CODATA 18th International Conference: Frontiers of Scientific and Technical Data (29 September - 3 October, 2002)

Prototype of TRC Integrated Information System for Physicochemical Properties of Organic Compounds: Evaluated Data, Models and Knowledge

Xinjian Yan,

Qian Dong, Xiangrong Hong, Robert D. Chirico, Michael Frenkel

Thermodynamics Research Center (TRC) National Institute of Standards and Technology

Introduction

- Requirement: Industrial and scientific developments require high quality data and models
- Key Point: High quality data system needs strong support from comprehensive knowledge base
- Aim: Develop a system with high quality data and models fully supported by domain knowledge

The Relationship between Data, Model and Knowledge





TRC Integrated Information System (TIIS) Structure

The Support of Knowledge to Data and Model Analysis



Data Background - TRC Databases

✤ Databases:

Source Database, Table Database, Density Database, Vapor Pressure Database, Ideal Gas Database, etc.

A Comprehensive Physicochemical Data System: Source Database contains more than 100 physical and chemical properties, over 2 million experimental records for 32,000 chemical systems (pure compounds, mixtures, and reaction systems)

Information for Recommended Data (RD)

Detailed information is crucial for a good understanding of data. The following information has been prepared for recommended data (also for experimental data).

- The uncertainty values of RD
- The number of data points used for obtaining RD
- The discreteness of the data used to process RD
- The description about the selection of RD
- The grade of RD

Data Processing for RD

- For compounds having multiple values, a weighted average method is used to obtain recommended data
- For compounds having only one or two values, the data are inspected by:
 - A. Theories and thermodynamics relationships
 B. Comparison with the values from models
 C. Comparison with other well characterized sources
 D. Similar compounds
- For doubtful data, original articles are reviewed

Criteria and Methods of Evaluating Models

The major problem in model evaluations is that very little attention is paid on the prediction abilities of models. The following factors have been considered in our evaluating and selecting of models for TIIS.

- Prediction ability
- Complexity of compounds used in developing and testing models
- Diversity of compounds used in developing and testing models
- Reliability of each parameter (how many and how well data were used in obtaining each parameter)
- Similarity analysis

Example of Prediction Ability of Models

WJ is a simple model with about 20 parameters, while MP is a model using 167 parameters. MP is better than WJ in correlating data, but not in predicting for new compounds.

Ranges of deviations (Dev, K) of the Tc predicted by WJ and MP group contribution models for the compounds having experimental data reported between 1996 and 2001, and number of compounds in each range.

	WJ	MP
Correlated result Total number Dev, K	467 6.3	434 5.7
Predicted result Total number Dev, K	48 10.2	42 13.0

Examination on the Reliability of Parameters

167 MP groups and estimated occurring frequency in compounds that have Tc data published before 1987 (N1) and up to now (N2).

Sumn	narv	NI	G	roups 40	NI	Group 36	ps NL 2	Gr	ups 0	
Sum	uar y			40	1	50	-		2	
			1987	Now		CH ₂	rC	1	0	4
			361	600		rC=	Ī	ī	Ō	Ó
Partl	Part2	Bond	NI	N2		rC=	ŇН	ī	Ō	Ō
C=	rC=	1	0	1		rC=	COOH	ī	Ō	Ō
ĒЊ	C*	ī	Ō	ō		CĪ	ĊŌ	ī	Ō	Ō
CH ₂	Ī	ī	Ō	ō		rC=	ΖŌ	ī	Ō	ī
rC	rrC=	ī	Õ	ō		ō	ČÕ	ī	Ō	ō
rĈ	rC=	ī	Ō	Ō		rC=	Br	ī	Ō	6
S	S	ī	Ō	ĩ		CO	CO	ī	ō	Ō
rCH	rrCH	ī	Õ	$\overline{2}$		rC=	mC=	ī	ŏ	ŏ
CH	CH=	ī	Ō	ō		H	CHO	ī	ō	Ō
ĊĦ	C=	ī	Ō	Ō		NH ₂	NH	ī	Ō	Ō
ĊĦ	ĊНО	ī	Ō	ō		H	COOH	ī	ī	ī
Ĉ	C=	ī	Ō	ī		CH ₂	F	ī	ī	ī
ĒЊ	Ē=	ī	Ō	ō		C	Ō	ī	ī	- 23
=C=	=0	2	Ō	Ō		ĒH*	ĒH*	3	ī	ī
$r\bar{C}$	rrC	ī	Ō	Ō		CH*	C *	3	ī	ī
C=	=C=	2	Ō	ō		rC=	<u>čoo</u>	ĩ	ī	3
ĒH=	CN	ī	Ō	ī		C=	ĊĪ –	ī	ī	4
CH=	COOo	ī	Ō	ō		Ē	COOH	ī	ī	i
CH=	COOH	ī	Õ	ŏ		ē	rC=	ī	ī	7
CH=	CHO	ī	Ō	ī		ĒН=	COO	ī	ī	2
CH=	F	ī	Õ	ō		CH=	ō	ī	ī	2
ČH=	rC=	ī	ŏ	ĩ		ČH=	Čl	ī	ī	ī
ČH=	Č*	ī	ŏ	õ		ČH.	Č*	ī	ĩ	ĩ
CH=	Č=	ī	Ň	ň		ČH-	СH=	ī	ĩ	ī
CH=	=C=	2	õ	ŏ		CH ₂	CH	î	ĩ	ĩ
Č*	$\tilde{C^{\star}}$	3	Ň	Ň		rC=	CN	ĩ	ĩ	ĩ
rO	rN=	ĩ	ň	ĭ		CH	CHO	î	î	ŝ
CH	Br	î	ň	ñ		CH.	CN	î	î	ĭ
ČH.	ĩ	î	ň	ň		ČH.	NO	î	î	î

Complexity of Organic Compounds - Definition

Group/ complexity	=1	>1		
СН	1	1	CH3-CH(CH3)-CH3	= 2
С	2	2		
C=C (double bond)	2	2		
=C= ` ´	2	2		
C*C (triple bond)	2	2		
F, CI, Br, I	3	5	2 (when groups >4)	
ĆN Ś	3	4		
Ν	3	4		
NC	3	4		
S	3	4		
SH	3	4		
СНО	4	10		
CO	4	10		
COO	4	10		
СООН	4	10		
N=	4	10		
NH	4	10		
NH=	4	10		
NH2	4	10		
NO2	4	10		
0	4	10		
ОН	4	10	OH-CH2-CH2-OH	= 18
SO	4	10		
SO2	4	10		
Ring / complexity	3	5	Including fused ring	
Terminale / complexity	6 (C-1)	2(C-2)	1 (C-2)	
reminals / complexity	0(0-1)	১ (८− ∠)	r (C-S)	
C atoms / complexity				
1-10 İ	11- 20	2	21-30 3	
31-40 4	41- 50	5	> 50 6	

Example of Complexity for Compounds Having Critical Temperature (Tc) Data

		CN	AC
Tc before	e 1996*	500	14
Tc after	1995**	100	21

CN - Compound Number; AC - Average Complexity

* 500 compounds having critical temperature reported before 1996.
** 100 new compounds reported between 1996 and 2001.

Example of Using the Information from Similar Compounds to Judge Uncertainty of the Value Estimated by Models

Property Estimation by Thermodynamic Models									
CASRN 7	5978		Calculate Info Print						Exit
Select a Property Select a Property normal freezing point rotical boiling point critical temperature critical pressure critical density/volume acentric factor standard state Gibbs energy of formation enthalpy change of atmospheric boiling enthalpy change of atmospheric melting ideal gas enthalpy of formation at 298 K	rmation boiling melting 298 K	Select a Model Joback Constantinou/G Marrero/Pardillo	iani				,		
 O ideal gas O vaporizat 	Gibbs energy of formation tion heat at normal boiling	on at 298 K g point	rmal boiling point t	oy Joback m	ethod = 387	7.3 (K)			
C Gibbs en C Enthalpy	ergy at 1 bar of ideal gas	Т	ne Information of S	imilar Molec	ules				_
 Entropy of O vapor pre 	or ideal gas essure	N	o CASRN	Similarity	PValue	MValue	PVError	Formula	
 Surface tension density of liquid heat capacity Thermal conductivity Viscosity Thermal diffusivity 		67641 563804 78933 108101 596220 5107879 7141797 8463821 9591786 0589388 106354 2123193 3110430	4 5 7 8 9 10 10 10 12 12 12	329.3 367.1 352.7 389.1 375 375.2 402.8 282.6 400.3 396.5 420.4 416.8 423.8	321.9 367.2 344.8 390.1 367.7 394.6 310.6 390.5 390.5 413.4 413.4 413.4	7.4 0.1 7.9 1 7.3 7.5 8.2 28 9.8 6 7 3.4 10.4	C3H60 C5H100 C4H80 C6H120 C5H100 C5H100 C6H100 C5H12 C6H120 C6H120 C7H140 C7H140 C7H140 C7H140		

Knowledge is the key to evaluate and understand scientific data as well as models

- Scientific experiment is a complicated process
- Experimental data tend to have uncertainty or error
- Evaluation of scientific data is extremely difficult, no way to guarantee their absolute correctness
- The true value of physicochemical property needs repeated experimental examination
- The above problems are also true for models

Domain Knowledge

- Thermophysics theory and concept
- Experimental and theoretical research methods
- Evaluation and comment on experimental data
- Compound physical and chemical characteristics
- Models (introduction, evaluation and comment)
- Molecular structure and interaction information
- Terminology
- Unit
- **♦**

Example about Knowledge and the Selection of Ethanol's Recommended Critical Temperature



Example about Knowledge and the Selection of Ethanol's Recommended Critical Temperature



Example of Knowledge Supporting System

🐂 TRC Knowle	edge Base	Display Records				_	
New	Save	List_All	Delete	Print	For	nt Exit	
Search Ke	ecord Number = 29	II I Domain K	nowledge of Therm	earch odynamics	Search Option	Description	•
Subject	Ethanol and Critic	al Point Measure	ement				
Description The critical temperatures of ethanol reported in the literature are divided into two distinct groups. The values reported by Young (1910), Griswold et al. (1943), Golik et al. (1955), Nozdrev (1956), McCracken et al. (1960), Efremov (1966), Marshall and Jones (1974), Hentze (1977), and Mousa (1907) all fall in the range from 516.0 to 516.7 K with uncertainties of 1 K or less, whereas the values reported by Mocharnyuk (1960), Skaates and Kay (1964), Ambrose et al. (1974), Wilson et al. (1984), and Rosenthal and Teja (1989) all lie between 513.9 and 514.2 K with uncertainties of less than 0.6 K. Similar behavior was observed in the case of the critical pressure, with values in the first group lying in the range of 6.325-6.391 MPa (reported uncertainties less than 0.001 MPa) and values in the second lying in the range of 6.129-6.148 MPa (uncertainties less than 0.02 MPa). Ethanol, like all alkanols, is very hygroscopic. Particularly, it forms an azeotrope with water. Hicks and Young (1975, Chem. Rev. 1975, 75, 119-175.) reported that critical property data for the ethanol + water system show that the influence of water is to increase the critical temperature and pressure of the mixture to a value above that for pure ethanol. Since the results of the second group of researchers), e in tre
Keyword	Ethanol, ethyl alcoh	al 64175					
Key Subject	compound	5, 51115		Idcode	9/4/2002 1:46:	50 PM	

Summary

- Uncertainty is everywhere
- Our knowledge on uncertainty is very limited
- Our awareness on uncertainty is low
- Knowledge is crucial to decrease the uncertainty
- For building a high quality information system, it is necessary to develop a strong ability for analyzing the uncertainty of data, models and text information