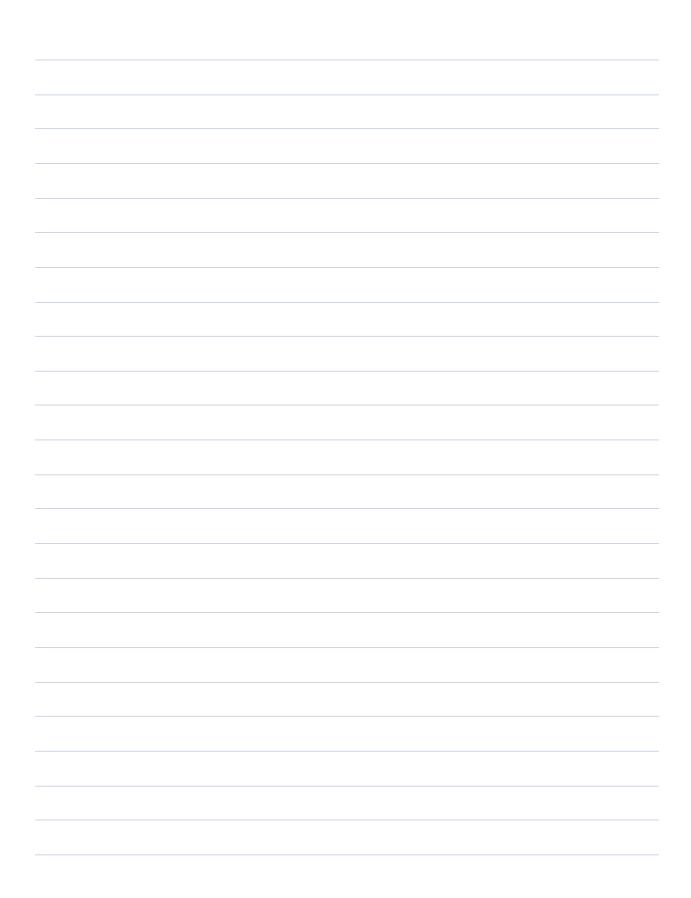
净资考高级信息系统
1事资考高级信息系统 要上
- 时代沿草中的信息系统
- 时代沿草中的温息系统 - 耐名分析与使能
一行为与实证分析 一不确定性信息 一须成展望
- 不确定性注意
一须成春望

一、准备。
一大数据时代
① 数据特征 (4V)
Volume 起规模, variety 落块体 (多质异构) value (纸价值密度) Velocity (流跃器)
2 NO BYTHE

2. 行复系程从开究	
l	



二、港式转逐	

三、技术的经济风雨

四、关联模式

1. 支持度与置信度

 $I=\{I,\dots,Im\}: 商品集合 , T=\{t,\dots tn\}, ti \in I: 事務据集$ 支持度 $Dsupp(X)=\frac{I|X|}{|T|}$,

11×11: 7中包含×的沿湖。 171: 总记储数

须受支持度1图值 d. Dsupp(x) ≥ d 刚积 ×为频繁项集

英觀则 $X \Rightarrow Y (X,Y \neq \emptyset, X,Y \leq I, X \cap Y = \emptyset)$ Dsupp $(X \Rightarrow Y) = D$ supp $(X \cup Y)$

The Denf
$$(X \Rightarrow Y) = \frac{||X \cup Y||}{||X||} = \frac{D \operatorname{supp}(XY)}{||X||}$$

Dsupp(XY) >2, Doonf (XY) > B, X=Y & 2. BAR.

2. 英联规则挖掘游荡 (Aprilor)

分数两大步

- ① 生成频繁项集 (Frequent itemset generation)
- ② 報文班共和 (Association rule generation)

门生放频强项集

Firs Dsup(x) = $\frac{||x||}{|T|} \ge \frac{||xY||}{|T|} = Dsupp(xY)$

如果一个集合不是频繁顶集,则其所有的超集均不是

- Grenerate frequent itemset LICK=1)
- Repeat until no new frequent itemsets are identified.

O Cardidate generation: generate Ck+1 from 4 连接 @ Prune candidate itemsets Ck+1 \$75 3 Court support for each itemset in Ca+1 (12/2014) @ Remove infrequent itemsets in Ca+1 and get Lk+1 Lz={AB, AC, AD, BD, EF, DF} 生成设造了一项集》 C3= {ABC, ABD, ACD} 证:按字到顺序排例, 当新k-1个项相同,等k个不同时,各 其分不考虑,如:ABE,EDF, ADF等 C3中的 BC, CD & L2, 二去掉 ABC, ACD 二 Co = {ABD} . 三后遍历数据集着 ABD 是否满足支持度部分、进而假 到 Lz ② 英秋 关联共和 松口:置待度对于从同一个预集中生成的规则而是是单调的 Doonf (ABC>D) > Doonf (AB>CD) > Doonf (A>BCD)