Unsupervised Learning

Spectral methods

- Eigenvector/singular vector decomposition (SVD)
- PCA, CCA

Reconstruction methods

• PCA, ICA, auto-encoders

Clustering and Probabilistic methods

- K-means
- Gaussian mixtures
- Latent Dirichlet Allocation (LDA)

SVD

Learning objectives SVD and 'thin SVD"

Matrix norms Generalized inverse Lyle Ungar

Eigenvectors (review)

- **A** $\mathbf{v}_i = \lambda_i \mathbf{v}_i$
- Eigen-decomposition of a symmetric matrix A (n x n)
 - A = VDV^T
- V: orthonormal, $V^T V = I (n \times n)$
 - Columns of V are the eigenvectors of A
- **D: diagonal** (n x n)
 - Diagonal elements of D are the eigenvalues of A
 - All non-negative if $\mathbf{A} = \mathbf{X}^T \mathbf{X}$
 - Reported in *decreasing* order of magnitude down the diagonal

We don't compute eigenvectors

- What symmetric matrix have we seen?
- In practice we rarely compute eigenvectors
 - Why not?

Singular Value Decomposition

- Singular value decomposition of matrix X (n x p)
 - X = UDV[⊤]
- U: orthonormal, $U^T U = I (n \times n)$
 - Columns of **U** are the *left singular vectors of* **X**
- **D: diagonal** (n x p)
 - Diagonal elements of **D** are the singular values of **X**
- ◆ V: orthonormal, V^TV = I (p x p)
 - Columns of **V** are the right singular vectors of **X**

SVD

Singular value decomposition of X: $X = UDV^{T}$



Let k = min(n,p). Then: $\mathbf{X} = \sum_{i=1}^{k} D_{ii} \boldsymbol{u}_{i} \boldsymbol{v}_{i}^{T}$

Since all $\boldsymbol{u}_i, \boldsymbol{v}_i$ are unit vectors, the importance of the i'th term in the sum is determined by the size of D_{ii} .

Review Questions

• $\mathbf{X}_{n^*p} = \mathbf{U} \mathbf{D} \mathbf{V}^T$

• What are the dimensions of U D and V?

- ♦ What are the eigenvectors of X^TX?
- ♦ What are the eigenvalues of X^TX?

Thin SVD – pick a smaller k

Singular value decomposition of X: $X = UDV^{T}$



Let k = min(n,p). Then: $\mathbf{X} = \sum_{i=1}^{k} D_{ii} \boldsymbol{u}_{i} \boldsymbol{v}_{i}^{T}$

Since all u_i , v_i are unit vectors, the importance of the i'th term in the sum is determined by the size of D_{ii} .

SVD and eigenvalues/eigenvectors

$X = UDV^T$, $X^T X = V(D^T D)V^T$

The columns v_1 ... v_p of V are the *eigenvectors* of the covariance matrix $X^T X$. Hence we can write

$$\boldsymbol{X}^{T}\boldsymbol{X} = \sum_{i=1}^{p} (D_{ii})^{2} \boldsymbol{v}_{i} \boldsymbol{v}_{i}^{T}$$

From before:

$$\boldsymbol{X} = \sum_{i=1}^{k} D_{ii} \boldsymbol{u}_{i} \boldsymbol{v}_{i}^{T}$$

k = min(n,p).

 D_{ii} are singular values of X, $(D_{ii})^2$ are eigenvalues of $X^T X$

Frobenius norm

Here: **A** is an arbitrary *m* x *n* matrix, what we often call **X**

How to measure the size of a matrix?

$$\|A\|_{ ext{F}} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{ ext{trace}(A^{\dagger}A)} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)}$$

Where σ_i are the singular values of A.
One can also use an L₁ norm ||A||₁ = ||σ||₁

Generalized Inverses

- Linear regression estimates w in y = Xw
- This uses a pseudo-inverse ("Moore-Penrose inverse")
 X⁺ of X, so
 - *w* = *X*[≁]*y*
- Thus far, we have done this by
 - $\mathbf{X}^{+} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}$

Generalized Inverses

We can also compute inverses using SVD
The idea:

 $X^{+} = (U D^{-1} V^{T})^{T} = V (D^{-1})^{T} U^{T}$

 You can't take the inverse of a rectangular matrix, but we can approximate it using the thin SVD
 X⁺ ~ V_k D_k⁻¹ U_k^T

Pseudo-inverse of X = U D V^T

- ♦ What are the dimensions of X⁺ = V D⁻¹ U^T
 ♦ What is X X_k⁺
 - $X X^+ = U D V^T V D^{-1} U^T$

Power Method

Power method for a square matrix A

- Write any $\mathbf{x} = \Sigma_i z_i \mathbf{v}_i$ where $z_i = \mathbf{v}_i^T \mathbf{x}$
- Then $Ax = A \Sigma_i z_i \mathbf{v}_i = \Sigma_i z_i A \mathbf{v}_i = \Sigma_i z_i \lambda_i \mathbf{v}_i$
- So AAAAx = A^4x = = $\Sigma_i z_i \lambda_i^4 \mathbf{v}_i$

Find the largest eigenvalue/eigenvector

• Project it off from x and repeat

• $\mathbf{x} := \mathbf{x} - (\mathbf{v}_1^T \mathbf{x}) \mathbf{v}_1$

Fast 'Randomized' SVD

- Generalizes the power method
- Input:
 - matrix A of size n × p,
 - the desired hidden state dimension k,
 - the number of "extra" singular vectors, I

 Simultaneously find all the largest singular values/vectors by alternately left and right multiplying by A

You are not required to know this

Randomized SVD for any matrix A

- 1. Generate a $(k + l) \times n$ random matrix Ω
- 2. Find the SVD $U_1D_1V_1^T$ of ΩA , and keep the k + l components of V_1 with the largest singular values
- 3. Find the SVD $U_2 D_2 V_2^T$ of AV_1 , and keep the 'largest' k + l components of U_2
- 4. Find the SVD $U_3D_3V_{final}^T$ of U_2^TA , and keep the 'largest' k components of V_{final}
- 5. Find the SVD $U_{final}D_4V_4^T$ of AV_{final} and keep the 'largest' k components of U_{final}

Output: The left and right singular vectors U_{final} , V_{final}^T **You are not required to know this**

What you should know

- Eigenvalues/vectors & singular values/vectors
- Eigenvectors as a basis
- Thin SVD
- Frobenius norm
- Pseudo ("Moore-Penrose") inverse
- Power method

To think about:

- What is an efficient way to do linear regression?
 - $\mathbf{w} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}}\mathbf{y}$
 - How does it scale with n and p?