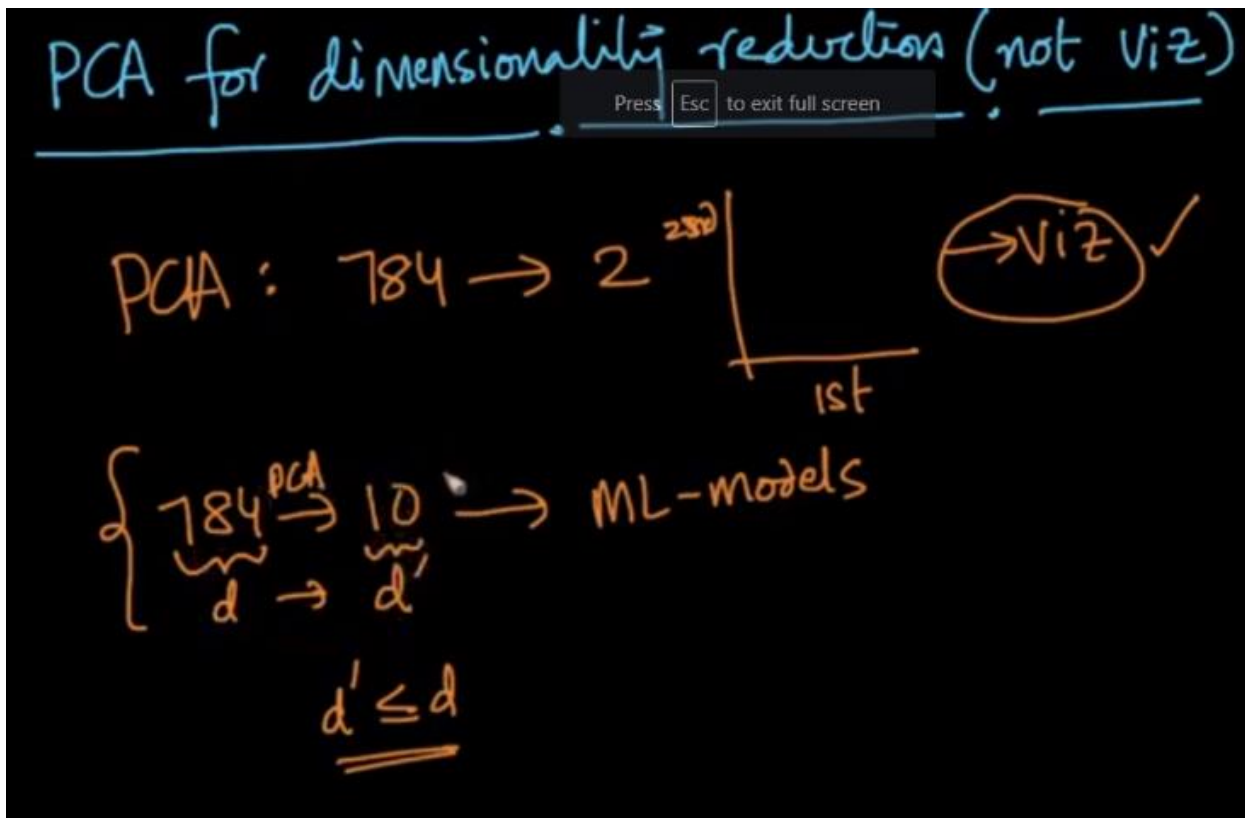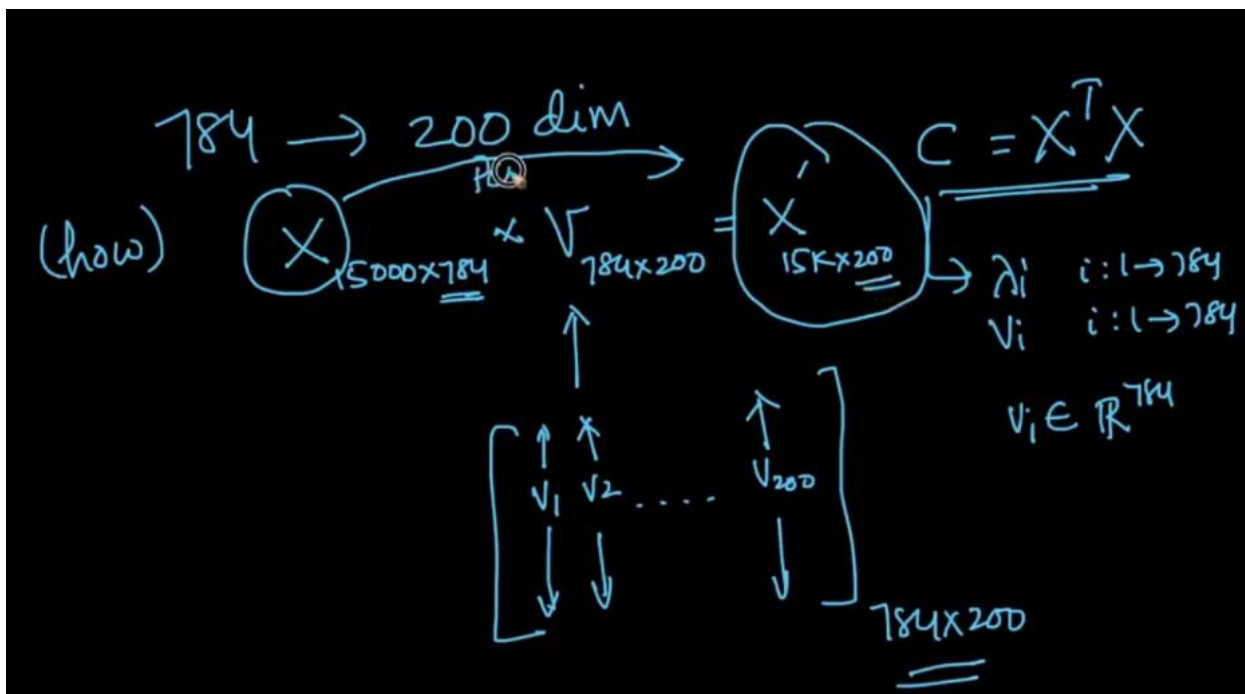This can be says as a summary of PCA, how we do it for dimensionality reduction.
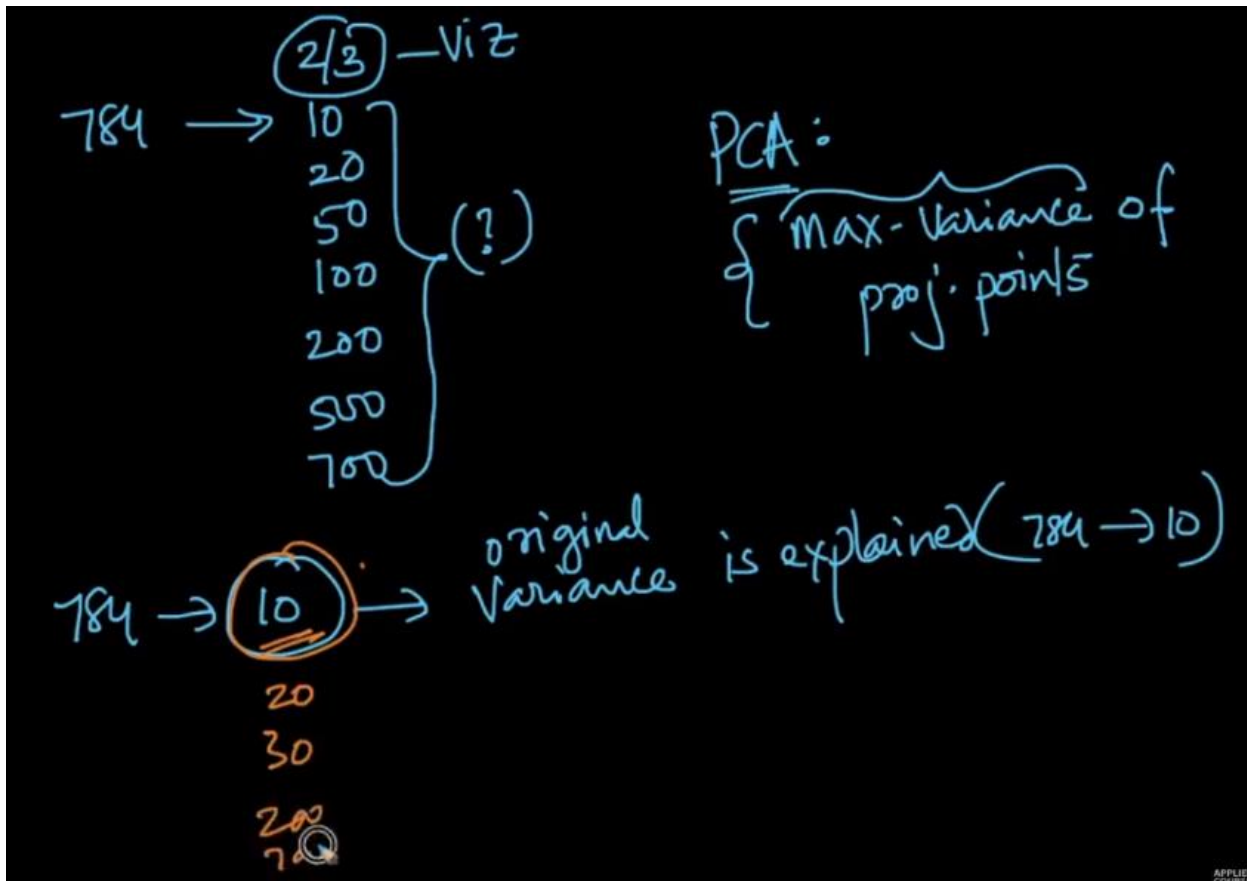
Our ultimate aim is to obtain features **d'** such that **d' <= d**, where d is original no. of features.

PCA for dimensionality reduction (not viz)

$$PCA: \quad 784 \longrightarrow 2^{\,250}$$

Press Esc to exit full screen

→Viz ✓

1st

$$\left\{ 784 \xrightarrow{PCA} 10 \longrightarrow ML\text{-}models \right.$$

$$d \longrightarrow d'$$

$$\underline{d' \leq d}$$

Let's say for mnist where we are obtaining 200 features from 784 features, here we have X of dimension 15000 * 784 and the eigen vectors matrix of dimension 784 * 200, which results the final matrix X' of dimension 15k * 200

$$784 \longrightarrow 200 \, dim$$

(how)

$$X_{15000 \times 784} \times V_{784 \times 200} = X_{15K \times 200}$$

$$C = X^T X$$

$$\rightarrow \lambda_i \quad i : 1 \rightarrow 784$$
$$V_i \quad i : 1 \rightarrow 784$$

$$V_i \in \mathbb{R}^{784}$$

$$\begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ V_1 & V_2 & \cdots & V_{200} \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}$$

$$784 \times 200$$

So our ultimate aim is to find the no of features for which the information retain percentage Is what we want which may b 75%, 99% and how it's calculated is given in below image.

# PCA:

$$C = X^T X$$
$784 \times 784$

$\lambda_i, V_i$ ✓

$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{784}$ ✓

$784 \rightarrow \boxed{10} \; dim$

Percentage of Variance explained in 10-dim

$$= \frac{\lambda_1 + \lambda_2 + \cdots + \lambda_{10}}{\sum_{i=1}^{784} \lambda_i}$$

$\boxed{0.2}$

20% of the total variance in 784-dim is explained in 10-dim

---

```
# PCA for dimensionality redcution (non-visualization)

pca.n_components = 784
pca_data = pca.fit_transform(sample_data)

percentage_var_explained = pca.explained_variance_ / np.sum(pca.explained_v

cum_var_explained = np.cumsum(percentage_var_explained)

# Plot the PCA spectrum
plt.figure(1, figsize=(6, 4))

plt.clf()
plt.plot(cum_var_explained, linewidth=2)
plt.axis('tight')
plt.grid()
plt.xlabel('n_components')
plt.ylabel('Cumulative_explained_variance')
plt.show()
```
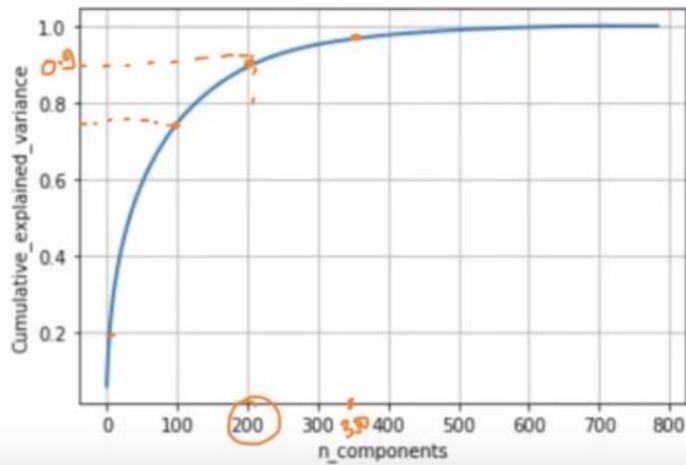
$\dfrac{\lambda_i}{\sum \lambda_i} \; \forall i$

$\dfrac{\lambda_1}{\sum \lambda_i}, \dfrac{\lambda_2}{\sum \lambda_i}, \dfrac{\lambda_3}{\sum \lambda_i} \cdots$

$\dfrac{\lambda_1}{\sum \lambda_i}, \dfrac{\lambda_1 + \lambda_2}{\sum \lambda_i}, \dfrac{\lambda_1 + \lambda_2 + \lambda_3}{\sum \lambda_i} \cdots$

# If we take 200-dimensions, approx. 90% of variance is expalined.



$784 \rightarrow 1n \quad (\sim 0.75)$

$\boxed{784} \rightarrow d' \quad (\sim 90\%)$

$PCA \rightarrow \boxed{200}$

$784 \rightarrow d' \quad (\sim 95\%)$

$\hookrightarrow \boxed{350}$

As we can see for 100 components we have 75%, for 200 we have 90% and for 300 components we have 95%.