

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

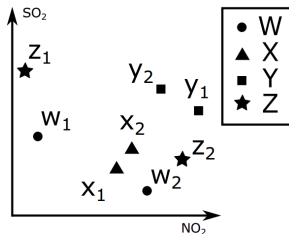
November 16, 2015



Probabilistic Skylines



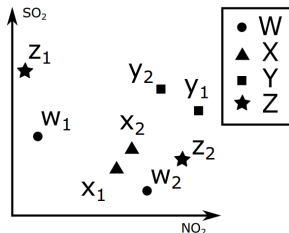
Object	Instance	NO ₂	SO ₂	Probability
W	w ₁	10	40	0.5
	w ₂	75	10	0.4
X	x ₁	55	20	0.2
	x ₂	65	30	0.2
Y	y ₁	95	60	0.8
	y ₂	80	70	0.2
Z	z ₁	5	80	0.5
	z ₂	90	25	0.5



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



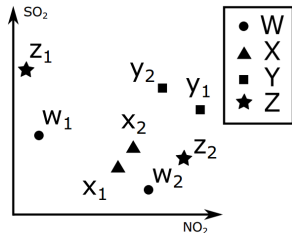
question

Which Object's skyline probability is larger than 0.6?

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 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 , $pSL(\mathbb{D}, T_p) = \{U \in \mathbb{D} | P_{sky}(U) \geq T_p\}$.

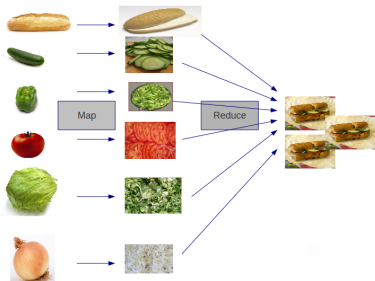
The discrete model:

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

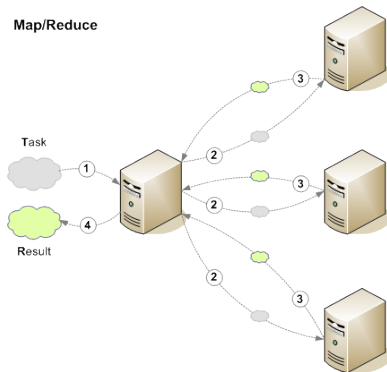
The continuous model:

$$P_{sky}(u_i) = \int_U U f(u) \times \prod_{V \in \mathbb{D}, V \neq U} (1 - \int_V V f(v) 1(v \prec u) dv) du$$

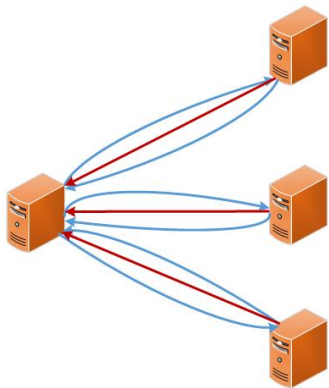
What is MapReduce?



Map/Reduce



PSMR: The State-of-the-art Algorithm



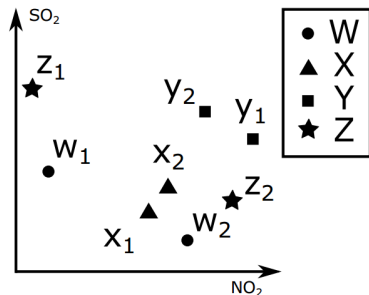
- 1 local computing candidate sets.
- 2 merge candidate sets, broadcast and local computing, reduce probabilities.

Early Pruning Techniques

Lemma (Zero-probability Filtering)

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

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



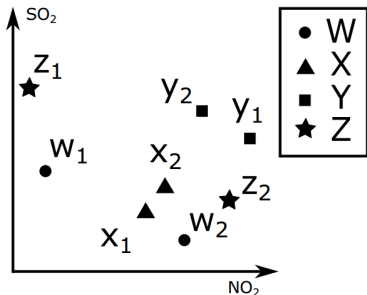
Example (Zero-probability Filtering)

we can see, x_1 dominate y_1 , if $P(x_1) = 1$, then $P_{sky}(y_1) = 0$

Early Pruning Techniques

Lemma (Upper-bound Filtering)

$$\beta(U, \mathbb{S}, R(u_i)) = \frac{\prod_{v \in \mathbb{S}} (1 - \sum_{v_j \in V, v_j \prec R(u_i).min} P(v_j))}{1 - \sum_{v_k \in U, v_k \prec R(u_k).min} P(v_k)}$$
$$up(u_i, U, \mathbb{S}, R(u_i)) = P(u_i) \times \beta(U, \mathbb{S}, R(u_i)).$$



Example (Upper-bound Filtering)

$$P_{sky}(y_1) = P(y_1)(1 - P(w_1) - P(w_2))(1 - P(x_1) - P(x_2))(1 - P(z_2))$$
$$\leq P(y_1)(1 - P(w_1)) = 0.4$$

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

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

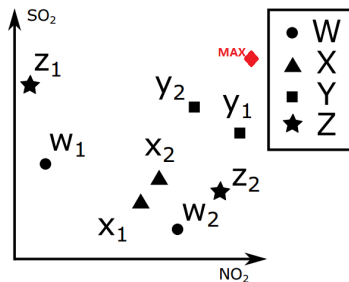
Early Pruning Techniques

Lemma (Dominance-Power Filtering)

$$DP(v_j) = \prod_{k=1}^d (b(k) - v_j(k)) = 0, b(k) = \max\{v_1(k), \dots, v_n(k)\}.$$

$$DP(V) = \sum_{v_j \in V} (P(v_j) \times DP(v_j)).$$

\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.



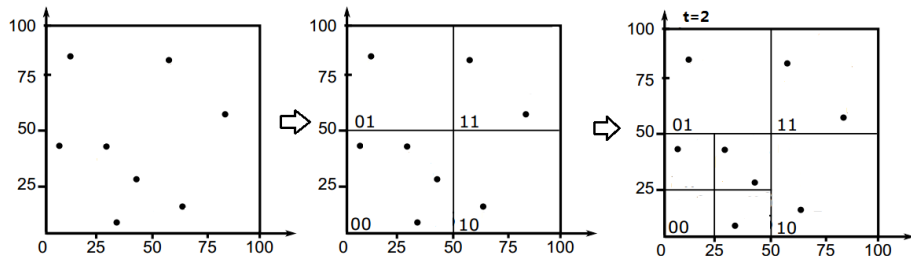
Example (Dominance-Power Filtering)

$$DP(y_1) = (100 - 95)(100 - 60) = 200$$

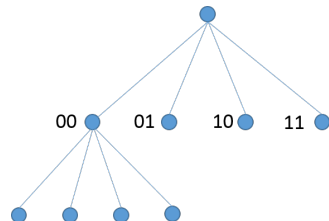
$$DP(y_2) = (100 - 80)(100 - 70) = 600$$

$$\begin{aligned} DP(Y) &= P(y_1) * DP(y_1) + P(y_2) * DP(y_2) \\ &= 0.8 * 200 + 0.2 * 600 = 280 \end{aligned}$$

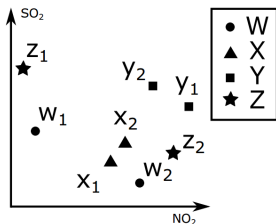
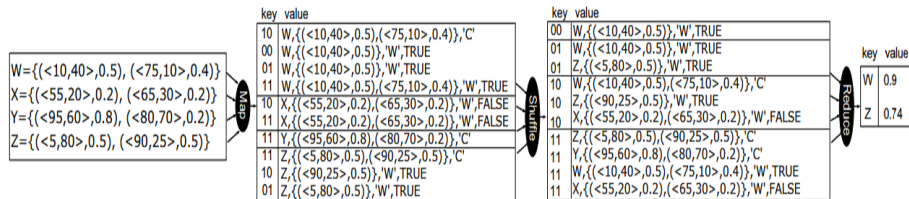
PSQtree for Pruning



- generate PSQtree
- traverse PSQtree for computing $P_{sky}(node.min)$
- zero-probability filtering
- upper-bound filtering
- partitioning objects by PSQtree(weakly dominate)



MapReduce Algorithms with PSQtree



Example (PS-QPF-MR)

- $T_p = 0.5$, map function is called with an uncertain object.
- for X , upper bound
 $node(10).P_{min}P(x1) + node(10).P_{min}P(x2) = 0.4 < T_p = 0.5$.
 X is not a skyline candidate.
- every instance of X , emit the key-value pairs
 $\langle 10, (\{(\langle 55, 20 \rangle, 0.2), (\langle 65, 30 \rangle, 0.2)\}, W, False) \rangle$ and $\langle 11, (\{(\langle 55, 20 \rangle, 0.2), (\langle 65, 30 \rangle, 0.2)\}, W, False) \rangle$ since $node(10)$ weakly dominates $node(10)$ and $node(11)$

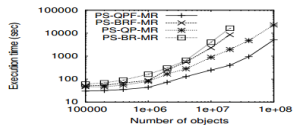
Experiments

- 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

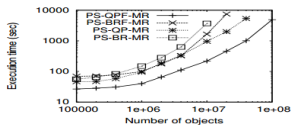
Algorithm	Description
PS-QP-MR	The algorithm with quadtree partitioning
PS-QPF-MR	The algorithm with quadtree partitioning and filtering
PS-BR-MR	The algorithm with random partitioning
PS-BRF-MR	The algorithm with random partitioning and filtering
PSMR	The state-of-the-art algorithm

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(ℓ)	$1 \sim 400$	40
No. of machines (t)	10~200	25

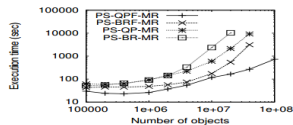
Experiments



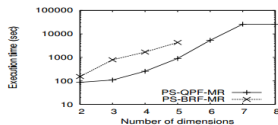
(a) ANTI



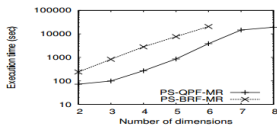
(b) IND



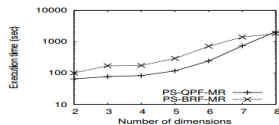
(c) COR



(a) ANTI

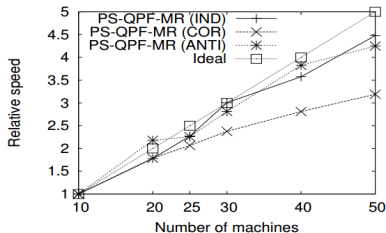


(b) IND

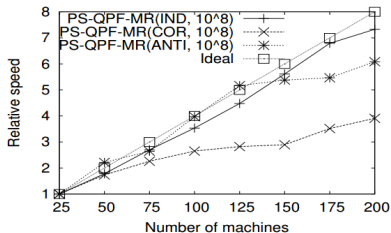


(c) COR

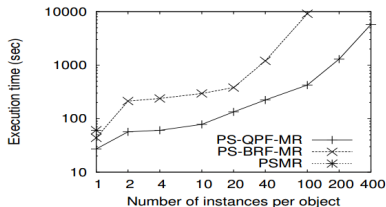
Experiments



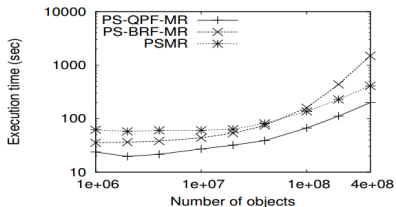
(a) With our cluster



(b) With Amazon EC2



(a) Varying ℓ



(b) Varying $|\mathbb{D}|$ when $\ell = 1$

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 techniques for optimization

Q & A