

# **Factorization Machines (FM) and Field-aware Factorization Machines (FFM)**

Hung-Hsuan Chen

<https://www.ncu.edu.tw/~hhchen/>

# Linear model

- Formula

$$y = \sum_{i=0}^d \theta_i x_i$$

- $y$ : target
  - $x_1, \dots, x_d$ : features
  - $x_0$ : bias
  - $\theta_0, \theta_1, \dots, \theta_d$ : parameters to learn
- The model can capture each feature's influence on the target

# Degree-2 polynomial model (poly2)

- Formula

$$y = \sum_{i=0}^d \theta_i x_i + \sum_{(j,k) \in C_2} \theta_{jk} x_j x_k$$

- $y$ : target
- $x_1, \dots, x_d$ : features
- $x_0$ : bias
- $\theta_0, \theta_1, \dots, \theta_d, \theta_{12}, \theta_{13}, \dots, \theta_{d-1,d}$ : parameters to learn
- $C_2$ : 2-combination of elements in  $[x_1, \dots, x_d]$

- The model can capture

- Each feature's influence on the target
- Each feature-pair's influence on the target

# Factorization machines (FM)

- Formula

$$y = \sum_{i=0}^d \theta_i x_i + \sum_{(j,k) \in C_2} \langle \mathbf{v}_j, \mathbf{v}_k \rangle x_j x_k$$

- $y$ : target
- $x_1, \dots, x_d$ : features
- $x_0$ : bias
- $\theta_0, \theta_1, \dots, \theta_d, \mathbf{v}_1, \mathbf{v}_d$ : parameters to learn, each  $\mathbf{v}_j$  is a vector of length  $\ell$
- $C_2$ : 2-combination of elements in  $[x_1, \dots, x_d]$

# Poly2 vs FM (# paras)

- If we have  $d$  features

- # parameters for poly-2:

$$(d + 1) + \binom{d}{2} = d + 1 + \frac{d(d - 1)}{2} \approx O(d^2)$$

- # parameters for FM (assuming the length of the vector  $\mathbf{v}_i$  is  $\ell$ :

$$(d + 1) + \ell d = 1 + (\ell + 1)d \approx O(\ell d)$$

- If  $d \gg \ell$ , FM has fewer parameters to learn

- FM is probably more appropriate when we have large but sparse features

# Example: ad classification

Country	Day	Ad type	Clicked?
USA	Thanksgiving	Movie	1
China	Chinese New Year	Game	0
China	Thanksgiving	Game	1

- Task: given features, predict click or not

# Example: ad classification (cont')

- Standard one-hot encoding

USA	China	Thanks giving	Chinee s new year	Movie	Game	Clicked ?
1	0	1	0	1	0	1
0	1	0	1	0	1	0
0	1	1	0	0	1	1

- Very large feature space
- Very sparse samples

# Example: ad classification (cont')

- Features might be more important in “pairs”
  - Country == “USA” and Day == “Thanksgiving”
  - Country == “China” and Day == “Chinese new year”
- If we create features for every pair of features
  - Number of features goes from  $d$  to  $\binom{d}{2}$
  - Samples: still sparse



# Gradients

$$\hat{y} = \theta_0 + \sum_{i=1}^d \theta_i x_i + \sum_{(j,k) \in C_2} \langle \mathbf{v}_j, \mathbf{v}_k \rangle x_j x_k$$

$$L := (y - \hat{y})^2$$

- $\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_i} = -2(y - \hat{y}) \frac{\partial \hat{y}}{\partial \theta_i}$

- $\frac{\partial \hat{y}}{\partial \theta_0} = 1$

- $\frac{\partial \hat{y}}{\partial \theta_i} = x_i$  (if  $1 \leq i \leq d$ )

- $\frac{\partial \hat{y}}{\partial v_{jf}} = v_{kf} x_j x_k$

# FFM example

Publisher (P)	Advertiser (A)	Gender (G)	Clicked?
ESPN	Nike	Male	Yes

- On ESPN, a male clicked an ad about Nike
- For FM:

$$\hat{y} = \mathbf{v}_{ESPN} \cdot \mathbf{v}_{Nike} + \mathbf{v}_{ESPN} \cdot \mathbf{v}_{Male} + \mathbf{v}_{Nike} \cdot \mathbf{v}_{Male}$$

- Every feature has one corresponding latent vector to learn
- E.g.,  $\mathbf{v}_{ESPN}$  is used to learn the effect with Nike ( $\mathbf{v}_{ESPN} \cdot \mathbf{v}_{Nike}$ ) and Male ( $\mathbf{v}_{ESPN} \cdot \mathbf{v}_{Male}$ )
- However, Nike and Male belong to different fields, the effects of  $\mathbf{v}_{ESPN}$  on ( $\mathbf{v}_{ESPN} \cdot \mathbf{v}_{Nike}$ ) and ( $\mathbf{v}_{ESPN} \cdot \mathbf{v}_{Male}$ ) could be different

# FFM example (cont')

Publisher (P)	Advertiser (A)	Gender (G)	Clicked?
ESPN	Nike	Male	Yes

- For FFM:

$$\hat{y} = \mathbf{v}_{ESPN,A} \cdot \mathbf{v}_{Nike,P} + \mathbf{v}_{ESPN,G} \cdot \mathbf{v}_{Male,P} + \mathbf{v}_{Nike,G} \cdot \mathbf{v}_{Male,A}$$

- Every feature has several corresponding latent vector to learn
- E.g., to learn the effect of (ESPN, Nike),  $\mathbf{v}_{ESPN,A}$  is used because Nike belongs to field Advertiser. However, to learn the effect of (ESPN, Male),  $\mathbf{v}_{ESPN,G}$  is used because Male belongs to field Gender

# Summary

- FM as an extension of linear model
- FM as a variation of poly2 model
- FFM as an extension of FM
- FM and FFM are especially useful when the features are large and sparse

# Quiz

- If we know that targets are influenced by each feature and each pair of features
  - Polynomial kernel is helpful in this case (True or false)
  - Linear SVM + poly is helpful in this case (True or false)