UNSUPERVISED LEARNING

OUTLIERS APPLICATIONS

Outliers

- Values outside the normal range
- Statistically, define in terms of deviation from mean or median
- Gaussian distribution number of standard deviations from mean
- Box and whisker plots outer and inner fences based on median, interquartile (IQR) value





Outliers and clustering

- Outliers are points that lie outside natural clusters
- K Means far away from all centroids
 - But outliers can distort the clustering process
- Density based clustering not connected to any core point
 - But density is applied uniformly



Outliers and density

Is the density of a point Similar to the density of to

- An outlier is less dense than its nearest neighbours
- But difference in density may be local
- A distance metric to eliminate o₂ could make all of C₁ outliers
- C₁ has 400 points, C₂ has 100 points
- Larger distance would make all of C₂ outliers with respect to C₁





Outliers and density

- For clustering, we defined a radius *Eps* and looked for *MinPts* heighbours within that ball
- Instead, fix *MinPts* and find smallest ball with that many neighbours
- Compare *radius(p)* with radius of its neighours
- A is an outlier because its radius is much more than that of its neighbours





Outliers and density

• Local outlier factor LOF(p)



- The smaller this ratio, the more likely that *p* is an outlier
- Comparison is local to neighbourhood, so this can deal with different densities across range of data







Outlin detection using mischire of 3 **X**2 Gammians Numeric data $(\mu_{1},\sigma_{1}),(\mu_{2},\sigma_{2}),(\mu_{3},\sigma_{3})$ > 30 for all Gams: as Owthing

Semi-supervised learning

(MNIST)

603

- Labelling training data is a bottleneck of supervised learning
- Handwritten digits 0,1,...,9
 - 1797 images
- Standard logistic regression model has 96.9% accuracy
- Suppose we take 50 random samples as training set
- Logistic regression gives 83.3%



Semi-supervised learning

- Instead of 50 random samples, 50 clusters using K means
- Use image nearest to each centroid as training set
 - 50 representative images
- Logistic regression accuracy jumps to 92.2%





Semi-supervised learning

- Propagate representative image label to entire cluster
- Logistic regression improves to 93.3%
- Propagage representive image label to only 20% items closest to centroid
- Logistic regression improves to 94%
- Only 50 actual labels used, about 5 per class!



- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours





- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours
- With 10 clusters, not much change







- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours
- With 10 clusters, not much change
- Same with 8







- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours
- With 10 clusters, not much change
- Same with 8
- At 6 colours, ladybug red goes







- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours
- With 10 clusters, not much change
- Same with 8
- At 6 colours, ladybug red goes
- 4 colours







- An image is a matrix of pixels
- Each pixel has (R,G,B) values
- K means clustering on these values merges colours
- With 10 clusters, not much change
- Same with 8
- At 6 colours, ladybug red goes
- 4 colours
- Finally 2 colours, flower and rest





Dimensionality reduction

Madhavan Mukund https://www.cmi.ac.in/~madhavan

Data Mining and Machine Learning August-December 2020

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Dimensionality reduction

Principal Component Anaylsis — transform *d*-dimensional input to *k*-dimensional input, preserving essential features

э

▶ ∢ ⊒

Dimensionality reduction

- Principal Component Anaylsis transform *d*-dimensional input to *k*-dimensional input, preserving essential features
- Example: PCA projection of blue points in 3D to black points in 2D



Dimensionality reduction

Unsupervised preprocessing technique — may make later steps easier, like simplifying classification boundaries

3

→ < ∃→

< A

Dimensionality reduction

- Unsupervised preprocessing technique may make later steps easier, like simplifying classification boundaries
- Swiss roll dataset: dimensionality reduction helps



Dimensionality reduction

- Unsupervised preprocessing technique may make later steps easier, like simplifying classification boundaries
- Swiss roll dataset: dimensionality reduction helps



Swiss roll dataset: dimensionality reduction does not help



Input matrix M, dimensions $n \times d$

Rows are items, columns are features



э

- Input matrix M, dimensions $n \times d$ \longrightarrow $n \times k$
 - Rows are items, columns are features
- Decompose M as $U D V^{\top}$
 - **D** is a $k \times k$ diagonal matrix, positive real entries
 - U is $n \times k$, V is $d \times k$
 - Columns of U, V are orthonormal unit vectors, mutually orthogonal

- Input matrix M, dimensions $n \times d$
 - Rows are items, columns are features
- Decompose M as UDV^{\top}
 - **D** is a $k \times k$ diagonal matrix, positive real entries
 - U is $n \times k$, V is $d \times k$
 - Columns of U, V are orthonormal unit vectors, mutually orthogonal
- Interpretation
 - Columns of V correspond to new abstract features



- Input matrix M, dimensions $n \times d$
 - Rows are items, columns are featurest
- Decompose M as UDV^{\top}
 - **D** is a $k \times k$ diagonal matrix, positive real entries
 - U is $n \times k$, V is $d \times k$
 - Columns of U, V are orthonormal unit vectors, mutually orthogonal
- Interpretation
 - Columns of V correspond to new abstract features
 - Rows of U describe decomposition of terms across features

- Input matrix M, dimensions $n \times d$
 - Rows are items, columns are features
- Decompose M as UDV^{\top}
 - **D** is a $k \times k$ diagonal matrix, positive real entries
 - U is $n \times k$, V is $d \times k$
 - Columns of U, V are orthonormal unit vectors, mutually orthogonal
- Interpretation
 - Columns of V correspond to new abstract features
 - Rows of U describe decomposition of terms across features
 - For columns u_i of U and v_i of V, $u_i \cdot v_i^{\top}$ is an $n \times d$ matrix, like M





- Input matrix M, dimensions $n \times d$
 - Rows are items, columns are features
- Decompose M as UDV^{\top}
 - **D** is a $k \times k$ diagonal matrix, positive real entries
 - U is $n \times k$, V is $d \times k$
 - Columns of U, V are orthonormal unit vectors, mutually orthogonal
- Interpretation
 - Columns of V correspond to new abstract features
 - Rows of U describe decomposition of terms across features
 - For columns u_i of U and v_i of V, $u_i \cdot v_i^{\top}$ is an $n \times d$ matrix, like M
 - **u**_i · \mathbf{v}_i^{\top} describes components of rows of M along direction \mathbf{v}_i

Unit vectors passing through the origin

3

▶ < ∃ ▶</p>

- Unit vectors passing through the origin
- Want to find "best" k singular vectors to represent feature space

э

< E

- Unit vectors passing through the origin
- Want to find "best" k singular vectors to represent feature space
- Suppose we project $a_i = (a_{i1}, a_{i2}, \dots, a_{id})$ onto v through origin



- Unit vectors passing through the origin
- Want to find "best" k singular vectors to represent feature space
- Suppose we project $a_i = (a_{i1}, a_{i2}, \dots, a_{id})$ onto v through origin



• Minimizing distance of a_i from v is equivalent to maximizing the projection of a_i onto v

Madhavan Muk	kund
--------------	------

- Unit vectors passing through the origin
- Want to find "best" k singular vectors to represent feature space
- Suppose we project $a_i = (a_{i1}, a_{i2}, \dots, a_{id})$ onto v through origin



- Minimizing distance of a_i from v is equivalent to maximizing the projection of a_i onto v
- Length of the projection is $a_i \cdot v$

Madhavan Mukund

• Sum of squares of lengths of projections of all rows in M onto $\mathbf{v} - |M\mathbf{v}|^2$

- 2

- Sum of squares of lengths of projections of all rows in M onto $\mathbf{v} |M\mathbf{v}|^2$
- First singular vector unit vector through origin that maximizes the sum of projections of all rows in M

 $oldsymbol{
u}_1 = rg\max_{|oldsymbol{
u}|=1} |Moldsymbol{
u}|$

э

▶ < ∃ ▶</p>

- Sum of squares of lengths of projections of all rows in M onto $\mathbf{v} |M\mathbf{v}|^2$
- First singular vector unit vector through origin that maximizes the sum of projections of all rows in M

$$oldsymbol{
u}_1 = rg\max_{|oldsymbol{
u}|=1} |Moldsymbol{
u}|$$

Second singular vector — unit vector through origin, perpendicular to v_1 , that maximizes the sum of projections of all rows in M

$$oldsymbol{v}_2 = rg\max_{oldsymbol{v} \perp oldsymbol{v}_1; \ |oldsymbol{v}| = 1} |Moldsymbol{v}|$$

- Sum of squares of lengths of projections of all rows in M onto $\mathbf{v} |M\mathbf{v}|^2$
- First singular vector unit vector through origin that maximizes the sum of projections of all rows in M

$$v_1 = \arg \max_{|v|=1} |Mv|$$
 $|Mv_1| > |Mv_2| > |Mv_3|$

Second singular vector — unit vector through origin, perpendicular to v_1 , that maximizes the sum of projections of all rows in M

$$oldsymbol{v}_2 = rg\max_{oldsymbol{v} \perp oldsymbol{v}_1; \ |oldsymbol{v}| = 1} |Moldsymbol{v}|$$

Third singular vector — unit vector through origin, perpendicular to v_1 , v_2 , that maximizes the sum of projections of all rows in M

$$\mathbf{v}_3 = \arg \max_{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2; \ |\mathbf{v}|=1} |M\mathbf{v}|$$

• With each singular vector \mathbf{v}_j , associated singular value is $\sigma_j = |M\mathbf{v}_j|$

< □ > < 向

3

- With each singular vector \mathbf{v}_j , associated singular value is $\sigma_j = |M\mathbf{v}_j|$
- Repeat *r* times till $\max_{\boldsymbol{\nu} \perp \boldsymbol{\nu}_1, \boldsymbol{\nu}_2, \dots, \boldsymbol{\nu}_r; \ |\boldsymbol{\nu}|=1} |M\boldsymbol{\nu}| = 0$
 - r turns out to be the rank of M
 - Vectors $\{v_1, v_2, \dots, v_r\}$ are orthonormal right singular vectors

- With each singular vector \mathbf{v}_i , associated singular value is $\sigma_i = |M\mathbf{v}_i|$
- Repeat r times till $\max_{\boldsymbol{\nu} \perp \boldsymbol{\nu}_1, \boldsymbol{\nu}_2, \dots, \boldsymbol{\nu}_r; \ |\boldsymbol{\nu}| = 1} |M\boldsymbol{\nu}| = 0$
 - r turns out to be the rank of M
 - Vectors $\{v_1, v_2, \dots, v_r\}$ are orthonormal right singular vectors
- Our greedy strategy provably produces "best-fit" dimension r subspace for M
 - Dimension r subspace that maximizes content of M projected onto it

3

- With each singular vector \mathbf{v}_j , associated singular value is $\sigma_j = |M\mathbf{v}_j|$
- Repeat *r* times till $\max_{\boldsymbol{v} \perp \boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_r; \ |\boldsymbol{v}|=1} |M\boldsymbol{v}| = 0$
 - r turns out to be the rank of M
 - Vectors $\{v_1, v_2, \dots, v_r\}$ are orthonormal right singular vectors
- Our greedy strategy provably produces "best-fit" dimension r subspace for M
 - Dimension r subspace that maximizes content of M projected onto it
- Corresponding left singular vectors are given by $\boldsymbol{u}_i = \frac{1}{\sigma_i} M \boldsymbol{v}_i$
- Can show that $\{u_1, u_2, \dots, u_r\}$ are also orthonormal

- *M*, dimension $n \times d$, of rank *r* uniquely decomposes as $M = UDV^{\top}$
 - $V = [v_1 \ v_2 \ \cdots \ v_r]$ are the right singular vectors
 - D is a diagonal matrix with $D[i, i] = \sigma_i$, the singular values
 - $U = [u_1 \ u_2 \ \cdots \ u_r]$ are the left singular vectors



■ *M* has rank *r*, SVD gives rank *r* decomposition

3

() < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < ()

- *M* has rank *r*, SVD gives rank *r* decomposition
- Singular values are non-increasing $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_r$

3

▶ < ∃ ▶</p>

- M has rank r, SVD gives rank r decomposition
- Singular values are non-increasing $\sigma_1 > \sigma_2 > \cdots > \sigma_r$
- Suppose we retain only k largest ones
- We have
 - Matrix of first k right singular vectors $V_k = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_k]$.
 - Corresponding singular values $\sigma_1, \sigma_2, \ldots, \sigma_k$
 - Matrix of k left singular vectors $U_k = [u_1 \ u_2 \ \cdots \ u_k]$

3

- *M* has rank *r*, SVD gives rank *r* decomposition
- Singular values are non-increasing $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_r$
- Suppose we retain only k largest ones
- We have
 - Matrix of first k right singular vectors $V_k = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_k]$,
 - Corresponding singular values $\sigma_1, \sigma_2, \ldots, \sigma_k$
 - Matrix of k left singular vectors $U_k = [u_1 \ u_2 \ \cdots \ u_k]$
- Let D_k be the $k \times k$ diagonal matrix with entries $\sigma_1, \sigma_2, \ldots, \sigma_k$
- Then $U_k D_k V_k^{\top}$ is the best fit rank-k approximation of M

3

- *M* has rank *r*, SVD gives rank *r* decomposition
- Singular values are non-increasing $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_r$
- Suppose we retain only k largest ones
- We have
 - Matrix of first k right singular vectors $V_k = [v_1 \ v_2 \ \cdots \ v_k]$,
 - Corresponding singular values $\sigma_1, \sigma_2, \ldots, \sigma_k$
 - Matrix of k left singular vectors $U_k = [u_1 \ u_2 \ \cdots \ u_k]$
- Let D_k be the $k \times k$ diagonal matrix with entries $\sigma_1, \sigma_2, \ldots, \sigma_k$
- Then $U_k D_k V_k^{\top}$ is the best fit rank-k approximation of M
- In other words, by truncating the SVD, we can focus on k most significant features implicit in M

Summary

- Singular Value Decomposition (SVD) finds best fit k-dimensional subspace for any matrix M
- Principal Component Analysis uses SVD for dimensionality reduction
- Unsupervised technique often helps simplify the problem, but may not
- SVD/PCA can only compress features that have a linear relationship
- More general techniques based on neural networks autoencoders