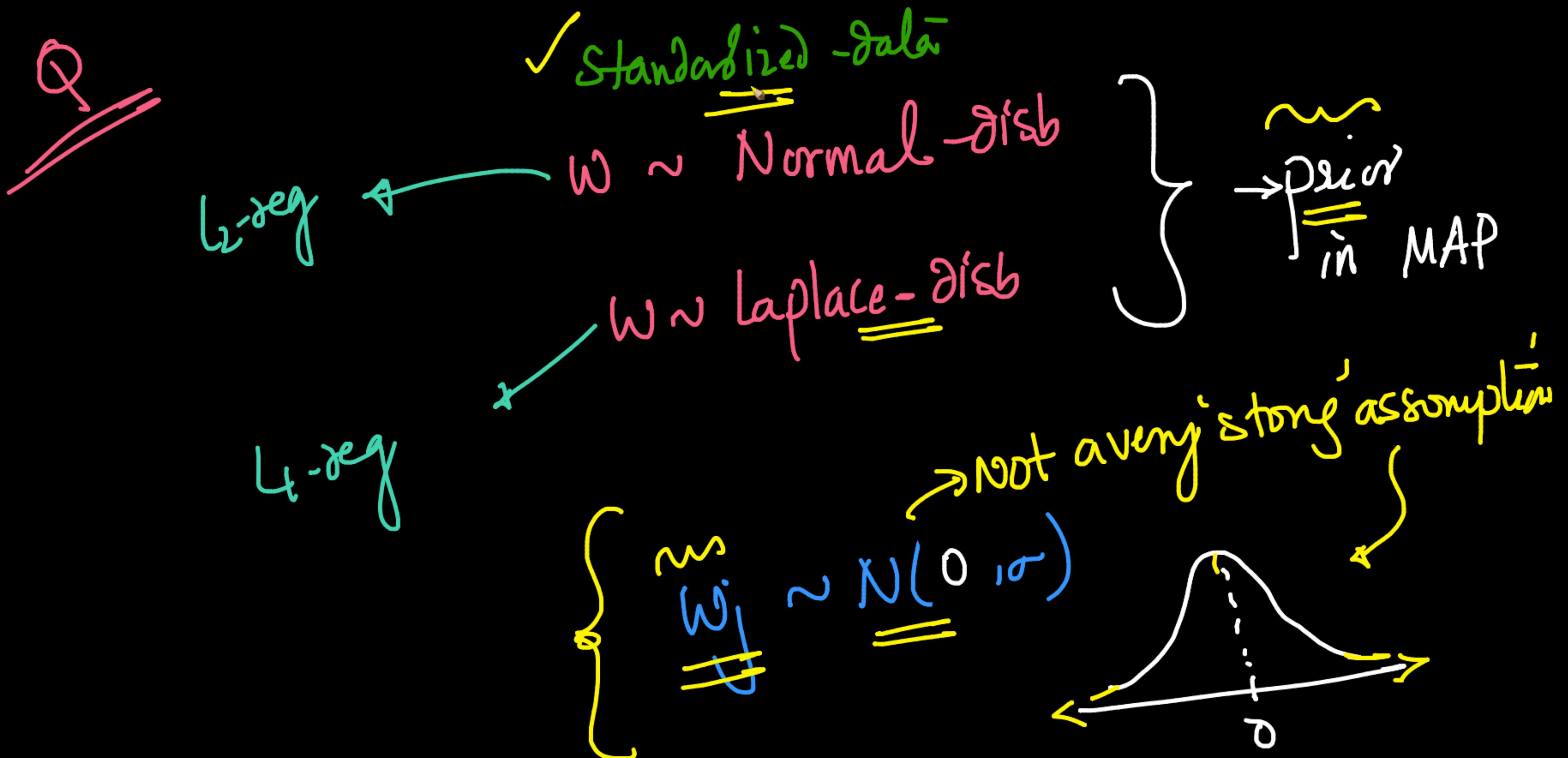


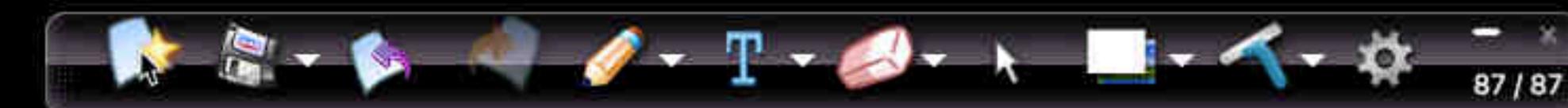
```
[ ] %%time
Z3 = UMAP(init='random', random_state=42).fit_transform(X)
```

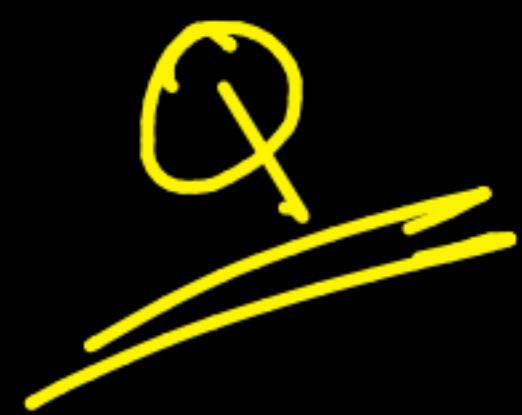
CPU times: user 8.39 s, sys: 30 ms, total: 8.42 s

Wall times: 8.27 s









tSNE

$x_i$ 's

① all

$P_{ij}$   $\rightarrow O(n^2)$

② Optimiza-pask

$\rightarrow$  multiple optimiza

$y_i$ 's  $\rightarrow Q_{ij}$   $O(n^2)$

KL-div  $O(n^2)$

$d \xrightarrow{\text{PCA}} d' \approx 2$

$$\text{minimize} \frac{1}{n} \sum_{i=1}^n \zeta_i + \lambda \|\mathbf{w}\|^2$$

subject to  $y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \zeta_i$  and  $\zeta_i \geq 0$ , for all  $i$ .

This is called the *primal* problem.

## Dual [edit]

By solving for the [Lagrangian dual](#) of the above problem, one obtains the simplified problem

$$\text{maximize } f(c_1 \dots c_n) = \sum_{i=1}^n c_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (\mathbf{x}_i^T \mathbf{x}_j) y_j c_j,$$

$$\text{subject to } \sum_{i=1}^n c_i y_i = 0, \text{ and } 0 \leq c_i \leq \frac{1}{2n\lambda} \text{ for all } i.$$

$O(n^2)$

This is called the *dual* problem. Since the dual maximization problem is a quadratic function of the  $c_i$  subject to linear constraints, it is efficiently solvable by [quadratic programming](#) algorithms.

Here, the variables  $c_i$  are defined such that

$$\mathbf{w} = \sum_{i=1}^n c_i y_i \mathbf{x}_i.$$

Moreover,  $c_i = 0$  exactly when  $\mathbf{x}_i$  lies on the correct side of the margin, and  $0 < c_i < (2n\lambda)^{-1}$  when  $\mathbf{x}_i$  lies on the margin's boundary. It follows that  $\mathbf{w}$  can be written as a linear combination of the support vectors.

The offset,  $b$ , can be recovered by finding an  $\mathbf{x}_i$  on the margin's boundary and solving

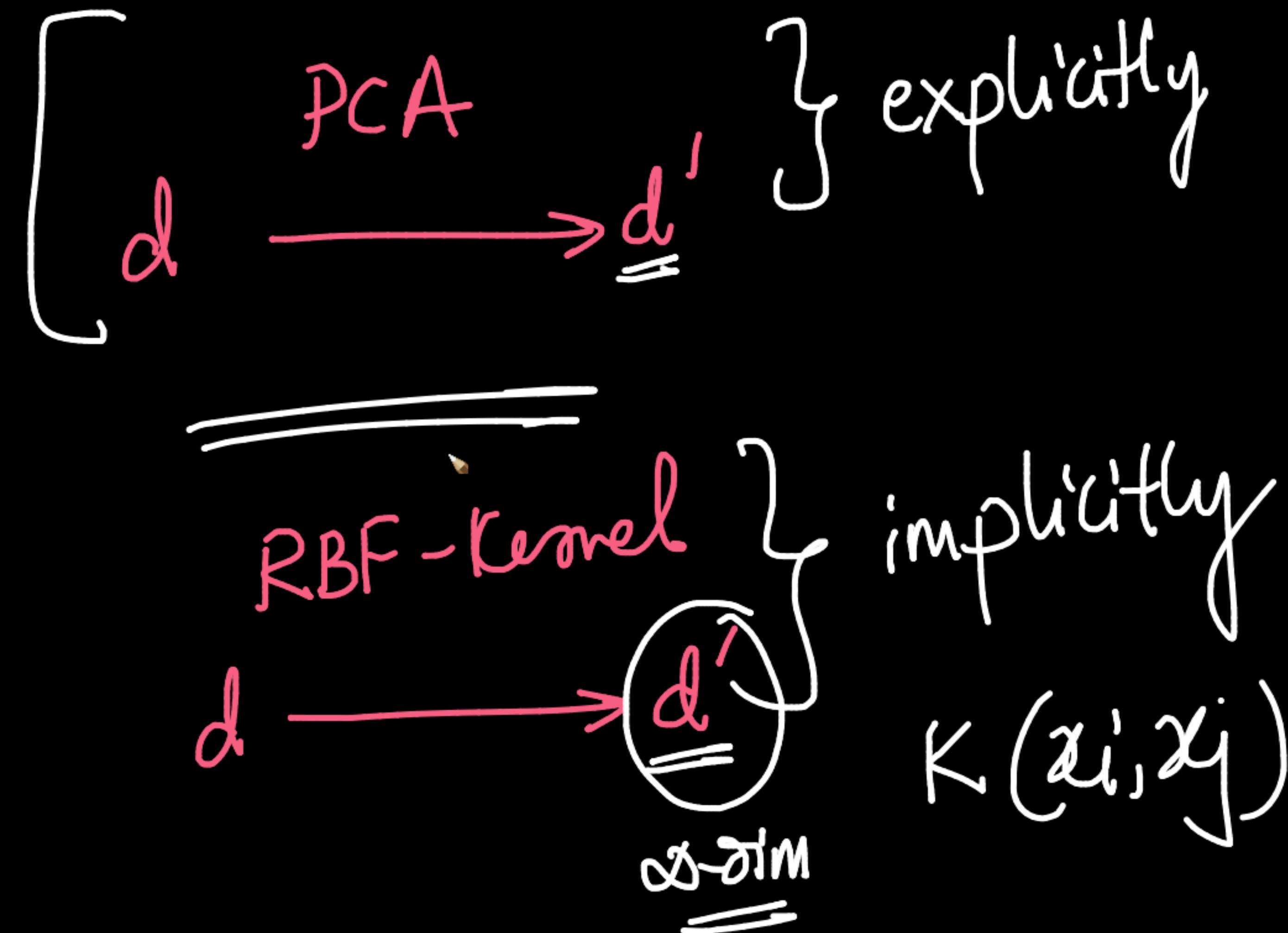
$$y_i(\mathbf{w}^T \mathbf{x}_i - b) = 1 \iff b = \mathbf{w}^T \mathbf{x}_i - y_i.$$

(Note that  $y_i^{-1} = y_i$  since  $y_i = \pm 1$ .)

## Kernel trick [edit]

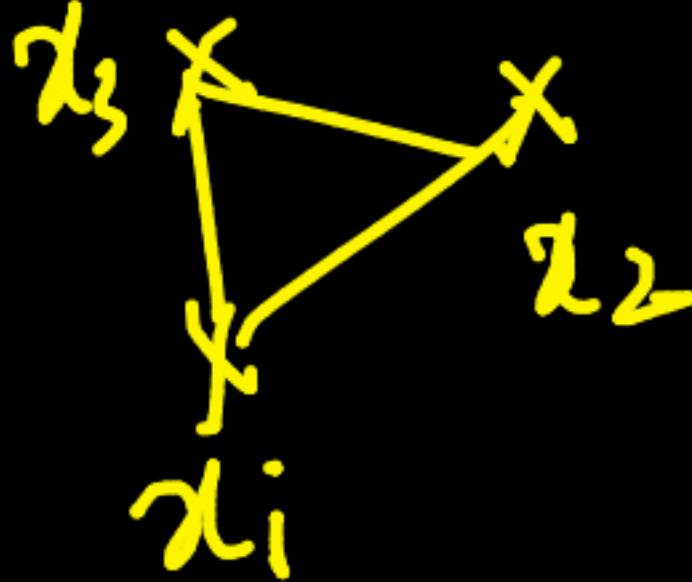
[Main article: Kernel methods](#)





Q

$x_1, x_2, \dots, x_{50}$



148dim

$x_i, x_j$

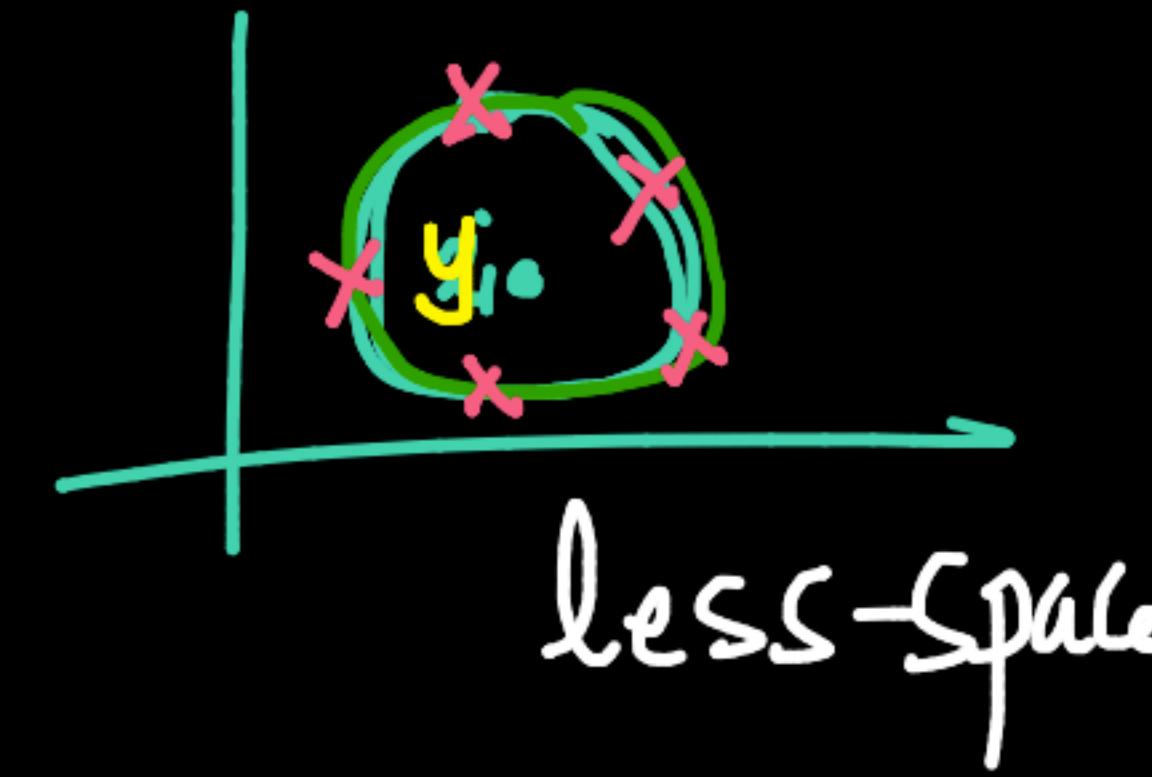
$tSNE$

Perplexity = 30  
(default)

3-dim



2-dim



1D



$$\text{minimize} \frac{1}{n} \sum_{i=1}^n \zeta_i + \lambda \|\mathbf{w}\|^2$$

subject to  $y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \zeta_i$  and  $\zeta_i \geq 0$ , for all  $i$ .

This is called the *primal* problem.

## Dual [edit]

By solving for the [Lagrangian dual](#) of the above problem, one obtains the simplified problem

$$\text{maximize } f(c_1 \dots c_n) = \sum_{i=1}^n c_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (\mathbf{x}_i^T \mathbf{x}_j) y_j c_j,$$

subject to  $\sum_{i=1}^n c_i y_i = 0$ , and  $0 \leq c_i \leq \frac{1}{2n\lambda}$  for all  $i$ .

This is called the *dual* problem. Since the dual maximization problem is a quadratic function of the  $c_i$  subject to linear constraints, it is efficiently solvable by [quadratic programming](#) algorithms.

Here, the variables  $c_i$  are defined such that

$$\mathbf{w} = \sum_{i=1}^n c_i y_i \mathbf{x}_i.$$

Moreover,  $c_i = 0$  exactly when  $\mathbf{x}_i$  lies on the correct side of the margin, and  $0 < c_i < (2n\lambda)^{-1}$  when  $\mathbf{x}_i$  lies on the margin's boundary. It follows that  $\mathbf{w}$  can be written as a linear combination of the support vectors.

The offset,  $b$ , can be recovered by finding an  $\mathbf{x}_i$  on the margin's boundary and solving

$$y_i(\mathbf{w}^T \mathbf{x}_i - b) = 1 \iff b = \mathbf{w}^T \mathbf{x}_i - y_i.$$

(Note that  $y_i^{-1} = y_i$  since  $y_i = \pm 1$ .)

## Kernel trick [edit]

Main article: [Kernel method](#)

