

Latent factor model & Learning algorithm

2013-5-30

Seung-hwan baek

FM model

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{j=1}^p w_j x_j + \sum_{j=1}^p \sum_{j'=j+1}^p x_j x_{j'} \sum_{f=1}^k v_{j,f} v_{j',f},$$

factorize Matrix V: use non-direct information

$$w_{j,j} \approx \langle \mathbf{v}_j, \mathbf{v}_{j'} \rangle = \sum_{f=1}^k v_{j,f} v_{j',f}$$

Lemma: $W = VV^T$ if k is chosen large enough

$$w_0 \in \mathbb{R}, \quad \mathbf{w} \in \mathbb{R}^p, \quad V \in \mathbb{R}^{p \times k}$$

Learning

$$\hat{y}(\mathbf{x}) = g_{\theta}(\mathbf{x}) + \theta h_{\theta}(\mathbf{x}) \quad \forall \theta \in \Theta,$$

$$h_{\theta}(\mathbf{x}) = \frac{\partial \hat{y}(\mathbf{x})}{\partial \theta} = \begin{cases} 1, & \text{if } \theta \text{ is } w_0 \\ x_l, & \text{if } \theta \text{ is } w_l \\ x_l \sum_{j \neq l} v_{j,f} x_j, & \text{if } \theta \text{ is } v_{l,f} \end{cases}$$

Gradient of predicted value

Object-Optimization

- Minimize the gap between true value and predicted value
- Loss function
 - Regression
 - Least square loss $l^{\text{LS}}(y_1, y_2) := (y_1 - y_2)^2$
 - Binary Classification
 - Sigmoid/logistic function $l^{\text{C}}(y_1, y_2) := -\ln \sigma(y_1 y_2)$
- To avoid overfitting
 - Regularization value per group
 $\lambda^0, \quad \lambda_{\pi}^w, \quad \lambda_{f,\pi}^v, \quad \forall \pi \in \{1, \dots, \Pi\}, \forall f \in \{1, \dots, k\},$

$$\text{OPT}(S) := \operatorname{argmin}_{\Theta} \sum_{(\mathbf{x}, y) \in S} l(\hat{y}(\mathbf{x}|\Theta), y)$$



$$\text{OPTREG}(S, \lambda) := \operatorname{argmin}_{\Theta} \left(\sum_{(\mathbf{x}, y) \in S} l(\hat{y}(\mathbf{x}|\Theta), y) + \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2 \right)$$

Probabilistic interpretation

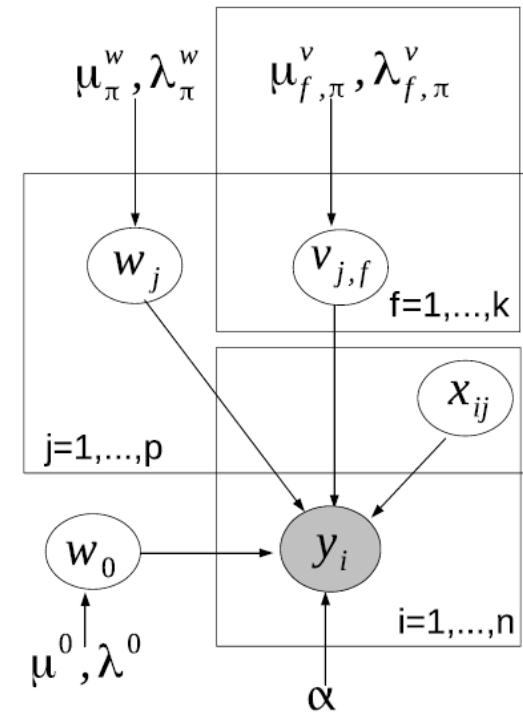
- Regression

$$y|\mathbf{x}, \Theta \sim \mathcal{N}(\hat{y}(\mathbf{x}, \Theta), 1/\alpha).$$

- Binary classification

- b: link function like logistic function

$$y|\mathbf{x}, \Theta \sim \text{Bernoulli}(b(\hat{y}(\mathbf{x}, \Theta))),$$



(a) Factorization machine (FM).

Gradients

- Loss function of gradient
 - Regression

$$\frac{\partial}{\partial \theta} l^{\text{LS}}(\hat{y}(\mathbf{x}|\Theta), y) = \frac{\partial}{\partial \theta} (\hat{y}(\mathbf{x}|\Theta) - y)^2 = 2 (\hat{y}(\mathbf{x}|\Theta) - y) \frac{\partial}{\partial \theta} \hat{y}(\mathbf{x}|\Theta),$$

- Binary classification

$$\frac{\partial}{\partial \theta} l^{\text{C}}(\hat{y}(\mathbf{x}|\Theta), y) = \frac{\partial}{\partial \theta} -\ln \sigma (\hat{y}(\mathbf{x}|\Theta) y) = (\sigma (\hat{y}(\mathbf{x}|\Theta) y) - 1) y \frac{\partial}{\partial \theta} \hat{y}(\mathbf{x}|\Theta).$$

$$\frac{\partial}{\partial \theta} \hat{y}(\mathbf{x}|\Theta) = h_{\theta}(\mathbf{x}).$$

Stochastic Gradient Descent(SGD)

- Object: Find Loss function's global minimum, but we can't
- At least local minimum
- Iterates over (\mathbf{x}, y) in S

$$\theta \leftarrow \theta - \eta \left(\frac{\partial}{\partial \theta} l(\hat{y}(\mathbf{x}), y) + 2 \lambda_{\theta} \theta \right),$$

ALGORITHM 1: Stochastic Gradient Descent (SGD)

Input: Training data S , regularization parameters λ , learning rate η , initialization σ

Output: Model parameters $\Theta = (w_0, \mathbf{w}, \mathbf{V})$

$w_0 \leftarrow 0; \mathbf{w} \leftarrow (0, \dots, 0); \mathbf{V} \sim \mathcal{N}(0, \sigma);$

repeat

for $(\mathbf{x}, y) \in S$ **do**

$w_0 \leftarrow w_0 - \eta \left(\frac{\partial}{\partial w_0} l(\hat{y}(\mathbf{x}|\Theta), y) + 2 \lambda^0 w_0 \right);$

for $i \in \{1, \dots, p\} \wedge x_i \neq 0$ **do**

$w_i \leftarrow w_i - \eta \left(\frac{\partial}{\partial w_i} l(\hat{y}(\mathbf{x}|\Theta), y) + 2 \lambda_{\pi(i)}^w w_i \right);$

for $f \in \{1, \dots, k\}$ **do**

$v_{i,f} \leftarrow v_{i,f} - \eta \left(\frac{\partial}{\partial v_{i,f}} l(\hat{y}(\mathbf{x}|\Theta), y) + 2 \lambda_{f,\pi(i)}^v v_{i,f} \right);$

end

end

end

until *stopping criterion is met*;

SGD hyperparameters

Learning rate η

Regularization λ

Initialization σ

ALGORITHM 1: Stochastic Gradient Descent (SGD)

Input: Training data S , regularization parameters λ , learning rate η , initialization σ

Output: Model parameters $\Theta = (w_0, \mathbf{w}, \mathbf{V})$

$w_0 \leftarrow 0$; $\mathbf{w} \leftarrow (0, \dots, 0)$; $\mathbf{V} \sim \mathcal{N}(0, \sigma)$;

repeat

for $(\mathbf{x}, y) \in S$ **do**

$w_0 \leftarrow w_0 - \eta \left(\frac{\partial}{\partial w_0} l(\hat{y}(\mathbf{x}|\Theta), y) + 2\lambda^0 w_0 \right)$;

for $i \in \{1, \dots, p\} \wedge x_i \neq 0$ **do**

$w_i \leftarrow w_i - \eta \left(\frac{\partial}{\partial w_i} l(\hat{y}(\mathbf{x}|\Theta), y) + 2\lambda_{\pi(i)}^w w_i \right)$;

for $f \in \{1, \dots, k\}$ **do**

$v_{i,f} \leftarrow v_{i,f} - \eta \left(\frac{\partial}{\partial v_{i,f}} l(\hat{y}(\mathbf{x}|\Theta), y) + 2\lambda_{f,\pi(i)}^v v_{i,f} \right)$;

end

end

end

until *stopping criterion is met*;

$$\theta \leftarrow \theta - \eta \left(\frac{\partial}{\partial \theta} l(\hat{y}(\mathbf{x}), y) + 2\lambda_{\theta} \theta \right),$$

Alternating Least-Squares/Coordinate Descent

- SGD: Minimizing the loss per training data
- ALS: Minimizing the loss per model parameter

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} \left(\sum_{(\mathbf{x}, y) \in S} (\hat{y}(\mathbf{x}|\Theta) - y)^2 + \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2 \right) \\ &= \operatorname{argmin}_{\theta} \left(\sum_{(\mathbf{x}, y) \in S} (g_{\theta}(\mathbf{x}|\Theta \setminus \{\theta\}) + \theta h_{\theta}(\mathbf{x}|\Theta \setminus \{\theta\}) - y)^2 + \sum_{\theta \in \Theta} \lambda_{\theta} \theta^2 \right) \\ &= \frac{\sum_{i=1}^n (y - g_{\theta}(\mathbf{x}_i|\Theta \setminus \{\theta\})) h_{\theta}(\mathbf{x}_i|\Theta \setminus \{\theta\})}{\sum_{i=1}^n h_{\theta}(\mathbf{x}_i)^2 + \lambda_{\theta}} \\ &= \frac{\theta \sum_{i=1}^n h_{\theta}^2(\mathbf{x}_i) + \sum_{i=1}^n h_{\theta}(\mathbf{x}_i) e_i}{\sum_{i=1}^n h_{\theta}(\mathbf{x}_i)^2 + \lambda_{\theta}}, \quad e_i := y_i - \hat{y}(\mathbf{x}_i|\Theta).\end{aligned}$$

ALS hyperparameters

Regularization λ

Initialization σ

ALGORITHM 2: Alternating least squares (ALS)

Input: Training data S , regularization parameters λ , initialization σ

Output: Model parameters $\Theta = (w_0, \mathbf{w}, \mathbf{V})$

$w_0 \leftarrow 0$; $\mathbf{w} \leftarrow (0, \dots, 0)$; $\mathbf{V} \sim \mathcal{N}(0, \sigma)$;

repeat

$\hat{\mathbf{y}} \leftarrow$ predict all cases S ;

$\mathbf{e} \leftarrow \mathbf{y} - \hat{\mathbf{y}}$;

$w_0 \leftarrow w_0^*$;

for $l \in \{1, \dots, p\}$ **do**

$w_l \leftarrow w_l^*$;

 update e ;

end

for $f \in \{1, \dots, k\}$ **do**

 init $q_{\cdot, f}$;

for $l \in \{1, \dots, p\}$ **do**

$v_{l, f} \leftarrow v_{l, f}^*$;

 update e, q ;

end

end

until stopping criterion is met;

$$\theta \leftarrow \theta^*.$$