

Dimensionality reduction

CISC 5800
Professor Daniel Leeds

Opening note on dimensional differences

Each dimension corresponds to a feature/measurement

Magnitude differences for each measurement (e.g., animals):

- x_1 – speed (mph) 0-100
- x_2 – weight (pounds) 10-1000
- x_3 – size (feet) 2-20



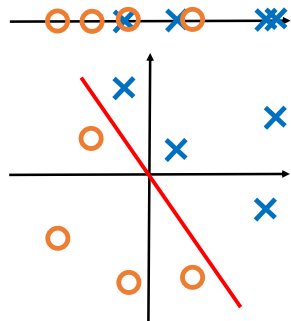
Problem for learning:

$$w_j \leftarrow w_j + \varepsilon x_j^i (y^i - g(w^T x^i)) - \frac{w_j}{\lambda}$$

Normalize: $r_1 = \frac{x_1 - \mu_1}{\sigma_1}$ or $r_1 = \frac{x_1 - \mu_1}{\max_i - \min_i}$

2

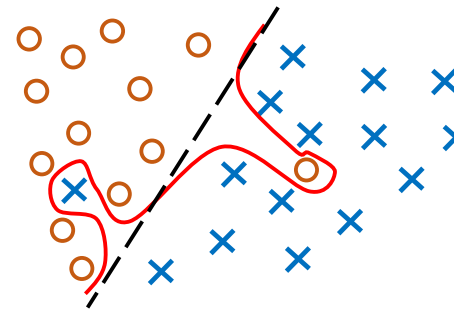
The benefits of extra dimensions



- Finds existing complex separations between classes

3

The risks of too-many dimensions

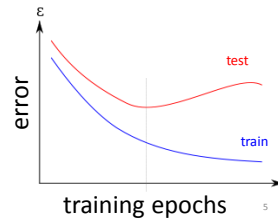


- High dimensions with kernels over-fit the outlier data
- Two dimensions ignore the outlier data

4

Training vs. testing

- **Training**: learn parameters from set of data in each class
- **Testing**: measure how often classifier correctly identifies new data
- More training reduces classifier error ϵ
 - More gradient ascent steps
 - More learned feature
- Too much training causes worse testing error – overfitting



Goal: High Performance, Few Parameters

- “Information criterion”: performance/parameter trade-off
- Variables to consider:
 - L likelihood of train data after learning
 - k number of parameters (e.g., number of features)
 - m number of points of training data
- Popular information criteria:
 - Akaike information criterion **AIC**: $\log(L) - k$
 - Bayesian information criterion **BIC**: $\log(L) - 0.5 k \log(m)$

6

Decreasing parameters

- Force parameter values to 0
 - L1 regularization
 - Support Vector selection
 - Feature selection/removal
- Consolidate feature space
 - Component analysis

8

Feature removal

- Start with feature set: $F = \{x_1, \dots, x_k\}$
 - Find classifier performance with set F : $\text{perform}(F)$
 - Loop
 - Find classifier performance for removing feature x_1, x_2, \dots, x_k :
 $\text{argmax}_i \text{perform}(F - x_i)$
 - Remove feature that causes least decrease in performance:
 $F = F - x_i$
- AIC**: $\log(L) - k$
BIC: $\log(L) - 0.5 k \log(m)$
- Repeat, using AIC or BIC as termination criterion

9

AIC testing: $\log(L)-k$

| Features | k (num features) | L (likelihood) | AIC |
|--------------------------------------|------------------|----------------|-------|
| F | 40 | 0.1 | -42.3 |
| F- $\{x_3\}$ | 39 | 0.03 | -41.5 |
| F- $\{x_3, x_{74}\}$ | 38 | 0.005 | -41.3 |
| F- $\{x_3, x_{24}, x_{32}\}$ | 37 | 0.001 | -40.9 |
| F- $\{x_3, x_{24}, x_{32}, x_{15}\}$ | 36 | 0.0001 | -41.2 |

10

Feature selection

AIC: $\log(L) - k$ **BIC:** $\log(L) - 0.5 k \log(m)$

- Find classifier performance for just set of 1 feature: $\text{argmax}_i \text{perform}(\{x_i\})$
- Add feature with highest performance: $F = \{x_i\}$
- Loop
 - Find classifier performance for adding one new feature: $\text{argmax}_i \text{perform}(F + \{x_i\})$
 - Add to F feature with highest performance increase: $F = F + \{x_i\}$

Repeat, using AIC or BIC as termination criterion

11

Capturing links between features

Document1 Document2 Document3

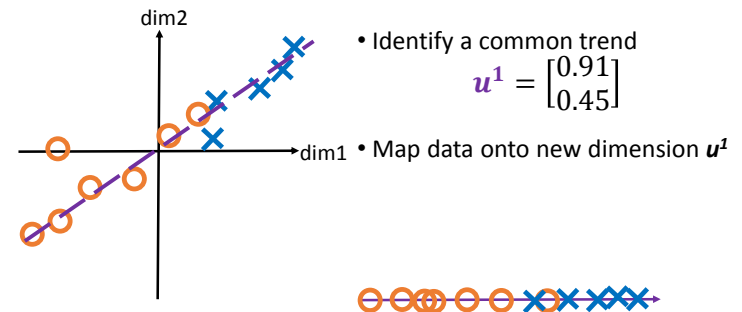
| | | | |
|--------|----|----|----|
| Wolf | 12 | 4 | 1 |
| Lion | 16 | 3 | 2 |
| Monkey | 5 | 11 | 4 |
| Sky | 7 | 3 | 14 |
| Tree | 2 | 8 | 5 |
| Cloud | 6 | 2 | 12 |
| ⋮ | ⋮ | ⋮ | ⋮ |

With large number of features,
some features x_j and x_k act similarly x_{wolf} & $x_{\text{lion}} \rightarrow u_{\text{predator}}$ x_{sky} & $x_{\text{cloud}} \rightarrow u_{\text{atmosphere}}$ Approximate $x^1 = \begin{bmatrix} x_1^1 \\ \vdots \\ x_N^1 \end{bmatrix}$ with $u^1 = \begin{bmatrix} u_1^1 \\ \vdots \\ u_{N'}^1 \end{bmatrix}$

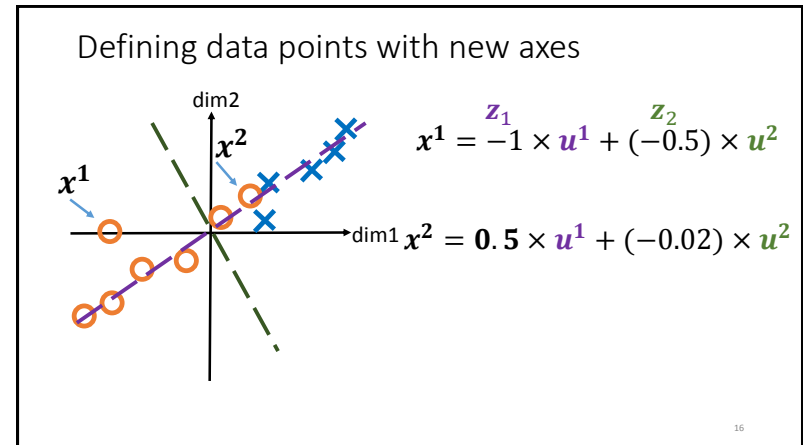
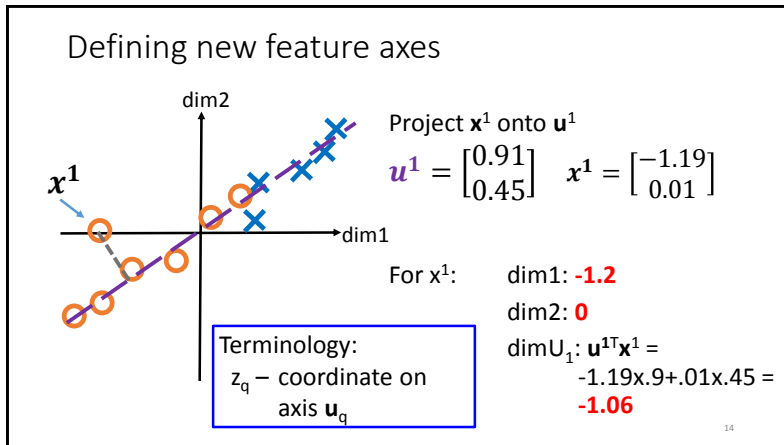
Automatically learn summary features

12

Defining new feature axes



13



Component analysis

Each data point \mathbf{x}^i in D can be reconstructed as sum of components \mathbf{u} :

- $\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$
- z_q^i is weight on q^{th} component to reconstruct data point \mathbf{x}^i

Image features

Image as grid of $n \times m$ pixels

Find representative component features as pixel patterns

Cartoon face example:

$$\mathbf{x}^1 \approx 1 \times \mathbf{u}^1 + 0 \times \mathbf{u}^2 + 1 \times \mathbf{u}^3 + 1 \times \mathbf{u}^4 + 0 \times \mathbf{u}^5$$

Add relevant face components
1000-pixel image becomes
6 co-efficients
Estimate is fairly close to
actual image

Component analysis

Each data point \mathbf{x}^i in D can be reconstructed as sum of components \mathbf{u} :

- $\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$
- z_q^i is weight on q^{th} component to reconstruct data point \mathbf{x}^i

Evaluating components

Components learned in order of descriptive power

Compute reconstruction error for all data by using first r components:

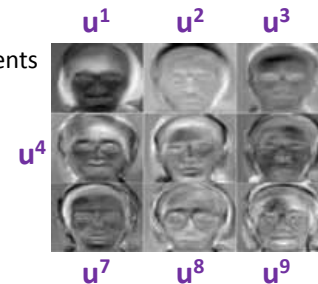
$$\text{error} = \sum_i \left(\sum_j \left(\mathbf{x}_j^i - \sum_{q=1}^r z_q^i \mathbf{u}_j^q \right)^2 \right)$$

Component analysis: examples

$$\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$$

“Eigenfaces” – learned from set of face images

\mathbf{u} : nine components



\mathbf{x}^4 : data reconstructed



Types of component analysis

Capture links between features as “components”

- Principal component analysis (PCA)
- Independent component analysis (ICA)
- Non-negative matrix factorization (NMF)

25

Principal component analysis (PCA)

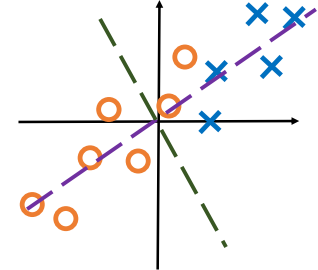
Describe every \mathbf{x}^i with small set of components $\mathbf{u}^{1:Q}$

Use same $\mathbf{u}^1, \dots, \mathbf{u}^T$ for all \mathbf{x}^i

All components orthogonal:

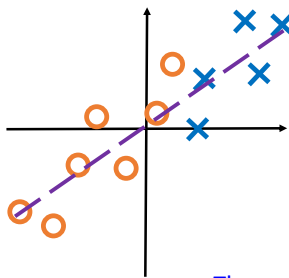
$$(\mathbf{u}^i)^T \mathbf{u}^j = 0 \quad \forall i \neq j$$

$$\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$$



26

Idea of learning in PCA



1. $D = \{\mathbf{x}^1, \dots, \mathbf{x}^n\}$, data 0-center

2. Component index: $q=1$

3. Loop

- Find direction of highest variance: \mathbf{u}^q

- Ensure $|\mathbf{u}^q| = 1$

- Remove \mathbf{u}^q from data:

$$D = \{\mathbf{x}^1 - z_q^1 \mathbf{u}^q, \dots, \mathbf{x}^n - z_q^n \mathbf{u}^q\}$$

$$(\mathbf{u}_i)^T \mathbf{u}_j = 0 \quad \forall i \neq j$$

Thus, we guarantee $z_j^i = \mathbf{u}_j^T \mathbf{x}^i$

27

Independent component analysis (ICA)

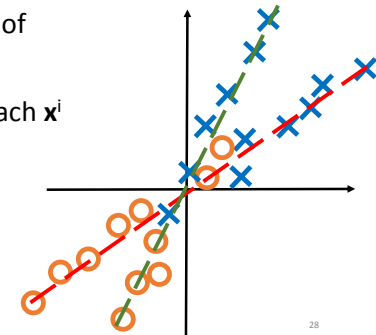
Describe every \mathbf{x}^i with small set of components $\mathbf{u}^{1:T}$

Can use **different** $\mathbf{u}^1, \dots, \mathbf{u}^Q$ for each \mathbf{x}^i

No orthogonality constraint:

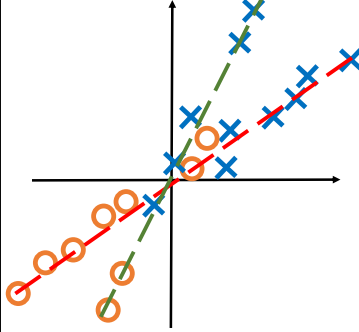
$$(\mathbf{u}^i)^T \mathbf{u}^j \neq 0 \quad \forall i \neq j$$

$$\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$$



28

Idea of learning ICA



1. $D = \{\mathbf{x}^1, \dots, \mathbf{x}^n\}$, data 0-center
2. Component index: $q=1$
3. Loop
 - Find next most common group across data points
 - Find component direct for group \mathbf{u}^q
 - Ensure $|\mathbf{u}^q| = 1$

We cannot guarantee $z_j^i = \mathbf{u}_j^T \mathbf{x}^i$

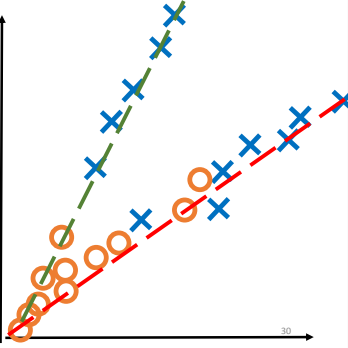
29

Non-negative matrix factorization (NMF)

Describe every \mathbf{x}^i with small set of components $\mathbf{u}^{1:T}$

All components and weights non-negative

$$\mathbf{u}^i \geq 0, z_q^i \geq 0 \quad \forall i, q$$

$$\mathbf{x}^i = \sum_{q=1}^T z_q^i \mathbf{u}^q$$


30

Comparing component analysis

PCA

- Represent each data point \mathbf{x} with low number of components
- All components used to reconstruct each data point \mathbf{x}

ICA / NMF

- Represent each data point \mathbf{x} with low number of components
- Subset of components used to reconstruct each data point \mathbf{x}

33