# Processing of Probabilistic Skyline Queries Using MapReduce

Qingfu Wen

School of Software, Tsinghua University

qingfu.wen@gmail.com

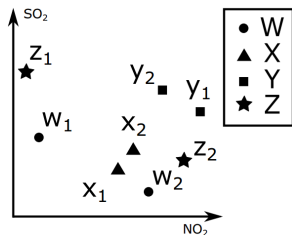Author: Yoonjae Park, Jun-Ki Min, Kyuseok Shim

November 14, 2015

# Probabilistic Skylines



| Object | Instance | NO$_2$ | SO$_2$ | Probability |
|--------|----------|--------|--------|-------------|
| W | $w_1$ | 10 | 40 | 0.5 |
|   | $w_2$ | 75 | 10 | 0.4 |
| X | $x_1$ | 55 | 20 | 0.2 |
|   | $x_2$ | 65 | 30 | 0.2 |
| Y | $y_1$ | 95 | 60 | 0.8 |
|   | $y_2$ | 80 | 70 | 0.2 |
| Z | $z_1$ | 5 | 80 | 0.5 |
|   | $z_2$ | 90 | 25 | 0.5 |



$$P_{sky}(y_1) = P(y_1)(1 - P(w_1) - P(w_2))(1 - P(x_1) - P(x_2))(1 - P(z_2))$$
$$= 0.024$$
$$P_{sky}(y_2) = 0.012$$
$$P_{sky}(Y) = P_{sky}(y_1) + P_{sky}(y_2) = 0.036$$
$$P_{sky}(W) = 0.9$$
$$P_{sky}(X) = 0.4$$
$$P_{sky}(Z) = 0.74$$

# Probabilistic Skylines

## Probabilistic Skyline Problem

For a set of uncertain objects $\mathbb{D}$ and a probability threshold $T_p$, the probabilistic skyline $pSL(\mathbb{D}, T_p)$, is the set of all objects whose skyline probabilities are at least $T_p$, $pSL(\mathbb{D}, T_p) = \{U \in \mathbb{D} | P_{sky}(U) \geq T_p\}$.
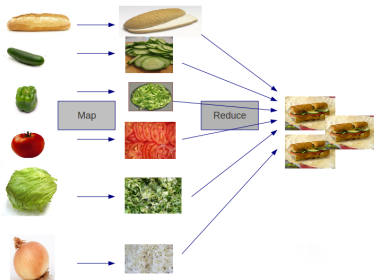
The discrete model:

$$P_{sky}(U) = \sum_{u_i \in U} P_{sky}(u_i) = \sum_{u_i \in U} (P(u_i) \times \prod_{V \in \mathbb{D}, V \neq U} (1 - \sum_{v_j \in V, v_j \prec u_i} P(v_j)))$$
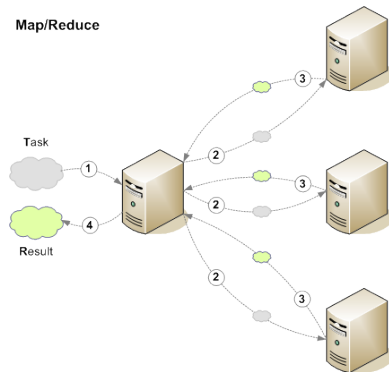
The continuous model:

$$P_{sky}(u_i) = \int_U Uf(u) \times \prod_{V \in \mathbb{D}, V \neq U} (1 - \int_V Vf(v) 1(v \prec u)\, dv)\, du$$
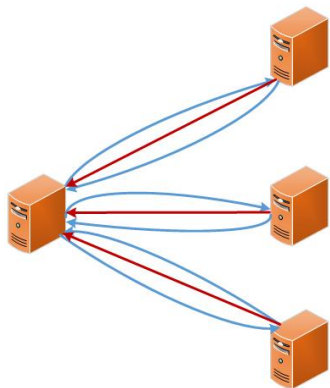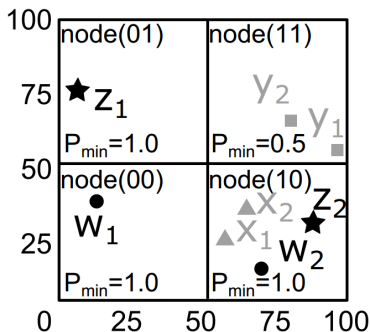
# What is MapReduce?



Map/Reduce

1. local computing candidate sets.
2. merge candidate sets, broadcast and local computing, reduce probabilities.

# Early Pruning Techniques

## Lemma (Zero-probability Filtering)

$$P_{sky}(U) = \sum_{u_i \in U} P(u_i) \times \prod_{V \in \mathbb{D}, V \neq U} (1 - \sum_{v_j \in V, v_j \prec u_i} P(v_j))$$

$$delete\ u_i\ if\ \prod_{V \in \mathbb{D}, V \neq U} (1 - \sum_{v_j \in V, v_j \prec u_i} P(v_j)) = 0.$$



## Example (Zero-probability Filtering)

$$P_{sky}(y_1) = P(y_1)(1 - P(w_1) - P(w_2))(1 - P(x_1) - P(x_2))(1 - P(z_2))$$
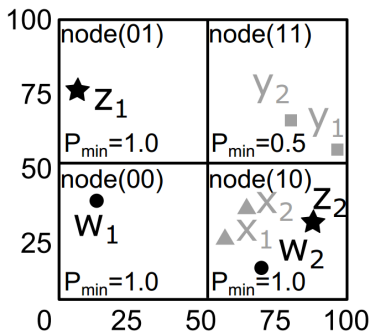$$\leq P(y_1)(1 - P(w_1)) = 0.4$$
$$P_{sky}(y_2) \leq 0.1$$
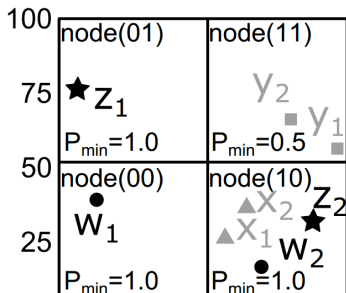$$P_{sky}(Y) = P_{sky}(y_1) + P_{sky}(y_2) \leq 0.5 \leq T_p$$

# Early Pruning Techniques

## Lemma (Upper-bound Filtering)

$$\beta(U, \mathbb{S}, R(u_i)) = \frac{\prod\limits_{V \in \mathbb{S}} (1 - \sum\limits_{v_j \in V, v_j \prec R(u_i).min} P(v_j))}{1 - \sum\limits_{v_k \in U, v_k \prec R(u_k).min} P(v_k)}$$

$$up(u_i, U, \mathbb{S}, R(u_i)) = P(u_i) \times \beta(U, \mathbb{S}, R(u_i)).$$



## Example (Upper-bound Filtering)

$P_{sky}(y_1) = P(y_1)(1 - P(w_1) - P(w_2))(1 - P(x_1) - P(x_2))(1 - P(z_2))$

$\qquad \leq P(y_1)(1 - P(w_1)) = 0.4$

$P_{sky}(y_2) \leq 0.1$

$P_{sky}(Y) = P_{sky}(y_1) + P_{sky}(y_2) \leq 0.5 \leq T_p$

# Early Pruning Techniques

**Lemma (Dominance-Power Filtering)**

$DP(v_j) = \prod_{1}^{d}(b(k) - v_j(k)) = 0, b(k) = max\{v_1(k), \cdots, v_n(k)\}$.
$DP(V) = \sum_{v_j \in V}(P(v_j) \times DP(v_j))$.
$\mathbb{F}$ is topK DP set, $\sum_{u_i \in U} P(u_i) \times \prod_{V \in \mathbb{F}, V \neq U}(1 - \sum_{v_j \in V, v_j \prec u_i} P(v_j)) < T_p$, U is
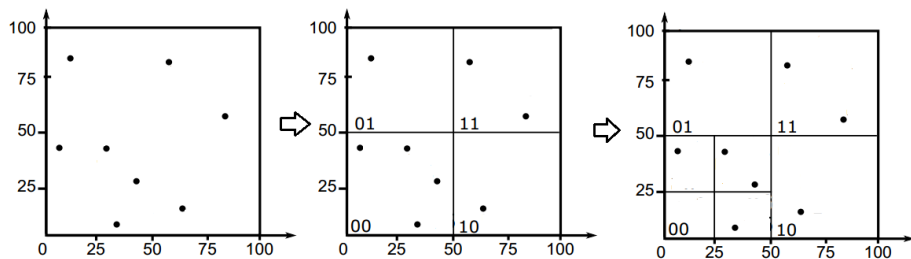not a probabilistic skyline Object.

Example (Dominance-Power Filtering)

$P_{sky}(y_1) = P(y_1)(1 - P(w_1) - P(w_2))(1 - P(x_1) - P(x_2))(1 - P(z_2))$
$\qquad \leq P(y_1)(1 - P(w_1)) = 0.4$
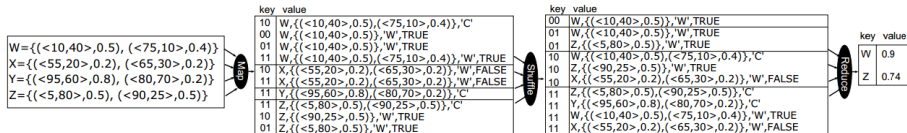$P_{sky}(y_2) \leq 0.1$
$P_{sky}(Y) = P_{sky}(y_1) + P_{sky}(y_2) \leq 0.5 \leq T_p$

# PSQtree for Pruning



- generate PSQtree using a random sample $\mathbb{S}$ of $\mathbb{D}$
- traverse PSQtree for computing $P_{sky}(node.min)$
- zero-probability filtering
- upper-bound filtering
- partitioning objects by PSQtree

**Function** PS-QPF-MR(𝔻, $T_p$, $\rho$)
𝔻: uncertain dataset, $T_p$: probability threshold, $\rho$: split threshold
**begin**
1. $\mathbb{S}$ = Sample(𝔻);
2. $PSQtree$ = GenQtree($\mathbb{S}$, $\rho$);
3. Broadcast $PSQtree$; Broadcast $T_p$;
4. $pSL$ = RunMapReduce(PS-QPFC-MR, 𝔻);
5. **return** $pSL$;
**end**

**Function** PS-QPF-MR.setup()
**begin**
1. $H$ = InitMinHeap(); $PSQtree$ = LoadPSQtree();
**end**

**Function** PS-QPFC-MR.map($U$)
$U$: an uncertain object
**begin**
1. $T_p$ = LoadThreshold();
2. $U'$ = ZeroProb($U$, $PSQtree$);
3. $upper$ = UpperBound($U'$, $PSQtree$);
4. $cand$ = FALSE;
5. **if** $upper \geq T_p$ **then**
6. $\quad cand$ = DP-Filter($U'$, $T_p$, $H$);
7. $\quad$ **if** $cand$ = **True** and $n_\ell = n(U'.max)$ **then** emit($n(U'.max)$, ($U'$, 'C'));
8. **for each** leaf node $n_\ell$ in $PSQtree$ **do**
9. $\quad$ **if** $cand$ = **True** and $n_\ell = n(U'.max)$ **then** continue;
10. $\quad I$ = NewList();
11. $\quad$ **for each** $u_i$ in $U'$ **do**
12. $\quad\quad$ **if** $n(u_i) \preceq n_\ell$ **then**
13. $\quad\quad\quad I$.add($u_i$));
14. $\quad$ emit($n_\ell$, ($I$, 'W', $cand$))
**end**

**Function** PS-QPFC-MR.reduce($n_\ell$, $L$)
**begin**
1. $(L_C, L_W^T, L_W^F)$ = SplitList($L$);
2. $T_p$ = LoadThreshold();
3. **for each** object $U$ in $L_C$ **do**
4. $\quad skyline\_prob$ = SkylineProb($U$, $L_C$, $L_W^T$, $L_W^F$);
5. $\quad$ **if** $skyline\_prob \geq T_p$ **then**
6. $\quad\quad$ emit( $U$, $skyline\_prob$ );
**end**

# MapReduce Algorithms with PSQtree

- Reducing network overhead by clustering
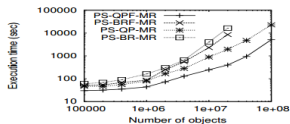- Workload balancing of reduce functions

## Experiments

- 50 machines with Intel i3 3.3GHz CPU and 4GB, Linux
- 200 machines with Intel Xeon 2.5GHz CPU and 3.75GB, Amazon EC2
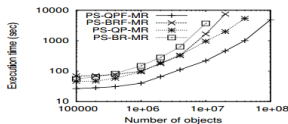- Java 1.6, Hadoop 1.2.1

| Algorithm | Description |
|-----------|-------------|
| PS-QP-MR | The algorithm with quadtree partitioning |
| PS-QPF-MR | The algorithm with quadtree partitioning and filtering |
| PS-BR-MR | The algorithm with random partitioning |
| PS-BRF-MR | The algorithm with random partitioning and filtering |
| PSMR | The state-of-the-art algorithm |

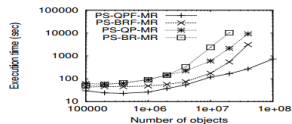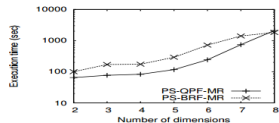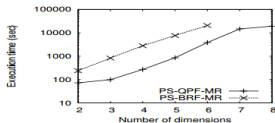| Parameter | Range | Default |
|-----------|-------|---------|
| No. of samples($\mathbb{S}$) | 1000~10,000 | 1000 for PS-QPF-MR 2000 for PS-QP-MR 10000 for PS-BRF-MR |
| No. of dominating objects($\mathbb{F}$) | 1000~10,000 | 100 for PS-QPF-MR 1000 for PS-BRF-MR |
| No. of objects($\mathbb{D}$) | $10^5 \sim 10^8$ | $10^7$ |
| No. of dimensions($d$) | $2 \sim 8$ | 4 |
| Probability threshold ($T_p$) | 0.1~0.6 | 0.3 |
| No. of inst. per object($\ell$) | $1 \sim 400$ | 40 |
| No. of machines ($t$) | 10~200 | 25 |

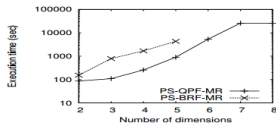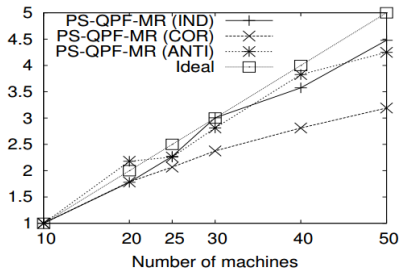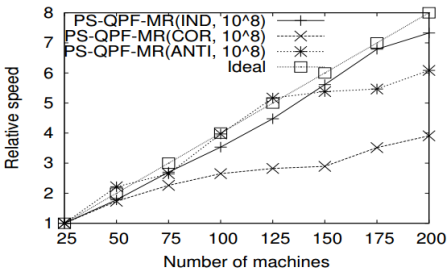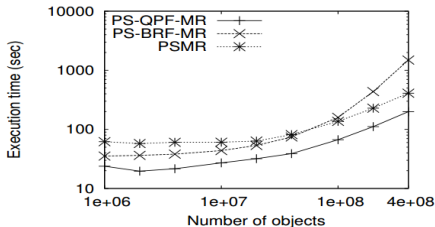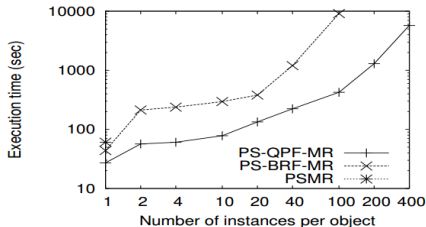(a) ANTI  (b) IND  (c) COR

(a) With our cluster

(b) With Amazon EC2

## Conclusion

- probabilistic skyline query for both discrete and continuous models
- zero-probability, the upper-bound, and dominance power filtering techniques
- using a PSQtree to distribute the instances of objects effectively
- a single MapReduce phase algorithm PS-QPF-MR and grouping optimization

Q & A