

1 InteractionArch

一共有三种 Interaction 备选：

1. Default: Returns the pairwise dot product of each sparse feature pair, the dot product of each sparse features with the output of the dense layer, and the dense layer itself (all concatenated).
2. DCN: Returns the output of a Deep Cross Net v2 <https://arxiv.org/pdf/2008.13535.pdf> with a low rank approximation for the weight matrix. The input and output sizes are the same for this interaction layer ($F \times D + D$).
3. Projection: Return $Y \times Z$ and the dense layer itself (all concatenated) where Y is the output of interaction branch 1 and Z is the output of interaction branch 2. Y and Z are of size $B \times (F_1 \times D)$ and $B \times (D \times F_2)$ respectively for some F_1 and F_2 .

2 Default Interaction

forward:

```
1 # dense_features.shape = [B, 128]
2 # sparse_features.shape = [B, 26, 128]
3 combined_values =
4     torch.cat((dense_features.unsqueeze(1), sparse_features),
5               dim=1)
6 # combined_values.shape = [B, 27, 128]
7
8 # [B, 27, 27] = [B, 27, 128] <bmm> [B, 128, 27]
9 interactions =
10     torch.bmm(combined_values,
11              torch.transpose(combined_values, 1, 2))
12 # [B, 351]
13 interactions_flat =
14     interactions[:, self.triu_indices[0], self.triu_indices[1]]
15
16 # concatenated_dense.shape = [B, 479]
17 concatenated_dense =
18     torch.cat((dense_features, interactions_flat), dim=1)
```

给 forward 形式化的表达。为了方便推导，我们假设 $B = 2$, $embedding_dim = 4$, $sparse_feature_size = 2$:

$$D \text{ for } dense[2, 4] = \begin{bmatrix} d_{0,0} & d_{0,1} & d_{0,2} & d_{0,3} \\ d_{1,0} & d_{1,1} & d_{1,2} & d_{1,3} \end{bmatrix}$$

$$S \text{ for } sparse[2, 2, 4] = \begin{bmatrix} [s_{0,0,0}, s_{0,0,1}, s_{0,0,2}, s_{0,0,3}] & [s_{0,1,0}, s_{0,1,1}, s_{0,1,2}, s_{0,1,3}] \\ [s_{1,0,0}, s_{1,0,1}, s_{1,0,2}, s_{1,0,3}] & [s_{1,1,0}, s_{1,1,1}, s_{1,1,2}, s_{1,1,3}] \end{bmatrix}$$

$$C \text{ for } Comb[2, 3, 4] = \begin{bmatrix} [d_{0,0}, d_{0,1}, d_{0,2}, d_{0,3}] & [s_{0,0,0}, s_{0,0,1}, s_{0,0,2}, s_{0,0,3}] & [s_{0,1,0}, s_{0,1,1}, s_{0,1,2}, s_{0,1,3}] \\ [d_{1,0}, d_{1,1}, d_{1,2}, d_{1,3}] & [s_{1,0,0}, s_{1,0,1}, s_{1,0,2}, s_{1,0,3}] & [s_{1,1,0}, s_{1,1,1}, s_{1,1,2}, s_{1,1,3}] \end{bmatrix}$$

$$C^T[2, 4, 3] = \begin{bmatrix} [d_{0,0}, s_{0,0,0}, s_{0,1,0}] & [d_{0,1}, s_{0,0,1}, s_{0,1,1}] & [d_{0,2}, s_{0,0,2}, s_{0,1,2}] & [d_{0,3}, s_{0,0,3}, s_{0,1,3}] \\ [d_{1,0}, s_{1,0,0}, s_{1,1,0}] & [d_{1,1}, s_{1,0,1}, s_{1,1,1}] & [d_{1,2}, s_{1,0,2}, s_{1,1,2}] & [d_{1,3}, s_{1,0,3}, s_{1,1,3}] \end{bmatrix}$$

$$flat(CC^T)|_{b=\{0,1\}} = \begin{bmatrix} \cdots & F_{b,0} & F_{b,1} \\ F_{b,0} & \cdots & F_{b,2} \\ F_{b,1} & F_{b,2} & \cdots \end{bmatrix}$$

其中：

$$\begin{aligned} F_{b,0} &= d_{b,0}s_{b,0,0} + d_{b,1}s_{b,0,1} + d_{b,2}s_{b,0,2} + d_{b,3}s_{b,0,3} \\ F_{b,1} &= d_{b,0}s_{b,1,0} + d_{b,1}s_{b,1,1} + d_{b,2}s_{b,1,2} + d_{b,3}s_{b,1,3} \\ F_{b,2} &= s_{b,0,0}s_{b,1,0} + s_{b,0,1}s_{b,1,1} + s_{b,0,2}s_{b,1,2} + s_{b,0,3}s_{b,1,3} \end{aligned}$$

$$Y[2, 7] = \begin{bmatrix} d_{0,0} & d_{0,1} & d_{0,2} & d_{0,3} & F_{0,0} & F_{0,1} & F_{0,2} \\ d_{1,0} & d_{1,1} & d_{1,2} & d_{1,3} & F_{1,0} & F_{1,1} & F_{1,2} \end{bmatrix}$$

backward：已知 $\frac{\partial L}{\partial Y}$ ，求 $\frac{\partial L}{\partial D}$ 和 $\frac{\partial L}{\partial S}$ 。

$$\frac{\partial L}{\partial Y} = \begin{bmatrix} \frac{\partial L}{\partial y_{0,0}} & \frac{\partial L}{\partial y_{0,1}} & \frac{\partial L}{\partial y_{0,2}} & \frac{\partial L}{\partial y_{0,3}} & \frac{\partial L}{\partial y_{0,4}} & \frac{\partial L}{\partial y_{0,5}} & \frac{\partial L}{\partial y_{0,6}} \\ \frac{\partial L}{\partial y_{1,0}} & \frac{\partial L}{\partial y_{1,1}} & \frac{\partial L}{\partial y_{1,2}} & \frac{\partial L}{\partial y_{1,3}} & \frac{\partial L}{\partial y_{1,4}} & \frac{\partial L}{\partial y_{1,5}} & \frac{\partial L}{\partial y_{1,6}} \end{bmatrix}, \frac{\partial L}{\partial D} = \begin{bmatrix} \frac{\partial L}{\partial d_{0,0}} & \frac{\partial L}{\partial d_{0,1}} & \frac{\partial L}{\partial d_{0,2}} & \frac{\partial L}{\partial d_{0,3}} \\ \frac{\partial L}{\partial d_{1,0}} & \frac{\partial L}{\partial d_{1,1}} & \frac{\partial L}{\partial d_{1,2}} & \frac{\partial L}{\partial d_{1,3}} \end{bmatrix}$$

以 $\frac{\partial L}{\partial d_{0,0}}$ 为例：

$$\frac{\partial L}{\partial d_{0,0}} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial d_{0,0}} = \frac{\partial L}{\partial Y} \begin{bmatrix} 1 & 0 & 0 & 0 & s_{0,0,0} & s_{0,1,0} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} = \frac{\partial L}{\partial y_{0,0}} + s_{0,0,0} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,0} \frac{\partial L}{\partial y_{0,5}}$$

类似的，可以计算出所有的项：

$$\begin{aligned} \frac{\partial L}{\partial d_{i,j}} &= \frac{\partial L}{\partial y_{i,j}} + s_{i,0,j} \frac{\partial L}{\partial y_{i,4}} + s_{i,1,j} \frac{\partial L}{\partial y_{i,5}} \\ \frac{\partial L}{\partial D} &= \frac{\partial L}{\partial Y}|_{[:,0:4]} + \begin{bmatrix} s_{0,0,0} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,0} \frac{\partial L}{\partial y_{0,5}} & s_{0,0,1} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,1} \frac{\partial L}{\partial y_{0,5}} & s_{0,0,2} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,2} \frac{\partial L}{\partial y_{0,5}} & s_{0,0,3} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,3} \frac{\partial L}{\partial y_{0,5}} \\ s_{1,0,0} \frac{\partial L}{\partial y_{1,4}} + s_{1,1,0} \frac{\partial L}{\partial y_{1,5}} & s_{1,0,1} \frac{\partial L}{\partial y_{1,4}} + s_{1,1,1} \frac{\partial L}{\partial y_{1,5}} & s_{1,0,2} \frac{\partial L}{\partial y_{1,4}} + s_{1,1,2} \frac{\partial L}{\partial y_{1,5}} & s_{1,0,3} \frac{\partial L}{\partial y_{1,4}} + s_{1,1,3} \frac{\partial L}{\partial y_{1,5}} \end{bmatrix} \\ &\quad + \begin{bmatrix} \frac{\partial L}{\partial y_{0,4}} & \frac{\partial L}{\partial y_{0,5}} \\ \frac{\partial L}{\partial y_{1,4}} & \frac{\partial L}{\partial y_{1,5}} \end{bmatrix} S \end{aligned}$$

注意中间这一项的 `shape` 是 `[2, 1, 2]`，左乘 `S` 之后得到一个 `shape = [2, 1, 4]`，需要把中间的这个维度 `squeeze` 掉再和第一项相加。因为 `Interaction` 是两个输入，因此还需要计算关于另一个输入的导数：

$$\frac{\partial L}{\partial S} = \begin{bmatrix} [\frac{\partial L}{\partial s_{0,0,0}}, \frac{\partial L}{\partial s_{0,0,1}}, \frac{\partial L}{\partial s_{0,0,2}}, \frac{\partial L}{\partial s_{0,0,3}}] & [\frac{\partial L}{\partial s_{0,1,0}}, \frac{\partial L}{\partial s_{0,1,1}}, \frac{\partial L}{\partial s_{0,1,2}}, \frac{\partial L}{\partial s_{0,1,3}}] \\ [\frac{\partial L}{\partial s_{1,0,0}}, \frac{\partial L}{\partial s_{1,0,1}}, \frac{\partial L}{\partial s_{1,0,2}}, \frac{\partial L}{\partial s_{1,0,3}}] & [\frac{\partial L}{\partial s_{1,1,0}}, \frac{\partial L}{\partial s_{1,1,1}}, \frac{\partial L}{\partial s_{1,1,2}}, \frac{\partial L}{\partial s_{1,1,3}}] \end{bmatrix}$$

以第一项为例：

$$\frac{\partial L}{\partial s_{0,0,0}} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial s_{0,0,0}} = \frac{\partial L}{\partial Y} \begin{bmatrix} 0 & 0 & 0 & 0 & d_{0,0} & 0 & s_{0,1,0} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} = d_{0,0} \frac{\partial L}{\partial y_{0,4}} + s_{0,1,0} \frac{\partial L}{\partial y_{0,6}}$$

由于 `S` 比较长，省略一些中间过程：

$$\frac{\partial L}{\partial S} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

其中：

$$\begin{aligned} A &= \begin{bmatrix} \frac{\partial L}{\partial y_{0,4}} & \frac{\partial L}{\partial y_{0,6}} \end{bmatrix} \begin{bmatrix} d_{0,0} & d_{0,1} & d_{0,2} & d_{0,3} \\ s_{0,1,0} & s_{0,1,1} & s_{0,1,2} & s_{0,1,3} \end{bmatrix}, B = \begin{bmatrix} \frac{\partial L}{\partial y_{0,5}} & \frac{\partial L}{\partial y_{0,6}} \end{bmatrix} \begin{bmatrix} d_{0,0} & d_{0,1} & d_{0,2} & d_{0,3} \\ s_{0,0,0} & s_{0,0,1} & s_{0,0,2} & s_{0,0,3} \end{bmatrix} \\ C &= \begin{bmatrix} \frac{\partial L}{\partial y_{1,4}} & \frac{\partial L}{\partial y_{1,6}} \end{bmatrix} \begin{bmatrix} d_{1,0} & d_{1,1} & d_{1,2} & d_{1,3} \\ s_{1,1,0} & s_{1,1,1} & s_{1,1,2} & s_{1,1,3} \end{bmatrix}, D = \begin{bmatrix} \frac{\partial L}{\partial y_{1,5}} & \frac{\partial L}{\partial y_{1,6}} \end{bmatrix} \begin{bmatrix} d_{1,0} & d_{1,1} & d_{1,2} & d_{1,3} \\ s_{1,0,0} & s_{1,0,1} & s_{1,0,2} & s_{1,0,3} \end{bmatrix} \end{aligned}$$

整理得：

3 DCN Interaction

forward:

```

1 # dense_features.shape = [B, 128]
2 # sparse_features.shape = [B, 26, 128]
3 # combined_values.shape = [B, 27, 128]
4 combined_values =
5     torch.cat((dense_features.unsqueeze(1), sparse_features),
6               dim=1)
7
8 # size B X (F*D + D)
9 concatenated_dense = crossnet(combined_values.reshape([B, -1]))

```

LowRankCrossNet (当前参考代码使用的是这个)

1. 初始化

- $N = in_features = (num_sparse_features + 1) * embedding_dim$
- num_layers
- $K = low_rank$
- 初始化 num_layers 个 w, v, b
- $w.shape = [in_features, low_rank]$
- $v.shape = [low_rank, in_features]$
- $b.shape = [in_features, 1]$

2. 计算

$$x_{l+1} = \underset{[B, N, 1]}{\hat{x}} \otimes \left[\underbrace{\underset{[N, K]}{w_l} * \begin{pmatrix} \underset{[K, N]}{v_l} * \underset{[B, N, 1]}{x_l} \\ [B, K, 1] \end{pmatrix}}_{[B, N, 1]} + \underset{[N, 1]}{b_l} \right] + \underset{[B, N, 1]}{x_l}$$

其中 $\hat{x} = input.unsqueeze(2), x_0 = \hat{x}, l \in [0, num_layers-1]$ 。⊗ 代表 element-wise multiplication, * 代表 matrix multiplication。

3. 最终输出 shape 为 $[B, in_features]$, 因此在使用这个 Interaction 的时候, OverArch 的第一层就是 $[in_features, 1024]$ 。

CrossNet (这个是论文新提出的方法)

1. 初始化

- $N = in_features = (num_sparse_features + 1) * embedding_dim$
- $num_layers = wait\ for\ input$
- 初始化 num_layers 个 w, v, b
- $w.shape = [in_features, in_features]$
- $b.shape = [in_features, 1]$

2. 计算

$$x_{l+1} = \underset{[B, N, 1]}{\hat{x}} \otimes \left[\underbrace{\underset{[N, N]}{w_l} * \underset{[B, N, 1]}{x_l}}_{[B, N, 1]} + \underset{[N, 1]}{b_l} \right] + \underset{[B, N, 1]}{x_l}$$

其中 $\hat{x} = input.unsqueeze(2), x_0 = \hat{x}, l \in [0, num_layers-1]$ 。⊗ 代表 element-wise multiplication, * 代表 matrix multiplication。

3. 最终输出 shape 为 $[B, in_features]$, 因此在使用这个 Interaction 的时候, OverArch 的第一层就是 $[in_features, 1024]$ 。

VectorCrossNet

1. 初始化

- $N = in_features = (num_sparse_features + 1) * embedding_dim$
- $num_layers = wait\ for\ input$

- 初始化 num_layers 个 w, v, b
- $w.shape = [in_features, 1]$
- $b.shape = [in_features, 1]$

2. 计算

$$x_{l+1} = \hat{x}_{[B,N,1]} * \left(\begin{matrix} w_l & \bullet & x_l \\ [N,1] & [B,N,1] \\ [B,1,1] \end{matrix} \right) + \begin{matrix} b_l \\ [N,1] \end{matrix} + \begin{matrix} x_l \\ [B,N,1] \end{matrix}$$

其中 $\hat{x} = input.unsqueeze(2)$, $x_0 = \hat{x}$, $l \in [0, num_layers - 1]$ 。• 代表 `tensordot`, * 代表 `matrix multiplication`。

LowRankMixtureCrossNet

1. 初始化

- $N = in_features = (num_sparse_features + 1) * embedding_dim$
- num_layers
- $K = num_experts$
- $L = low_rank$
- $activation = for\ example\ relu$
- 初始化 num_layers 个 u, v, b, c
- $u.shape = [num_experts, in_features, low_rank]$
- $v.shape = [num_experts, low_rank, in_features]$
- $b.shape = [in_features, 1]$
- $c.shape = [num_experts, low_rank, low_rank]$
- $num_experts$ 个线性层 $g = Linear(N, 1, bias = False)$, 对所有的 LowRankMixtureCrossNet 层共用。

2. 计算

- $\hat{x} = input.unsqueeze(2), [B, N, 1]$
- $x_l[B, N]$ 分别经过 $num_experts$ 个线性层, 每层输出一个 $[B, 1]$, stack 到一起组成 $gating_l[B, K, 1]$
- 用 $f(x)$ 表示激活函数, 计算

$$expert_l[e] = \hat{x} * [u_l[e] * f(c_l[e] * f(v_l[e] * x_l)) + b_l], e \in [0, num_experts - 1]$$

- 此时 $expert_l[e].shape = [B, N, 1]$, squeeze 掉最后一维之后 $expert_l[e].shape = [B, N]$ 。
- 把所有的 $expert_l[e]$ 叠起来, $experts_l.shape = [B, N, K]$:

$$experts_l = stack([expert_l[e] \ e \ from \ 0 \ to \ num_experts], 2)$$

- 再计算

$$x_{l+1} = experts_l * softmax(gating_l) + x_l$$

- 以上这是一层 LowRankMixtureCrossNet, 每层都是一样的, 首尾相接, 最后输出的是一个 $[B, N, 1]$, 后续处理与其他的 CrossNet 类似。

4 Projection Interaction

简单描述就是来自 SparseArch 和 DenseArch 的 feature 先拼在一起, 然后分别过两个 mlp, 然后再 bmm, 再 concatenate。forward:

```

1 # combined_values.shape = [B, 27, 128]
2 # dense_features.shape = [B, 128], D=128
3 combined_values = torch.cat(
4     (dense_features.unsqueeze(1), sparse_features), dim=1)
5
6 interaction_branch1_out = self.interaction_branch1(
7     torch.reshape(combined_values, (B, -1))) # [B, 256]
8

```

```

9  interaction_branch2_out = self.interaction_branch2(
10      torch.reshape(combined_values, (B, -1))) # [B, 512]
11
12  interactions = torch.bmm(
13      interaction_branch1_out.reshape([B, -1, D]), # [B, 2, 128]
14      interaction_branch2_out.reshape([B, D, -1]), # [B, 128, 4]
15  )
16
17  # interactions_flat.shape = [B, 8]
18  interactions_flat = torch.reshape(interactions, (B, -1))
19  # concatenated_dense.shape = [B, 136]
20  concatenated_dense = torch.cat((dense_features, interactions_flat),
21                                  dim=1)

```

一些关于初始化参数的细节

- $num_sparse_features = 26$ 。
- $interaction_branch1 = DenseArch\{-?\}$ ，需要给定超参来指定，例如 D_1, D_2, \dots, D 。
- $interaction_branch2 = DenseArch\{-?\}$ ，需要给定超参来指定，例如 E_1, E_2, \dots, E 。
- 这两个 DenseArch 最后一层的输出 size D 和 E 必须都是 $embedding_dim$ 的整数倍，但 D 和 E 可以不同。
- 设 $projected_dim_X = interaction_branchX_layer_sizes[-1] // embedding_dim$ ，则使用 ProjectionInteraction 是的 OverArch 的第一层的输入 size 就是 $embedding_dim + projected_dim_1 * projected_dim_2 = 128 + 2 * 4 = 136$
- $interaction_branch$ 的具体结构和层数是超参控制的。

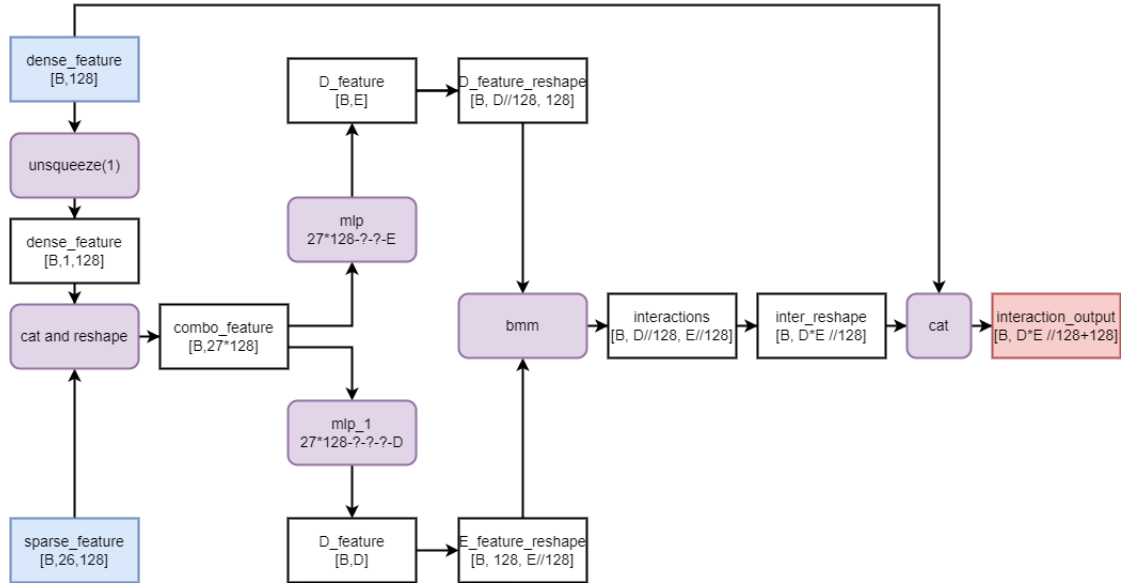


Figure 1: Projection Interaction