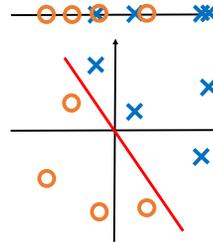# Dimensionality reduction

CISC 5800
Professor Daniel Leeds
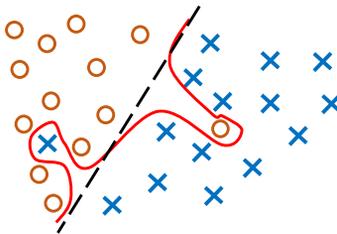
---

## The benefits of extra dimensions



- **Finds existing complex separations between classes**

---

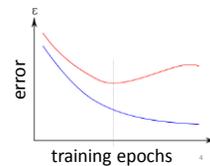## The risks of too-many dimensions



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

---

## 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 $\varepsilon$
  - 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)

---

## Decreasing parameters

- Force parameter values to 0
  - L1 regularization
  - Support Vector selection
  - Feature selection/removal

- Consolidate feature space
  - Component analysis

## Feature removal

- Start with feature set: F={x$_1$, ..., x$_k$}
- Find classifier performance with set F: perform(F)
- Loop
  - Find classifier performance for removing feature x$_1$, x$_2$, ..., x$_k$: argmax$_i$ perform(F-x$_i$)
  - Remove feature that causes least decrease in performance: F=F-x$_i$

Repeat, using AIC or BIC as termination criterion
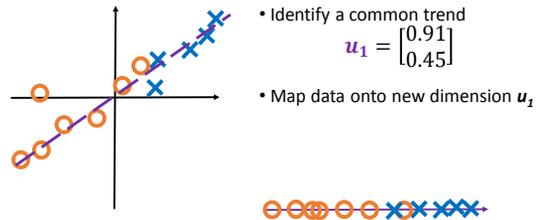
**AIC**: log(L) - k
**BIC**: log(L) - 0.5 k log(m)

## AIC testing: log(L)-k

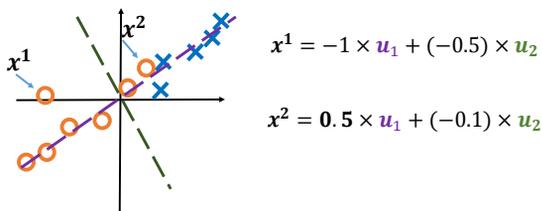| Features | k (num features) | L (likelihood) | AIC |
|---|---|---|---|
| F | 40 | 0.1 | -42.3 |
| F-{x$_3$} | 39 | 0.03 | -41.5 |
| F-{x$_3$,x$_{24}$} | 38 | 0.005 | -41.3 |
| F-{x$_3$,x$_{24}$,x$_{32}$} | 37 | 0.001 | -40.9 |
| F-{x$_3$,x$_{24}$,x$_{32}$,x$_{15}$} | 36 | 0.0001 | **-41.2** |

## Feature selection

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

Repeat, using AIC or BIC as termination criterion

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

## Defining new feature axes

- Identify a common trend
$$u_1 = \begin{bmatrix} 0.91 \\ 0.45 \end{bmatrix}$$
- Map data onto new dimension $u_1$

## Defining data points with new axes

$$x^1 = -1 \times u_1 + (-0.5) \times u_2$$
$$x^2 = 0.5 \times u_1 + (-0.1) \times u_2$$

## Component analysis

Each data point $x^i$ in D can be reconstructed as sum of components $u$:
- $x^i = \sum_{q=1}^{T} z_q^i u_q$
- $z_q^i$ is weight on q$^{th}$ component to reconstruct data point $x^i$

## Component analysis: examples

**Components**          **Data**



2
-1
0

14

## Component analysis: examples

"Eigenfaces" – learned from set of face images

**u**: nine components



$x^i$: data reconstructed

15

## Types of component analysis

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

Learn new axes from data sets: common "components"

- Principal component analysis (PCA):
  - Best reconstruction of each data point $x^i$ with first $t$ components
  - Each component perpendicular to all others: $(u_i)^T u_j = 0 \ \forall i \neq j$

- Independent component analysis (ICA):
  - Minimize number of components to describe each $x^i$
  - Can focus on different components for different $x^i$

- Non-negative matrix factorization (NMF):
  - All data $x^i$ non-negative
  - All components and weights non-negative $u_j \geq 0, \ z_q^i \geq 0 \ \forall i,q$

16

## Principle component analysis (PCA)



Start with
- D = $\{x^1,...,x^n\}$ , data 0-center
- Component index: q=1

Loop
- Find direction of highest variance: $u_q$
  - Ensure $|u_q| = 1$
- Remove $u_q$ from data:
  $$D = \left\{ x^1 - (x^1)^T u_q, \cdots, x^n - (x^n)^T u_q \right\}$$

We **require** $(u_i)^T u_j = 0 \ \forall i \neq j$

Thus, we guarantee $z_j^i = u_j^T x^i$

17

## Independent component analysis (ICA)



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

Find group(s) for each data point

Find direction for each group $u_q$
- Ensure $|u_q| = 1$

We do **not** require $(u_i)^T u_j = 0 \ \forall i \neq j$
Thus, we cannot guarantee $z_j^i = u_j^T x^i$

18

## Evaluating components

- Components learned in order of descriptive power

- Compute reconstruction error for all data by using first v components:

$$error = \sum_i \left( \sum_j \left( x_j^i - \sum_{q=1}^{v} a_q^i u_{q,j} \right)^2 \right)$$

19