16 Matrix Sketching

The singular value decomposition (SVD) can be interpreted as finding the most dominant *directions* in an $(n \times d)$ matrix A (or n points in \mathbb{R}^d). Typically n > d. It is typically easy to call a built in version of the SVD in many programming languages

$$[U, S, V] = \operatorname{svd}(A)$$

where $U = [u_1, \ldots, u_n]$, $S = \operatorname{diag}(\sigma_1, \ldots, \sigma_d)$, and $V = [v_1, \ldots, v_d]$. Then $A = USV^T$ and in particular $A = \sum_{j=1}^d \sigma_j u_j v_j^T$. To approximate A we just use the first k components to find $A_k = \sum_{j=1}^k \sigma_j u_j v_j^T = U_k S_k V_k^T$ where $U_k = [u_1, \ldots, u_k]$, $S_k = \operatorname{diag}(\sigma_1, \ldots, \sigma_k)$, and $V_k = [v_1, \ldots, v_k]^T$. Then the vectors v_j (starting with smaller indexes) provide the best subspace representation of A.

But, although SVD has been *heavily* optimized on data sets that fit in memory (via LAPACK, found in Matlab; Python in SciPy; C++ in Armadillo, and just about every other language), it can sometimes be improved. The traditional SVD takes $O(\min\{nd^2, n^2d\})$ time to compute, which can be prohibitive for large n and/or d. Here we highlight two of these ways:

- to provide better interpretability of each v_i .
- to be more efficient on enormous scale, in a stream, or in distributed settings.

16.0.1 Covariance Matrix Summation

The first regime we focus on is when n is extremely large, but d is moderate. For instance n=100 million, and d=100. The a simple approach in a stream is to make one pass using d^2 space, and just maintain the sum of outer products $C=\sum_{i=1}a_ia_i^T$, the $d\times d$ covariance matrix of A exactly.

Algorithm 16.0.1 Summed Covariance

```
Set C all zeros (d \times d) matrix.

for rows (i.e. points) a_i \in A do

C = C + a_i a_i^T

return C
```

We have that at any point i in the stream, where $A_i = [a_1; a_2; \dots, a_i]$, the maintained matrix C is precisely $C = A_i A_i^T$. Thus the eigenvectors of C are the right singular vectors of C, and the eigenvalues of C are the squared singular values of C. This only requires d^2 space, and nd^2 total time, and incurs no error.

We can choose the top k eigenvectors of C as V_k , and on a second pass of the data, project all vectors on a_i onto V_k to obtain the best k-dimensional embedding of the dataset.

16.1 Frequent Directions

The next regime assumes that n is extremely large (say n=100 million), but that d is also uncomfortably large (say d=100 thousand), and our goal is something like a best rank k-approximation with $k\approx 10$. So $k\ll d\ll n$. In this regime perhaps d is so large that d^2 space is too much, but something close to dk space and ndk time is reasonable. We will not be able to solve things exactly in the streaming setting under these constraints, but we can provide a provable approximation with slightly more space and time.

This approach, called *Frequent Directions*, can be viewed as an extension of the Misra-Gries sketch. We will consider a matrix A one row (one point $a_i \in \mathbb{R}^d$) at a time. We will maintain a matrix B that is $2\ell \times d$,

that is it only has 2ℓ rows (directions). We say a row of B is *empty* if it contains all 0s. We maintain that one row of B is always empty at the end of each round (this will always include the last row $B_{2\ell}$).

We initialize with the first $2\ell-1$ rows a_i of A as B, again with the last row B_ℓ left as all zeros. Then on each new row, we put a_i in the empty row of B. We set $[U,S,V]=\operatorname{svd}(B)$. Now examine $S=\operatorname{diag}(\sigma_1,\ldots,\sigma_{2\ell})$, which is a length 2ℓ diagonal matrix. If $\sigma_{2\ell}=0$ (then a_i is in the subspace of B), do nothing. Otherwise subtract $\delta=\sigma_\ell^2$ from each (squared) entry in S, that is $\sigma_j'=\sqrt{\max\{0,\sigma_j^2-\delta\}}$ and in general $S'=\operatorname{diag}(\sqrt{\sigma_1^2-\delta},\sqrt{\sigma_2^2-\delta},\ldots,\sqrt{\sigma_{\ell-1}^2-\delta},0,\ldots,0)$.

Now we set $B = S'V^T$. Notice, that since S' only has non-zero elements in the first $\ell - 1$ entries on the diagonal, then B is at most rank $\ell - 1$ and we can then treat V and B as if the ℓ th row does not exist.

Algorithm 16.1.1 Frequent Directions

```
Set B all zeros (2\ell \times d) matrix. 

for rows (i.e. points) a_i \in A do

Insert a_i into a zero-valued row of B

if (B has no zero-valued rows) then
[U,S,V] = \operatorname{svd}(B)
\operatorname{Set} \delta_i = \sigma_\ell^2 \qquad \qquad \text{# the $\ell$th entry of $S$}
\operatorname{Set} S' = \operatorname{diag}\left(\sqrt{\sigma_1^2 - \delta}, \sqrt{\sigma_2^2 - \delta}, \ldots, \sqrt{\sigma_{\ell-1}^2 - \delta}, 0, \ldots, 0\right).
\operatorname{Set} B = S'V^T \qquad \qquad \text{# the last rows of $B$ will again be all zeros}
\mathbf{return} B
```

The result of Algorithm 16.1.1 is a matrix B such that for any (direction) unit vector $x \in \mathbb{R}^d$

$$0 \le ||Ax||^2 - ||Bx||^2 \le ||A - A_k||_F^2 / (\ell - k)$$

and

$$||A - A\Pi_{B_k}||_F^2 \le \frac{\ell}{\ell - k} ||A - A_k||_F^2,$$

for any $k < \ell$, including when k = 0. So setting $\ell = 1/\varepsilon$, then in any direction in \mathbb{R}^d , the squared mass in that direction is preserved up to $\varepsilon \|A\|_F^2$ (that is, ε times the total squared mass) using the first bound. And in the second bound if we set $\ell = \lceil k/\varepsilon + k \rceil$ then we have $\|A - A\Pi_{B_k}\|_F^2 \le (1+\varepsilon)\|A - A_k\|_F^2$. Recall that $\|A\|_F = \sqrt{\sum_{a_i \in A} \|a_i\|^2}$.

• Why does this work?

Just like with Misra-Greis, when some mass is deleted from one counter it is deleted from all ℓ counters, and none can be negative. So here when one direction has its (squared) mass decreased, at least ℓ directions (with non-zero squared mass) are decreased by the same amount. So no direction can have more than $1/\ell$ fraction of the total squared mass $\|A\|_F^2$ decreased from it.

Finally, since squared mass can be summed independently along any set of **orthogonal** directions, we can subtract each of them without affecting others.

• Why do we use the svd?

The SVD defines the true axis of the ellipse associated with the norm of B at each step. If we shrink along an basis (or even a set of non-orthogonal vectors) we will warp the ball, and we will not be able to ensure that each direction of B shrinks in squared norm by at most δ_i .

- Did we need to use the svd? (its expensive, right)?
 The cost is amortized. We only call the svd once every \(\ell \) steps, so at most \(O(n/\ell) \) times. Since each call takes \(O(d\ell^2) \) time, the total cost is \(O(nd\ell) \), or only \(\ell \) times as long as reading the matrix. It is also possible to call approximate versions of the SVD. This allows versions which have runtime depending on the number of non-zeros in the input matrix. This makes a big difference for very sparse word count or recommendation system matrices.
- What happened to U in the svd output?

 The matrix U just related the main directions to each of the n points (rows) in A. But we don't want to keep around the space for this. In this application, we only care about the directions or subspace that best represents the points; e.g. PCA only cares about the right singular vectors.

The Frequent Direction algorithm calls an SVD operation multiple times, so what it is useful to bound its runtime to ensure it does not take much longer than a single decomposition of A, which would take time proportional to nd^2 . In this case, the matrix B which has its SVD computed is only $2\ell \times d$, so this operation can be performed in time proportional to only $d\ell^2$. We do this step every ℓ rows, so it occurs n/ℓ times. In total this takes time proportional to $n/\ell \cdot d\ell^2 = nd\ell$. So when $\ell \ll d$, (recall $\ell \approx k/\varepsilon$) this is a substantial decrease in running time. Moreover, it only requires $2\ell d$ space, so it can perform these operations without reading data from disk more than once.

16.2 Row Sampling

We next move to a regime where n and d are again both large, and so might be k. But a runtime of ndk may be too large – that is we can read the data, but maybe a factor of k times reading the data is also large. The next algorithms have runtime slightly more than nd, they are almost as fast as reading the data. In particular, if there are $\operatorname{nnz}(A)$ non-zero entries in a very sparse matrix, and we only need to read these entries, then the runtime is almost proportional to $\operatorname{nnz}(A)$.

The goal is to approximate A up to the accuracy of A_k . But in A_k the directions v_i are linear combinations of features.

- What is a linear combination of genes?
- What is a linear combination of typical grocery purchases?

Instead our goal is to choose V so that the columns of V are also columns of A.

For each row of $a_i \in A$, set $w_i = ||a_i||^2$. Then select $t = (k/\varepsilon)^2 \cdot \log(1/\delta)$ rows of A, each proportional to w_i (recall weighted random sampling in a stream). Let R be the "stacking" of these rows.

These t rows will jointly act in place of V_k^T . However since V was orthogonal, then the columns $v_i, v_j \in V_k$ were orthogonal. This is not the case for R, we need to orthogonalize R. Let $\Pi_R = R^T (RR^T)^{-1} R$ be the projection matrix for R, so that $A_R = A\Pi_R$ describes the *projection* of A onto the subspace of the directions spanned by R. Now

$$||A - A\Pi_R||_F \le ||A - A_k||_F + \varepsilon ||A||_F$$

with probability at least $1 - \delta$.

Why did we not just choose the t rows of A with the largest w_j values?
 Some may point along the same "direction" and would be repetitive. This should remind you of the choice to run k-means++ versus the Gonzalez algorithm for greedy point-assignment clustering.

Instructor: Jeff M. Phillips, U. of Utah

• Why did we not factor out the directions we already picked?
We could, but this allows us to run this in a streaming setting. (See next approach)

- But AΠ_R could be rank t, can we get it rank k ≪ t?
 Yes, you can take its best rank k approximation [Π_RA]_k and about the same bounds hold, you may need to increase t slightly.
- Can we get a better error bound? Yes. First take SVD [U, S, V] = svd(A) and let U_k be the top k left singular vectors. Let $U_k(i)$ be the ith row of U_k . Now the leverage score of data point a_i is $\ell_i = \|U_k(i)\|^2$. Using the leverage scores as weights $w_i = \ell_i$ allows one to achieve stronger bounds

$$||A - A\Pi_R||_F \le (1 + \varepsilon)||A - A_k||_F.$$

But this requires us to first take the SVD (or other time-consuming procedures), so its is harder to do in a stream; although it is possible to get good enough approximations of the leverage scores. In many cases, these approaches do not seem to provide tangible benefits over the faster $||a_i||^2$ -weighted sampling.

Can we also sample columns this way?
 Yes. All tricks can be run on A^T the same way (in fact most of the literature talks about sampling columns instead of rows). And, both approaches can be combined. This is known as the CUR-decomposition of A.

A significant downside of these row sampling approaches is that the $(1/\varepsilon^2)$ coefficient can be quite large for a small error tolerance. If $\varepsilon=0.01$, meaning 1% error, then this part of the coefficient alone is $10{,}000$. In practice, the results may be better, but for guarantees, this may only work on very enormous matrices.

16.3 Projection and Count Sketch

Stronger guarantees can be obtained through random projection-based sketches (see next lecture). The starting point is a random projection matrix $S \in \mathbb{R}^{\ell \times n}$ that maps A to a $\ell \times d$ matrix B = SA. As with the random projections approach, each element $S_{i,j}$ of S is drawn iid $S_{i,j} \sim \mathbb{N} \cdot \sqrt{n/\ell}$ from a normal distribution. That is, each row s_i of S is a n-dimensional Gaussian random variable $\sim \mathcal{G}_n$, properly normalized.

Using $\ell \approx k/\varepsilon$ columns in S yields the rank-k approximation result, with a few linear algebra steps. First let $V \in \mathbb{R}^{d \times d}$ be the orthogonal basis defined by the right singular vectors of B = SA. Then let $[AV]_k$ be the best rank-k approximation of AV, found using the SVD. Ultimately, we have

$$||A - [AV]_k V^T||_F \le (1 + \varepsilon)||A - A_k||_F.$$

A factorized form (like the output of the SVD) can be computed for the product $[AV]_kV^T$ in time which scales linearly in the matrix size nd, and polynomially in $\ell=k/\varepsilon$.

Moreover, this preserves a stronger bound, called an oblivious subspace embedding, using $\ell \approx d/\varepsilon^2$ so, for all $x \in \mathbb{R}^d$

$$(1-\varepsilon) \le \frac{\|Ax\|}{\|Bx\|} \le (1+\varepsilon).$$

This is a very strong bound that also ensures that given a matrix A of d-dimensional explanatory variables, and a vector b of dependent variables, then the result of linear regression on SA and Sb provides a $(1 \pm \varepsilon)$ approximation to the result on the full A and b.

Count Sketch Hashing for Sparse Matrices

By increasing the size of the sketch to $\ell \approx k^2 + k/\varepsilon$ for the rank-k approximation result, or to $\ell \approx d^2/\varepsilon^2$ for the oblivious subspace embedding result, then a faster *count sketch* based approach can be used. In this approach S has each column S_j as all 0s, except for one randomly chosen entry (this can be viewed as a hash to a row of S) that is either S1 or S2 at random. This works just like a count sketch but for matrices.

The runtime of these count-sketch approaches becomes proportional to nnz(A), the number of non-zeros in A. For very sparse data, such as those generated from bag-of-word approaches, this is as fast as only reading the few relevant entries of the data. Each row is now randomly accumulated onto one row of the output sketch B instead of onto all rows as when using the Gaussian random variables approach. However, this sketch B is not as interpretable as the row selection methods, and in practice for the same space often works a bit worse (due to extra factors necessary for the concentration of measure bounds to kick in) than the Frequent Directions approach.

Instructor: Jeff M. Phillips, U. of Utah