

Information theory and databases

Concept of information (through an example)

Information content of data streams, information rate

(Extension: physical entropy and information)

Databases: concept and examples



Concept of information (through an example)

Intuitive concept:

"informare" (Lat.) : „to give form to the mind", or to teach, *instruct somebody*

Thus: „We can only change our minds, when we receive **information**."

Or:

„a type of input to an organism or designed device" : Ecology, sensory input
(Smell of food → movement of animal)

Or:

„information is any type of pattern that influences the formation or transformation of other patterns."
(RNA sequence → Protein structure)



Transmitting information – information coding

in general

Information source

Which event occurred from a set of possibilities?

encoding

Encoding: We represent **events** with **NUMBERS**

Transmission channel

decoding

Decoding: We reconstruct **events** from **NUMBERS**

Information receiver
destination

(news)

Transmitting information – information coding

in general

Information source

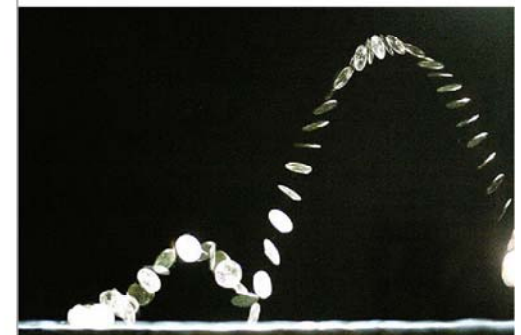
encoding

Transmission channel

decoding

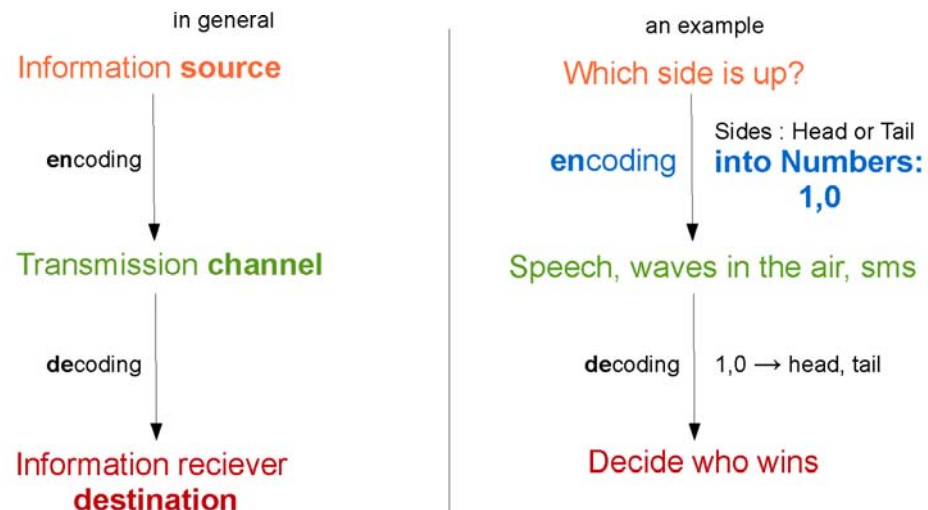
Information receiver
destination

an example
Tossing a dime

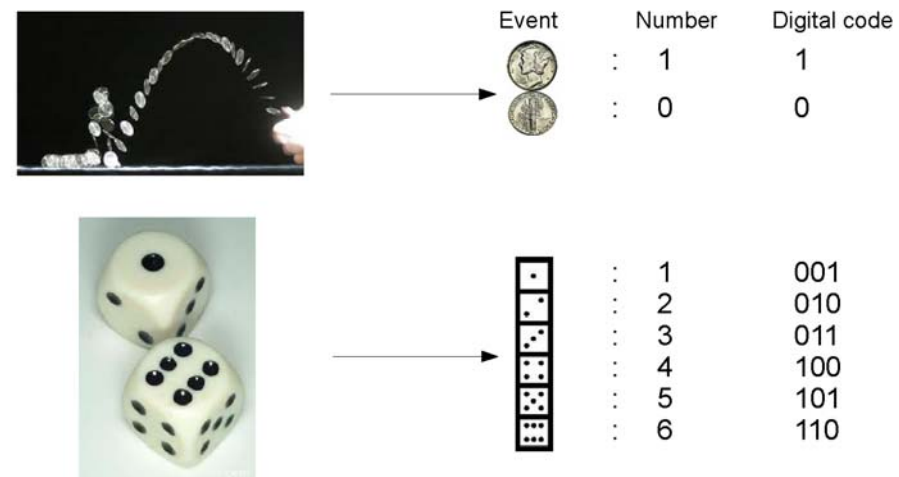


Head or Tail?

Transmitting information – information coding



Transmitting information – digital coding





How many **bits** we need?
Bit: **binary digit**
0 or 1







Transmitting information – digital coding

How many **bits** we need?





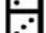



Bit: **binary digit**

0 or 1

Event	Number	Digital code	Bits needed
	: 1	1	1
	: 0	0	

Event	Number	Digital code	Bits needed
	: 1	001	3
	: 2	010	
	: 3	011	
	: 4	100	
	: 5	101	
	: 6	110	

Transmitting information – coding efficiency

Event	Number	Digital code	Bits needed	Maximum number of events
	: 1	1	1	2
	: 0	0		
	: 1	001	3	8
	: 2	010		
	: 3	011		
	: 4	100		
	: 5	101		
	: 6	110		
	7	111		
	0	000		

Here we only have 6 events,
but could encode 8 in 3 bits!

Transmitting information – coding *efficiency*

Event	Number	Digital code	Bits needed	Maximum number of events
	: 1	001	3	8
	: 2	010		
	: 3	011		
	: 4	100		
	: 5	101		
	: 6	110		
	7	111		
	0	000		

Here we only have 6 events,
but could encode 8 in 3 bits!

A better encoding:

$\{X_1 X_2 X_3\}$ group 3 events together
Classic coding
3x3 bits = 9 bits

Transmitting information – coding *efficiency*

Event	Number	Digital code	Bits needed	Maximum number of events
	: 1	001	3	8
	: 2	010		
	: 3	011		
	: 4	100		
	: 5	101		
	: 6	110		
	7	111		
	0	000		

Here we only have 6 events,
but could encode 8 in 3 bits!







A better encoding:

$\{X_1 X_2 X_3\}$ group 3 events together : number of possibilities = $6^3 = 216$
Classic coding
3x3 bits = 9 bits
 $\xrightarrow{\text{1 bit less!}}$
 $256 = 2^8$
 It is possible to encode a group of
any 3 events in 8 bits

Transmitting information – information content







Information content = how many bits do we *minimally* need to encode
(This also gives the encoding efficiency limit)

Transmitting information – measure of information – example of two dices

Fair	P_i	probability	code example	bits needed	$p^*(\text{number of bits needed})$
	1/6	0,17	000	3	0,5
	1/6	0,17	001	3	0,5
	1/6	0,17	010	3	0,5
	1/6	0,17	011	3	0,5
	1/6	0,17	100	3	0,5
	1/6	0,17	101	3	0,5

Expected number of bits needed: 3 $\mu = \sum (x_i * p_i)$



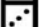



Loaded
dice

	1/2	0,5	0	1	0,5
	1/4	0,25	10	2	0,5
	1/8	0,13	110	3	0,38
	1/16	0,06	1110	4	0,25
	1/32	0,03	11110	5	0,16
	1/32	0,03	11111	5	0,16

We can encode more efficiently here, for example such:

Expected number of bits needed: 1,94

Transmitting information – measure of information – example of two dices

Fair	p_i	probability	code example	bits needed	$p \cdot (\text{number of bits needed})$
	1/6	0,17	000	3	0,5
	1/6	0,17	001	3	0,5
	1/6	0,17	010	3	0,5
	1/6	0,17	011	3	0,5
	1/6	0,17	100	3	0,5
	1/6	0,17	101	3	0,5







Here we do NOT
Expect anything

Maximal uncertainty

Expected number of bits needed: 3

Gained information is proportional to the number of bits needed

Loaded dice

	1/2	0,5	0	1	0,5	Here we <i>expect</i>
	1/4	0,25	10	2	0,5	„one“ (most probable)
	1/8	0,13	110	3	0,38	or „two“
	1/16	0,06	1110	4	0,25	
	1/32	0,03	11110	5	0,16	
	1/32	0,03	11111	5	0,16	

On average
we *learn less*

Expected number of bits needed: 1,94

Here the overall, average information content is less.

Transmitting information – measure of information

How should be information content **mathematically** specified? (Shannon 1948)

1.: H should be *continuous* in the p_i (small change in $p_i \rightarrow$ small change in H)

2.: *Unlikely events carry a high information content:*

H should be in some way inverse proportional to p

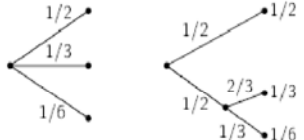
If all the p_i are equal, ($p_i = 1/n$)

then H should be a monotonic increasing function of n .

With equally likely events there is more choice, or uncertainty, when there are more possible events.

3.: *Branching Choices:*

If a choice can be broken down into two successive choices, the original H should be the weighted sum of the individual values of H .

$$H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) = H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2} \cdot H\left(\frac{2}{3}, \frac{1}{3}\right)$$


Transmitting information – measure of information

Shannon : define measure as: $H = p \cdot \log_2 \left(\frac{1}{p} \right)$

\log_2 : 2-base logarithm

Examples:

$$\log_2(2) = 1$$

$$\log_2(4) = 2$$

$$\log_2(8) = 3$$

Transmitting information – measure of information

Shannon

$$H = p \cdot \log_2 \left(\frac{1}{p} \right) \quad [\text{bit}]$$

If we have multiple events in the set, then it is a sum for every possible event:

$$H = \sum_i p_i \cdot \log_2 \left(\frac{1}{p_i} \right) = \sum_i -p_i \cdot \log_2 p_i$$

other log-bases:
 $\log_e(\ln)$: [nat]
 $\log_{10}(\lg)$: [ban]

measure of information - entropy

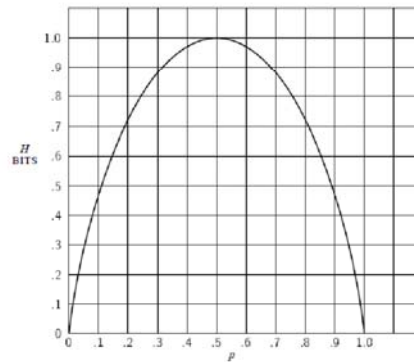
Dime tossing



p



q = 1-p



$$H = \sum_i -p_i \cdot \log_2 p_i = -p \cdot \log_2 p - q \cdot \log_2 q = -p \cdot \log_2 p - (1-p) \cdot \log_2 (1-p)$$

measure of information - entropy

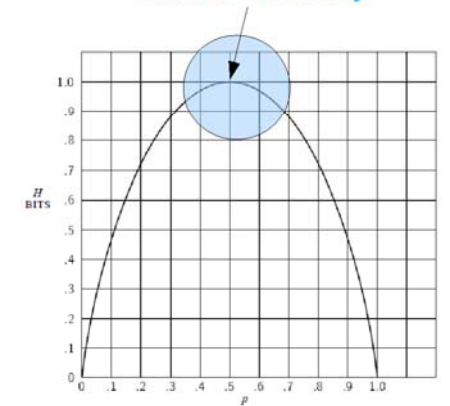
Dime tossing



p



q = 1-p



$$H = \sum_i -p_i \cdot \log_2 p_i = -p \cdot \log_2 p - q \cdot \log_2 q = -p \cdot \log_2 p - (1-p) \cdot \log_2 (1-p)$$

measure of information - entropy

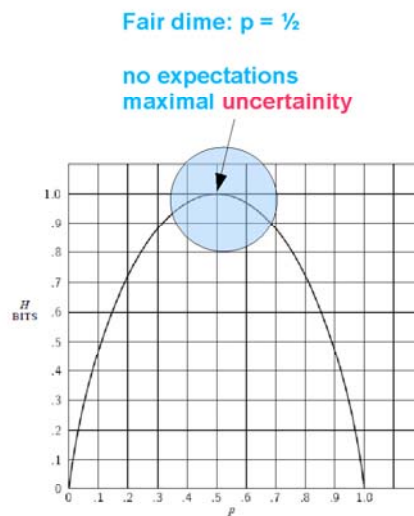
Dime tossing



p



q = 1-p



H has another name: **Shannon-entropy**

H has a **maximum** when we know nothing in advance (all p_i -s are equal, $p_i = 1/n$)

Expected outcomes are maximized: each state is equally probable



Physical entropy (S) has a maximum if the number of microstates is maximal.

measure of information - entropy

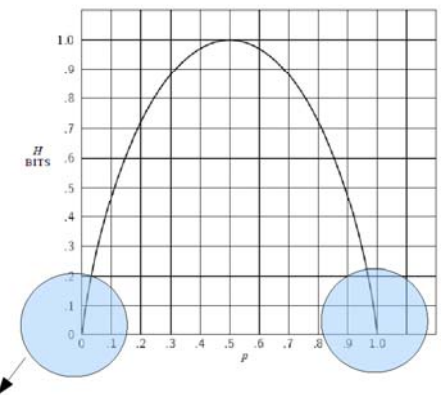
Dime tossing



p



q = 1-p



H has another name: **Shannon-entropy**

H vanishes **ONLY** if we are absolutely certain of the outcome: $p=0$ or $p=1$



Physical entropy (S) vanishes **ONLY** if there is exactly 1 microstate

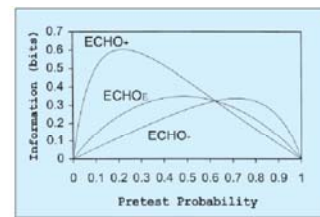
Examples of usage in medicine

Bayes-theorem based methods:

The amount of information gained by performing a diagnostic test can be quantified by calculating the relative entropy between the posttest and pretest probability distributions

Application:

- Diagnostic tests
- expert systems



a: pretest probability
b: post test probability

$$D(b||a) = \sum_{i=1}^n b_i \log_2(b_i/a_i)$$

Testing Situation	Pretest Probability of Disease	Test Operating Characteristics: Sensitivity/Specificity	Test Result	Posttest Probability of Disease	Information Gained
Breast cancer screening with mammography	0.01	0.75/0.94	Positive	0.11	0.25 bits
			Negative	0.003	0.006 bits
Mammography given palpable breast mass	0.2	0.80/0.90	Positive	0.67	0.74 bits
			Negative	0.05	0.13 bits
Screening for HIV with antibody test	0.001	0.96/0.998	Positive	0.33	2.4 bits
			Negative	0.00001	0.001 bits
Presence of tonsillar exudate in diagnosing infection with group A streptococci	0.1	0.45/0.84	Positive	0.24	0.11 bits
			Negative	0.07	0.01 bits
Colon cancer screening by fecal occult blood testing	0.005	0.40/0.90	Positive	0.02	0.02 bits
			Negative	0.003	0.0005 bits

Databases

Databases store information:

Databases are used for: storage, structuring and **extraction of information** gathered previously.

It is hard to **extract** or modify information stored on paper

FOSTER CITY EYE CARE - OPTOMETRIC CENTER
PATIENT HISTORY QUESTIONNAIRE

Last name: _____ First name: _____ Mr. ☐ Ms. ☐ Miss ☐ Mx. ☐

Address: _____ (H) _____ (Cell) _____

Telephone (W): _____ (H) _____ (Cell) _____

SSN: _____ Date of Birth: _____ Age: _____

Occupation: _____ Computer Hours Per Day: _____

Emergency contact/telephone no.: _____

Date of last eye exam: _____ Dilated: _____ Today's Date: _____

Hobbies or Sports: _____

Primary reason for today's exam: _____

MEDICAL INFORMATION

What is your general health: _____

Do you have any problems with any of these systems? (please circle all that apply)

Gastrointestinal	Y/N	Nervous	Y/N	Eyes	Y/N
Ear/Nose/Throat	Y/N	Genitourinary	Y/N	Mental	Y/N
Cardiovascular	Y/N	Musculoskeletal	Y/N	Endocrine (glands)	Y/N
Respiratory	Y/N	Insomniac (skin)	Y/N	Blood/Lymph	Y/N
				Allergic/immunologic	Y/N
				Pregnant or nursing	Y/N

Please explain: _____

Please answer all that apply:

Diabetes	Y/N	Type	Date of diagnosis
Allergies	Y/N	Allergic to what?	What happens?
Medication allergy	Y/N	What happens?	Headaches
Other health problems	Y/N		HIV/AIDS
Current medication(s)	Y/N		
Have you had any operations?	Y/N	Kind?	When?
Do you use cigarettes/tobacco?	Y/N	Alcohol?	Other substance(s)?
Name of family doctor			Date of last visit
Date of last tetanus shot			

FAMILY HISTORY

High blood pressure	Y/N	Relation	Myocardial degeneration	Y/N	Relation
Diabetes	Y/N	Relation	Retinal detachment	Y/N	Relation
Glaucoma	Y/N	Relation	Cataracts	Y/N	Relation
Other eye condition(s)	Y/N	What kind?	Relation		

PERSONAL EYE INFORMATION

Have you had an eye operation?	Y/N	Type	Date
Have you had an eye injury?	Y/N	Kind	Date
Do you have glaucoma?	Y/N	Cataracts?	Y/N
Do you wear glasses?	Y/N	Dry eyes?	Y/N
Other eye problems?	Y/N	What kind?	Blurred vision?
Do you wear glasses?	Y/N	Contact lenses?	Y/N
Additional information		Are you interested in new contact lenses?	Y/N

Whom may we thank for referring you? _____

Doctor's initials: _____

Databases

Databases store information:

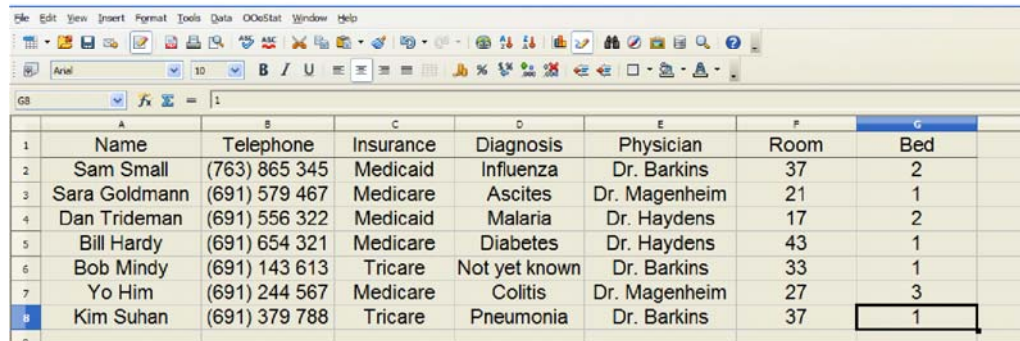
Databases are used for: storage, structuring and extraction of **information** gathered previously.

Databases – storing information

Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
Sam Small						

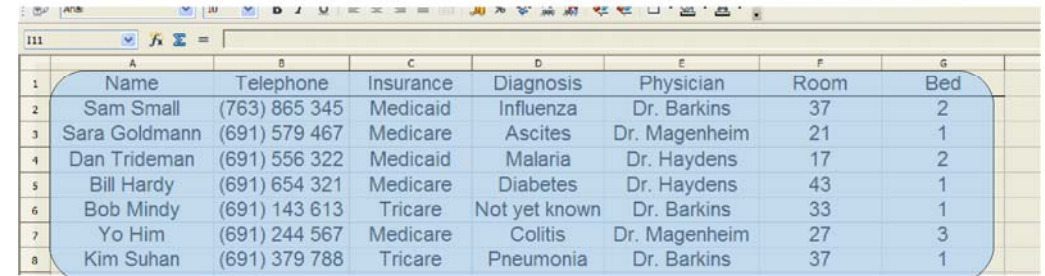
Instead of paper, one could start typing the data into a spreadsheet (Excel, OpenOffice, etc)

Databases – storing information



	A	B	C	D	E	F	G
	Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
2	Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
3	Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
4	Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
5	Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
6	Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
7	Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
8	Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

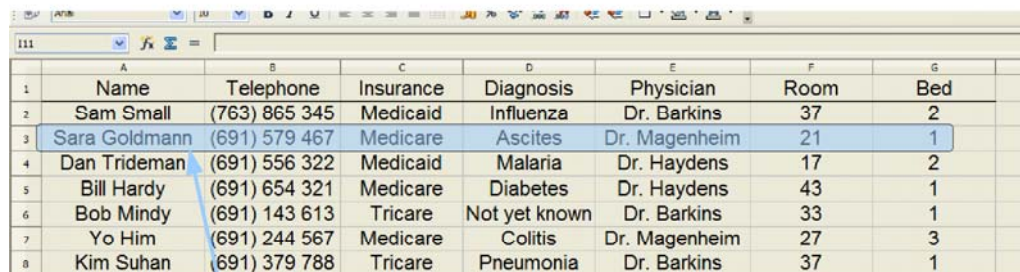
Databases – storing information



	A	B	C	D	E	F	G
	Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
2	Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
3	Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
4	Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
5	Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
6	Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
7	Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
8	Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

Table : ordered set of data (information)

Databases – storing information



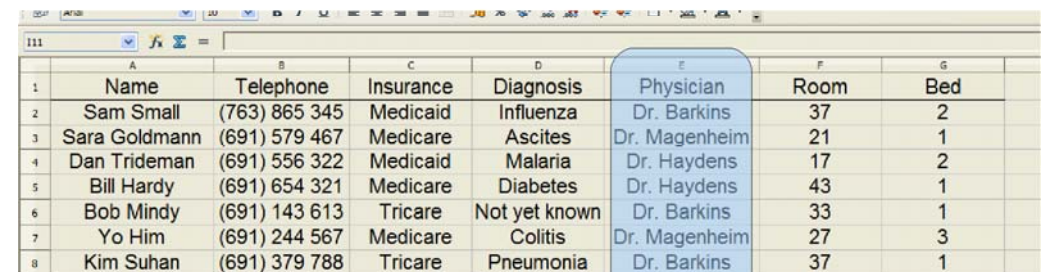
	A	B	C	D	E	F	G
	Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
2	Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
3	Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
4	Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
5	Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
6	Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
7	Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
8	Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

Record : Information grouped together
(one ROW in a Table)

Each row is a selected set of data

Every row has the same structure

Databases – storing information



	A	B	C	D	E	F	G
	Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
2	Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
3	Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
4	Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
5	Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
6	Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
7	Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
8	Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

Column: data type

Databases – manipulating information

Sorting data

Name	Insurance	Diagnosis	Physician	Room	Bed
Sam Small	Medicaid	Influenza	Dr. Barkins	37	2
Sara Goldmann	Medicare	Ascites	Dr. Magenheim	21	1
Dan Trideman	Medicaid	Malaria	Dr. Haydens	17	2
Bill Hardy	Medicare	Diabetes	Dr. Haydens	43	1
Yo Him	Medicare	Colitis	Dr. Magenheim	27	3
Kim Suhan	Tricare	Pneumonia	Dr. Barkins	37	1

Databases – manipulating information

Name	Telephone	Insurance	Diagnosis	Physician	Room	Bed
Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

Databases – retrieving information

Telephone	Insurance	Diagnosis	Physician	Room	Bed
763) 865 345	Medicaid	Influenza	Dr. Barkins	37	2
691) 556 322	Medicaid	Malaria	Dr. Haydens	17	2
691) 654 321	Medicare	Diabetes	Dr. Haydens	43	1
691) 579 467	Medicare	Ascites	Dr. Magenheim	21	1
691) 244 567	Medicare	Colitis	Dr. Magenheim	27	3
691) 143 613	Tricare	Not yet known	Dr. Barkins	33	1
691) 379 788	Tricare	Pneumonia	Dr. Barkins	37	1

Databases – problems with simple methods

Name	Telephone	Insurance	Diagnosis	Physician	Medication	Medication	Room	Bed
Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	Aspiryn		37	2
Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	Halo fantrine		17	2
Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	Insulin		43	1
Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	Triamterene	spironolactone	21	1
Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	sulfasalazine		27	3
Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins			33	1
Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	Aspiryn	Augmentin	37	1

Records do not have the same size

Waste of space

Adding new data types tedious

Inconsistency : is a field empty by error?

Databases – problems with simple methods

	A	B	C	D	E	F	G	H	I
	Name	Telephone	Insurance	Diagnosis	Physician	Medication	Medication	Room	Bed
1	Sam Small	(763) 865 345	Medicaid	Influenza	Dr. Barkins	Aspiryn		37	2
2	Dan Trideman	(691) 556 322	Medicaid	Malaria	Dr. Haydens	Halofantrine		17	2
3	Bill Hardy	(691) 654 321	Medicare	Diabetes	Dr. Haydens	Insulin		43	1
4	Sara Goldmann	(691) 579 467	Medicare	Ascites	Dr. Magenheim	Triamterene	spironolactone	21	1
5	Yo Him	(691) 244 567	Medicare	Colitis	Dr. Magenheim	sulfasalazine		27	3
6	Bob Mindy	(691) 143 613	Tricare	Not yet known	Dr. Barkins			33	1
7	Kim Suhan	(691) 379 788	Tricare	Pneumonia	Dr. Barkins	Aspiryn	Augmentin	37	1

Entering the same data multiple times:

Typos

Redundancy

Later change almost impossible – too many items

...

Databases – SQL

Relational database:

we store data together with a logic.

Data types have logical connections.

Every data is stored only once,
but can be used multiple times.

A Relational Model of Data for Large Shared Data Banks

E. F. Codd
IBM Research Laboratory, San Jose, California

Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A grouping service which supplies such information is not a satisfactory solution. Activities of users at terminals and user application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed. Changes in data representation will often be needed as a result of changes in query, update, and report traffic and natural growth in the types of stored information.

Existing nonrelational, formatted data systems provide users with tree-structured files or slightly more general network models of the data. In Section 1, inadequacies of these models are discussed. A model based on *n*-ary relations, a normal form for data base relations, and the concept of a universal data sublanguage are introduced. In Section 2, certain operations on relations (other than logical inference) are discussed and applied to the problem of redundancy and consistency in the user's model.

KEY WORDS AND PHRASES: data bank, data base, data structure, data organization, inventories of data, networks of data, relations, denormality, redundancy, consistency, association, join, ordered language, products, setoids, country, data integrity
CR CATEGORIES: 3.70, 3.75, 3.76, 4.00, 4.02, 4.20

The relational view (or model) of data described in Section 1 appears to be superior in several respects to the graph or network model [3, 4] presently in vogue for nonrelational systems. It provides a means of describing data with its natural structure only—that is, without superimposing any additional structure for machine representation purposes. Accordingly, it provides a basis for a high level data language which will yield maximal independence between programs on the one hand and machine representation and organization of data on the other.

A further advantage of the relational view is that it forms a sound basis for treating denormality, redundancy, and consistency of relations—these are discussed in Section 2. The network model, on the other hand, has spawned a number of confusion, not the least of which is mistaking the derivation of connections for the derivation of relations (see remarks in Section 2 on the "connection trap"). Finally, the relational view permits a clearer evaluation of the scope and logical limitations of present formatted data systems, and also the relative merits (from a logical standpoint) of competing representations of data within a single system. Examples of this clearer perspective are cited in various parts of this paper. Implementations of systems to support the relational model are not discussed.

1.2. DATA DESCRIPTION IN PRESENT SYSTEMS
The provision of data description tables in recently developed information systems represents a major advance toward the goal of data independence [5, 6, 7]. Such tables facilitate changing certain characteristics of the data representation stored in a data bank. However, the variety of data representation characteristics which can be changed without logically requiring some application programs to be modified is still quite limited. Further, the model of data with which users interact is still cluttered with representational properties, particularly in regard to the representation of collections of data (as opposed to individual items). Three of the principal kinds of data dependencies which still need to be removed are: ordering dependencies, indexing dependencies, and access path dependencies. In some systems these dependencies are not clearly separable from one another.

1.2.1. Ordering Dependencies. Elements of data in a data bank may be stored in a variety of ways, some involving no concern for ordering, some permitting each element to participate in one ordering only, others permitting each element to participate in several orderings. Let us consider these existing systems which either require or permit data elements to be stored in at least one total ordering which is closely associated with the hardware-determined ordering of addresses. For example, the records of a file concerning parts might be stored in ascending order by part serial number. Such systems normally permit application programs to assume that the order of presentation of records from such a file is identical to (or is a subordering of) the

Databases

clinical example:
part of a relational database:
each record has its own KEY.

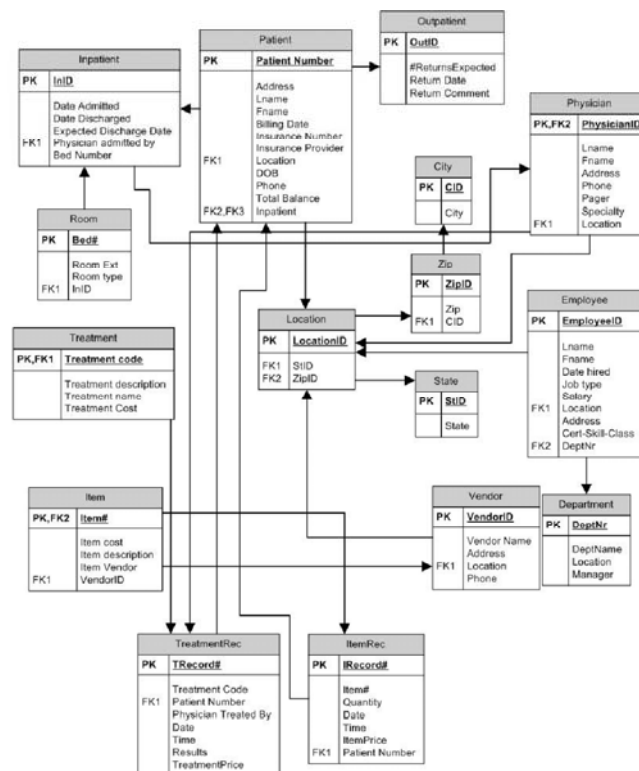
Two types of keys:

Primary Key (PK)

Foreign Key (FK)

FK-s join the tables together,
and enable re-use of data.

The structure of the database
(layout/content of tables)
represents the logic behind
the workflow for which the
database is being designed.



Extension material:

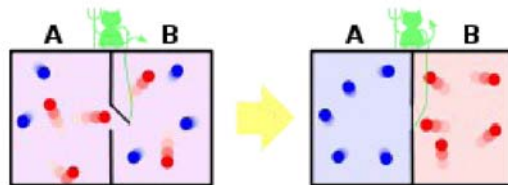
information entropy and physical entropy

(these last few slides are for those who wish to have a starting point for understanding why the term „entropy“ is appropriate also for information, and what it has to do with the concept of entropy in physics. It may be a good idea to re-read this part after you study entropy in physics...)

Information entropy and physical entropy

„ In an isolated system, entropy never decreases.” Second Law of Thermodynamics

The Maxwell demon

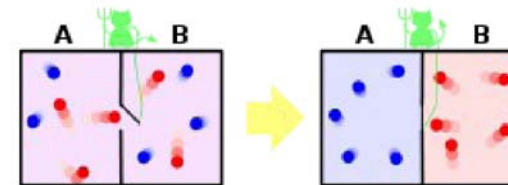


Temperature of A **decreases**, B **increases** → Violation of the Second Law ?

Information entropy and physical entropy

„ In an isolated system, entropy never decreases.” Second Law of Thermodynamics

The Maxwell demon

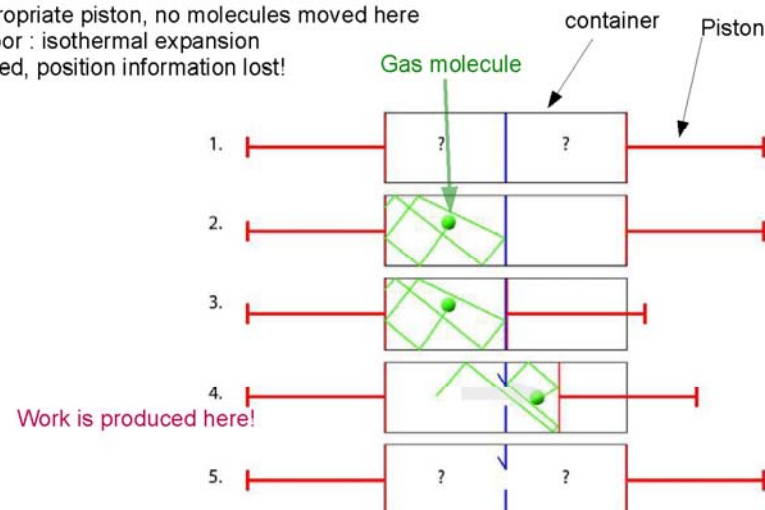


Temperature of A **decreases**, B **increases** → Violation of Law II. ?

Solution: NO, since the demon interacts with the system, it must be considered. The demon acquires **information**, and this changes its state!

Information entropy and physical entropy

1. : molecule's position unknown
2. : measure position, information = 1 bit
3. : move appropriate piston, no molecules moved here
4. : release door : isothermal expansion
5. : door opened, position information lost!



Information entropy and physical entropy

1. : molecule's position unknown
2. : measure position, information = 1 bit
3. : move appropriate piston, no molecules moved here
4. : **release door : isothermal expansion**
5. : door opened, position information lost!

Isothermal expansion:

$$W_{A \rightarrow B} = NkT \ln \left(\frac{V_A}{V_B} \right)$$

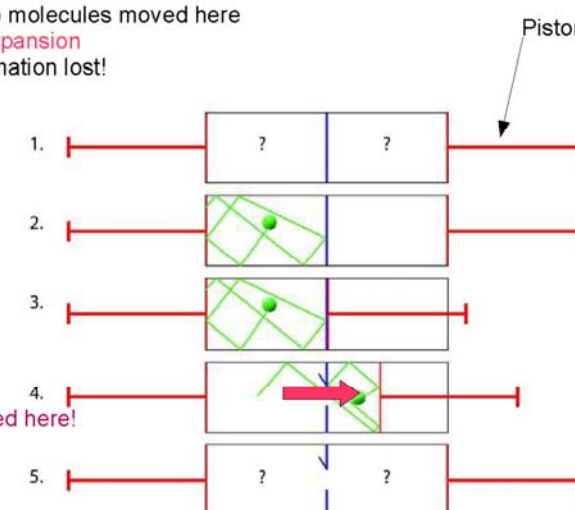
In this case:

$$N=1$$

$$V_A/V_B = 2$$

Hence

$$W = kT \ln(2) \text{ Work is produced here!}$$



Information entropy and physical entropy

1. : molecule's position unknown
2. : measure position, information = 1 bit
3. : move appropriate piston, no molecules moved here
4. : release door : isothermal expansion
5. : door opened, position information lost!

Isothermal expansion:

$$W_{A \rightarrow B} = NkT \ln \left(\frac{V_A}{V_B} \right)$$

In this case:

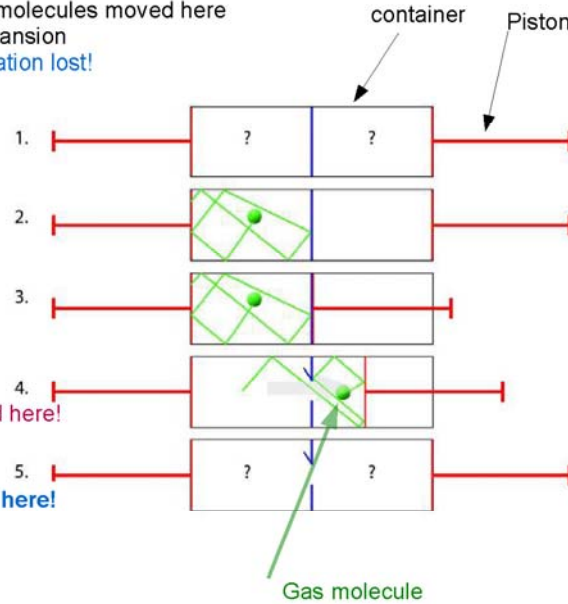
$$N=1$$

$$V_A/V_B = 2$$

Hence

$$W = kT \ln(2) \text{ Work is produced here!}$$

Information is lost here!



Information entropy and physical entropy

Leo Szilárd:

From Law II. taking into account that $W = T\Delta S$

$$W_{\text{produced by piston}} = W_{\text{loss of information}}$$

$$T\Delta S_{\text{inf}} = kT \ln 2$$

$$\Delta S_{\text{1bit}} = k \ln 2$$

Erasing 1 bit of information increases physical entropy by $k \ln 2$

(Landauer 1971, logically irreversible processes, eg. AND-gate produce entropy → heating of circuits!)

