

## 4. non-linear classification

### 4.1 polynomial function

↳ transforms your input into a new set of features by including all possible combinations up to a certain degree  $d$

↳ they include

1. original features  $(n_1, n_2)$
2. power of features  $(n_1^2, n_2^3)$
3. interaction terms  $(n_1 n_2, n_1^2 n_2)$

↳ say your original data is  $n = [n_1, n_2] \in \mathbb{R}^2$  & you perform a polynomial transformation of degree  $d=2$  the new features are:

1. degree 0: a constant term (ex. 1)
2. degree 1: original features  $(n_1, n_2)$
3. degree 2: squares  $(n_1^2, n_2^2)$  & interactions terms  $(n_1 n_2)$

the new feature vector  $\phi(n) \Rightarrow$  includes all combinations of powers of  $n_1$  &  $n_2$  where total power  $\leq d$

$$\phi(n) = [1, n_1, n_2, n_1^2, n_2^2, n_1 n_2]$$

$\Downarrow_{1 \leq 2} \quad \Downarrow_{1 \leq 2} \quad \Downarrow_{2 \leq 2} \quad \Downarrow_{2 \leq 2} \quad \Downarrow_{2 \leq 2}$

now  $n \in \mathbb{R}^2$  is transformed into  $\phi(n) \in \mathbb{R}^6$

to apply factorial normalization for each term:

degree 0 (constant term)

- term = 1

- degree  $\Rightarrow d=0$

$j_1 \in j_2$  represent the features  $u_1, u_2$   
 $j_3$  is there to hold the degree when  $j_1 = 0 = j_2, j_3 = 2$  not used

$$\sqrt{\frac{d!}{j_1! \cdot j_2! \cdot j_3!}} \cdot \sqrt{\frac{20!}{0! \cdot 0! \cdot 2!}} = \sqrt{1} = 1$$

so here, the degree is 0 but the transformation is of degree 2, so we put the 2 in  $j_3$  to ensure  $u_1 \in u_2$  are of degree 0. Same with the rest

degree 1 (linear terms)

- term =  $u_1, u_2$

- degree  $\Rightarrow d=1$

$$\text{for } u_1 \Rightarrow j_1=1, j_2=0, j_3=1$$

$$\text{for } u_2 \Rightarrow j_1=0, j_2=1, j_3=1$$

$$u_1 \Rightarrow \sqrt{\frac{2!}{1! \cdot 0! \cdot 1!}} = \sqrt{2}, \quad u_2 \Rightarrow \sqrt{\frac{2!}{0! \cdot 1! \cdot 1!}} = \sqrt{2}$$

degree 2 (quadratic terms)

- term =  $u_1^2, u_2^2, u_1 u_2$

- degree  $\Rightarrow d=2$

$$\text{for } u_1^2 \Rightarrow j_1=2, j_2=0, j_3=0$$

$$\text{for } u_2^2 \Rightarrow j_1=0, j_2=2, j_3=0$$

$$\text{for } u_1 u_2 \Rightarrow j_1=1, j_2=1, j_3=0$$

$$u_1^2 \Rightarrow \sqrt{\frac{2!}{2! \cdot 0! \cdot 0!}} = 1$$

$$u_2^2 \Rightarrow \sqrt{\frac{2!}{0! \cdot 1! \cdot 1!}} = 1$$

$$u_1 u_2 = \sqrt{\frac{2!}{1! \cdot 1! \cdot 0!}} = \sqrt{2}$$

the classifier now has the following form

$$h(\phi(u); \theta, \theta_0) = \text{sign}(\theta^T \phi(u) + \theta_0)$$

$$= \text{sign}(\theta_1 u_1 + \theta_2 u_2 + \theta_3 u_1^2 + \dots + \theta_d u_d + \theta_0)$$

the size of the vector  $\phi(u) =$

$$\binom{n+d}{d} = \frac{(n+d)!}{d! n!}$$

## 4.2 the Kernel trick

↳ allows for dot product calculation

$\phi(u)^T \phi(u')$  in higher dimensional feature space  
without explicitly computing  $\phi(u)$

$$K(u, u') = \phi(u)^T \phi(u')$$

↳ Kernel functions

1. additivity: sum of 2 valid Kernels is a valid Kernel

2. scalar multiplication: valid Kernel  $\cdot$  +ve scalar = valid Kernel

3. product of Kernels: product of 2 valid Kernels is a valid Kernel

4. exponentiation

### 4.3 Kernel least square

$$\min_{\beta} \|y - \phi(x)\beta\|_2^2$$

$$\beta = (\phi(x)^T \phi(x))^{-1} \phi(x)^T y$$

### Limitations

1. complex patterns
2. overfitting
3. curse of dimensionality
4. Kernel & parameter tuning