

## Idea

- Goal:** Build *least squares regression models* over training datasets defined by arbitrary join queries over databases.
- Observation:** Joins entail a *high degree of redundancy* in both computation and data representation, which is not required for an end-to-end solution to learning over joins.
- Solution:** **F** uses gradient descent to learn the model parameters in one pass over a factorized database view.

## Recap on Linear Regression

The model:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n = \sum_{i=0}^n \theta_i x_i, \text{ where } x_0 = 1$$

Least Squares Objective Function gives a measure of the error between the actual value  $y^{(i)}$  and the model  $h_{\theta}(x^{(i)})$ :

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

To minimise the error, we calculate the gradient and update the parameters iteratively as follows:

$$\begin{aligned} \forall 0 \leq j \leq n : \theta_j &:= \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta) \\ &:= \theta_j - \alpha \sum_{i=1}^m \left( \sum_{k=0}^n \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}. \end{aligned}$$

By letting  $\theta_0 = -1$ , the gradient becomes:

$$\sum_{i=1}^m \left( \sum_{k=0}^n \theta_k x_k^{(i)} \right) x_j^{(i)}$$

## Highlights of Our Solution

- Decouple data-dependent (aggregate) computation from data-independent parameter convergence:
$$\forall 0 \leq j \leq n : S_j = \sum_{k=0}^n \theta_k \times \text{Cofactor}[k, j]$$

where  $\text{Cofactor}[k, j] = \sum_{i=1}^m x_k^{(i)} x_j^{(i)}$

- Compute the cofactor matrix in one pass over the factorized join.
- Time and space complexity:  $O(|D|^{fhtw(Q)})$ , where  $fhtw(Q)$  is the fractional hypertree width of the query hypergraph.
- Principles also applicable to polynomial regression, factorized machines, and various regularizers.

## Factorized Join Example

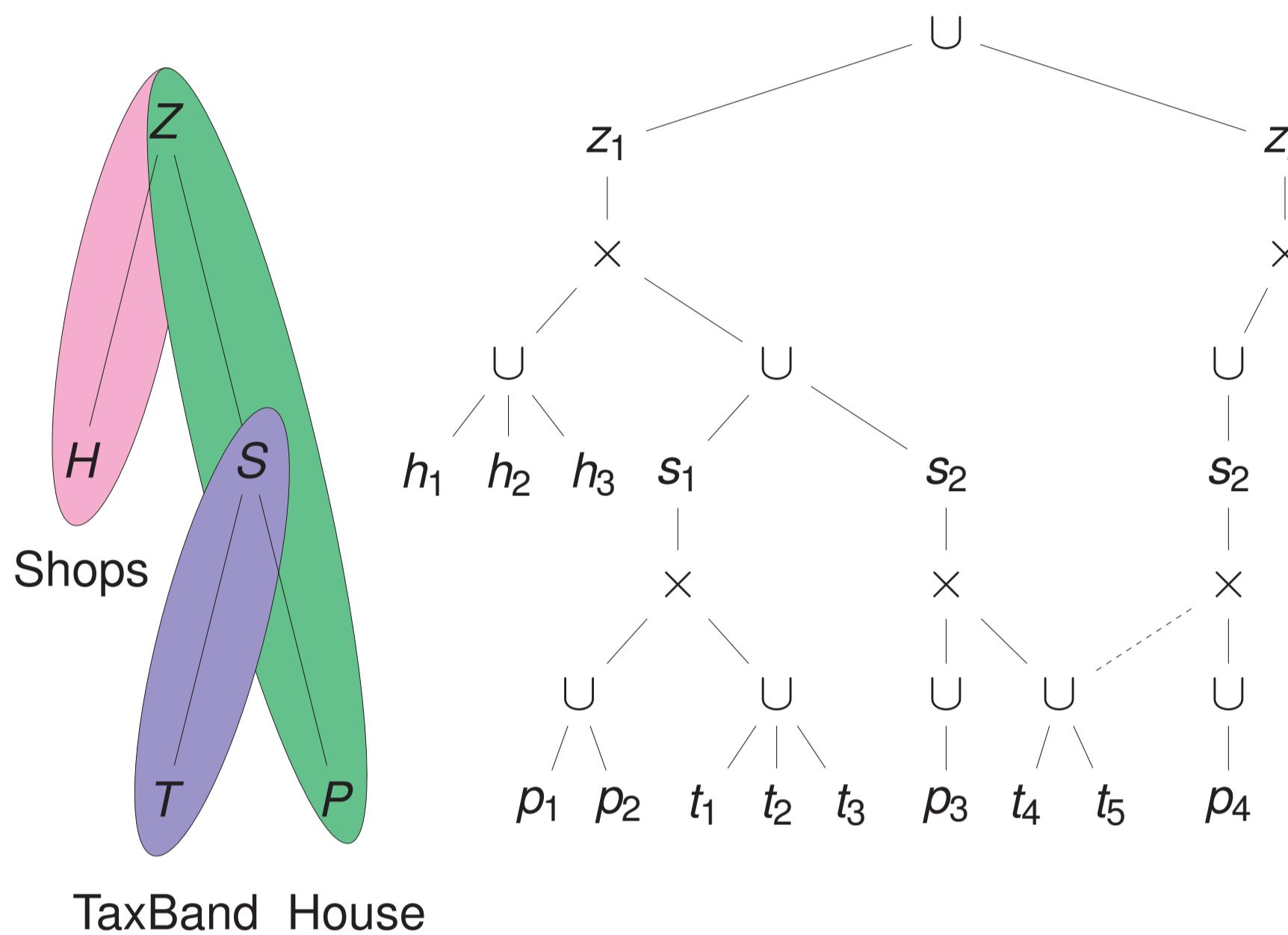
- Natural join of relations House, Shop, TaxBand.
- Redundancy in flat join:  $z_1$  occurs in 24 tuples,  $h_1$  to  $h_3$  occur in eight tuples each and are paired with the same combination of values for  $P$ ,  $T$  and  $S$ .
- Avoid to explicitly materialise the local products.

| Shops     | House         | Shops $\bowtie$ House $\bowtie$ TaxBand |
|-----------|---------------|---|
| Z H       | Z S P         | Z H S P T                               |
| $z_1 h_1$ | $z_1 s_1 p_1$ | $z_1 h_1 s_1 p_1 t_1$                   |
| $z_1 h_2$ | $z_1 s_1 p_2$ | $z_1 h_1 s_1 p_1 t_2$                   |
| $z_1 h_3$ | $z_1 s_2 p_3$ | $z_1 h_1 s_1 p_1 t_3$                   |
| $z_2 h_4$ | $z_2 s_2 p_4$ | $z_1 h_1 s_1 p_2 t_1$                   |
|           |               | $z_1 h_1 s_1 p_2 t_2$                   |
|           |               | $z_1 h_1 s_2 p_3 t_3$                   |
|           |               | $z_1 h_1 s_2 p_3 t_4$                   |
|           |               | $z_2 h_4 s_2 p_4 t_5$                   |
|           |               | $z_2 h_4 s_2 p_4 t_5$                   |

| TaxBand   |   |
|-----------|---|
| S         | T |
| $s_1 t_1$ |   |
| $s_1 t_2$ |   |
| $s_1 t_3$ |   |
| $s_2 t_4$ |   |
| $s_2 t_5$ |   |

the above for  $h_2$  and  $h_3$



(a) The three relations of database **D** and natural join  $Q(\mathbf{D})$ .

(b) F-tree **F**.

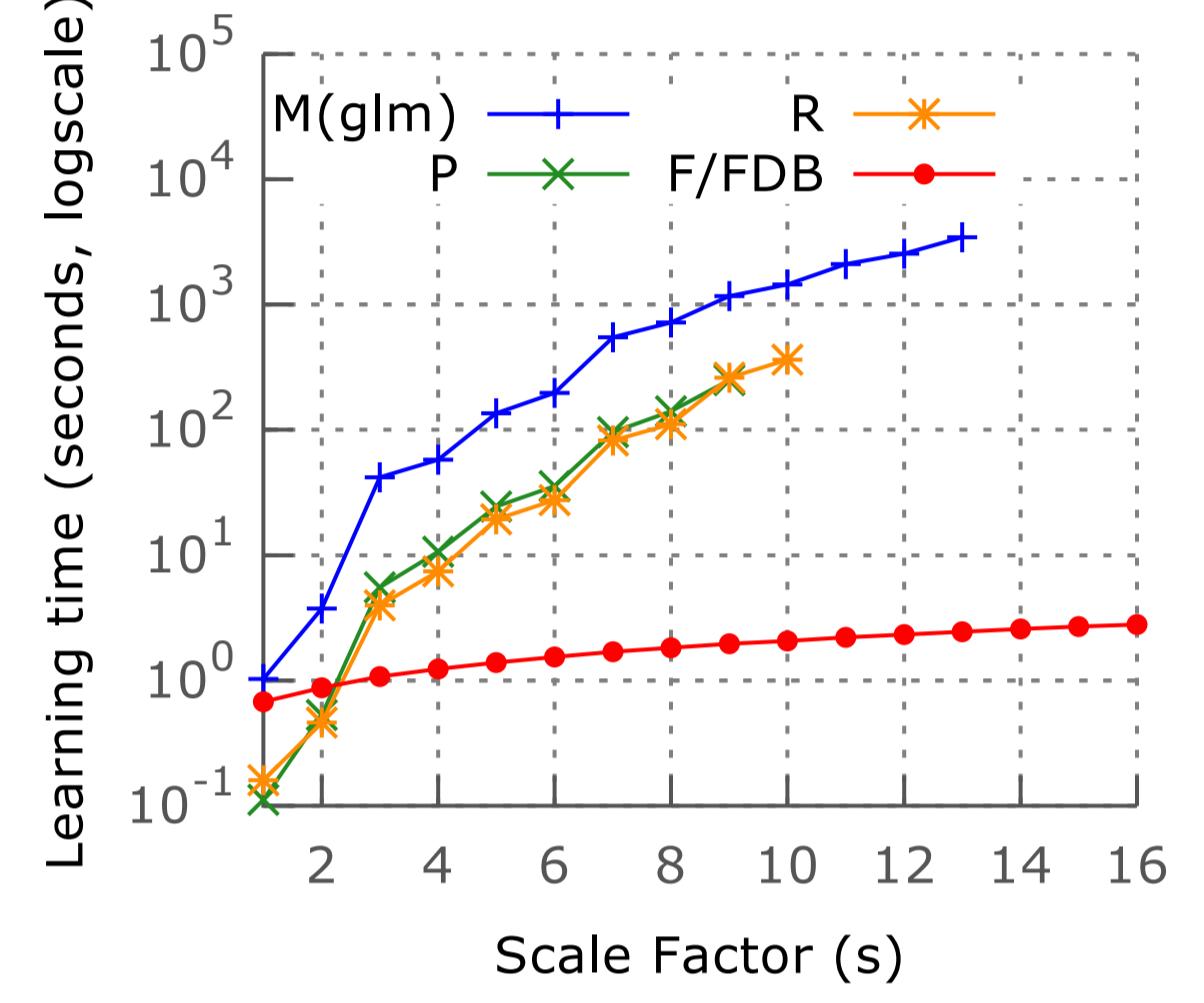
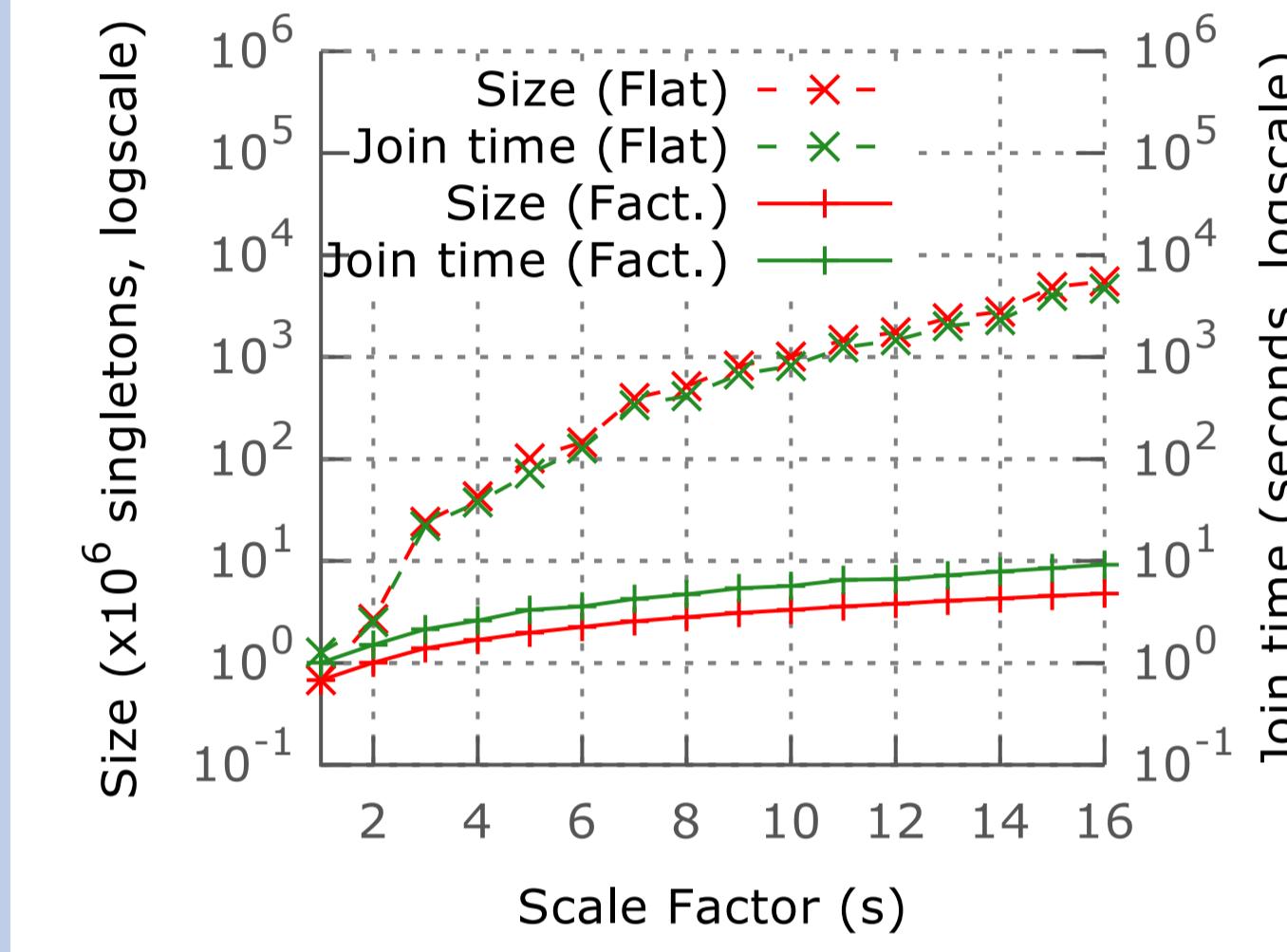
(c) Factorization  $F(\mathbf{D})$  of  $Q(\mathbf{D})$  over **F**.

## F: Flavors and Competitors

- F/FDB:** Cofactors computed in one pass over the materialized factorized join.
  - F:** Factorized join and cofactor computation intermixed.
  - F/SQL:** SQL-encoding of **F**, intertwining joins and cofactors in one query.
- Competitors: **R** (QR decomp.), **Python StatsModels** (ols), **MADlib** (ols, glm).

## Experiments on Synthetic Data

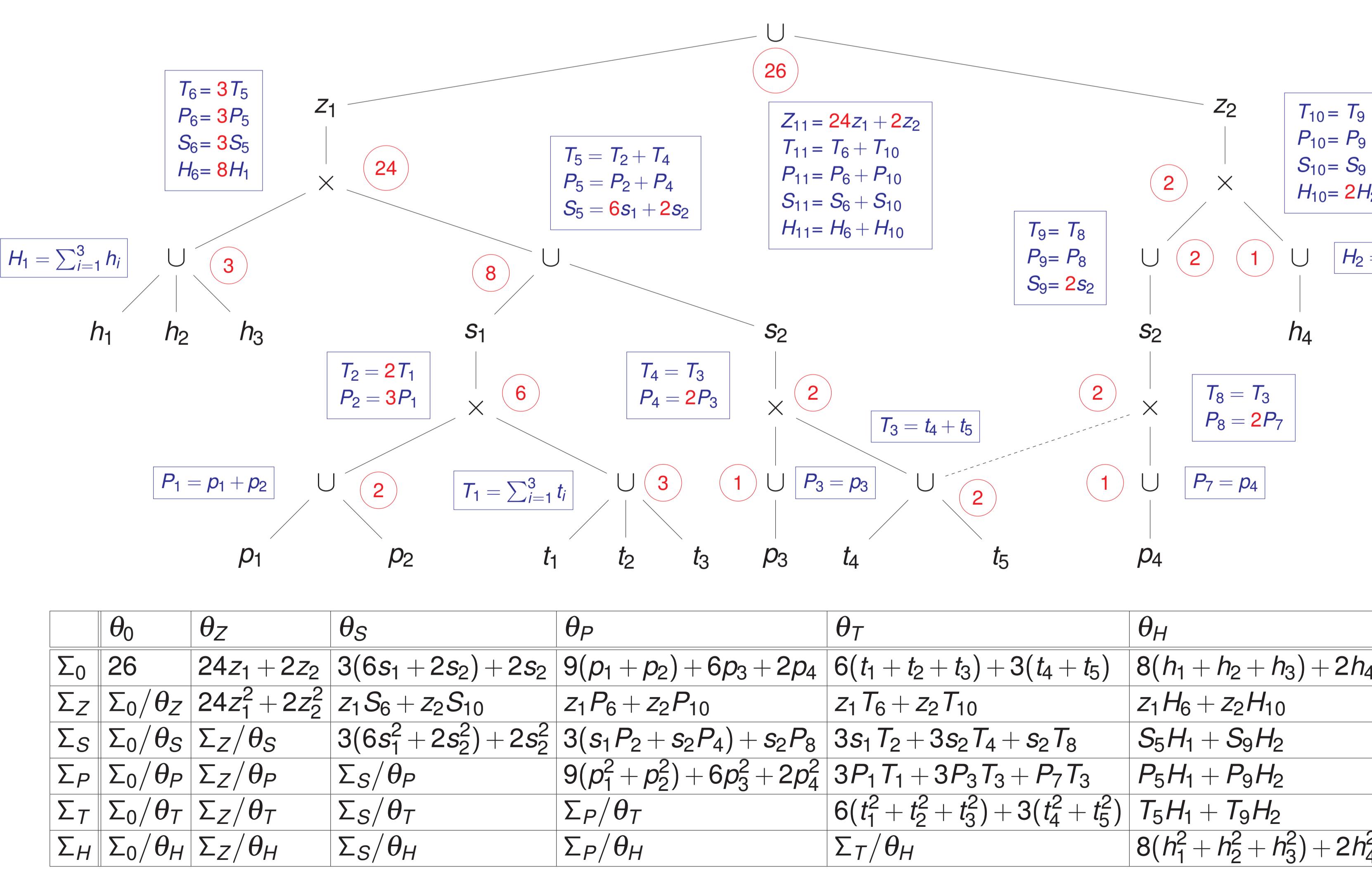
- (left) Factorized join: Compression ratio and speedup.
- (right) Learning: Speedup.



## F Computes the Cofactors of Model Parameters in One Pass over the Factorized Join

Reformulate the sum-aggregates defining cofactors to save computation. **This rewriting is already performed by the factorization of the join!**

- Compute **occurrence count** of each value and **weighted sums** of each attribute at each node.
- Complete incrementally the cofactor matrix.



## Experiments on Real Data

|                      | US retailer | LastFM (1)    | LastFM (2)  | MovieLens   |
|----------------------|-------------|---------------|-------------|-------------|
| # parameters         | 31          | 6             | 10          | 27          |
| Join Size            | Factorized  | 97,134,675    | 376,402     | 315,818     |
| Flat                 |             | 2,585,046,352 | 369,986,292 | 590,793,800 |
| Compression          |             | 26.61×        | 982.86×     | 1870.68×    |
| Join Time            | Factorized  | 36.03         | 4.79        | 9.94        |
| Flat                 |             | 249.41        | 54.25       | 61.33       |
| <b>F and M</b>       |             | 0             | 0           | 0           |
| Import R             |             | 1189.12*      | 155.91      | 276.77      |
| P                    |             | 1164.40*      | 179.16      | 328.97      |
| <b>F</b>             |             | 9.69          | 0.53        | 0.89        |
| Learn M (glm)        |             | 2671.88       | 572.88      | 746.50      |
| R                    |             | 810.66*       | 268.04      | 466.52      |
| P                    |             | 1199.50*      | 35.74       | 148.84      |
| <b>F</b>             |             | 16.29         | 0.11        | 0.25        |
| <b>F/FDB</b>         |             | 45.72         | 5.32        | 10.83       |
| <b>F/SQL</b>         |             | 108.81        | 0.58        | 2.00        |
| Total Time           | M (ols)     | 680.60        | 152.37      | 196.60      |
|                      | M (glm)     | 2921.29       | 627.13      | 807.83      |
|                      | R           | 2249.19*      | 478.20      | 804.62      |
|                      | P           | 2613.31*      | 269.15      | 539.14      |
| <b>F vs. M (ols)</b> |             | 41.78×        | 1385.18×    | 786.40×     |
| <b>F vs. M (glm)</b> |             | 179.33×       | 5701.18×    | 3231.32×    |
| <b>F vs. R</b>       |             | 138.07×       | 4347.27×    | 3218.48×    |
| <b>F vs. P</b>       |             | 57.16×        | 50.59×      | 49.78×      |