

5-17

5-18

5-19

6 Feature dimension reduction

6-1

6-2

6-3

6-4

6.1 Feature selection

Given : • $F_d = \{x_1, \dots, x_d\}$ → set of all d features
• desired number of features $\bar{d} < d$

Feature selection problem:

Find an optimum subset

$$F_{\bar{d}, \text{opt}} = \{x_{i_1}, \dots, x_{i_{\bar{d}}}\} \subset F_d$$

Goal: min a certain cost function $J(F)$

$$J(F_{\bar{d}, \text{opt}}) \leq J(F_d) \quad \forall F_d$$

- filter approach : J ? → not depending on task at hand
- wrapper approach : J = classification error rate

How to find $F_{\bar{d}, \text{opt}}$?

(a) exhaustive search:

try all ${}^d C_{\bar{d}} = \frac{d!}{\bar{d}!(d-\bar{d})!}$ possible subsets

eg $d = 300, \bar{d} = 50 \quad {}^{300}C_{50} \approx 3 \times 10^{57}$

impossible

(b) single feature ranking

- find the best single feature x_{b_1}

$$J(x_{b_1}) \leq J(x_i) \quad \forall x_i \in F_d$$

\Rightarrow d single-feature classification

- find the 2nd best single feature x_{b_2}

$$J(x_{b_2}) \leq J(x_i) \quad \forall x_i \in F_d \setminus \{x_{b_1}\}$$

\Rightarrow $d-1$ single feature classification

\vdots

- find the \bar{d} -th best single feature $x_{b_{\bar{d}}}$

Does it work well?

- many features may be redundant, there may be dependencies between them, no new information is gained

- bad single features may combine to form a better multiple feature eg XOR

$$\begin{array}{c|c|c}
 x_2 & + & - \\
 \hline
 + & 0 & 0 \\
 - & 0 & 0 \\
 \hline
 \end{array}$$

x_1 : bad feature
 x_2 : " "
 (x_1, x_2) : perfect

6-5

(1) heuristic sequential algorithms

6-6

6-7

6-8

Pro/cons:

+) Features keeping their physical meanings

+) only \bar{d} features needs to be calculated

(-) a special case of feature transform, not optimum

6.2 Feature transform

$$\underline{x} \in \mathbb{R}^d \xrightarrow{\underline{\Phi}(\cdot)} \bar{\underline{x}} = \underbrace{\underline{\Phi}(\underline{x})}_{\substack{\uparrow \\ \bar{d}}} \in \mathbb{R}^{\bar{d}} \quad \bar{d} < d$$

How to choose $\underline{\Phi}$?

Pro/cons:

(+) more general, better performance than purely feature selection for the same \bar{d} because feature selection is a special case of $\underline{\Phi}(\cdot)$

$$\text{eg } \bar{\underline{x}} = \begin{bmatrix} x_2 \\ x_4 \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{\underline{\Phi}(\underline{x})} \begin{bmatrix} \underline{x} \end{bmatrix}$$

(-) compute all d features

(-) loss of physical meaning of \underline{x}

optimum $\phi(\cdot)$: difficult to find

a simple $\phi(\cdot)$

• linear transform: $\underline{\bar{x}} = \underline{W}^T \underline{x}$

• affine transform: $\underline{\bar{x}} = \underline{W}^T (\underline{x} - \underline{\mu}) + \underline{\mu}$

\Rightarrow only optimize \underline{W} (and $\underline{\mu}$)

Principal Component Analysis (PCA):

N samples $\underline{x}_n \in \mathbb{R}^d$ $1 \leq n \leq N$

no class labels, unsupervised

• $\underline{W} = [\underline{v}_1 \dots \underline{v}_{\bar{d}}] \in \mathbb{R}^{d \times \bar{d}}$ $\bar{d} < d$

• $\underline{\bar{x}} = \underline{W}^T \underline{x} \in \mathbb{R}^{\bar{d}}$

• Reconstruction: $\hat{\underline{x}} = \underline{W} \cdot \underline{\bar{x}} = \underline{W} \underline{W}^T \underline{x} \in \mathbb{R}^d$

Optimization criterion:

$$\min \sum_{n=1}^N \left\| \underline{x}_n - \underline{W} \underbrace{\underline{W}^T \underline{x}_n}_{\text{compression}} \right\|^2$$

s. t.

$$\underline{W}^T \underline{W} = \underline{I}$$

i. e. orthonormal

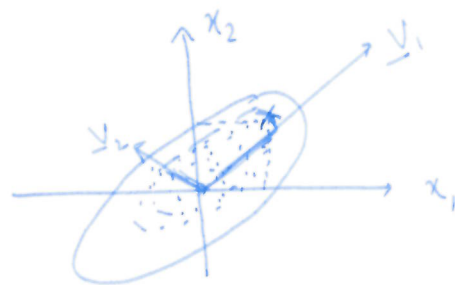
$\underbrace{\text{compression: } \mathbb{R}^d \rightarrow \mathbb{R}^{\bar{d}}}$
 $\underbrace{\text{reconstruction: } \mathbb{R}^{\bar{d}} \rightarrow \mathbb{R}^d}$

reconstruction error

6-9

6-11

Illustration of PCA:



project $(x_1, x_2) \rightarrow (y_1, y_2)$

PCA: keep \bar{d} coordinates $v_i^T x$ of x along

the most dominant directions v_i

$\Rightarrow \bar{d}$ principal components of x

Why DCT over PCA \rightarrow ~~depends on data~~, DCT is independent of input whereas PCA is more computationally costly

Non-linear extension of PCA \rightarrow auto encoder \rightarrow compression done by encoder,
decompression done by decoder.