

Singular Value Decomposition

Szymon Rusinkiewicz

COS 302, Fall 2020



Singular Value Decomposition (SVD)

- Matrix *decomposition* that reveals structure
- Useful for:
 - Inverses, pseudoinverses
 - Stable least-squares, even for unconstrained problems
 - Matrix similarity and approximation
 - Dimensionality reduction and PCA
 - Orthogonalization
 - Constrained least squares and multidimensional scaling



Let's look at
motivaion for these

Condition Number

- $\text{cond}(\mathbf{A})$ is function of \mathbf{A}
- $\text{cond}(\mathbf{A}) \geq 1$, bigger is bad
- Measures how change in input propagates to output:

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(\mathbf{A}) \frac{\|\Delta \mathbf{A}\|}{\|\mathbf{A}\|}$$

- E.g., if $\text{cond}(\mathbf{A}) = 451$ then can lose $\log(451) = 2.65$ digits of accuracy in x , compared to precision of \mathbf{A}
- For matrices with real eigenvalues, $\text{cond}(\mathbf{A}) = |\lambda_{\max}| / |\lambda_{\min}|$

Normal Equations are Bad

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}$$

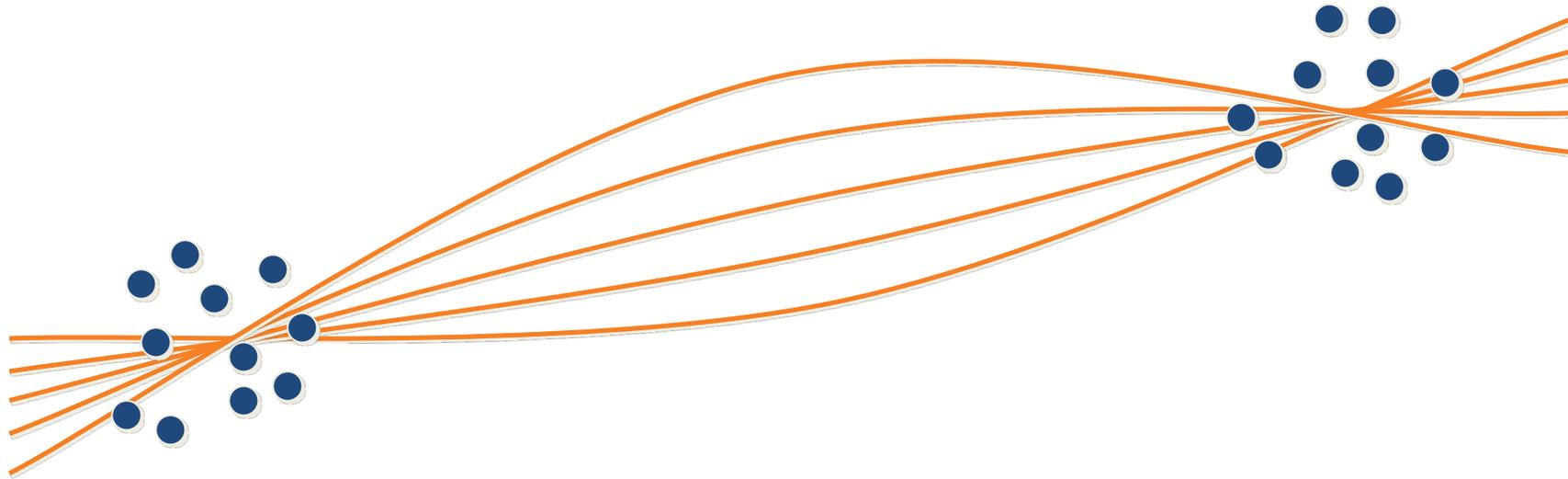
- Least squares using normal equations involves solving $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$
- $\text{cond}(\mathbf{A}^T \mathbf{A}) = [\text{cond}(\mathbf{A})]^2$
- E.g., if $\text{cond}(\mathbf{A}) = 451$ then can lose $\log(451^2) = 5.3$ digits of accuracy, compared to precision of \mathbf{A}

Underconstrained Least Squares

- What if you have fewer data points than parameters in your function?
 - Intuitively, can't do standard least squares
 - Solution takes the form $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$
 - When \mathbf{A} has more columns than rows, $\mathbf{A}^T \mathbf{A}$ is singular: can't take its inverse, etc.

Underconstrained Least Squares

- More subtle version: more data points than unknowns, but data poorly constrains function
- Example: fitting to $y = ax^2 + bx + c$



Underconstrained Least Squares

- Problem: if problem very close to singular, roundoff error can have a huge effect
 - Even on “well-determined” values!
- Can detect this:
 - Uncertainty proportional to covariance $\mathbf{C} = (\mathbf{A}^T\mathbf{A})^{-1}$
 - In other words, unstable if $\mathbf{A}^T\mathbf{A}$ has small values
 - More precisely, care if $\mathbf{x}^T(\mathbf{A}^T\mathbf{A})\mathbf{x}$ is small for any \mathbf{x}
- Idea: if part of solution unstable, set answer to 0
 - Avoid corrupting good parts of answer

Singular Value Decomposition (SVD)

- Handy mathematical technique that has application to many problems
- Given any $m \times n$ matrix A , algorithm to find matrices U , V , and W with:

$$A = UWV^T$$

U is $m \times m$ and **orthonormal**

W is $m \times n$ and **zero except main diagonal**

V is $n \times n$ and **orthonormal**

- Won't derive algorithm – treat as black box (e.g., `numpy.linalg.svd`)

“Full” SVD

$$\begin{bmatrix} & & \\ & \mathbf{A} & \\ & & \end{bmatrix} = \begin{bmatrix} & & \\ & \mathbf{U} & \\ & & \end{bmatrix} \begin{bmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} & & \\ & \mathbf{V} & \\ & & \end{bmatrix}^T$$

$m \times n$ $m \times m$ $m \times n$ $n \times n$

```
u,w,vt = numpy.linalg.svd(a)
```

SVD

- Handwavy explanation: rotate to a basis where all the scaling and stretching of \mathbf{A} is along coordinate axes
 - Should remind you of eigendecomposition (which would have $\mathbf{U} = \mathbf{V}$)
- The w_i are called the **singular values** of \mathbf{A}
- If \mathbf{A} is singular, some of the w_i will be 0
- In general $\text{rank}(\mathbf{A}) = \text{number of nonzero } w_i$
- SVD is mostly unique (up to permutation of singular values, or if some w_i are equal)
 - The w_i are conventionally returned in sorted order, largest to smallest

Singular Value Decomposition (SVD)

- If $m > n$, only n nonzero rows in \mathbf{W} , many useless columns in \mathbf{U}
- If $n > m$, only m nonzero columns in \mathbf{W} , many useless columns in \mathbf{V}
 - Define “compact” or “reduced” versions that omit all those zeroes

“Compact” SVD, if $m > n$

$$\begin{array}{c} \left[\begin{array}{c} A \\ \end{array} \right] \\ m \times n \end{array} = \begin{array}{c} \left[\begin{array}{c} U \\ \end{array} \right] \\ m \times n \end{array} \begin{array}{c} \left[\begin{array}{ccc} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{array} \right] \\ n \times n \end{array} \begin{array}{c} \left[\begin{array}{c} V \\ \end{array} \right]^T \\ n \times n \end{array}$$

```
u,w,vt = numpy.linalg.svd(a, full_matrices=False)
```

“Compact” SVD, if $n > m$

$$\begin{array}{c} \left[\begin{array}{c} \mathbf{A} \\ \end{array} \right] = \left[\begin{array}{c} \mathbf{U} \\ \end{array} \right] \left[\begin{array}{ccc} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_m \end{array} \right] \left[\begin{array}{c} \mathbf{V} \\ \end{array} \right]^T \\ \begin{array}{c} m \times n \\ \\ m \times m \\ \\ m \times m \\ \\ m \times n \end{array} \end{array}$$

```
u,w,vt = numpy.linalg.svd(a, full_matrices=False)
```

SVD and Inverses

- Why is SVD so useful?
- Application #1: inverses
- $\mathbf{A}^{-1} = (\mathbf{V}^T)^{-1} \mathbf{W}^{-1} \mathbf{U}^{-1} = \mathbf{V} \mathbf{W}^{-1} \mathbf{U}^T$
 - Using fact that inverse = transpose for orthogonal matrices
 - Since \mathbf{W} is diagonal, \mathbf{W}^{-1} also diagonal with reciprocals of entries of \mathbf{W}

SVD and the Pseudoinverse

- $\mathbf{A}^{-1} = (\mathbf{V}^T)^{-1} \mathbf{W}^{-1} \mathbf{U}^{-1} = \mathbf{V} \mathbf{W}^{-1} \mathbf{U}^T$
- This fails when some w_i are 0
 - It's *supposed* to fail – singular matrix
 - Happens when rectangular \mathbf{A} is **rank deficient**
- Pseudoinverse \mathbf{A}^+ : if $w_i = 0$, set $1/w_i$ to 0 (!!)
 - “Closest” matrix to inverse
 - Defined for all (even non-square, singular, etc.) matrices
 - Equal to $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ if $\mathbf{A}^T \mathbf{A}$ invertible

SVD and Least Squares

- Solving $\mathbf{Ax}=\mathbf{b}$ by least squares:
- $\mathbf{A}^T\mathbf{Ax} = \mathbf{A}^T\mathbf{b} \rightarrow \mathbf{x} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$
- Replace with \mathbf{A}^+ : $\mathbf{x} = \mathbf{A}^+\mathbf{b}$
 - Compute pseudoinverse using SVD
- Lets you see if data is singular ($< n$ nonzero singular values)
- Singular values tell you how stable the solution will be
 - Condition number = ratio of largest to smallest singular values
- For better stability, set $1/w_i$ to 0 if w_i is small (even if not exactly 0)
 - Accuracy / stability tradeoff? Not if that component was underconstrained...

SVD and Matrix Similarity

- One common definition for the norm of a matrix is the Frobenius norm:

$$\|\mathbf{A}\|_F = \sum_i \sum_j a_{ij}^2$$

- Frobenius norm can be computed from SVD

$$\|\mathbf{A}\|_F = \sum_i w_i^2$$

- Euclidean (spectral) norm can also be computed:

$$\|\mathbf{A}\|_2 = \{\max |\lambda| : \lambda \in \sigma(\mathbf{A})\}$$

- So changes to a matrix can be evaluated by looking at changes to singular values

SVD and Matrix Similarity

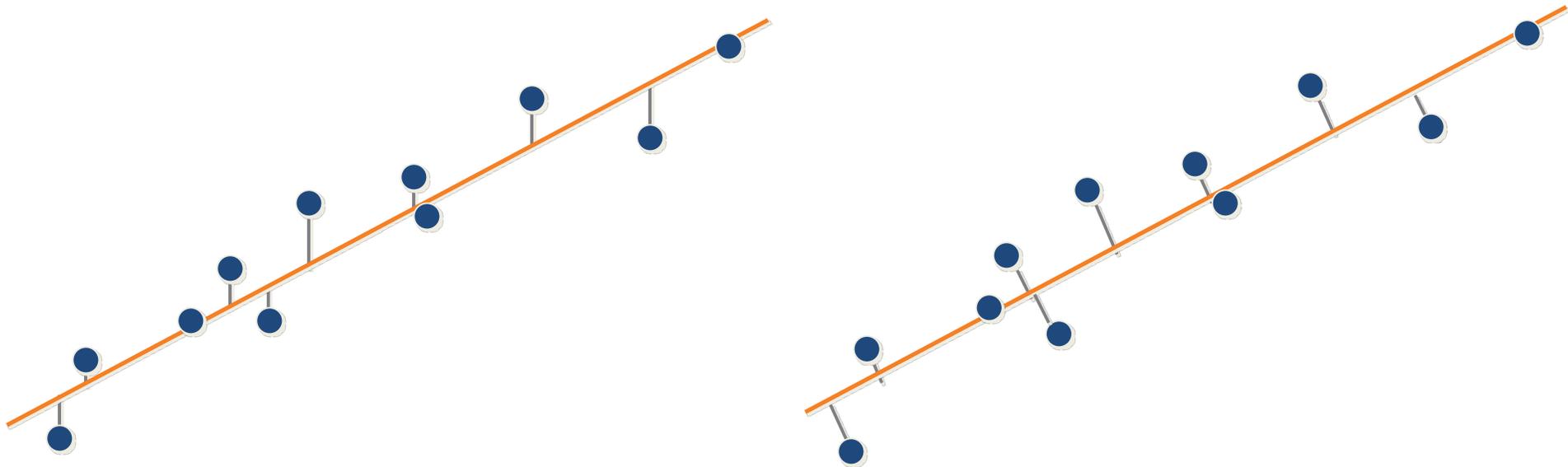
- Suppose you want to find best rank- k approximation to A
- Answer: set all but the largest k singular values to zero
- Can form compact representation by eliminating columns of U and V corresponding to zeroed w_i

SVD and Orthogonalization

- U and V are orthonormal, all stretching and scaling in W
- The matrix UV^T is the “closest” orthonormal matrix to A
 - Yet another useful application of the matrix-approximation properties of SVD
 - Much more stable numerically than Gram-Schmidt orthogonalization

Total Least Squares

- One final least squares application
- Fitting a line: vertical vs. perpendicular error



Total Least Squares

- Distance from point to line:

$$d_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \cdot \vec{n} - a$$

where \mathbf{n} is normal vector to line, a is a constant

- Minimize:

$$\chi^2 = \sum_i d_i^2 = \sum_i \left[\begin{pmatrix} x_i \\ y_i \end{pmatrix} \cdot \vec{n} - a \right]^2$$

Total Least Squares

- First, let's pretend we know \mathbf{n} , solve for a

$$\chi^2 = \sum_i \left[\begin{pmatrix} x_i \\ y_i \end{pmatrix} \cdot \vec{n} - a \right]^2$$

$$a = \frac{1}{m} \sum_i \begin{pmatrix} x_i \\ y_i \end{pmatrix} \cdot \vec{n}$$

- Then

$$d_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix} \cdot \vec{n} - a = \begin{pmatrix} x_i - \frac{\Sigma x_i}{m} \\ y_i - \frac{\Sigma y_i}{m} \end{pmatrix} \cdot \vec{n}$$

Total Least Squares

- So, let's define

$$\begin{pmatrix} \tilde{x}_i \\ \tilde{y}_i \end{pmatrix} = \begin{pmatrix} x_i - \frac{\Sigma x_i}{m} \\ y_i - \frac{\Sigma y_i}{m} \end{pmatrix}$$

and minimize

$$\sum_i \left[\begin{pmatrix} \tilde{x}_i \\ \tilde{y}_i \end{pmatrix} \cdot \vec{n} \right]^2$$

Total Least Squares

- Write as linear system

$$\begin{pmatrix} \tilde{x}_1 & \tilde{y}_1 \\ \tilde{x}_2 & \tilde{y}_2 \\ \tilde{x}_3 & \tilde{y}_3 \\ \vdots & \vdots \end{pmatrix} \begin{pmatrix} n_x \\ n_y \end{pmatrix} = \vec{0}$$

- Have $\mathbf{A}\mathbf{n}=\mathbf{0}$
 - Problem: lots of \mathbf{n} are solutions, including $\mathbf{n}=\mathbf{0}$
 - Standard least squares will, in fact, return $\mathbf{n}=\mathbf{0}$

Constrained Optimization

- Solution: constrain \mathbf{n} to be unit length
- So, try to minimize $\|\mathbf{A}\mathbf{n}\|^2$ subject to $\|\mathbf{n}\|^2 = 1$

$$\|\mathbf{A}\vec{n}\|^2 = (\mathbf{A}\vec{n})^T (\mathbf{A}\vec{n}) = \vec{n}^T \mathbf{A}^T \mathbf{A} \vec{n}$$

- Expand in eigenvectors \mathbf{e}_i of $\mathbf{A}^T \mathbf{A}$:

$$\vec{n} = \mu_1 \mathbf{e}_1 + \mu_2 \mathbf{e}_2$$

$$\vec{n}^T (\mathbf{A}^T \mathbf{A}) \vec{n} = \lambda_1 \mu_1^2 + \lambda_2 \mu_2^2$$

$$\|\vec{n}\|^2 = \mu_1^2 + \mu_2^2$$

where the λ_i are eigenvalues of $\mathbf{A}^T \mathbf{A}$

Constrained Optimization

- To minimize $\lambda_1\mu_1^2 + \lambda_2\mu_2^2$ subject to $\mu_1^2 + \mu_2^2 = 1$
set $\mu_{\min} = 1$, all other $\mu_i = 0$
- That is, \mathbf{n} is eigenvector of $\mathbf{A}^T\mathbf{A}$ with
the smallest corresponding eigenvalue

SVD and Eigenvectors

- Let $\mathbf{A} = \mathbf{U}\mathbf{W}\mathbf{V}^T$, and let \mathbf{x}_i be i^{th} column of \mathbf{V}

- Consider $\mathbf{A}^T\mathbf{A} \mathbf{x}_i$:

$$\mathbf{A}^T\mathbf{A} \mathbf{x}_i = \mathbf{V}\mathbf{W}^T\mathbf{U}^T\mathbf{U}\mathbf{W}\mathbf{V}^T \mathbf{x}_i = \mathbf{V}\mathbf{W}^2\mathbf{V}^T \mathbf{x}_i = \mathbf{V}\mathbf{W}^2 \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{V} \begin{pmatrix} 0 \\ \vdots \\ w_i^2 \\ \vdots \\ 0 \end{pmatrix} = w_i^2 \mathbf{x}_i$$

- So elements of \mathbf{W} are sqrt(eigenvalues) and columns of \mathbf{V} are eigenvectors of $\mathbf{A}^T\mathbf{A}$

Constrained Optimization

- To minimize $\lambda_1\mu_1^2 + \lambda_2\mu_2^2$ subject to $\mu_1^2 + \mu_2^2 = 1$
set $\mu_{\min} = 1$, all other $\mu_i = 0$
- That is, \mathbf{n} is eigenvector of $\mathbf{A}^T\mathbf{A}$ with the smallest corresponding eigenvalue
- That is, \mathbf{n} is column of \mathbf{V} corresponding to smallest singular value
 - Provides a solution to the total least squares problem
 - Also very related to PCA – next time