

Command Line Interface, Stochastic SVD*

1 Overview.

Stochastic SVD method in Mahout produces reduced rank Singular Value Decomposition output in its strict mathematical definition:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top,$$

i. e. it creates outputs for matrices \mathbf{U} , \mathbf{V} and $\mathbf{\Sigma}$, each of which may be requested individually. The desired rank of decomposition, henceforth denoted as $k \in \mathbb{N}_1$, is a parameter of the algorithm. The singular values inside diagonal matrix $\mathbf{\Sigma}$ satisfy $\sigma_{i+1} \leq \sigma_i \ \forall i \in [1, k-1]$, i.e. sorted from biggest to smallest. Cases of rank deficiency $\text{rank}(\mathbf{A}) < k$ are handled by producing 0s in singular value positions once deficiency takes place.

Single space for comparing row-items and column-items. On top of it, there's an option to present decomposition output in a form of

$$\mathbf{A} = \mathbf{U}_\sigma \mathbf{V}_\sigma^\top, \tag{1}$$

where one can request $\mathbf{U}_\sigma = \mathbf{U}\mathbf{\Sigma}^{0.5}$ instead of \mathbf{U} (but not both), $\mathbf{V}_\sigma = \mathbf{V}\mathbf{\Sigma}^{0.5}$ instead of \mathbf{V} (but not both). Here, notation $\mathbf{\Sigma}^{0.5}$ implies diagonal matrix containing square roots of the singular values:

$$\mathbf{\Sigma}^{0.5} = \begin{pmatrix} \sqrt{\sigma_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sqrt{\sigma_k} \end{pmatrix}.$$

*Dmitriy Lyubimov, dlyubimov at apache dot org

Original singular values Σ are still produced and saved regardless.

This option is a nod to a common need of comparing actors represented by both input rows and input columns in a common space. E.g. if LSI is performed such that rows are documents and columns are terms then it is possible to compare documents and terms (ether existing or fold in new ones) in one common space and perform similarity measurement between a document and a term, rather than computing just term2term or document2document similarities.

Folding in new observations. It is probably worth mentioning the operation of “folding in” new observations in context of this method, since it is often a basis for incremental methods.

If $\tilde{\mathbf{c}}_r$ ($\tilde{\mathbf{c}}_c$) is a new row (column) observation in addition to original input \mathbf{A} , then correspondent “new” row vectors of $\tilde{\mathbf{U}}$ ($\tilde{\mathbf{V}}$) can be obtained as

$$\begin{aligned}\tilde{\mathbf{u}} &= \Sigma^{-1} \mathbf{V}^\top \tilde{\mathbf{c}}_r, \\ \tilde{\mathbf{v}} &= \Sigma^{-1} \mathbf{U}^\top \tilde{\mathbf{c}}_c.\end{aligned}$$

Similarly, for the case (1) folding in new observations into rows of $\tilde{\mathbf{U}}_\sigma$ ($\tilde{\mathbf{V}}_\sigma$) would look like

$$\begin{aligned}\tilde{\mathbf{u}}_\sigma &= \mathbf{V}_\sigma^\top \tilde{\mathbf{c}}_r, \\ \tilde{\mathbf{v}}_\sigma &= \mathbf{U}_\sigma^\top \tilde{\mathbf{c}}_c.\end{aligned}$$

Thus, new rows can be added to matrices denoted as $\tilde{\mathbf{U}}$ ($\tilde{\mathbf{V}}$) corresponding to new observations as new observations become available, i.e. incrementally. Given that new observations are usually moderately sparse vectors, it might be feasible to do fold-in in real time or almost real time, assuming proper fast row-wise indexing of \mathbf{U} (\mathbf{V}) exists (e.g. using a batch request to an HBase table containing rows of \mathbf{U} (\mathbf{V})). However, since operation of folding in new observations doesn’t change original decomposition and its spaces, such new observations cannot be considered ’training’ examples. Typically, from time to time accumulated new observations can be added to original input \mathbf{A} and the whole decomposition can be recomputed again.

Common applications for SVD include Latent Semantic Analysis (LSA), Principal Component Analysis (PCA), dimensionality reduction and others.

2 File formats

Input \mathbf{A} , as well as outputs \mathbf{U} (\mathbf{U}_σ), \mathbf{V} (\mathbf{V}_σ), are in Mahout’s Distributed Row Matrix format, i.e. set of sequence files where value is of `VectorWritable`

type. As far as keys are concerned, rows of **A** may be keyed (identified) by any **Writable** (for as long as it is instantiable thru a default constructor). That, among other things, means that this method can be applied directly on the output of **seq2sparse** where keys are of **Text** type¹.

Definition of output **U** (**U_σ**) is identical to definition of the input matrix **A**, and the keys of corresponding rows in **A** are copied to corresponding rows of output **U** (**U_σ**).

Definition of output **V** (**V_σ**) is always sequence file(s) of (**IntWritable**, **VectorWritable**) where key corresponds to a column index of the input **A**.

Output of **Σ** is encoded by a single output file with a single vector value (**VectorWritable**) with main diagonal entries of **Σ** aka singular values ($\sigma_1 \dots \sigma_k$).

3 Usage²

```
mahout ssvd <options>
```

Options.

-k, --rank <int-value> (required): the requested SVD rank (minimum number of singular values and dimensions in U, V matrices). The value of $k + p$ directly (as in $O([k + p]^3)$) impacts running time and memory requirements. ***k+p=500 is probably more than reasonable.*** Typically $k + p$ is taken within range 20...200.

-p, --oversampling <int-value> (optional, default 15): stochastic SVD oversampling. p doesn't seem to have to be very significant. If power iterations ($q > 0$) are used then p perhaps could be kept quite low, not to exceed 10% of k .

-q, --powerIter <int-value> (optional, default 0): number of power iterations to perform. This helps fighting data noise and improve precision significantly more than just increasing p . Each additional power iteration adds 2 more steps (map/reduce + map-only). Experimental data suggests using $q = 1$ is already producing quite good results which are hard to much improve upon.

¹(TODO: re-verify)

²As of Mahout 0.6 trunk.

12/13/2011 adjusted for MAHOUT-922 changes.

-r, --blockHeight <int-value> (optional, default 10,000): the number of rows of source matrix for block computations during $\mathbf{Y} = \mathbf{QR}$ decomposition. Taller blocking causes more memory use but produces less blocks and therefore somewhat better running times. The most optimal mode from the running time point of view should be 1 block per 1 mapper. *This cannot be less than $k+p$.*

-oh, --outerProdBlockHeight <int-value> (optional, default 30,000): the block height during $\mathbf{B} = \mathbf{Q}^\top \mathbf{A}$ operation³⁴.

-abth, --abtBlockHeight <int-value> (optional, default 200,000): the block height during $\mathbf{Y}_i = \mathbf{AB}_i^\top$ multiplication³⁴.

-s, --minSplitSize <int-value> (optional, default: use Hadoop's default): minimum split size to use in mappers reading \mathbf{A} input.⁵

--computeU <true|false> (optional, default true). Request computation of the U matrix

--computeV <true|false> (optional, default true). Request computation of the V matrix

--vHalfSigma <true|false> (optional, default: false): compute $\mathbf{V}_\sigma = \mathbf{V}\Sigma^{0.5}$ instead of \mathbf{V} (see overview for explanation).

--uHalfSigma <true|false> (optional, default: false): compute $\mathbf{U}_\sigma = \mathbf{U}\Sigma^{0.5}$ instead of \mathbf{U} .

--reduceTasks <int-value> optional. The number of reducers to use (where applicable): depends on the size of the hadoop cluster. At this point it could also be overwritten by a standard hadoop property using -D option¹. ***Probably always needs to be specified as by default Hadoop would set it to 1, which is certainly far below the cluster capacity.*** Recommended value for this option $\sim 95\%$

³With extreme sparse matrices increasing this parameter may lead to better performance by reducing computational pressure on the shuffle and sort and grouping sparse records together.

⁴ Watch for GC thrashing and swap. Values too high may cause GC thrashing and/or swapping, both of which are capable of bringing job performance down to a halt. Don't starve jobs for memory. Defaults are believed to work well for -Xmx800mb or above per child task.

⁵As of this day, I haven't heard of a case where somebody would actually have to use this option and actually increase split size and how it has played out. So this option is experimental.

Since in this version projection block formation happens in mappers, for a sufficiently wide input matrix the algorithm may not be able to read minimum $k + p$ rows and form a block of minimum height required, so in that case the job would bail out at the very first mapping step. If this happens, one of the recourses available is to force increase in the MapReduce split size using SequenceFileInputFormat.setMinSplitSize() property. Increasing this significantly over HDFS block size may result in network IO to mappers. Another caveat is that one sometimes does not want too many mappers because it may in fact increase time of the computation. Consequently, this option should probably be left alone unless one has significant amount of mappers (as in thousands of map tasks) at which point reducing amount of mappers may actually improve the throughput (just a guesstimate at this point).

or ~190% of available reducer capacity to allow for opportunistic executions.

Standard Mahout options.

`--input <glob>` HDFS glob specification where the DistributedRowMatrix input to be found.

`--output <hdfs-dir>` non-existent hdfs directory where to output \mathbf{U} , \mathbf{V} and $\mathbf{\Sigma}$ (singular values) files.

`--tempDir <temp-dir>` temporary dir where to store intermediate files (cleaned up upon normal completion). This is a standard Mahout optional parameter.

`-ow, --overwrite` overwrite output if exists.

4 Embedded use

It is possible to instantiate and use `SSVDSolver` class in embedded fashion in Hadoop-enabled applications. This class would have getter and setter methods for each option available via command line. See javadoc for details.