

Index

A		
ACID	6, 9	
Additional storage cost	52	
Aggregated data	9	
Algebraic aggregation function	48	
Application layer	5	
Approximation	80, 133	
Arrays	33	
B		
Bandwidth	92	
Base data	8	
Base table	10	
Bitmap indexing	27	
Equality encoded	29	
Range based encoded	31	
Standard encoding	27	
Blocksize	16, 91	
Bottom-up structures	51	
Bulk loading	20	
C		
Cache	15	
Cell-tree	26	
Compression of Bitmaps	33	
Cplex	39	
D		
Data specific parameters	89	
Data Warehouse	7	
Database layer	5	
Database management systems	1	
Datacube	12	
DBMS	1	
Decision support system	6	
Denormalized schemas	10	
Dimension table	10	
Disk specific parameters	91	
Distributive		48
Additive aggregation function ..		50
Aggregation function		48
Non-additive aggregation function		50
DSS		6
F		
Fanout		23, 52
Fractal dimension		69
Fractal model		69
G		
Galaxy schema		11
Generic tree structures		26
Grid model		64
GUI		5
H		
HG-tree		24
Hierarchical attributes		11
Hilbert curve		24
Hilbert R-tree		24
Holistic aggregation function		48
K		
KD-tree		22
KDB-tree		22
L		
Latency time		92
Lattice		12
Leaf node		43
Locking		13
Long transactions		6
M		
Main memory		16
Margin		36

- Materialized views 11
- Mechanics of disks 16
- Memory hierarchy 15
- Metadata 8
- MIP 35
- Mixed integer problem 35
- MOLAP 34
- MOPS 41
- N**
- Nearest neighbor 19
- Non-volatile 7
- Normally distributed data 79
- O**
- OLAP 8
- OLTP 6
- Optimal index structure 35
- Overlaps 23, 56
- P**
- Packed R-tree 24
- Partial range query 19
- Partial sums 34
- PISA model 71
- Point query 18, 52
- Q**
- Quadtree 21
- Query box
 dimensions 90
 size 90
- Query specific parameters 90
- R**
- R-tree 22
- R^* -tree 23
- R^+ -tree 23
- Random disk access 17
- Range query 19, 52
- RDBMS 9
- Read-mostly environment 1
- ROLAP 9
- S**
- S-Plus 97
- Scale factor 91
- Secondary memory 16
- Sector 16
- Sequential disk access 17
- Skewed data 81
- Skewed distributed data 77
- Snowflake Schema 11
- SR-tree 26
- SS-tree 24
- Star Schema 10
- STR-tree 24
- Subcube 12
- Subject oriented 7
- SUM model 66
- System specific parameters 91
- T**
- Tertiary memory 16
- Time variant 7
- Top-down 20
- Transactions 5
- U**
- UB-tree 26
- Uniformly distributed data 76, 81
- V**
- View selection problem 12
- View update problem 13
- X**
- X-tree 24
- Z**
- Z-Ordering 26

A. List of Symbols

a	length of interval A	page 71
A	set of all index entries	page 18
a_j	name of j th attribute	page 18
A_j	set of extensions of attribute a_j	page 18
b	blocksize	page 16
B	set of distinct blocksizes	page 92
B_{dir}	maximum fanout of directory node	page 23
b_{dir}	minimum fanout of directory node	page 23
B_{equal}	expected number of bitmap vectors which are read for an equality encoded bitmap index for processing a range query	page 31
b_j	number of bitmap vectors in j 's dimension	page 30
b_{ji}	the i th bitmap vector on the j th attribute	page 30
B_{leaf}	maximum fanout of leaf node	page 23
b_{leaf}	minimum fanout of leaf node	page 23
B_{range}	expected number of bitmap vectors which are read for a range encoded bitmap index for processing a range query	page 33
$Border(q)$	number of blocks that have to be accessed when an R_a^* -tree is used: $Inter(q) - Contain(q)$	page 46
bw	bandwidth	page 92
BW	set of distinct bandwidths	page 92
c	cardinality of data space	page 90

C	set of distinct cardinalities of data space	page 92
c_j	cardinality of data space in j th dimension	page 18
$Contain(q)$	number of leaf nodes contained in query box q	page 46
cp	child pointer	page 23
d	number of dimensions	page 18
D	set of distinct number of dimensions	page 92
$d_a(x)$	density function of position of interval A	page 72
$d_b(y)$	density function of position of interval B	page 72
d_f	fractal dimension	page 69
$data_{entry}$	entry of a data node	page 48
dir_{entry}	entry of a directory node	page 47
DW	data warehouse	page 8
e	configuration vector $(d, t, c, qs, qd, b, sf, bw, t_l)$	page 92
E	set of all configurations: $\{(d, t, c, qs, qd, b, sf, bw, t_l) (d, t, c, qs, qd, b, sf, bw, t_l)$ $\in D \times N_t \times C \times Q_s \times Q_d \times B \times SF \times BW \times T_l,$ $(qd \leq d) \wedge (c \leq t)\}$	page 92
$f(u, v)$	additional space when switching from fanout of u to fanout of v	page 53
$f(x, y)$	characteristic function to decide if two intervals intersect	page 72
$g(x, y)$	characteristic function to decide if one intervals is contained in the other	page 73
G	input set for creation of classification tree	page 97
h	height of tree	page 51
$h_1(a, b)$	PISA model: probability that A intersects B	page 72
$h_2(a, b)$	PISA model: probability that A contains B	page 73
i	index	
I	multi dimensional interval	page 18

$Inter(q)$	number of leafs intersecting query box q	page 46
j	index on the number of dimensions	
k	number of classes of the user defined density function	page 75
K	upper bound for additional space	page 53
l_{ik}	lower border of block i in dimension k	page 36
m	total number of bitmap vectors that are stored by an bitmap index in all dimensions	page 95
m_j	number of bitmap vectors that are stored by an bitmap index in the j th dimension	page 30
M	boxes for fractal dimension model	page 69
M_f	boxes filled for fractal dimension model	page 69
n	number of leaf nodes	page 36
N	set $\{1, \dots, n\}$	page 69
n_i	number of nodes on level i	
	n_0 : number of leaf nodes	
	n_1 : number of inner nodes on level 1	
	n_h : number of root nodes ($n_h = 1$)	
n_u	number of leaf nodes of structure with fanout u	page 53
n_v	number of leaf nodes of structure with fanout v	page 53
N_t	set with distinct number of tuples	page 92
O	$O_1 \times \dots \times O_d = \{0, \dots, c_1 - 1\} \times \dots \times \{0, \dots, c_d - 1\}$	page 18
O_j	$\{0, \dots, c_j - 1\}$	page 18
OP	set of operations	page 8
P	set of d -dimensional points (tuples)	page 36
p_i	SUM model: probability that rectangle i intersects query box q	page 66
q	vector with size of query box $q = (q_1, \dots, q_d)$ (length in each dimension)	page 19

q_j	length of query box in j th dimension	page 19
qd	query box dimensions	page 90
Q_d	set with distinct query box dimensions	page 92
qs	query box size	page 90
Q_s	set of distinct query box sizes	page 92
R	Relation $R(a_1, \dots, a_n, s)$	page 10
r_{ij}	size of leaf node i in dimension j	page 66
\tilde{r}	average length of rectangle in SUM model and PISA model (1-case)	page 74
r'	average length of rectangle in FRACTAL model	page 70
\bar{r}	average length of rectangle in GRID model	page 65
s_{dir}	size of a directory entry	page 46
s_{min}	select the index of the structure with the minimum value	page 97
s	number of index structures which are compared	page 93
s_0	number of slice for approximation of PISA model	page 133
S	data space $[0, 1)^d$	page 64
sf	scale factor	page 91
SF	set with distinct scale factors	page 92
t	number of tuples	page 36
T	$\{1, \dots, t\}$	page 36
tid	tuple identifier	page 18
t_r	time for random block access to secondary memory	page 16
t_s	time for sequential block access to secondary memory	page 16
$t_i : E \rightarrow \mathbb{R}^+$	expected time for processing a range queries with index structure i	page 93
	$t_1 : E \rightarrow \mathbb{R}^+$ R -tree without aggregated data	page 94
	$t_2 : E \rightarrow \mathbb{R}^+$ R -tree with aggregated data	page 94

	$t_3 : E \rightarrow \mathbb{R}^+$ equality encoded bitmap index	page 95
	$t_4 : E \rightarrow \mathbb{R}^+$ range encoded bitmap index	page 96
<i>TDB</i>	target database	page 8
u_l	user defined density function for distribution of sizes of rectangles	page 75
u_{ik}	upper border of block i in dimension k	page 36
v	blocks per bitmap vector	page 95
x	leftmost point of interval A	page 72
y	leftmost point of interval B	page 72
Y_e	percentage error of modeled values for R^*/R_a^* -tree	page 80

B. Approximation of PISA Model

The integral in Equation 6.14 on page 72 and the integral in Equation 6.17 on page 73 cannot be computed for every $d_a(x)$ and $d_b(y)$ analytically. For these cases we apply approximation methods to compute the integrals numerically. Appendix B presents approximation method for computing the integral in Equation 6.14. Equation 6.17 is approximated similarly. We present the computation of the integral over the gray shaded area in Figure B.1. Functions u and l describe this area:

$$u(x) = \begin{cases} x + a & : 0 \leq x < 1 - a - b \\ 1 - b & : 1 - a - b \leq x \leq 1 - a \end{cases}$$

$$l(x) = \begin{cases} 0 & : 0 \leq x < b \\ x - b & : b \leq x \leq 1 - a \end{cases}$$

Figure B.1 shows the area divided in $s_0 = 14$ slices. In computations presented in this thesis the x -interval $[0, 1 - a)$ is divided in $s_0 = 200$ equidistant slices. Integrals over two rectangles are calculated for each interval. One sum of rectangles calculates an upper bound (U) of the real value and another sum of rectangles calculates a lower bound (L) of the real result:

$$U(a, b) = \sum_{x \in \{y | y = (1-a) \frac{i}{s_0} \wedge i \in \{0, \dots, s_0-1\}\}} \left(\Phi_{\mu, \sigma} \left(x + \frac{1}{s_0}, u_{a,b} \left(x + \frac{1}{s_0} \right) \right) \right) \quad (\text{B.1})$$

$$- \Phi_{\mu, \sigma} \left(x, u_{a,b} \left(x + \frac{1}{s_0} \right) \right) \quad (\text{B.2})$$

$$- \Phi_{\mu, \sigma} \left(x + \frac{1}{s_0}, l_{a,b} \left(x + \frac{1}{s_0} \right) \right) \quad (\text{B.3})$$

$$+ \Phi_{\mu, \sigma} (x, l_{a,b}(x))$$

$$L(a, b) = \sum_{x \in \{y | y = (1-a) \frac{i}{s_0} \wedge i \in \{0, \dots, s_0-1\}\}} \left(\Phi_{\mu, \sigma} \left(x + \frac{1}{s_0}, u_{a,b}(x) \right) \right) \quad (\text{B.4})$$

$$- \Phi_{\mu, \sigma} (x, u_{a,b}(x)) \quad (\text{B.5})$$

$$- \Phi_{\mu, \sigma} \left(x + \frac{1}{s_0}, l_{a,b} \left(x + \frac{1}{s_0} \right) \right) \quad (\text{B.6})$$

$$+ \Phi_{\mu, \sigma} \left(x, l_{a,b} \left(x + \frac{1}{s_0} \right) \right)$$

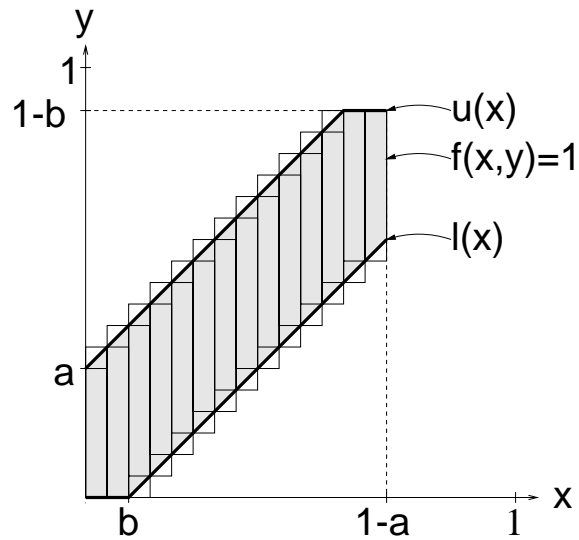


Figure B.1.: Approximation of the integral

The arithmetic average of these two sums is the approximation of the integral:

$$h_2(a, b) = \frac{1}{2}(U(a, b) + L(a, b)) \quad (\text{B.7})$$

The equations shown here apply the distribution function $\Phi_{\mu, \sigma}(x)$ of normal distributed data as an example. Other distribution functions can be applied as well.

Zusammenfassung der Ergebnisse

In dieser Arbeit wird untersucht, welche Indexstrukturen Anfragen in typischen Data Warehouse Systemen effizient unterstützen. Indexstrukturen, seit mehr als zwanzig Jahren Forschungsgegenstand im Datenbankbereich, wurden in der Vergangenheit für transaktionsorientierte Systeme optimiert. Ein Kennzeichen dieser Systeme ist die effiziente Unterstützung von Einfüge-, Änderungs- und Löschoptionen auf einzelnen Datensätzen.

Typische Operationen in Data Warehouse Systemen sind dagegen komplexe Anfragen auf großen relativ statischen Datenmengen. Aufgrund dieser veränderten Anforderungen müssen Datenbankmanagementsysteme, die für Data Warehouses eingesetzt werden, andere Techniken nutzen, um komplexe Anfragen effizient zu unterstützen.

Zunächst wird ein Ansatz untersucht, der mit Hilfe eines gemischt ganzzahligen Optimierungsproblems eine *optimale* Indexstruktur berechnet. Da die Kosten für die Berechnung dieser *optimalen* Indexstruktur mit der Anzahl der zu indizierenden Datensätze exponentiell steigen, wird in anschließenden Teilen der Arbeit heuristischen Ansätzen nachgegangen, die mit der Größe der zu indizierenden Datensätze skalieren.

Ein Ansatz erweitert auf Bäumen basierende Indexstrukturen um aggregierte Daten in den inneren Knoten. Experimentell wird gezeigt, daß mit Hilfe der materialisierten Zwischenergebnisse in den inneren Knoten Bereichsanfragen auf aggregierten Daten wesentlich schneller bearbeitet werden.

Um das Leistungsverhalten von Indexstrukturen mit und ohne materialisierte Zwischenergebnisse zu untersuchen, wird das PISA Modell (*Performance of Index Structures with and without Aggregated Data*) entwickelt. In diesem Modell wird die Verteilung der Daten und die Verteilung der Anfragen berücksichtigt. Das PISA Modell wird an gleich-, schief- und normalverteilte Datensätze angepaßt. Experimentell wird gezeigt, daß das PISA Modell mit einer höheren Präzision als die bisher aus der Literatur bekannten Modelle arbeitet.

Die Leistung von Indexstrukturen hängt von unterschiedlichen Parametern ab. In dieser Arbeit werden zwei Techniken vorgestellt, die abhängig von einer bestimmten Menge von Parametern Indexstrukturen vergleichen. Mit Hilfe von Klassifikationsbäumen wird z. B. gezeigt, daß die Blockgröße die relative Leistung weniger beeinflußt als andere Parameter. Ein weiteres Ergebnis ist, daß Bitmap-Indexstrukturen von den Verbesserungen neuerer Sekundärspeicher stärker profitieren als heute übliche auf Bäumen basierende Indexstrukturen. Bitmap-Indexierungstechniken bieten noch ein großes Potential für weitere Leistungssteigerungen im Datenbankbereich.