

DLRM

DeepSync, South Korea

0. 기존의 personalization and recommendation

1. Recommendation system
 - content filtering (expert-user)
 - past user behavior
 - neighborhood method
 - latent factor method
2. Predictive analytics
 - statistical model (ex: linear regression, logistic regression)
 - deep networks (-> embedding)
 - latent factor method

1. DLRM (Deep Learning Recommendation Model)

앞에 소개한 방법 중 몇 개 + mlp + interacting dense feature

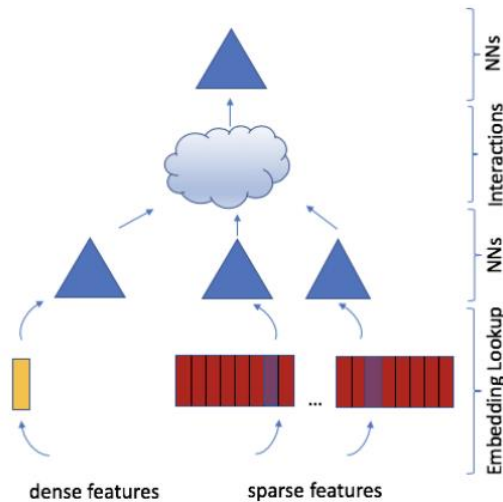
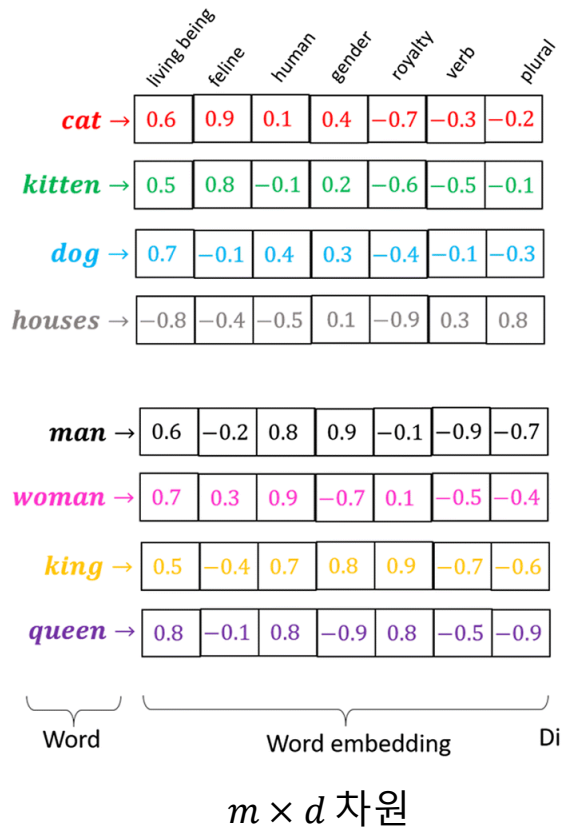


Figure 1: A deep learning recommendation model

2. Model Architecture: embedding

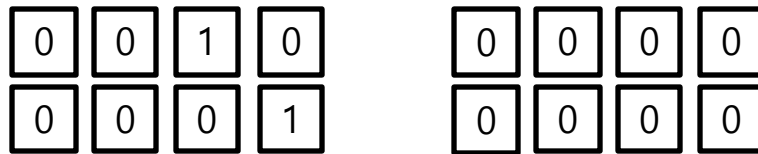


만약 dog 차원을 알고 싶다면?



$1 \times m$ 차원

만약 dog, houses 차원을 알고 싶다면?



$2 \times m$ 차원

$$w_i^T = e_i^T W$$

2. Model Architecture: latent factor method

Product vector: $w (\mathbb{R}^d)$

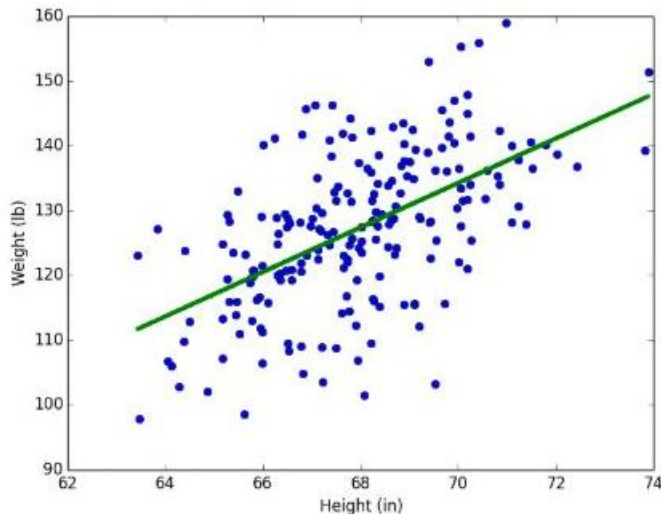
User vector: $v (\mathbb{R}^d)$

Dot product: $w_i^T \cdot v_j$ (i번째 상품과 j번째 사용자)

Ground truth: r_{ij}

Goal:
$$\min \sum_{(i,j) \in \mathcal{S}} r_{ij} - w_i^T v_j$$

2. Model Architecture: Factorization Machine



Linear regression: $y = w_1x_1 + w_0$

Multiple linear regression:

$$y = w_1x_1 + w_2x_2 + w_0$$

Multiple linear regression with interaction terms:

$$y = w_1x_1 + w_2x_2 + m_3(x_1 \times x_2) + w_0$$

$$y = w_1x_1 + w_2x_2 + \langle v_1 \times v_2 \rangle (x_1 \times x_2) + w_0$$

N term

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

$$\hat{y} = b + \mathbf{w}^T \mathbf{x} + \mathbf{x}^T \text{upper}(\mathbf{V}\mathbf{V}^T) \mathbf{x}$$

2. Model Architecture: Factorization Machine

Feature vector x																				Target y		
$x^{(1)}$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...	5	$y^{(1)}$
$x^{(2)}$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...	3	$y^{(2)}$
$x^{(3)}$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...	1	$y^{(2)}$
$x^{(4)}$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...	4	$y^{(3)}$
$x^{(5)}$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...	5	$y^{(4)}$
$x^{(6)}$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...	1	$y^{(5)}$
$x^{(7)}$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...	5	$y^{(6)}$
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...		
	User				Movie					Other Movies rated						Last Movie rated						

2. Model Architecture: Factorization Machine

장점1 : linear complexity

$$\begin{aligned}
 & O(kn^2) \quad \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j - \frac{1}{2} \sum_{i=1}^n \langle \mathbf{v}_i, \mathbf{v}_i \rangle x_i x_i \\
 &= \frac{1}{2} \left(\sum_{i=1}^n \sum_{j=1}^n \sum_{f=1}^k v_{i,f} v_{j,f} x_i x_j - \sum_{i=1}^n \sum_{f=1}^k v_{i,f} v_{i,f} x_i x_i \right) \\
 &= \frac{1}{2} \sum_{f=1}^k \left(\left(\sum_{i=1}^n v_{i,f} x_i \right) \left(\sum_{j=1}^n v_{j,f} x_j \right) - \sum_{i=1}^n v_{i,f}^2 x_i^2 \right) \\
 & \quad \downarrow \\
 & O(kn) \quad \boxed{= \frac{1}{2} \sum_{f=1}^k \left(\left(\sum_{i=1}^n v_{i,f} x_i \right)^2 - \sum_{i=1}^n v_{i,f}^2 x_i^2 \right)}
 \end{aligned}$$

2. Model Architecture: Factorization Machine

장점2: sparse한 상황에서 interaction을 예상 가능

Task: 사람 Alice의 Star Trek의 rating을 예측 하고 싶음

- no interaction ($w_{A,ST} = 0$)

- Bob과 Charlie는 star wars를 봤고 비슷한 평을 남김
($\langle v_B, v_{SW} \rangle$ and $\langle v_C, v_{SW} \rangle$ are similar)

- Alice는 Charlie와 different factor vector를 보임
($\langle v_A, v_{SW} \rangle$ and $\langle v_C, v_{SW} \rangle$ are not similar)

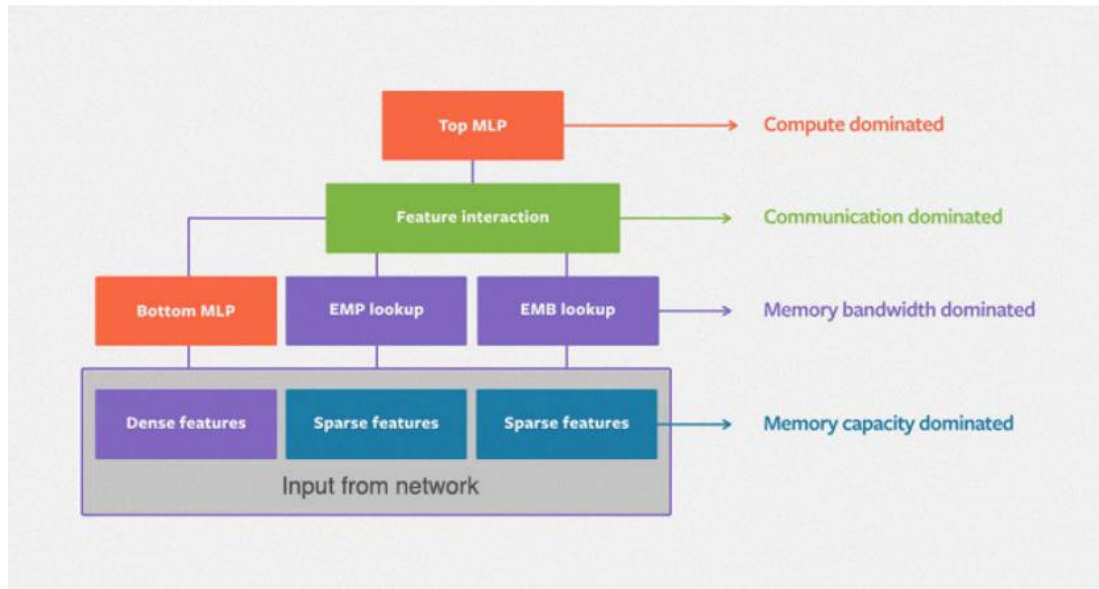
- Bob은 star wars와 star trek의 평을 비슷하게 함
($\langle v_B, v_{SW} \rangle$ and $\langle v_B, v_{ST} \rangle$ are similar)

$\Rightarrow (\langle v_A, v_{ST} \rangle$ and $\langle v_A, v_{SW} \rangle$ are similar)

2. Model Architecture: MLP

$$\hat{y} = W_k \sigma(W_{k-1} \sigma(\dots \sigma(W_1 \mathbf{x} + \mathbf{b}_1) \dots) + \mathbf{b}_{k-1}) + \mathbf{b}_k$$

2. Model Architecture and Parallelism



1. Categorical features -> emb
2. Continuous features -> MLP
3. Second-order interaction
4. 3의 결과 -> Top MLP, Dense features로
5. Sigmoid function -> $T = \{+1, -1\}$

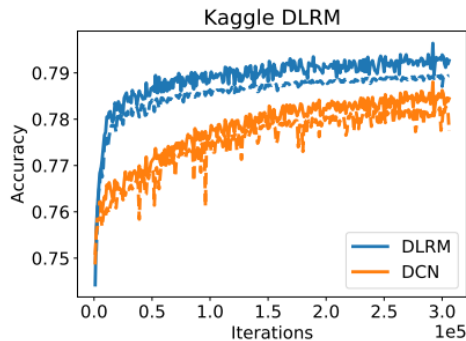
1. Model parallel (size of emb-> large parameter)
2. Data parallel (smaller parameter, large computation)
- 3.
4. Data parallel

3. experiments

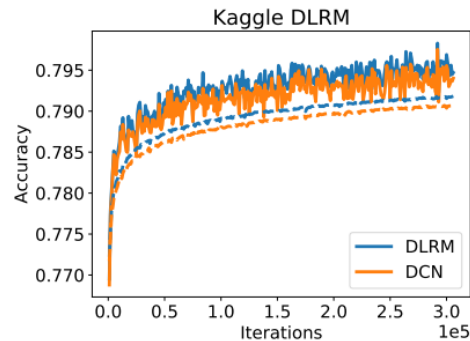
- Kaggle data set
- DCN과 비교

- DLRM
 1. bottom MLP: 512, 256, 64
 2. top MLP: 512, 256
 3. emb dimension: 16

- DCN
 1. 6 cross layer
 2. deep network: 512, 256
 3. emb dimension: 16



(a) SGD



(b) Adagrad

3. experiments

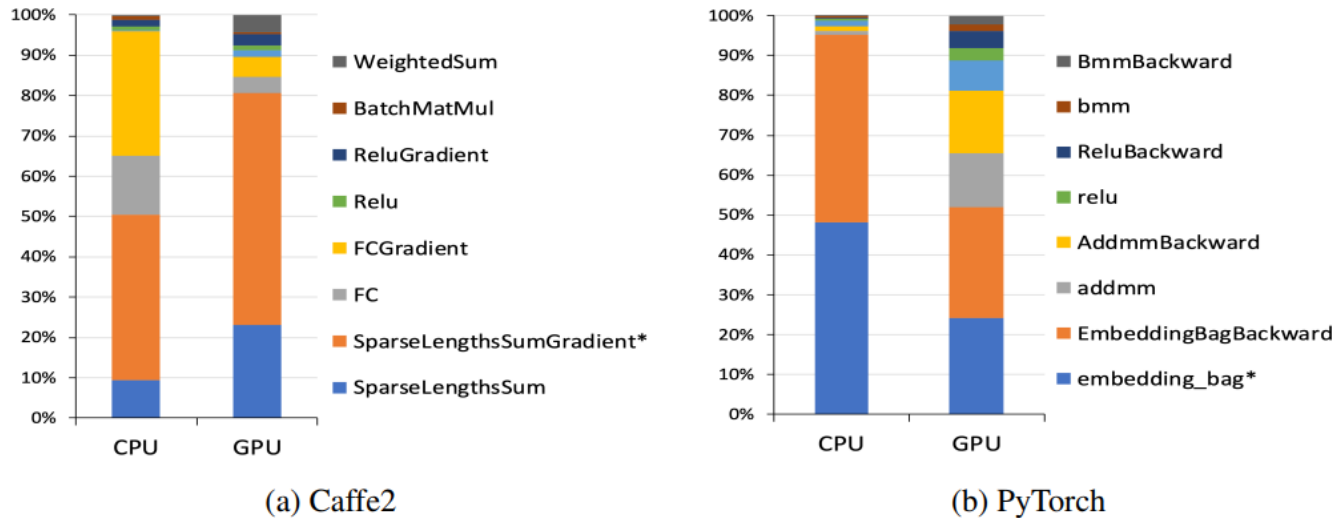


Figure 6: Profiling of a sample DLRM on a single socket/device

reference

[\[추천시스템\]\[paper review\]\[구현\] Factorization Machines \(tistory.com\)](#)

[DLRM: An advanced, open source deep learning recommendation model \(meta.com\)](#)

[paper.dvi \(ntu.edu.tw\)](#)

Thank you for your time