

### Data Mining for Data Modelers

Asoka Diggs

Data Scientist, Intel IT

March 5, 2015



### Legal Notices

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

\* Other names and brands may be claimed as the property of others.

Copyright ° 2014, Intel Corporation. All rights reserved.

### 2014 Intel IT Vital Statistics

#### >6,300 IT employees 59 global IT sites

#### >98,000 Intel employees<sup>1</sup> 168 Intel sites in 65 Countries

### 59 Data Centers

(91 Data Centers in 2010) 80% of servers virtualized (42% virtualized in 2010)

### >147,000+ Devices

100% of laptops encrypted 100% of laptops with SSDs >43,200 handheld devices 57 mobile applications developed

Source: Information provided by Intel IT as of Jan 2014 <sup>1</sup>Total employee count does not include wholly owned subsidiaries that Intel IT does not directly support,



Copyright © 2014, Intel Corporation. All rights reserved.



Introduction and assumptions Data Mining demo Break Overview of Data Mining concepts

Questions throughout Resources and links for further study



# My Background

### DBA starting in late 90's

Various data management roles since then

- Entity-Relationship modeling
- ETL development
- Data integration
- Data architecture
- MS Predictive Analytics started Fall 2011
- 9/11 classes completed

Designed and teach an internal Intro to Data Mining class





## Assumptions

Audience is composed of experienced Data Management Professionals

You're interested in Predictive Analytics / Data Mining but haven't studied it

You're already convinced of the value of Predictive Analytics / Data Mining

• We won't particularly be covering the value here

# What you will learn

Class is designed to teach a functional mental model of data mining Definition: Mental Model

- A mental model is an explanation of someone's thought process about how something works in the real world
- acts as information filter causes selective perception or perception of only selected parts of information

With the mental model, you'll be making your own connections between data mining and skills / experience you already have

Goal: Begin learning how to frame business problems as data mining problems

- Real goal: realize that you can learn to frame business problems as data mining problems with minimal effort
  - But more effort than a 3 hour demo version

http://en.wikipedia.org/wiki/Mental\_model

### Demo

Decision Tree and

- Iris data set
- eReader Adoption data set





### Effort to achieve skill level



intel

#### IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved.

### Deliverable Expectations





# A variation on <u>CRiSP-DM</u>



Copyright © 2014, Intel Corporation. All rights reserve

IT@Intel

intel



#### IT@Intel

inte

# Machine Learning (ML) input and output

Any ML algo will have certain kinds of data it can take as input And will produce certain kinds of data as output

- Example: Logistic Regression requires numeric input (including binary variables), and generates a binary (1/0) output
- Data prep activity converting data from an unusable form into a usable form

The nature of the data mining problem will shape the kind of answer needed (binary, category, number, ..)

• Which will start filtering the available techniques

Further filter the techniques based on the form of the data that is available

And we will talk more about this



# Supervised and Unsupervised learning

Machine Learning and predictive modeling falls into two buckets...

- Supervised or Directed (Classification)
  - Presence of a response or outcome variable
  - Finds relationships in the data
  - ~75% of the techniques
  - Generally more interesting and useful biz q's can be addressed

- Unsupervised or Undirected (Clustering)
  - Absence of a response or outcome variable
  - Finds patterns in the data
  - Not as many techniques
  - More easy to find patterns that aren't useful

Think of response variable as "human input" into the model training and validation process



class	sepallength	sepalwidth 🔺	petallength	petalwidth
Iris-versicolor	5	2	3.500	1
Iris-versicolor	6	2.200	4	1
Iris-versicolor	6.200	2.200	4.500	1.500
Iris-virginica	6	2.200	5	1.500
Iris-setosa	4.500	2.300	1.300	0.300
Iris-versicolor	5.500	2.300	4	1.300
Iris-versicolor	6.300	2.300	4.400	1.300
trie versiceler	F	2 200	2 200	4



#### Training Data

Build Analytic Model

**Model Validation** 

#### Scoring Data

Scored Data



			TRUE		
		Iris-	Iris-	Iris-	class
		setosa	versicolor	virginica	precision
	lris- setosa	50	0	0	100.00%
	Iris-	0	47	6	88.68%
PREDICTED	Iris-	0	3	44	93.62%
	class	100.00%	94.00%	88.00%	
	Tecall				

sepallength	sepalwidth 🔻	petallength	petalwidth
5.700	4.400	1.500	0.400
5.500	4.200	1.400	0.200
5.200	4.100	1.500	0.100
5.800	4	1.200	0.200
5.400	3.900	1.700	0.400
5.400	3.900	1.300	0.400
5.700	3.800	1.700	0.300
5 100	3 800	1 500	0.300

class	sepallength	sepalwi 🔻	petallength	petalwidth
Iris-setosa	5.700	4.400	1.500	0.400
Iris-setosa	5.500	4.200	1.400	0.200
Iris-setosa	5.200	4.100	1.500	0.100
Iris-setosa	5.800	4	1.200	0.200
Iris-setosa	5.400	3.900	1.700	0.400
Iris-setosa	5.400	3.900	1.300	0.400
Iris-setosa	5.700	3.800	1.700	0.300
tric cotoco	5 100	2 000	1 500	0.200

#### IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved.

ínte

# Variable modeling types (aka data types):

- Discrete / Categorical / Qualitative (alphanumeric/text)
  - Nominal no meaning in order
    - Distinctness ( = != ); Example: Name, favorite color
    - Binary/Binomial variables Represents presence or absence of something
      - » Not necessarily Black/White; Example: gender
  - Ordinal order makes sense
    - Order ( < > ); Example: Order ID, Pay Grade, low/medium/high
- Continuous / Quantitative (numbers only)
  - Interval sum/difference meaningful, but product/ratio not
    - Addition (+ -); Example: temperature in Fahrenheit, time of day AM/PM
  - Ratio ratios make sense in addition to sum/difference
    - Multiplication ( \* / ); Example: absolute time, absolute temp, height, weight, wavelength

### Grain: Another way to think about data

Level of detail - What is an observation (row), an observation of?

- Think of the grain in a piece of wood
  - Example: Cumulative GPA vs. per course grade-point (transcript)

Has proven useful in reporting, data structure, formulating ML problems, etc..

- Can you identify the grain of an answer to the business question?
- What is the grain of the source data you are analyzing?
- Is the grain of the (data mining) analysis of the source data compatible with the grain of the business question?
  - Example: GPA is not a good grain to analyze "fairness" across faculty

### Don't mix grain in a data set

Or at least only mix grain consciously and intentionally

If you ever have an itchy feeling that something about the data is wrong, start with the grain

Can you even identify what the grain is?

# Data Terminology

### Data Model

- Usually NOT an ERD (rather, a data mining model)
- I am trying to use ER Model and Analytic Model
  Observation
- Row, example, case, tuple, sample, member

### Variable

• **Column**, attribute, field, feature, parameter, dimension, metric, measure

### Roles played by variables

Label (dependent), attribute (independent), id

Emphasis in data mining is on denormalized data

Tuple		ATT	ribu	ITES		LABEL
(Row)	X1	X2	Х3		Xn	Y
1						
2						
3						
4						
5						
6						
•						
•						
Μ						



# Modeling Terminology

### Data Set

- Training. Used to train the data (mining) model
  - The response is a known value
- Test. An out of sample (observations not part of the training set) used to validate the results of the model
  - The response is a known value
- Scoring. A set of observations where the response or outcome is unknown. We use the model created by the data mining method and the training set to Score or create values for the response variable
  - After scoring, the response is predicted
    - NOT determined





# Decision Tree

Output / response / label: categorical (or continuous)

Input / observation / attribute: categorical or continuous

Model: (example below)

- Very friendly for human interpretation
- <u>http://en.wikipedia.org/wiki/Decision\_t</u>
  <u>ree\_learning</u>
- <u>http://en.wikipedia.org/wiki/C4.5\_algo</u>
  <u>rithm</u>



A tree showing survival of passengers on the Titanic ("sibsp" is the number of spouses or siblings aboard). The figures under the leaves show the probability of survival and the percentage of observations in the leaf.



# Model Validation and Testing

<b>≺</b>	ginal Set	
Training		Testing
Training	Validation	Testing
FIGURE 3. Two and three way splitting		
<b>`raining set</b> – a set used for learning and	d estimating para	meters of the model.
Alidetion set - a set used to evaluate the	o modol uquallu	for model coloction

**Testing set** - a set of examples used to assess the predictive performance of the model.

### **Citation**

# Validation approaches

### Most important part of model building

- How do you know the model you've built is reasonable?
  - This answers the technical question
  - Still need to answer the business question

### 70/30

- 70% of training data as an in-sample model building / training data set
- 30% of training data as an out-of-sample testing or validation set
  - Both training and testing data have a known value for the response variable
  - (Y has an actual value for all observations)
- Scoring data the Y value is unknown (but we'd like to have an idea of what it is likely to be)

# Confusion Matrix

### We'll look at the model validation results together

- More reading on Precision and Recall
- http://en.wikipedia.org/wiki/Precision\_and\_recall

"In simple terms, high **precision** means that an algorithm returned substantially more relevant results than irrelevant, while high **recall** means that an algorithm returned most of the relevant results."

- Quality
- Quantity

# Logistic Regression

### Output / response: binary (1/0)

Also polynomial; not covered in this book
 Input / observation: numeric

- Can use categorical data after converting to numeric
  - Discretization
  - Binarization or <u>Dummy variables</u>

### Model: logistic equation

http://en.wikipedia.org/wiki/Logistic\_regression

Conceptually:

- Event occurred / didn't occur
- We care / don't care
- True / false
- Black / not-Black

Absence of proof, is not proof of absence



### Resources



Copyright © 2014, Intel Corporation. All rights reserved



# Self-Study Resources

Textbook: Data Mining for the Masses, by Matthew North Method:

- Download the textbook
- Follow the step-by-step instructions from the beginning to the end of the book
- By the time you finish the last chapter, you will have:
  - Analyzed 4 data sets using different supervised learning methods
  - Analyzed 2 more data sets using unsupervised learning methods
  - Used rudimentary text analytics on a data set
  - Incorporated model validation into at least 1 model

Our demo today has made use of 2 of these data sets and is drawn from this book

# Additional Resources

Broadening Access to Advanced Analytics in the Enterprise, by Asoka Diggs & Christy Foulger

- And the subject of my Tuesday afternoon presentation "Democratization of Data Analytics"
- Short video series on Text Analytics, by Neil McGuigan

Data Science Specialization through Coursera

• For those looking for an Intermediate level of knowledge of data science

<u>The Signal and the Noise: Why So Many Predictions Fail – but Some Don't</u>, by Nate Silver

Contact me at asoka.diggs@intel.com



### IT@Intel Sharing Intel IT Best Practices With the World



# Learn more about Intel IT's Initiatives at www.intel.com/IT

#### IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved





# Questions?



# Thank You

1	ra	INI	ng	Da	Ita
			<u> </u>		

class	sepallength	sepalwidth 🔺	petallength	petalwidth
Iris-versicolor	5	2	3.500	1
Iris-versicolor	6	2.200	4	1
Iris-versicolor	6.200	2.200	4.500	1.500
Iris-virginica	6	2.200	5	1.500
Iris-setosa	4.500	2.300	1.300	0.300
Iris-versicolor	5.500	2.300	4	1.300
Iris-versicolor	6.300	2.300	4.400	1.300
Iris-versicolor	5	2 300	3 300	1

IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved

31

ínte



**Training Data** 

#### **Build Analytic Model**

IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved

32

íntel

	class	sepallength	sepalwi	dth 🔬 🛛 petall	ength	n petalwidth		
Ir	ris-versicolor	5	2	3.500		1		
Ir	ris-versicolor	6	2.200	4		1		
Ir	ris-versicolor	6.200	2.200	4.500		1.500		petallength
Ir	ris-virginica	6	2.200	5		1.500		
Ir	ris-setosa	4.500	2.300	1.300		0.300	>2	450 ≤ 2.450
Ir	ris-versicolor	5.500	2.300	4		1.300		
Ir	ris-versicolor	6.300	2.300	4.400		1.300	petalwidth	1
Ir	ris-versicolor	5	2 300	3 300		1		
						Iris-1	> 5.350	≤ 5.350 Iris-ver
						TRUE		
				Iris-	- Ir	ris-	Iris-	class
				setosa	V	ersicolor	virginica	precision
		lris- setos	а	5	0	0	0	100.00%
F	PREDICTE	lris- D versio	color		0	47	6	88.68%
		lris- virgin	ica		0	3	44	93.62%
		class recall		100.009	%	94.00%	88.00%	

IT@Intel

Copyright © 2014, Intel Corporation. All rights reserved



class	sepalle	ngth sep	alwidth 🔺	petallen	gth petalwid	th
s-versicolor	5	2		3.500	1	
-versicolor	6	2.20	00	4	1	
s-versicolor	6.200	2.20	00	4.500	1.500	petallengt
s-virginica	6	2.20	00	5	1.500	
s-setosa	4.500	2.30	00	1.300	0.300	> 2.450
is-versicolor	5.500	2.30	00	4	1.300	
is-versicolor	6.300	2.30	00	4.400	1.300	petalwidth
is-versicolor	5	2 30	00	3 300	1	
			TRUE		Iris-virgin	ica petallength > 5.350
		Iris-	Iris-	Iris-	class	Iris-virginica
1	Iris-	setosa	versicolor	virginica	precision	
5	setosa	50		) 0	100.00%	
	lris- versicolor	0	47	' 6	88.68%	
N	Iris- virginica	0	) =	3 44	93.62%	
	class recall	100.00%	94.00%	<mark>i 88.00%</mark>		
Į.	lecall					
sepalle	ength	sepa	lwidth	v p	etallengt	h petalwidth
5.700		4.400	)	1.	500	0.400
5.500		4.200	)	1.	400	0.200
5.200		4.100	)	1.	500	0.100
5.800		4		1.	200	0.200
5.400		3.900	)	1.	700	0.400

0.400

0.300

1.300

1.700

IT@Intel

5.400

5.700

3.900

3.800

(intel)

class	sepall	length	sepalw	vidth 🔺	petalle	ngth	petal	width						
ris-versicolor	5		2		3.500		1							
ris-versicolor	6		2.200		4		1							
ris-versicolor	6.200		2.200		4.500		1.500	(		petallengt	h			
is-virginica	6		2.200		5		1.500			/				
is-setosa	4.500		2.300		1.300		0.300	`		> 2.450	≤ 2.450		- 1	
is-versicolor	5.500		2.300		4		1.300			F	A			Training Data
is-versicolor	6.300		2.300		4.400		1.300		P	etalwidth	Iris-se	tosa		Training Data
							Iris-vi	> 1.75	50	≤ 1.750 petallength				Build Analytic Mod
•	Iris-	lris- seto	sa ve	TRUE is- ersicolor	Iris- virginica	clas a pre	ss cision	Iris-virgir	hica	- 5.350	s 5.350 Iris-versicolo			Model Validation
	setosa Iris-		0	47		5 8	8.68%							
PREDICTED	versicolo Iris- virginica		0	3	44	4 9	<mark>3.62%</mark>							Scoring Data
	class	100	.00%	94.00%	88.00%	6		sepallen	igth	sepalwidth	petallength	petalwidth	n	
	recall							5.700		4.400	1.500	0.400	_ (	
								5.500		4.200	1.400	0.200		Scored Data
								5.200		4.100	1.500	0.100	_	Cooled Bala
								5.800		4	1.200	0.200		
								5.400		3.900	1.700	0.400		
								5.400		3.900	1.300	0.400		
								5.700		3.800	1.700	0.300		
class	S	epall	ength	sepa	alwi		petal	length	р	etalwidth				
Iris-setos	sa 5.	700		4.40	0		1.500		0.4	100				
Iris-setos	sa 5.	500		4.20	0		1.400		0.2	200				
Iris-setos	sa 5.	200		4.10	0		1.500		0.1	100				
Iris-setos	sa 5.	800		4			1.200		0.2	200				
Iris-setos	sa 5.	400		3.90	0		1.700		0.4	400				
Iris-setos	sa 5.	400		3.90	0		1.300		0.4	100				
Iris-setos	sa 5.	700		3.80	0		1.700		0.3	300				

#### IT@Intel

íntel

