

MAHOUT-817

PCA options for SSVD

working notes

November 29, 2011

1 Mean of rows

1.1 Recap of SSVD flow.

Modified SSVD Algorithm. Given an $m \times n$ matrix \mathbf{A} , a target rank $k \in \mathbb{N}_1$, an oversampling parameter $p \in \mathbb{N}_1$, and the number of additional power iterations $q \in \mathbb{N}_0$, this procedure computes an $m \times (k + p)$ SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ (some notations are adjusted):

1. Create seed for random $n \times (k + p)$ matrix $\mathbf{\Omega}$. The seed defines matrix $\mathbf{\Omega}$ using Gaussian unit vectors per one of suggestions in [?].
2. $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$, $\mathbf{Y} \in \mathbb{R}^{m \times (k+p)}$.
3. Column-orthonormalize $\mathbf{Y} \rightarrow \mathbf{Q}$ by computing thin decomposition $\mathbf{Y} = \mathbf{Q}\mathbf{R}$. Also, $\mathbf{Q} \in \mathbb{R}^{m \times (k+p)}$, $\mathbf{R} \in \mathbb{R}^{(k+p) \times (k+p)}$. I denote this as $\mathbf{Q} = \text{qr}(\mathbf{Y}) \cdot \mathbf{Q}$.
4. $\mathbf{B}_0 = \mathbf{Q}^\top \mathbf{A} : \mathbf{B} \in \mathbb{R}^{(k+p) \times n}$. (Another way is $\mathbf{R}^{-1} \mathbf{Y}^\top \mathbf{A}$, depending on whether we believe if size of \mathbf{A} less than size of \mathbf{Q}).
5. If $q > 0$ repeat: for $i = 1..q$: $\mathbf{B}_i^\top = \mathbf{A}^\top \text{qr}(\mathbf{A}\mathbf{B}_{i-1}^\top) \cdot \mathbf{Q}$ (power iterations step)
6. Compute Eigensolution of a small Hermitian $\mathbf{B}_q \mathbf{B}_q^\top = \hat{\mathbf{U}} \mathbf{\Lambda} \hat{\mathbf{U}}^\top$. $\mathbf{B}_q \mathbf{B}_q^\top \in \mathbb{R}^{(k+p) \times (k+p)}$.
7. Singular values $\mathbf{\Sigma} = \mathbf{\Lambda}^{0.5}$, or, in other words, $s_i = \sqrt{\sigma_i}$.
8. If needed, compute $\mathbf{U} = \mathbf{Q} \hat{\mathbf{U}}$.
9. If needed, compute $\mathbf{V} = \mathbf{B}_q^\top \hat{\mathbf{U}} \mathbf{\Sigma}^{-1}$. Another way is $\mathbf{V} = \mathbf{A}^\top \mathbf{U} \mathbf{\Sigma}^{-1}$.

1.2 \mathbf{B}_0 pipeline mods

This option considers that data points are rows in the $m \times n$ input matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_m \end{pmatrix}$$

Mean of rows is n-vector

$$\begin{aligned} \boldsymbol{\xi} &= \begin{pmatrix} \xi_1 \\ \xi_2 \\ \vdots \\ \xi_n \end{pmatrix} \\ &= \frac{1}{m} \sum_i^m \mathbf{a}_i. \end{aligned}$$

Let $\tilde{\mathbf{A}}$ be \mathbf{A} with the mean subtracted.

$$\tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{a}_1 - \boldsymbol{\xi} \\ \mathbf{a}_2 - \boldsymbol{\xi} \\ \vdots \\ \mathbf{a}_m - \boldsymbol{\xi} \end{pmatrix}.$$

We denote $m \times n$ mean matrix

$$\boldsymbol{\Xi} = \begin{pmatrix} \boldsymbol{\xi} \\ \boldsymbol{\xi} \\ \vdots \\ \boldsymbol{\xi} \end{pmatrix}$$

\mathbf{B}_0 pipeline starts with notion that since $\tilde{\mathbf{A}}$ is dense, its mutliplications are very expensive. Hence, we factorize \mathbf{Y} as

$$\begin{aligned} \mathbf{Y} &= \tilde{\mathbf{A}}\boldsymbol{\Omega} \\ &= \mathbf{A}\boldsymbol{\Omega} - \boldsymbol{\Xi}\boldsymbol{\Omega} \end{aligned}$$

Current \mathbf{B}_0 pipeline already takes care of $\mathbf{A}\boldsymbol{\Omega}$, but the term $\boldsymbol{\Xi}\boldsymbol{\Omega}$ will need more work.

The term $\Xi\Omega$ will have identical rows $\xi\Omega$ so we need to precompute just one dense n-vector $\xi\Omega$. This computation is very expensive since matrix Ω is dense (potentially several orders of magnitude bigger than input \mathbf{A}) and the median ξ is dense as well, even that we don't actually have to materialize any of Ω . *Question is whether we could just ignore it since $\mathbb{E}(\xi\Omega) = 0$.* Alternatively, we could just brute-force it by creating a separate distributed computation of this over n .

⇐ Outstanding issue!!!

Moving onto \mathbf{B} and $\mathbf{B}\mathbf{B}^\top$. Here and on we assume $\mathbf{B} \equiv \mathbf{B}_0$ and omit the index for compactness.

$$\mathbf{B} = \mathbf{Q}^\top \tilde{\mathbf{A}} \quad (1)$$

$$= \mathbf{Q}^\top \mathbf{A} - \mathbf{Q}^\top \Xi. \quad (2)$$

Again, current pipeline takes care of $\mathbf{Q}^\top \mathbf{A}$ but product $\mathbf{Q}^\top \Xi$ would need more work.

Let $\mathbf{W} = \mathbf{Q}^\top \Xi$.

We see that all columns of \mathbf{W} are identical, and, more specifically*,

$$\begin{aligned} \mathbf{W}_{*,i} &= \mathbf{w} \\ &= (\mathbf{Q}^\top \Xi)_{*,i} \\ &= \left[\sum_{i=1}^m \mathbf{Q}_{i,*} \right] \circ \xi \\ &= \mathbf{s}_Q \circ \xi \quad \forall i \in [1, n], \end{aligned}$$

where $\mathbf{s}_Q = \sum_{i=1}^m \mathbf{Q}_{i,*}$ is sum of all rows of \mathbf{Q} .

Since B_0 pipeline computes $\mathbf{Q}^\top \mathbf{A}$ column-wise over columns of \mathbf{Q} and \mathbf{A} , the first thought is that (2) can be computed column-wise as well with computation seeded by the \mathbf{w} vector.

One problem with our first thought is that the \mathbf{s}_Q term is not yet known at the time of formation of \mathbf{B} columns because formation of final \mathbf{Q} blocks happens in the same distributed map task that produces initial $\mathbf{Q}^\top \mathbf{A}$ blocks. Hence, the sum of \mathbf{Q} rows at that moment would not be available. But we probably can fix our output later at the time when \mathbf{s}_Q would already have been known.

*Let also $\mathbf{a} \circ \mathbf{b} = \begin{pmatrix} a_1 b_1 \\ a_2 b_2 \\ \vdots \\ a_k b_k \end{pmatrix}$ to be a notation for element-wise vector product (Hadamard

product?).

Let $\mathbf{b}_i = \mathbf{B}_{*,i}$, $\tilde{\mathbf{b}}_i = (\mathbf{Q}^\top \mathbf{A})_{*,i}$. Then correction for \mathbf{B} output would be

$$\mathbf{b}_i = \tilde{\mathbf{b}}_i - \mathbf{w}. \quad (3)$$

Moving on to $\mathbf{B}\mathbf{B}^\top$:

$$\mathbf{B}\mathbf{B}^\top = \sum_i^n \mathbf{b}_i \mathbf{b}_i^\top$$

$$\begin{aligned} \mathbf{b}_i \mathbf{b}_i^\top &= (\tilde{\mathbf{b}}_i - \mathbf{w})(\tilde{\mathbf{b}}_i - \mathbf{w})^\top \\ &= \tilde{\mathbf{b}}_i \tilde{\mathbf{b}}_i^\top - \tilde{\mathbf{b}}_i \mathbf{w}^\top - \mathbf{w} \tilde{\mathbf{b}}_i^\top - \mathbf{w} \mathbf{w}^\top \\ &= \tilde{\mathbf{b}}_i \tilde{\mathbf{b}}_i^\top - \tilde{\mathbf{b}}_i \mathbf{w}^\top - (\tilde{\mathbf{b}}_i \mathbf{w}^\top)^\top + \mathbf{w} \mathbf{w}^\top. \end{aligned}$$

$$\mathbf{B}\mathbf{B}^\top = \sum_i^n \tilde{\mathbf{b}}_i \tilde{\mathbf{b}}_i^\top \quad (4)$$

$$- \sum_i^n \left[\tilde{\mathbf{b}}_i \mathbf{w}^\top + (\tilde{\mathbf{b}}_i \mathbf{w}^\top)^\top \right] \quad (5)$$

$$+ n \cdot \mathbf{w} \mathbf{w}^\top. \quad (6)$$

Let $k \times k$ matrix $\mathbf{C} = \sum_i^n \tilde{\mathbf{b}}_i \mathbf{w}^\top$, and then we can rewrite (5) as

$$\mathbf{B}\mathbf{B}^\top = \sum_i^n \tilde{\mathbf{b}}_i \tilde{\mathbf{b}}_i^\top - \mathbf{C} - \mathbf{C}^\top + n \cdot \mathbf{w} \mathbf{w}^\top.$$

So we can compute $\tilde{\mathbf{B}} = \sum_i \tilde{\mathbf{b}}_i \tilde{\mathbf{b}}_i^\top$ right away, that's what Bt-job does. We also can add $n \cdot \mathbf{w} \mathbf{w}^\top$ in front end before we do eigendecomposition since it is a tiny matrix and at that point \mathbf{w} is already known. The task boils down to computing small $(k+p) \times (k+p)$ matrix \mathbf{C} and then subtracting $[\mathbf{C} + \mathbf{C}^\top]$ in the front end as well. Note that

$$\begin{aligned} \mathbf{C} &= \sum_i^n \tilde{\mathbf{b}}_i \mathbf{w}^\top \\ &= \left(\sum_i^n \tilde{\mathbf{b}}_i \right) \mathbf{w}^\top \\ &= \mathbf{s}_{\tilde{B}} \mathbf{w}^\top. \end{aligned}$$

In this case, $\mathbf{s}_{\tilde{B}} = \sum_i^n \tilde{\mathbf{b}}_i$ can be output by Bt job as well. Hence \mathbf{C} can be computed as an outer product of two small k-vectors in the front end as well.

PCA would be primarily interested in \mathbf{V} or \mathbf{V}_σ output of the decomposition in order to fold in new items back into PCA space, so we need to correct \mathbf{V} job as well in this case to fix output of Bt-job per (3).

1.3 Power Iterations (aka \mathbf{B}_i pipeline) additions

Power iterations pipeline produces $\mathbf{B}_i^\top = \tilde{\mathbf{A}}^\top \text{qr}(\tilde{\mathbf{A}}\mathbf{B}_{i-1}^\top) \cdot \mathbf{Q}$. Similarly to versions of \mathbf{B} , each iteration would produce corrective vector \mathbf{w}_{i-1} .

First, we need to amend power iteration work flow to fix output of previous Bt-job on the fly with \mathbf{w}_{i-1} to reconstruct correct \mathbf{B}_{i-1} similarly to what is done in the \mathbf{V} per (3):

$$\mathbf{B}_{i-1} = \tilde{\mathbf{B}}_{i-1} - \mathbf{W}_{i-1}.$$

Second, again, $\tilde{\mathbf{A}}$ multipliers are a problem because they would be dense and perhaps should be decomposed in a way similar to \mathbf{B}_0 pipeline.

=====> to be ctd. <=====

~~Another note is that we run eigendecomposition only after the last iteration so the term $\mathbf{s}_{\tilde{\mathbf{B}}}$ needs to be computed only during the last iteration.~~