

How to use PCA when the feature dimension is so huge then svd or eig will be so slow and out of memory?

Rasool Fakoor *

Problem Definition

Assume:

m: #data points(training examples).

n: #features.

A: training examples set which is $A_{n \times m}$.

In PCA algorithm, to calculate the Σ :

$$\Sigma = AA^T \quad (1)$$

Problem: The Σ will be $n \times n$. If $n \gg 64k$ then Σ is **HUGE**.

Clever solution when ($m \ll 64K$)

Instead of calculating $\Sigma = AA^T$, lets calculate $L = A^T A$. The size of L would be $m \times m$. As a result, now svd or eig deals with much smaller matrix to calculate the eigenvalues and eigenvectors. However, the question here is what is the relationship between L 's eigenvalues and eigenvectors and Σ 's eigenvalues and eigenvectors?

Proof

If v is eigenvector of L then Av is eigenvector of Σ :

$$\begin{aligned} Lv &= \gamma v \\ A^T A &= \gamma v \\ A(A^T A) &= A(\gamma v) = \gamma Av \\ (AA^T)Av &= \gamma(Av) \\ \Sigma(Av) &= \gamma(Av) \end{aligned}$$

Summary

When face with above problem, follow the following steps:

1. $L = A^T A$
2. $[v, u, s] = \text{svd}(L)$
3. $v' = Av$

Reference

- <http://ranger.uta.edu/~heng/CSE6363.html>

*<https://sites.google.com/site/rfakoor/>