



Algorithm Optimization for Data Mining in Oil and Gas Exploration and Development

油气勘探开发常用数据挖掘算法优选

李大伟

中国石油勘探开发研究院



OUTLINE

- 1. Background**
- 2. Big data in petroleum industry**
- 3. Data mining in petroleum industry**
- 4. Case study**
- 5. Conclusions**



1 Background

(1) Self-introduction

Dawei Li (English name: David Lee)

- Petroleum geologist
- Got doctoral degree from China University of Geosciences (1996)
- Worked as a post-doctor in PetroChina (2001-2003)
- Managers of some large IT systems (2004-Present)
- Published about 50 technical papers and 3 books





PetroChina

1 Background

(2) E&P informatization

As the upstream of petroleum industry, E&P plays a very important role in the whole business flow of oil companies. Most data (about 90%, seismic, logging, etc.) of petroleum industry are in E&P.

By E&P informatization, management and application levels of E&P data can be improved greatly, and the efficiency of E&P research etc. can also be advanced efficiently.

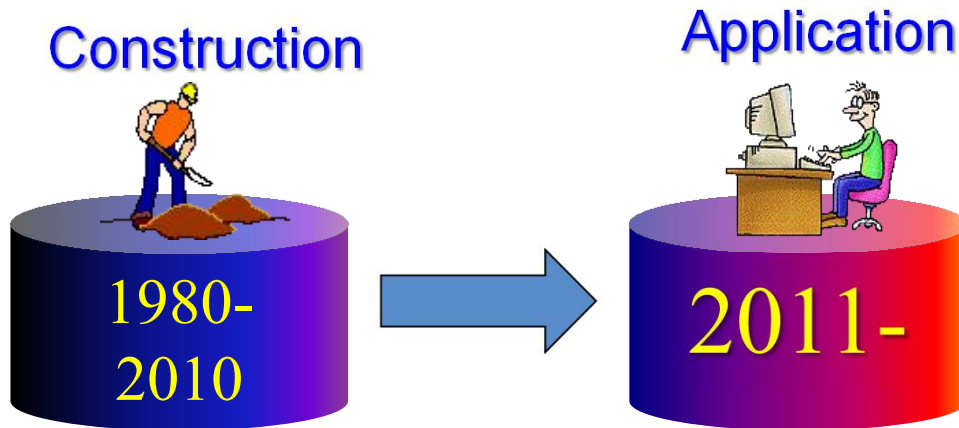
Thus E&P informatization is very important to oil companies.



PetroChina

1 Background

(3) Informatization of PetroChina



Guided by four 5-years' IT plans (2000 → 2005 → 2010 → 2015 →), PetroChina informatization has changed from construction stage to application stage since 2010.

For example:

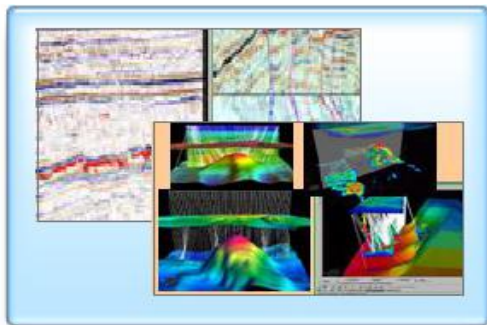
An oilfield of PetroChina: It has built a large E&P data center, including 18 types of E&P databases, about 100000 wells, 40 TB data, providing 80% data and more than 30% efficiency increase for E&P research.



OUTLINE

1. Background
- 2. Big data in petroleum industry**
- 3. Data mining in petroleum industry**
- 4. Case study**
- 5. Conclusions**

2 Big data in petroleum industry



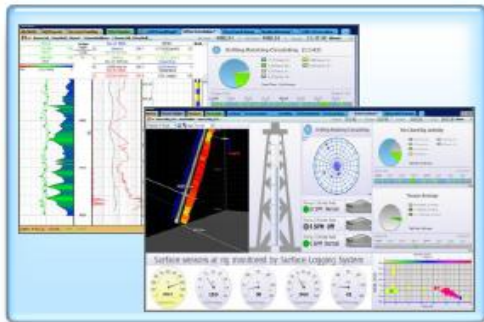
Seismic and geology modeling data



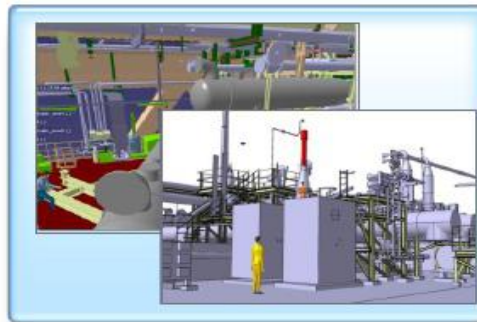
Production data



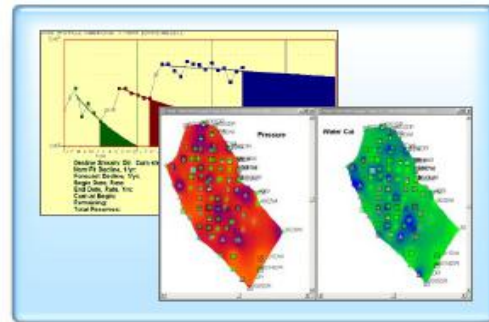
Monitoring data



Well logging data



Refining data

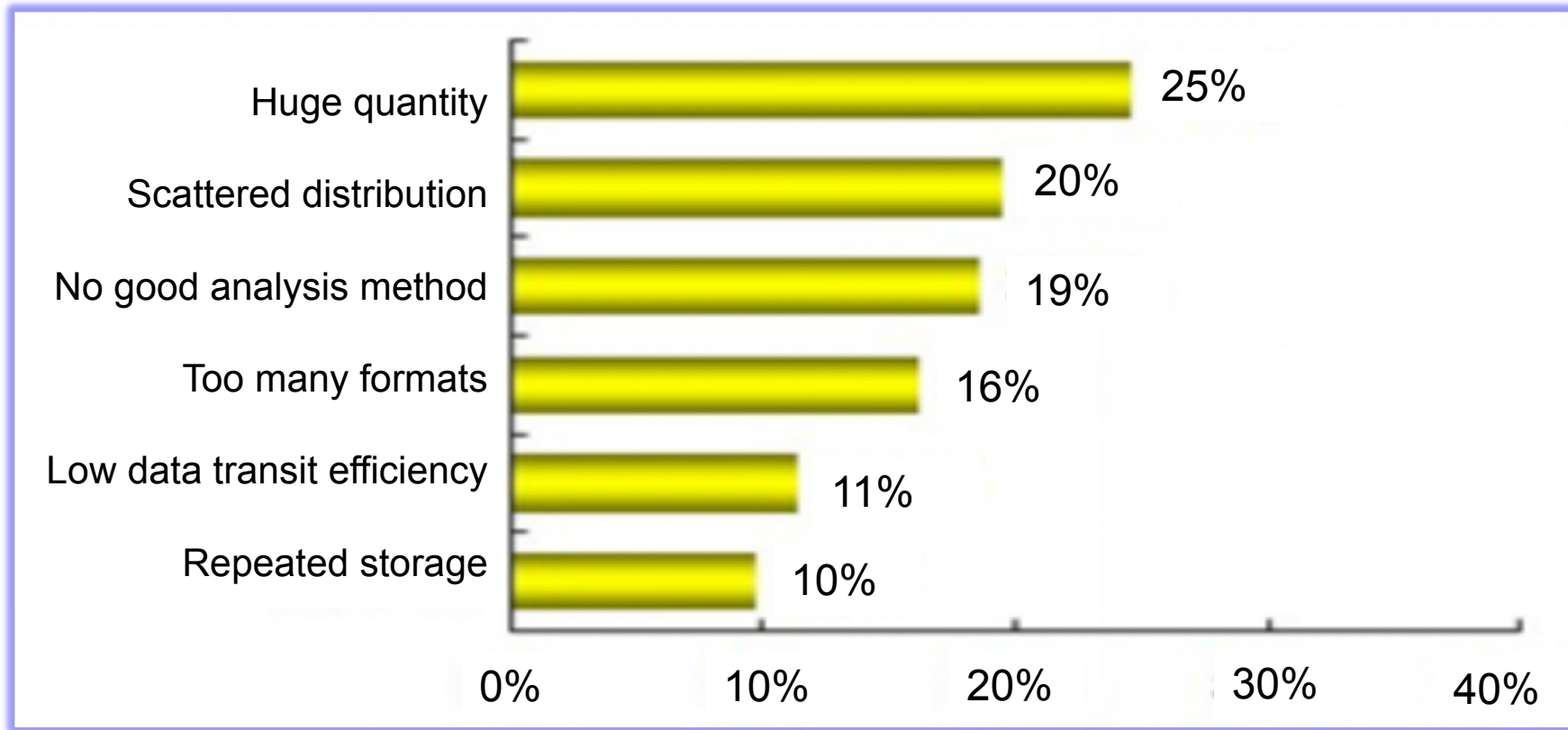


Numerical simulation data

Major data sources in petroleum industry



2 Big data in petroleum industry

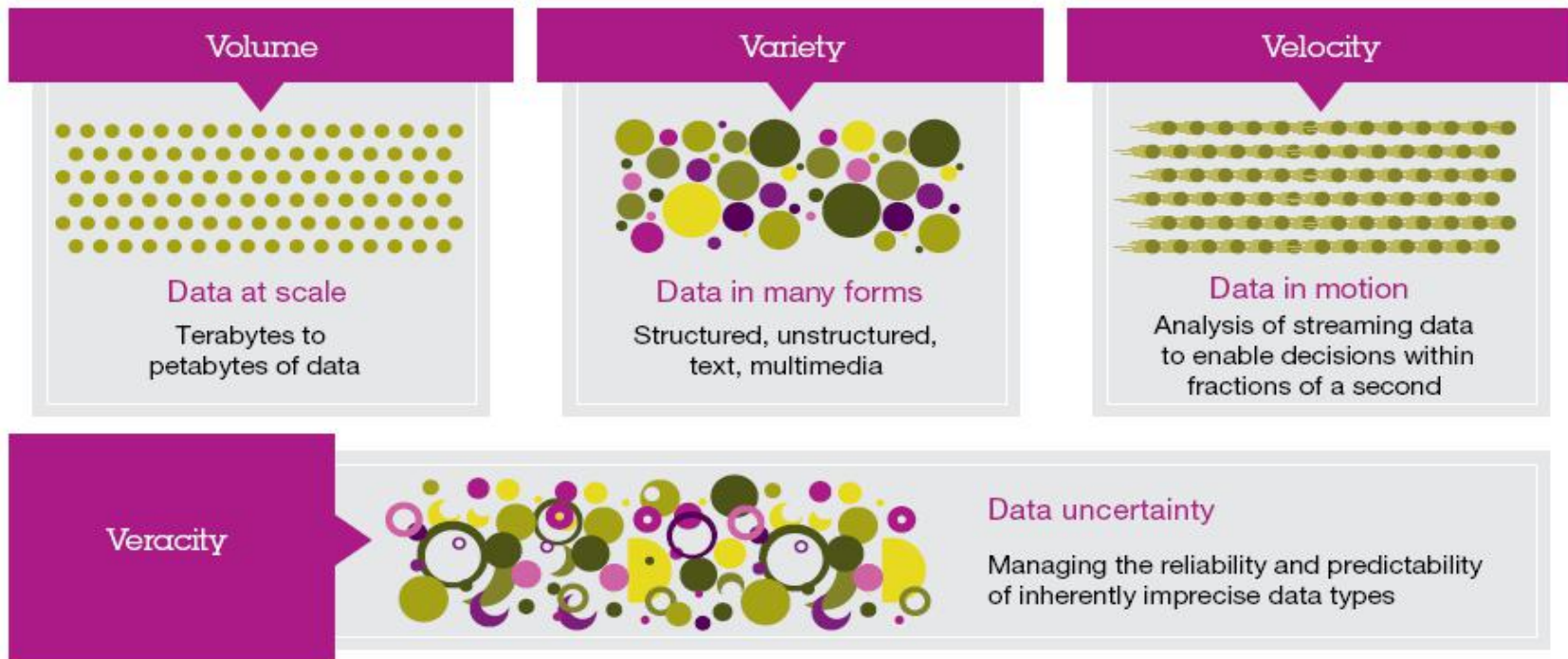


Major data problems in petroleum industry



2 Big data in petroleum industry

According to the “4Vs” of big data by IBM, petroleum industry has already entered “Big Data” era.



Example: a 3D seismic survey data in a block of Africa owned by PetroChina has 40 000 GB.



PetroChina

2 Big data in petroleum industry

● PetroChina

- 70 large IT systems
- About 2000 TB data in DBs
- Applications: data query, download, simple calculations and reports etc.
- The deep value of data has not been fully utilized



A 勘探开发与 管道项目	B 炼油化工与 销售项目	C 服务与支持 项目	D ERP项目	E 综合管理 项目	F 基础设施 项目	G 组织与保障 项目
A1.勘探与生产 业务数据管理系统	B1.炼油与化工 生产系统	C1.客户数据系统 和供应链管理	D1.ERP系统 与设备管理	E1.综合信息 系统	F1.广域网改进	G1.信息资源 管理
A2.油气生产 数据管理系统	B2.炼化物料与 生产系统	C2.设备管理 系统	D2.勘探与生产 ERP系统	E2.应急管理 系统	F2.局域网改进	G2.信息技术 标准制定
A3.勘探生产 管理系统	B3.客户关系 管理系统	C3.项目管理 系统	D3.炼油与生产 ERP系统	E3.企业信息 门户系统	F3.网络接入 系统	G3.数据资源 管理
A4.地理信息 系统	B4.加油站管理 系统	C4.设备维修 系统	D4.炼油与化工 ERP系统	E4.数据仓库 系统	F4.数据中心 建设	G4.帮助热线 建设
A5.采油与地面工 程运行管理系统		C5.物流管理 系统	D5.油田服务 ERP系统	E5.办公管理 系统	F5.企业信息 系统管理	G5.信息技术 培训中心
A6.数字盆地 系统		C6.发电供电 系统	D6.油田ERP系统	E6.档案管理 系统	F6.电子邮件 系统	G6.信息技术 培训中心
A7.工程技术生产 运行管理系统		C7.工程项目 管理系统	D7.工程ERP系统	E7.档案管理 系统	F7.网络信息 系统	
A8.勘探与生产服 务调度系统		C8.装备制造设计 与生产管理系统	D8.人力资源 管理系统			
A9.管道完整性 管理系统		C9.在用系统 推广				

12th five-year IT plan of PetroChina (after confidential treatment)



2 Big data in petroleum industry

- We have invested much in petroleum informatization.
- How to fully utilize the value of the IT systems and data assets has become a critical problem for oil companies.

There are many solutions for this problem.

Data mining (DM) is one of the good solutions.



PetroChina

OUTLINE

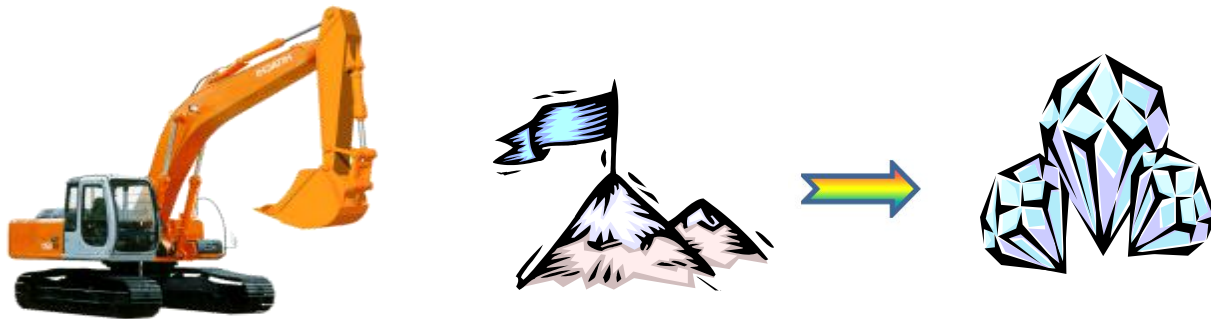
1. Background
2. Big data in petroleum industry
- 3. Data mining in petroleum industry**
- 4. Case study**
- 5. Conclusions**



3 DM in petroleum industry

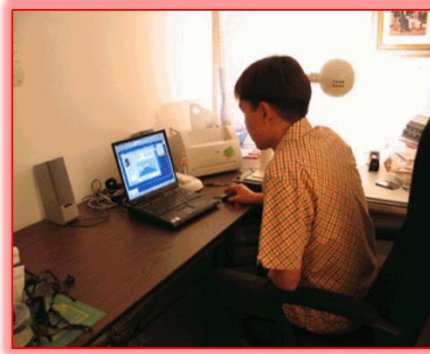
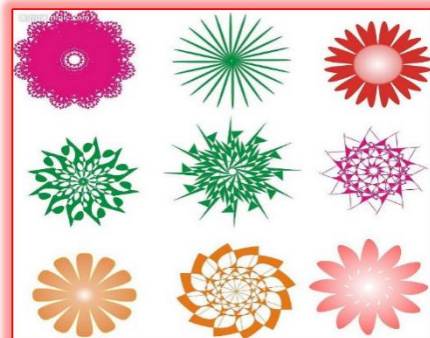
Data mining (DM) is the computing process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems ([Wikipedia](#)).

数据挖掘是在大型数据集中发现模式的计算过程, 涉及机器学习、统计和数据库系统。



3 DM in petroleum industry

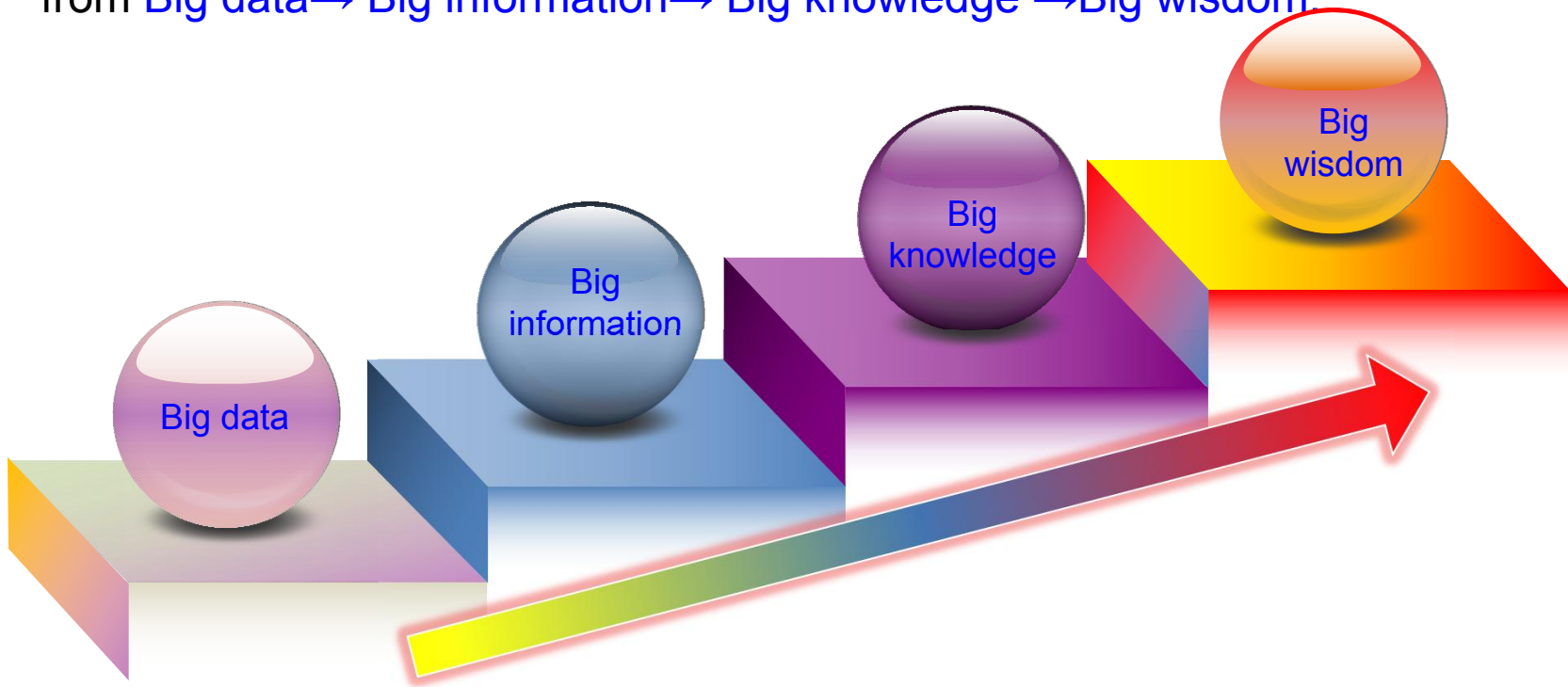
DM objects: Pure data, txt, graph, spatial data, audio data, video data, web data etc.





3 DM in petroleum industry

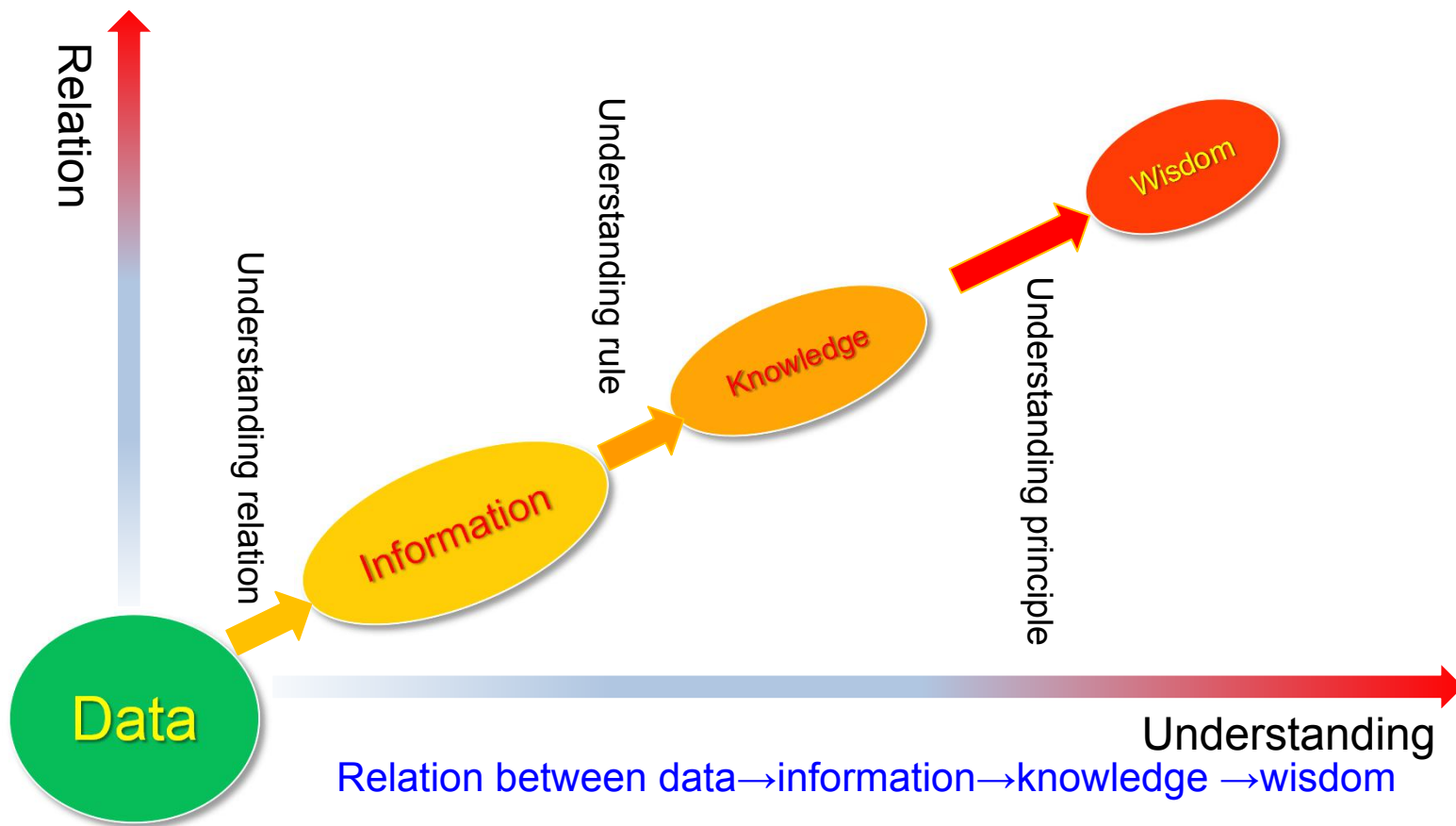
DM can fully utilize the deep value of data assets, and achieve leaps from **Big data**→ **Big information**→ **Big knowledge** →**Big wisdom**





PetroChina

3 DM in petroleum industry



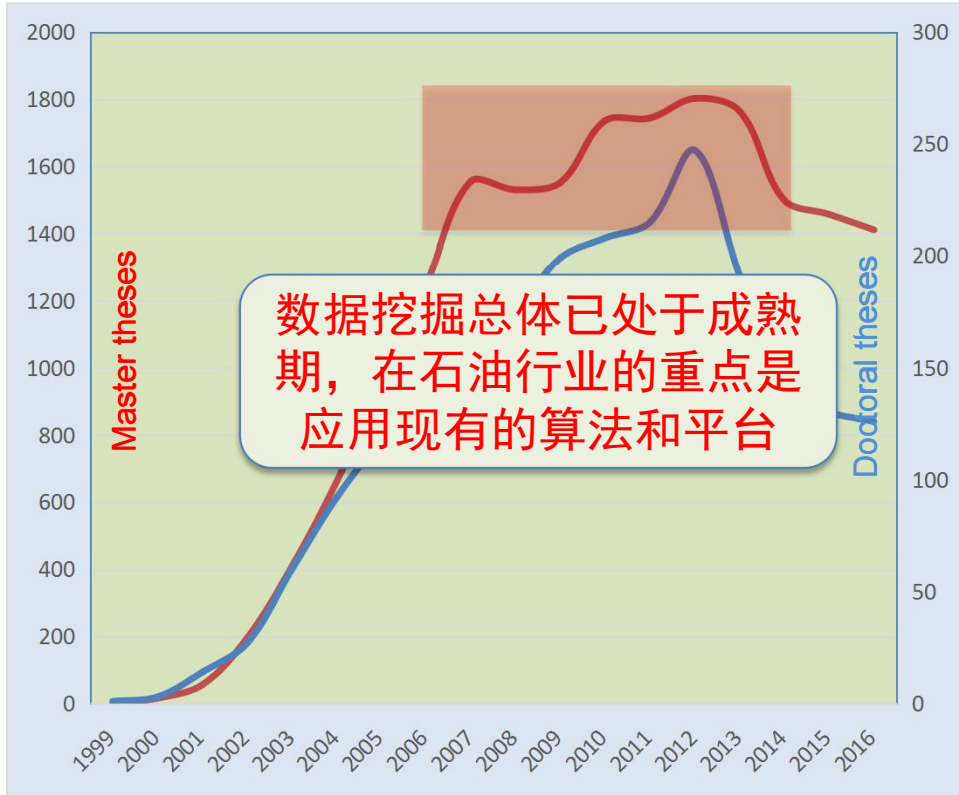


PetroChina

3 DM in petroleum industry



Statistics of published papers in various sectors related to data mining in Chinese (Dawei Li, 2016)



Statistics of master and doctoral theses related to data mining in Chinese



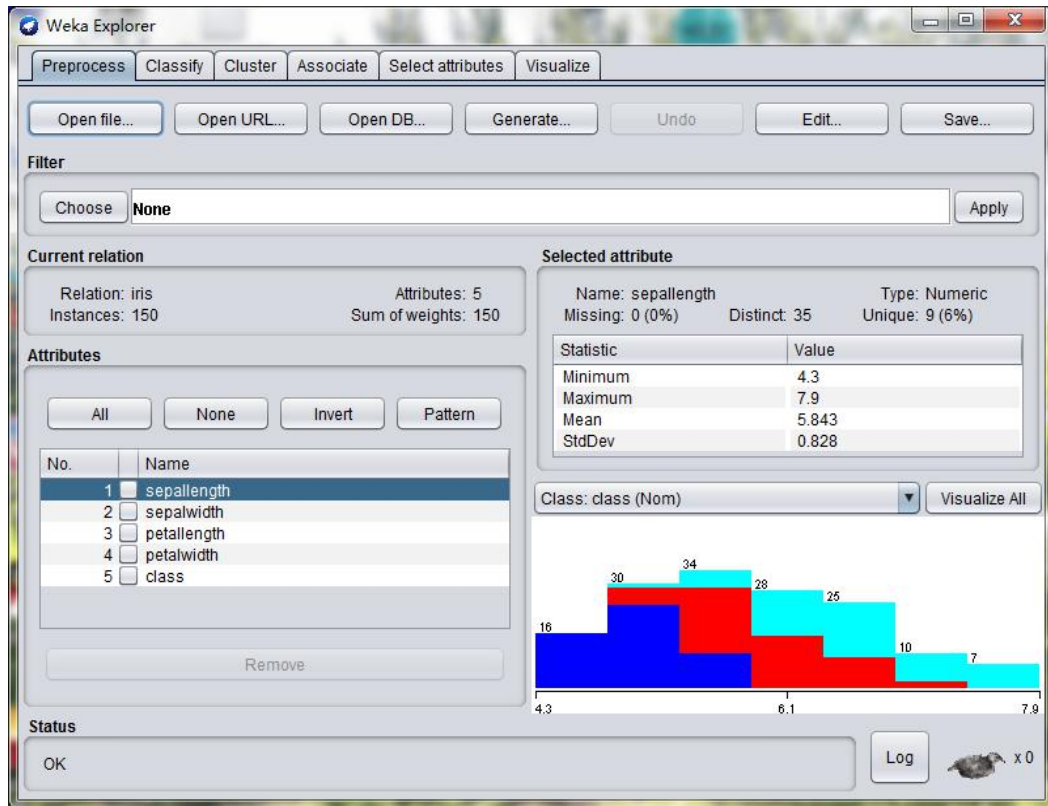
PetroChina

3 DM in petroleum industry

常用数据挖掘平台: Weka, SPSS, Rapid Miner, Matlab, TipDM...

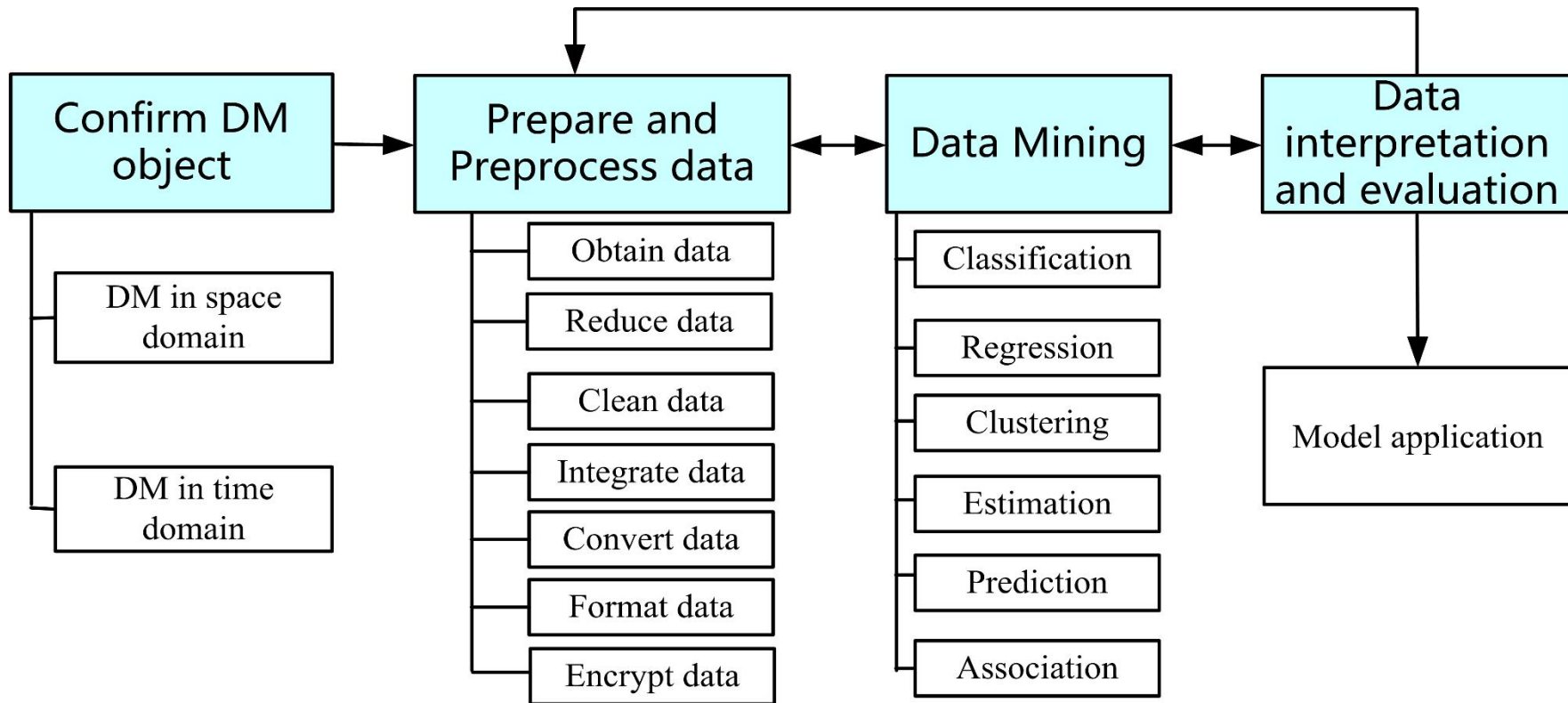


Weka (3.9.1版本) 中, 共3大类、约130种不同的算法, 其中关联算法6种、聚类算法12种、分类算法110种, 预处理算法约79种。



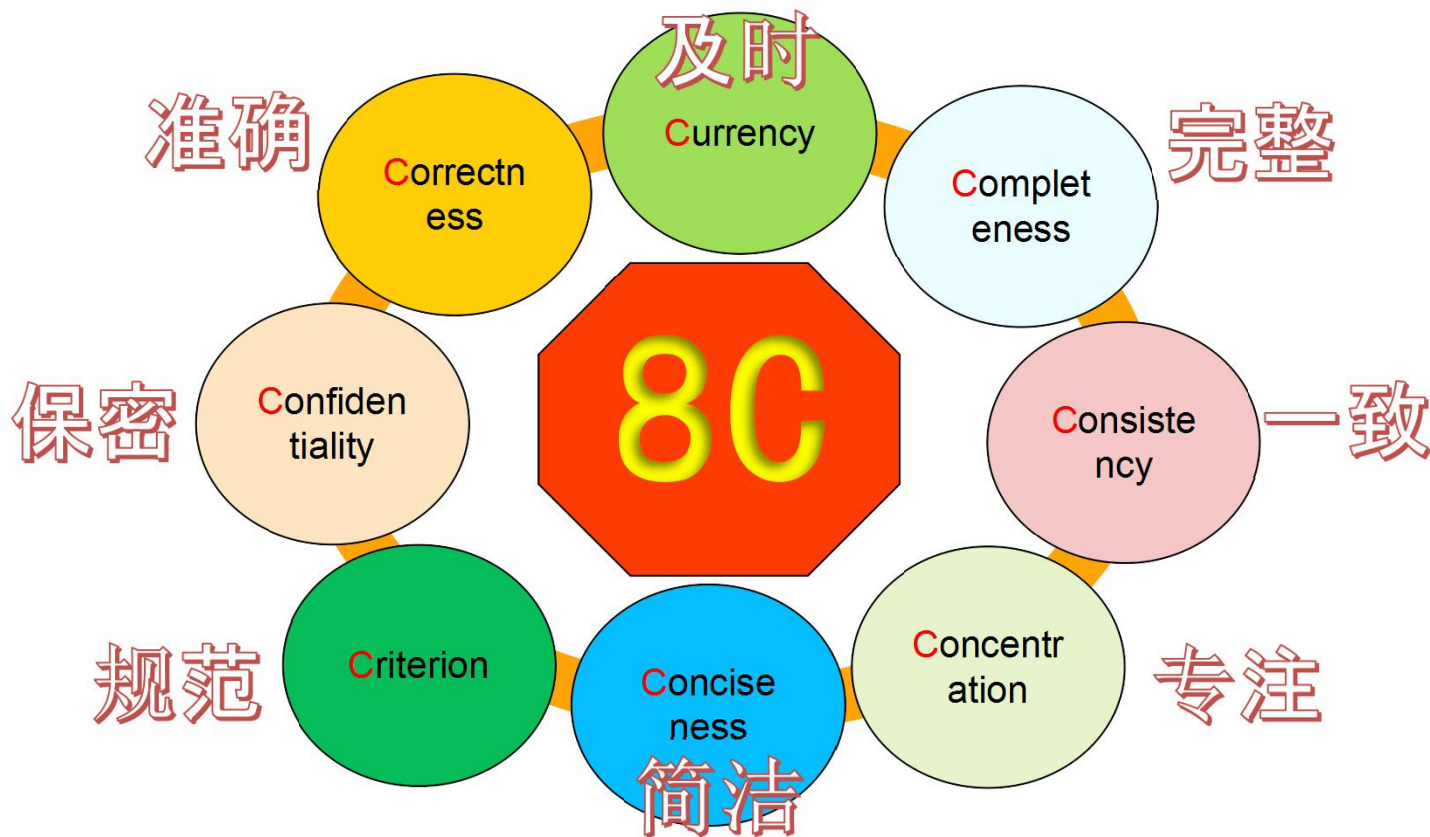


3 DM in petroleum industry



Flow chart of data mining

3 DM in petroleum industry



"8C" standard of data mining (Dawei Li, 2017)



3 DM in petroleum industry

Common DM algorithms

Regression:

- Multiple regression analysis (MRA)
- Error back-propagation neural network (BPNN)
- Regression of support vector machine (R -SVM)

Classification:

- Classification of support vector machine (C -SVM)
- Naïve Bayesian (NBAY)
- Bayesian discrimination (BAYD)



OUTLINE

1. Background
2. Big data in petroleum industry
3. Data mining in petroleum industry
- 4. Case study**
- 5. Conclusions**



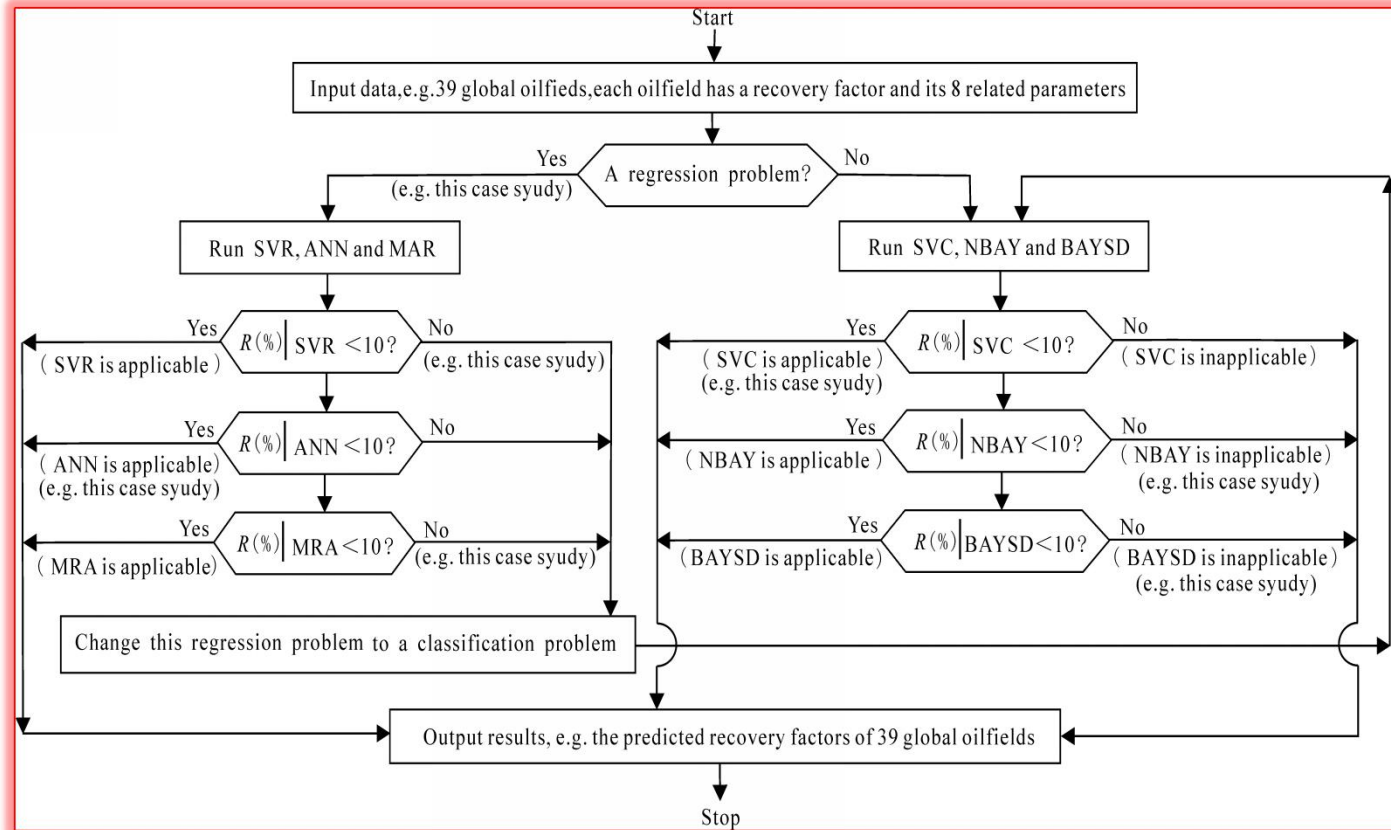
4 Case study

Based on a commercial global oil and gas field database bought from C&C (an American company), we analyzed the key affecting attributes for recovery factor and the applicability of common data mining algorithms.

Database name	Oil and gas field data	Static reservoir data	Dynamic reservoir data
Field numbers	55	377	423
Record numbers	1079	1446	699



4 Case study



Flow chart of data mining of case study

4 Case study

Major attributes of oilfield database from C&C

Oilfield code	Total production years	Current production stage code	Current production well numbers	Original in-place oil equivalent (MMBOE)	EUR oil equivalent (MMBOE)	Oil recovery factor (%)
Field 1	63	2	465	3200	311.7	12.03
Field 2	21	5	82	146	24	39.04
Field 3	46	3	193	1850	109	10
Field 4	58	5	63	405	126.8	32.59
Field 5	63	4	570	36840	1452.4	6.7
.....
Field 1076	9	2	28	75	6353	43.6
Field 1077	73	2	6150	2984	37452	33.91
Field 1078	31	3	145	66.7	1960	42
Field 1079	25	1	643	94.1	11308	14.77

Current production stage code: 1--primary rejuvenating, 2--secondary peak or plateau, 3--secondary decline, 4-- secondary mature, 5-- secondary rejuvenating, 6--tertiary peak or plateau

4 Case study

Input data for recovery factor in 39 selected oilfields

PetroChina

Sample type	Sample No.	Field No.	8 parameters related to y								y*	
			X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	RF ^b	RFC ^c
Learning samples	1	F001	63	1	465	792	3200	385	311.7	16555	12.03	5
	2	F002	21	3	82	148	146	57	24	3884	10	5
	3	F003	46	5	193	308	1850	185	109	7646	32.59	3
	4	F004	58	3	63	288	405	132	126.8	2535	6.7	5
	5	F005	63	2	570	800	36840	2470	1452.4	103000	58.97	1
	6	F006	37	3	47	96	3900	2300	1830	63500	46	2
	7	F007	34	2	201	730	25000	11500	4394	447052	32.8	3
	8	F008	53	4	7	186	25	8.2	8.1	19	15.08	5
	9	F009	21	2	122	178	325	49	25	6353	6.19	5
	10	F010	59	1	2002	3632	13436	982	267.3	71470	37.18	3
	11	F011	44	3	153	430	347	129	122	4196	42.74	2
	12	F012	61	3	185	632	15667	7355	3650	335000	33.51	3
	13	F013	41	3	93	326	6900	3580	3697	211868	30.62	3
	14	F014	37	1	3200	4600	967	324	270	20580	36.26	3
	15	F015	54	1	145	350	60520	22530	7222	500000	24.39	4
	16	F016	37	4	179	437	353	128	116.6	1910	18.13	5
	17	F017	50	2	782	1218	41000	10000	7083	350000	20	5
x_1 = total production years, x_2 = current production stage code, x_3 = current producing well numbers, x_4 = total well numbers, x_5 = original in-place oil equivalent (MMBOE), x_6 = EUR (Estimated Ultimate Recovery) oil equivalent (MMBOE), x_7 = production of cumulative oil equivalent (MMBOE), x_8 = production rate of current oil equivalent (BOEPD)], y^* = RF or y^* = RFC												
Prediction samples	35	F035	74	3	2221	6000	2800	1481	1340	26676	32.14	3
	36	F036	20	1	23	31	549	55	31	5900	(39.04)	(3)
	37	F037	74	1	2613	4121	28437	7565	5055	246441	(51.88)	(1)
	38	F038	44	4	33	73	140	45	41.8	198	(26.6)	(4)
	39	F039	25	1	327	643	637.3	94.1	59.8	11308	(14.77)	(5)



4 Case study

Classification of recovery factor

Recovery Factor	RF (recovery factor) (%)	RFC (recovery factor classification) (integer)
Very high recovery factor	>50	1
High recovery factor	$40 < RF \leq 50$	2
Intermediate recovery factor	$30 < RF \leq 40$	3
Low recovery factor	$20 < RF \leq 30$	4
Very low recovery factor	≤ 20	5



4 Case study

Prediction results of recovery factor regression in 39 global oilfields

Sample type	Sample No.	Recovery factor						
		y^*						
			R-SVR		ANN		MRA	
			y	R(%)	y	R(%)	y	R(%)
Learning samples	1	12.03	30.7558	155.659	11.6	3.88	18.7	55.1
	2	10	33.2312	232.312	9.35	6.49	35.9	259
	3	32.59	32.5433	0.143	32.1	1.67	28.9	11.4
	4	6.7	31.0992	364.167	6.19	7.61	-3.6	154
	5	58.97	32.4701	44.938	58.1	1.56	33.2	43.7
	6	46	31.862	30.735	47.2	2.52	42.5	7.62
	7	32.8	32.3528	1.363	32.9	0.169	30.1	8.21
	8	15.08	30.8125	104.327	15.2	0.629	22.6	49.7
	9	6.19	30.8843	398.939	6.19	0.0000308	16.6	168
	10	37.18	32.4168	12.811	36.4	2.05	29.3	21.3
	11	42.74	33.1068	22.539	44.4	3.87	44.2	3.34
	12	33.51	30.9235	7.719	35.0	4.53	29.7	11.4
	13	30.62	32.1243	4.913	29.5	3.60	31.9	4.23
	14	36.26	32.2356	11.099	35.6	1.88	30.7	15.4
	15	24.39	31.7768	30.286	22.4	8.10	21.8	10.6
	16	18.13	30.6217	68.901	19.5	7.55	20.6	13.7
	17	20	30.3865	51.933	16.3	18.7	18.8	6.24
	18	34.89	33.8271	3.046	32.7	6.23	40.0	14.8
	19	65.85	34.7358	47.25	67.6	2.59	69.1	4.99
	20	25.64	31.1359	21.435	26.2	1.98	20.5	19.9
	21	15.6	30.3802	94.745	17.6	12.8	18.8	20.2
	22	14.9	30.5754	105.204	15.7	5.62	19.0	27.4
	23	32.68	32.3982	0.862	32.1	1.88	29.9	8.39
	24	79.05	34.8936	55.859	79.1	0.000029	77.7	1.68
	25	35.36	32.3736	8.446	33.1	6.30	30.3	14.4
	26	34.37	32.5481	5.301	34.3	0.149	30.3	11.8
	27	63.5	33.5591	47.151	63.0	0.871	44.4	30.1
	28	11.12	30.8219	177.175	12.1	8.89	19.9	78.7
	29	13.91	33.568	141.323	12.7	8.77	40.0	188
	30	33.33	33.23	0.3	36.0	8.04	41.7	25.1
	31	37.81	32.4452	14.189	37.1	1.93	37.8	0.0841
	32	32.09	32.19	0.312	31.7	1.32	30.0	6.63
	33	52.89	33.1232	37.373	52.9	0.0276	38.6	27.1
	34	10.02	30.396	203.353	12.0	20.1	18.5	84.9
	35	32.14	32.2773	0.427	31.0	3.49	30.1	6.26
Prediction samples	36	39.04	32.2065	17.504	24.8	36.5	28.9	25.9
	37	51.88	32.7885	36.799	44.3	14.7	39.3	24.2
	38	26.6	31.6331	18.921	26.5	0.276	27.0	1.49
	39	14.77	30.447	106.141	12.9	12.6	19.8	3.3.8



4 Case study

30/33

PetroChina

Prediction results of recovery factor classification in 39 global oilfields

Sample type	Sample No.	γ^*	RFC					
			C-SVC		NBAY		BAYSD	
			γ	R(%)	γ	R(%)	γ	R(%)
Learning samples	1	5	5	0	5	0	5	0
	2	5	5	0	5	0	3	40
	3	3	3	0	5	66.7	5	66.7
	4	5	5	0	4	20	5	0
	5	1	1	0	1	0	3	200
	6	2	2	0	2	0	2	0
	7	3	3	0	3	0	3	0
	8	5	5	0	5	0	5	0
	9	5	5	0	5	0	5	0
	10	3	3	0	5	66.7	5	66.7
	11	2	2	0	2	0	2	0
	12	3	3	0	1	66.7	3	0
	13	3	3	0	3	0	3	0
	14	3	3	0	3	0	3	0
	15	4	4	0	4	0	4	0
	16	5	5	0	5	0	5	0
	17	5	5	0	5	0	5	0
	18	3	3	0	5	66.7	3	0
	19	1	1	0	1	0	1	0
	20	4	4	0	5	25	4	0
	21	5	5	0	5	0	5	0
	22	5	5	0	5	0	5	0
	23	3	3	0	5	66.7	5	66.7
	24	1	1	0	1	0	1	0
	25	3	3	0	3	0	3	0
	26	3	3	0	3	0	3	0
	27	1	1	0	5	400	3	200
	28	5	5	0	5	0	5	0
	29	5	5	0	5	0	3	40
	30	3	3	0	3	0	3	0
	31	3	3	0	3	0	3	0
	32	3	3	0	3	0	3	0
	33	1	1	0	1	0	3	200
	34	5	5	0	5	0	5	0
	35	3	3	0	3	0	3	0
Prediction samples	36	3	3	0	5	66.7	5	66.7
	37	1	1	0	3	200	5	400
	38	4	4	0	3	25	4	0
	39	5	5	0	5	0	5	0



Summary of the case study with 39 global oilfields

Problem type	Algorithm	Mean absolute relative residual error	Time consuming on PC	Conclusion
		R(%) (平均相对剩余误差的绝对值)		
Regression	SVR	68.9	3 s	Inapplicable
	ANN	5.89	30 s	Applicable
	MRA	38.4	<1 s	Inapplicable
Classification	SVC	0	5 s	Applicable
	NBAY	24.7	<1 s	Inapplicable
	BAYSD	34.5	1 s	Inapplicable

1. This program ran on a PC with configuration: HP Z230, Windows 7 (64 bit), Intel E3-1231, 3.40 GHz, 16GB RAM
2. If $R > 10\%$, it is inapplicable.



PetroChina

OUTLINE

1. Background
2. Big data in petroleum industry
3. Data mining in petroleum industry
4. Case study
- 5. Conclusions**



Conclusions

- Petroleum industry has entered “Big Data” era;
- Data mining (DM) is one of the good solutions to fully utilize the value of the IT systems and data assets;
- For the case study, the best algorithm is ANN for recovery factor regression, while the optimum algorithm is SVC for recovery factor classification;
- The best algorithms are different for various data sets.



Thank You !

Questions?



Email: leedw@petrochina.com.cn

Wechat: 13671142720