

Big Data Dimension Reduction using PCA

Xuelong Wang

November 28, 2017

- 1 Introduction
- 2 PCA Regression
- 3 Sufficient Statistics
- 4 PCA on Big Data

Challenge of Big Data

① Memory Barrier

- The size of data is too large to load into memory
- More specifically, n is way more large than p

② The computation time

- It could be very time consuming if only use single core or cluster

Solution

① Sufficient statistics

- PCA regression only uses sufficient statistics
- Sufficient statistics can be calculated by scanning the data row-by-row

② Parallel computation

- Multiple threads
- Map-Reduced structure

Basic idea of PCA

Singular Value Decomposition

$$X_s = UDV^T, \text{ where } x_{ij,s} = \frac{x_{ij} - \bar{x}}{s_j}$$

$U = (u_1, \dots, u_r)$ is a n by r orthogonal matrix

$D = \text{diag}(d_1, \dots, d_r)$ is a r by r diagonal matrix

$V = (v_1, \dots, v_r)$ is a p by r orthogonal matrix

Basic idea of PCA

Principle Component and Loading

$$X_s = \underbrace{\begin{bmatrix} d_1 u_1 & \dots & d_r u_r \end{bmatrix}}_{\text{PCs}} \underbrace{\begin{bmatrix} v_1^T \\ \vdots \\ v_r^T \end{bmatrix}}_{\text{Loading}}$$

- $PC_j = d_j \mathbf{u}_j = X \mathbf{v}_j$ is the j th principle component
- The sample variance of PC_j is d_j^2 / n

Basic idea of PCA

Reduced matrix $X_{s,k}$

$$X_{s,k} = \sum_{j=1}^k d_j \mathbf{u}_j \mathbf{v}_k^T = U_k D_k V_k^T, \quad \text{Its Variation} \quad \sum_{j=1}^k d_j^2 / n.$$

Its proportion of the total variation is

$$\lambda_k = \frac{\sum_{j=1}^k d_j^2}{\sum_{j=1}^r d_j^2}$$

- If a small k such that $\lambda_k \approx 1$, we can use $U_k D_k$ in the follow up analysis

Follow-up analysis

The PCA approach is applied to a linear regression

Model

$$y = \mathbb{1}_n \alpha_s + U_k D_K \beta_{s,k} + \epsilon_{s,k},$$

Where $\epsilon_{s,k} \sim N(0, \sigma_{s,k}^2 \mathbb{I}_n)$

LSE and their variance

$$\hat{\alpha}_s = \bar{Y}, \quad \hat{\beta}_{s,k} = D_k^{-1} U_k^T y \quad (\text{PCs are Orthogonal})$$

$$\mathbb{V}(\hat{\sigma}_{s,k}^2) = [y^T (\underbrace{\mathbb{I}_n - \mathbb{J}_n/n - U_k U_k^T}_{\mathbb{I}_n - P_k}) y] / (n - k)$$

Sufficient Statistics

Factorization Theorem

$$f(x_1, x_2, \dots, x_n; \theta) = \phi[u(x_1, \dots, x_n); \theta] h(x_1, \dots, x_n)$$

- $u(x_1, \dots, x_n)$ is the sufficient statistics for θ
- If θ is a vector, then $u(x_1, \dots, x_n)$, the Joint Sufficient Statistics, will be also a vector.

Sufficient Statistics

model

$$y = \mathbb{1}_n \alpha + X\beta + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I_n)$$

Log-Likelihood function

$$\begin{aligned} \log \{f(y|x, \alpha, \beta, \sigma)\} &= \frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 \\ &\quad - \frac{1}{2\sigma^2} (c_{yy} - 2\alpha c_y + 2c_{xy}\beta + n\alpha^2 + 2\alpha c_x^T \beta + \beta^T C_{xx} \beta) \end{aligned}$$

- Define: $\mathcal{C}(y, X) = (c_0, c_{yy}, c_y, \mathbf{c}_{xy}, \mathbf{c}_x, C_{xx})$
- Note that $\ell(\alpha, \beta, \sigma)$ only depends on $\mathcal{C}(y, X)$

More about $\mathcal{C}(y, X)$

Elements of $\mathcal{C}(y, X)$

$$c_0 = n, \quad c_{yy} = \sum_{i=1}^n Y_i^2, \quad c_y = \sum_{i=1}^n Y_i,$$

$$c_{xy} = \sum_{i=1}^n Y_i x_i, \quad c_x = \sum_{i=1}^n x_i, \quad C_{xx} = \sum_{i=1}^n x_i x_i^T$$

- Notice that $\mathcal{C}(y, X)$ is the Joint sufficient statistics for (α, β, σ)
- All terms are in the summation format, so it can be calculated by reading the data row-by-row

Computation of $\mathcal{C}(y, X)$

Algorithm 1 Computation of $\mathcal{C}(y, X)$ Based on A Single Processor

Input: row-by-row of the data

Output: $\mathcal{C}(y, X)$

- 1: **procedure** ALGORITHM FOR $(c_0, c_{yy}, c_y, \mathbf{c}_{xy}, \mathbf{c}_x, \mathbf{C}_{xx})$
 - 2: Let $c_0, c_{yy}, c_y, \mathbf{c}_{xy}, \mathbf{c}_x$, and \mathbf{C}_{xx} be values, vectors, and matrix, respectively, all equal to zero
 - 3: **for** the i th row of the data **do** update $c_0 = c_0 + 1$,
 $c_{yy} = s_{yy} + Y_i^2$, $c_y = c_y + Y_i$, $\mathbf{c}_{xy} = \mathbf{c}_{xy} + Y_i \mathbf{x}_i$, $\mathbf{c}_x = \mathbf{c}_x + \mathbf{x}_i$, and $\mathbf{C}_{xx} = \mathbf{C}_{xx} + \mathbf{x}_i \mathbf{x}_i^T$ until the last row is scanned
 - 4: **end for**
 - 5: Output
 - 6: **end procedure**
-

- $\mathcal{O}((p+1)^2)$ memory size
- $\mathcal{O}(n(p+1)^2)$ floating operations

Sufficient Statistics based on X_S

Statistics affected by standardization

$$C_{S,XX} = X_S^T X_S, \quad c_{S,XY} = X_S^T y, \quad c_{S,X} = \mathbb{1}_n^T X_S / n$$

- It can be proved that, those statistics can be computed directly from $\mathcal{C}(y, X)$

Key step

$$c_{S,k,xy} = V_k^T c_{S,xy} = V_k^T V D U^T y = D_k U_k^T y \Rightarrow U_k^T y = D_k^{-1} c_{S,k,xy}$$

Back to PCA regression

$$X_s = UDV^T$$

Since n is large, it's not available to calculate U . However, we can only use V and D to get the estimated coefficients and variance.

PCA regression based on sufficient statistics

- ① D, V can be calculated by $C_{s,xx} = X_s X_s^T = VD^2V^T$
- ② $U_k^T y = D_k^{-1} c_{s,k,xy}$, so $\hat{\beta}_{s,k} = D_k^{-2} c_{s,k,xy}$
- ③ $\mathbb{V}(\hat{\sigma}_{s,k}^2) = [c_{yy} - c_{yy}/n - c_{s,k,xy}^T D_k^{-2} c_{s,k,xy}]/(n - k)$

Parallel Computation with Distributed Systems

Algorithm 2 Computation of $\mathcal{C}(\mathbf{y}, \mathbf{X})$ in MapReduce

Input: row-by-row of individual sub-data sets

Output: $\mathcal{C}(\mathbf{y}, \mathbf{X})$

- 1: **procedure** PARALLEL COMPUTATION BASED ON (17)
- 2: **Map tasks:** compute $\mathcal{C}(\mathbf{y}_i, \mathbf{X}_i)$ using Step 3 of Algorithm 1 for $i = 1, \dots, K$ individually
- 3: **Reduce task:** let $\mathcal{C}(\mathbf{y}, \mathbf{X}) = \sum_{i=1}^K \mathcal{C}(\mathbf{y}_i, \mathbf{X}_i)$
- 4: **Output**
- 5: **end procedure**

Thank you