

advanced issues in data preparation and modeling

Carlos Soares

(partly using materials kindly provided by José Luís Borges, from Han, Kamber & Pei, from Moreira, Carvalho & Horvath and from Ceja)

plan

- advanced issues in data preparation
 - data reduction
- ... and modeling
 - dealing with unbalanced classes
- final reflections on data quality

- advanced issues in data preparation
 - data reduction
 - context
 - attribute aggregation
 - feature selection
- ... and modeling
 - dealing with unbalanced classes
- final reflections on data quality

data reduction

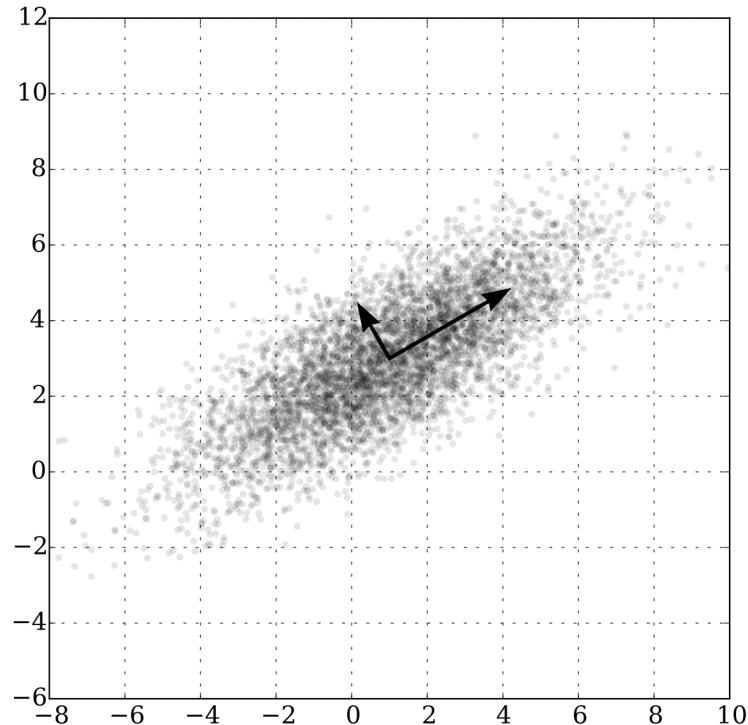
- obtain a reduced representation of the data set that is
 - much smaller in volume
- ... producing the **same analytical results**
 - or almost the same
- ... improved visualization of data
- ... with more interpretable models
- ... much faster

dimensionality reduction: data is cursed!

- curse of dimensionality
 - when dimensionality increases, data becomes increasingly sparse
 - density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
 - the possible combinations of subspaces will grow exponentially
- ... number of data points required for robust patterns grows exponentially with number of attributes

two approaches

attribute aggregation

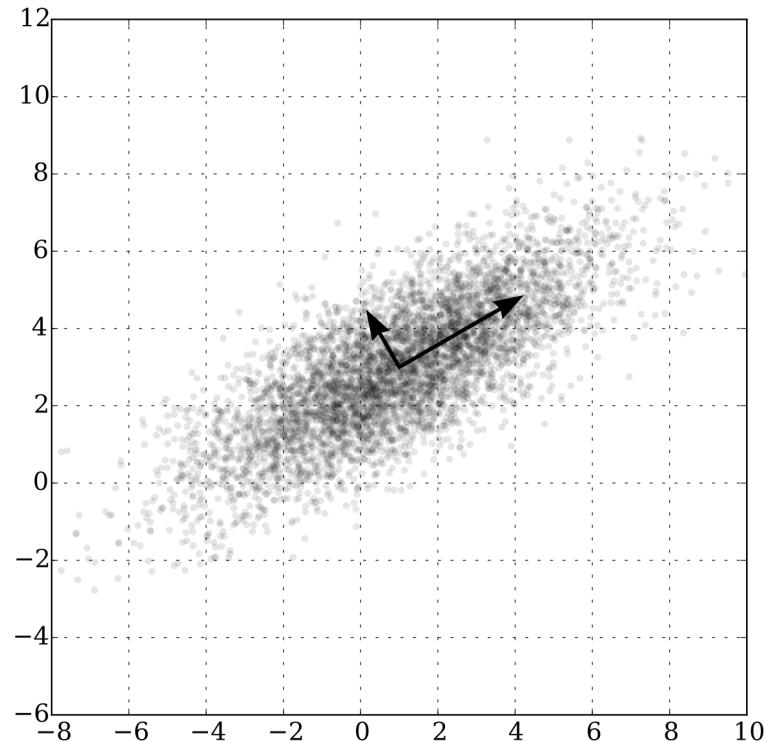


feature selection



attribute aggregation: PCA

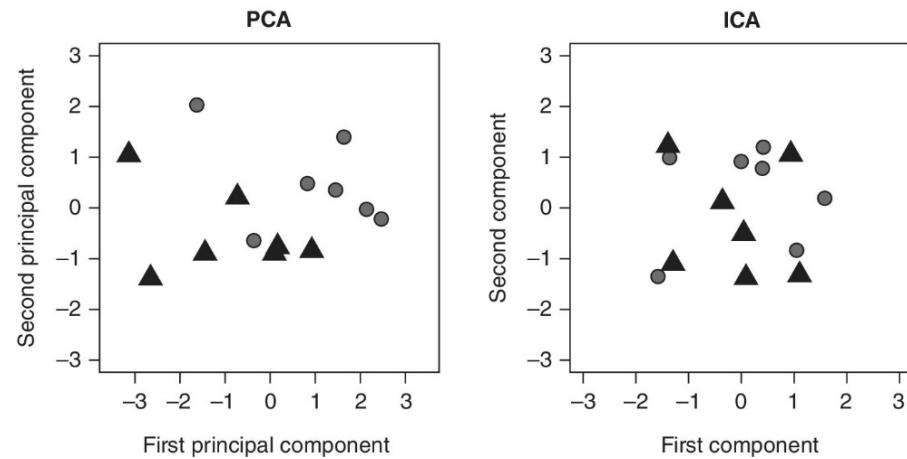
- Principal Component Analysis
 - n new features
 - linear combinations of existing n attributes
 - orthogonal to each other
 - $k \ll n$ explain most of the variance in the data



By Nicoguaro - Own work, CC BY 4.0, <https://commons.wikimedia.org/w/index.php?curid=46871195>

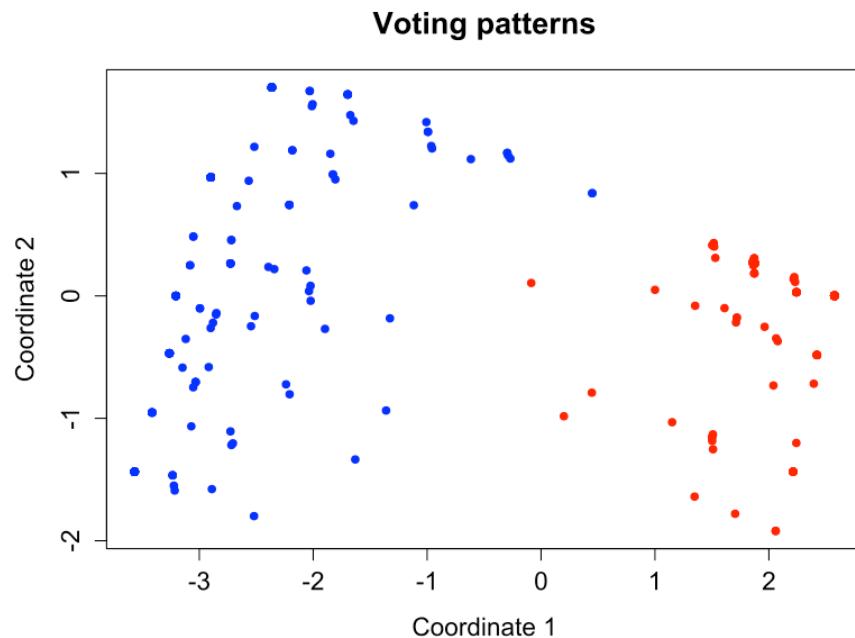
attribute aggregation: ICA vs PCA

- both create linear combinations of the attributes
- ICA assumes the original attributes are statistically independent
- ... reduces higher order statistics
 - e.g. kurtosis
- ... does not rank components



attribute aggregation: multidimensional scaling

- linear projection of a data set
- uses the distances between pairs of objects
 - not the values of the attributes of the objects
- particularly suitable when it is difficult to extract relevant features to represent the objects



https://en.wikipedia.org/wiki/Multidimensional_scaling

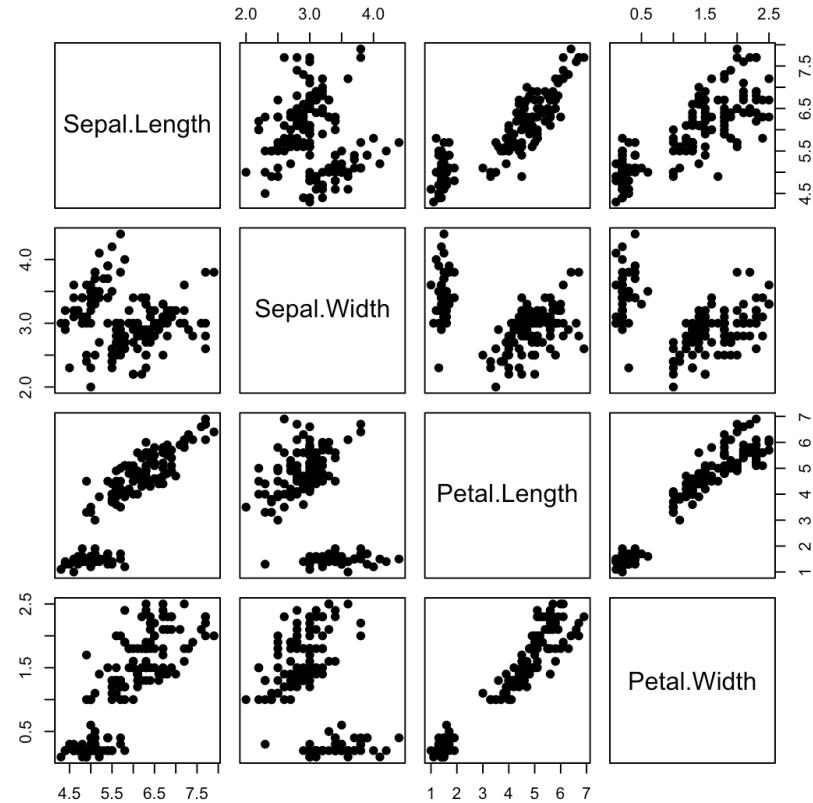
- advanced issues in data preparation
 - data reduction
 - context
 - attribute aggregation
 - feature selection
- ... and modeling
 - dealing with unbalanced classes
- final reflections on data quality

feature selection: eliminate...

- redundant attributes
 - duplicate much or all of the information contained in one or more other attributes
 - e.g. purchase price of a product and the amount of sales tax paid
- irrelevant attributes
 - contain no useful information
 - e.g. students' ID is often irrelevant to the task of predicting students' GPA

feature selection: filter methods

- 2 attributes
 - remove redundant attributes
- 1 attribute vs target
 - identify relevant variables



feature selection: wrapper methods (1/4)

Feature Selection

Full Feature Set



<https://medium.com/@mehulved1503/feature-selection-and-feature-extraction-in-machine-learning-an-overview-57891c595e96>

feature selection: wrapper methods (2/4)

Feature Selection

Full Feature Set



<https://medium.com/@mehulved1503/feature-selection-and-feature-extraction-in-machine-learning-an-overview-57891c595e96>

feature selection: wrapper methods (3/4)

Feature Selection

Full Feature Set



<https://medium.com/@mehulved1503/feature-selection-and-feature-extraction-in-machine-learning-an-overview-57891c595e96>

feature selection: wrapper methods (4/4)

Feature Selection

Full Feature Set



Identify Useful Features



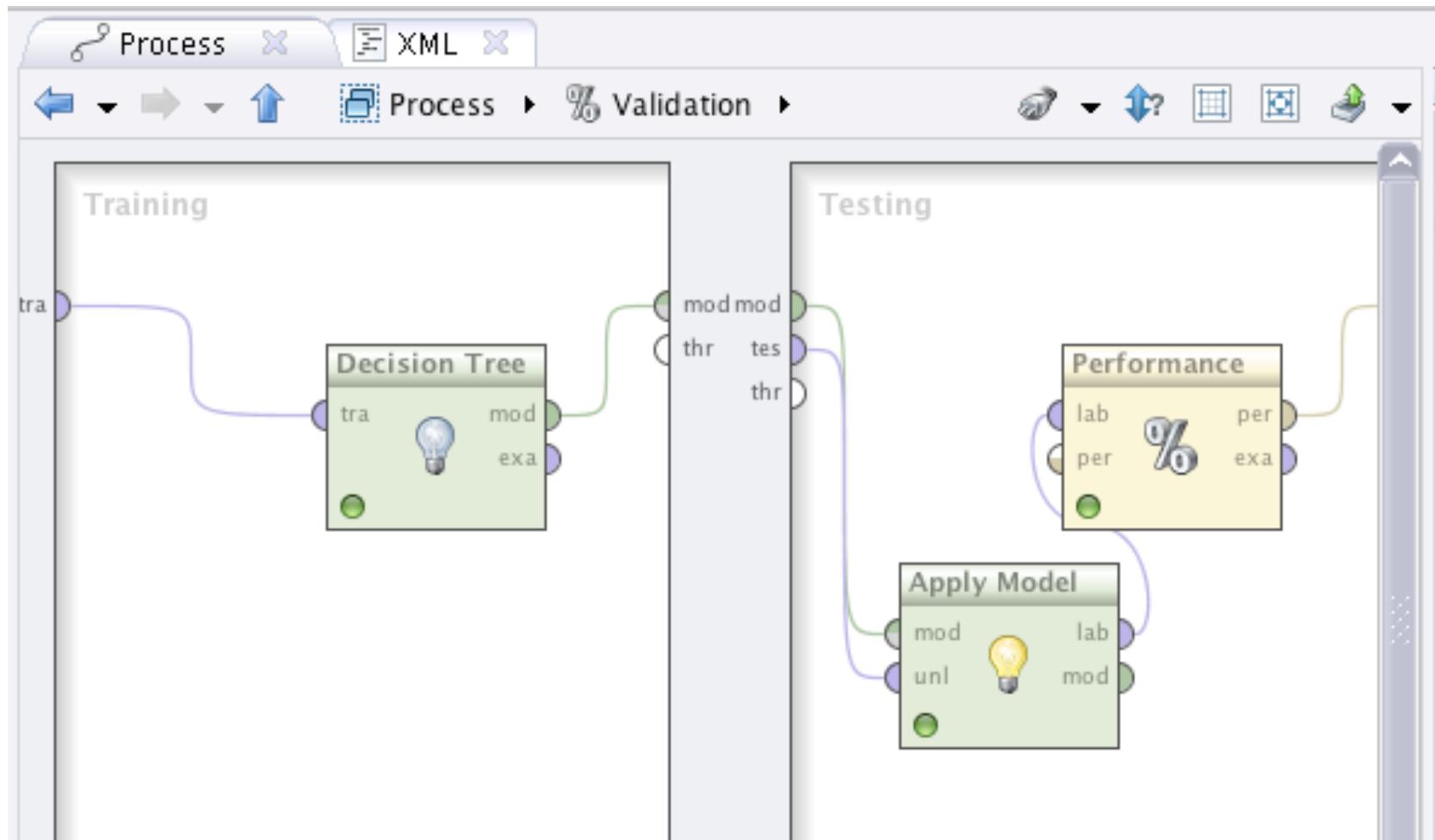
Selected Feature Set



<https://medium.com/@mehulved1503/feature-selection-and-feature-extraction-in-machine-learning-an-overview-57891c595e96>

- advanced issues in data preparation
 - data reduction
- ... and modeling
 - dealing with unbalanced classes
 - context
 - resampling
 - create synthetic data with SMOTE
 - cost sensitive learning
- final reflections on data quality

you developed a model for the competition...



accuracy = 99.9%!!



... wth?!

 InClass Prediction Competition

To loan or not to loan - that is the question

Practical assignment of the Machine Learning course at U.Porto

58 teams · 19 days to go

Overview Data Code Discussion **Leaderboard** Rules Team Host My Submissions **Submit Predictions** ...

[Public Leaderboard](#) [Private Leaderboard](#)

This leaderboard is calculated with approximately 50% of the test data. [Raw Data](#) [Refresh](#)

The final results will be based on the other 50%, so the final standings may be different.

#	Team Name	Notebook	Team Members	Score	Entries	Last
58	G45 - auntdulce		  	0.47160	2	17h

all predictions = negative class!

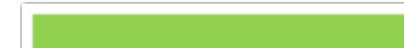


**maybe I should have done the
exploratory data analysis**

(as advised by the lecturers...)



negatives



positives

so what?

		Predicted class	
		Yes	No
Actual class	Yes	TP: True positive	FN: False negative
	No	FP: False positive	TN: True negative

- ML methods usually minimize FP+FN
 - ... or, at the very least give the same weight to both types of errors
- ... but potentially FP >> FN
 - i.e. quality of the model more affected by FP
- ... so algorithm effectively minimizes FP!
- ... and there's an easy model for that
 - prediction = N

- advanced issues in data preparation
 - data reduction
- ... and modeling
 - dealing with unbalanced classes
 - context
 - resampling
 - create synthetic data with SMOTE
 - cost sensitive learning
- final reflections on data quality

learning with class imbalance

- collect more data
 - difficult in many domains
- resample existing data
 - delete data from the majority class
 - duplicate data from the minority class
- create synthetic data
 - e.g. SMOTE
- adapt your learning algorithm
 - e.g. cost sensitive learning

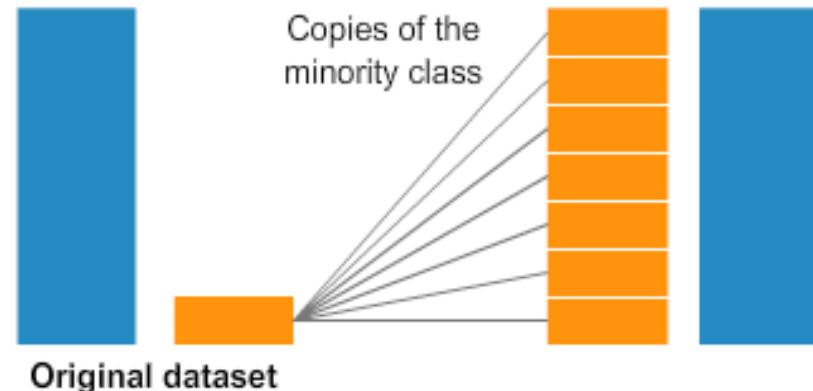
resampling

Undersampling



<https://www.kaggle.com/rafjaa/resampling-strategies-for-imbalanced-datasets#t1>

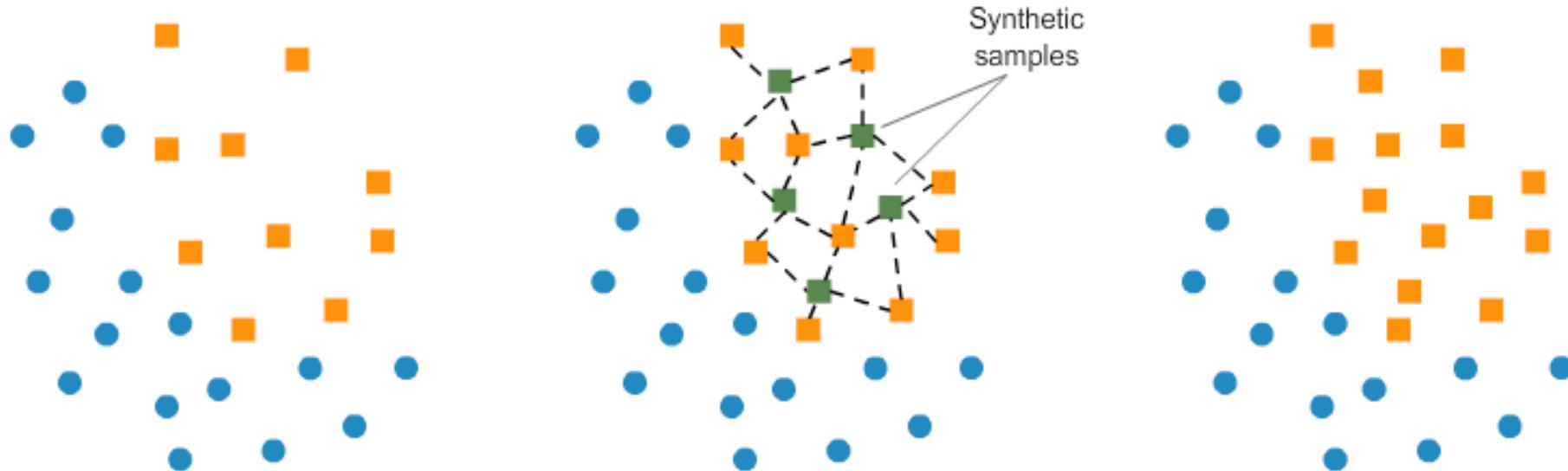
Oversampling



- ... but possible loss of information
- ... but fixed boundaries
- ... and danger of overfitting

furthermore, what is the best ratio?

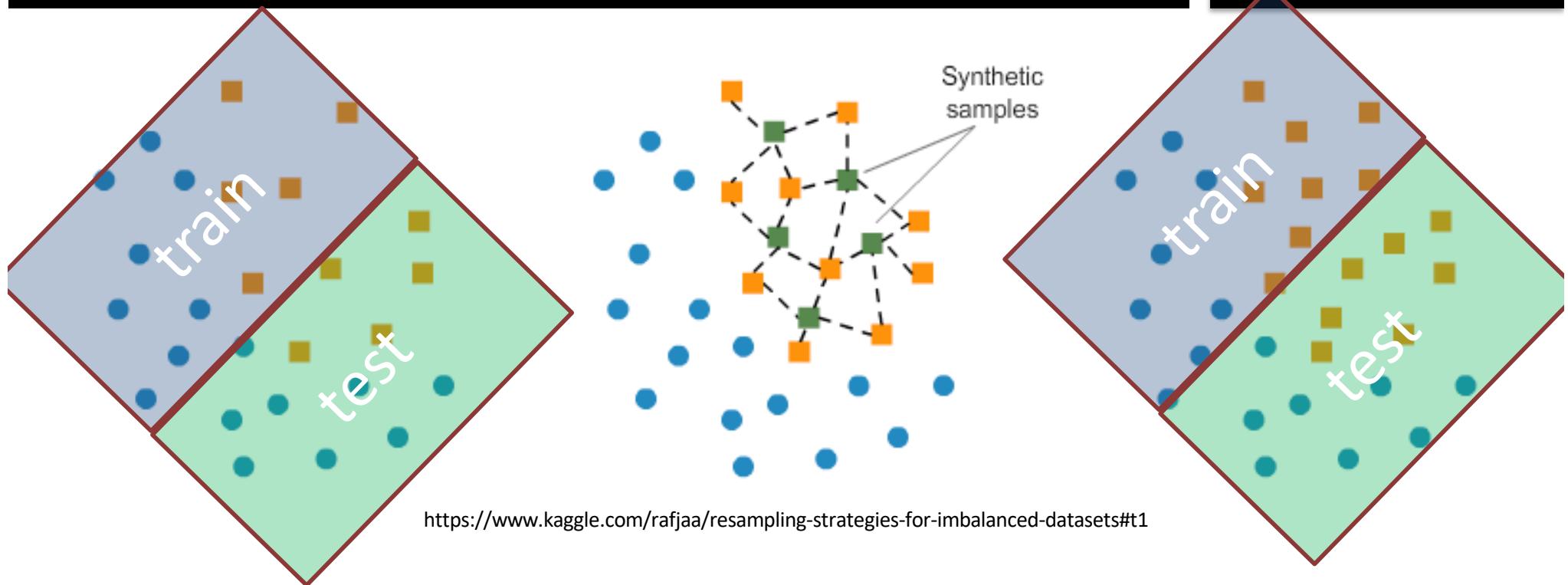
SMOTE (Synthetic Minority Over-sampling Technique)



<https://www.kaggle.com/rafaa/resampling-strategies-for-imbalanced-datasets#t1>

- ... but possibility of inadequate boundaries
 - what happens if minority observations are too far apart?
- ... and danger of overfitting

SMOTE + lack of basic statistical knowledge?



which one?
(or both?)

- advanced issues in data preparation
 - data reduction
- ... and modeling
 - dealing with unbalanced classes
 - context
 - resampling
 - create synthetic data with SMOTE
 - cost sensitive learning
- final reflections on data quality

the cost of errors

		Predicted class	
		Yes	No
Actual class	Yes	TP: True positive	FN: False negative
	No	FP: False positive	TN: True negative

- FP and FN errors often incur different costs
 - medical diagnostic
 - loan decisions
 - marketing campaigns
 - fraud detection in bank transactions
 - fault detection in machines
- ... but ML methods still usually minimize FP+FN

error vs missclassification costs: medical diagnosis example

unnecessary suffering + more expensive
procedures + eventually death
e.g. 100

		Predicted class	
		Yes	No
Actual class	Yes	10	5
	No	5	80

unnecessary exams & anxiety
e.g. 5

- error = $10 / 100 = 10\%$
- missclassification costs
 $= 5 \times 100 + 5 \times 5 = 525$
- ... per patient
– 5.25

cost-sensitive learning

- simple methods
 - resampling according to costs
 - weighting according to costs
 - basically, the same thing
- complex methods
 - e.g. metacost

1. create bootstrap replicates of training data
2. learn model from each replicate
3. relabel examples

$$\operatorname{argmin}_i \sum_j P(j|x)C(i,j)$$

- $C(i | j)$ = cost of mistaking j by i
 - $P(j | x)$ = class probability of x by voting
4. learn model on relabelled data
-
- independent of algorithm

- advanced issues in data preparation
 - data reduction
- ... and modeling
 - dealing with unbalanced classes
- final reflections on data quality

coming next...

- data
- cleaning
 - data quality: a data scientist's worst nightmare
 - quality issues
 - can we do better?
- integration
- reduction
- transformation and discretization
- challenges

data quality: multidimensional view

- accuracy
 - correct or wrong, accurate or not
- completeness
 - not recorded, unavailable, ...
- consistency
 - some modified but some not, dangling, ...
- timeliness
 - timely update?
- believability
 - how trustable is the data and its sources?
- interpretability
 - how easily the data can be understood?

a data scientist's worst nightmare

no worries: our data is clean!



no worries: our data is clean!

(1/6)

we have a data warehouse

- DWs are built with a different purpose
 - descriptive
 - aggregated data
 - typically



<https://pixabay.com/en/archive-bookcase-boxes-business-1850170/>

no worries: our data is clean! (2/6)

our IS was just revamped

- how far was analytics taken into account in the process?



<https://pixabay.com/en/confused-muddled-illogical-880735/>

no worries: our data is clean! (3/6)

we had a major data cleanup

- how was project success measured?



<https://pixabay.com/en/clean-dare-cleaning-1706439/>

no worries: our data is clean! (4/6)

our data is collected automatically

- ... and we all know that machines never break

A problem has been detected and windows has been shut down to prevent damage to your computer.

The problem seems to be caused by the following file: SPCMDCON.SYS

PAGE_FAULT_IN_NONPAGED_AREA

If this is the first time you've seen this Stop error screen, restart your computer. If this screen appears again, follow these steps:

Check to make sure any new hardware or software is properly installed. If this is a new installation, ask your hardware or software manufacturer for any Windows updates you might need.

If problems continue, disable or remove any newly installed hardware or software. Disable BIOS memory options such as caching or shadowing. If you need to use Safe Mode to remove or disable components, restart your computer, press F8 to select Advanced Startup Options, and then select Safe Mode.

Technical information:

*** STOP: 0x00000050 (0xFD3094C2,0x00000001,0xFBFE7617,0x00000000)

*** SPCMDCON.SYS - Address FBFE7617 base at FBFE5000, Datestamp 3d6dd67c

https://en.wikipedia.org/wiki/Blue_Screen_of_Death#/media/File:Windows_XP_BSOD.png

no worries: our data is clean! (5/6)

our data collection is human-error proof

- “Data errors, uh, find a way”



<http://knowyourmeme.com/memes/life-uh-finds-a-way>

no worries: our data is clean! (6/6)

tell us what you need: we have everything

- should be read as

“if something goes wrong,
it’s your fault”
- ... often associated with

“you do magic, right?”



[https://en.wikipedia.org/wiki/Marvin_\(character\)#/media/File:Marvin_\(HHGG\).jpg](https://en.wikipedia.org/wiki/Marvin_(character)#/media/File:Marvin_(HHGG).jpg)

human resources

- remember when IT director was not a C-level job?



Fair use, <https://en.wikipedia.org/w/index.php?curid=24782741>

can we do any better? (2/3)

analytics at the core of IS development

- requirements analysis includes analytics
- analytics components built in the same process as the rest of the functionalities
- we don't mind that systems still do old-fashioned tasks
 - sell stuff, pay salaries, etc.

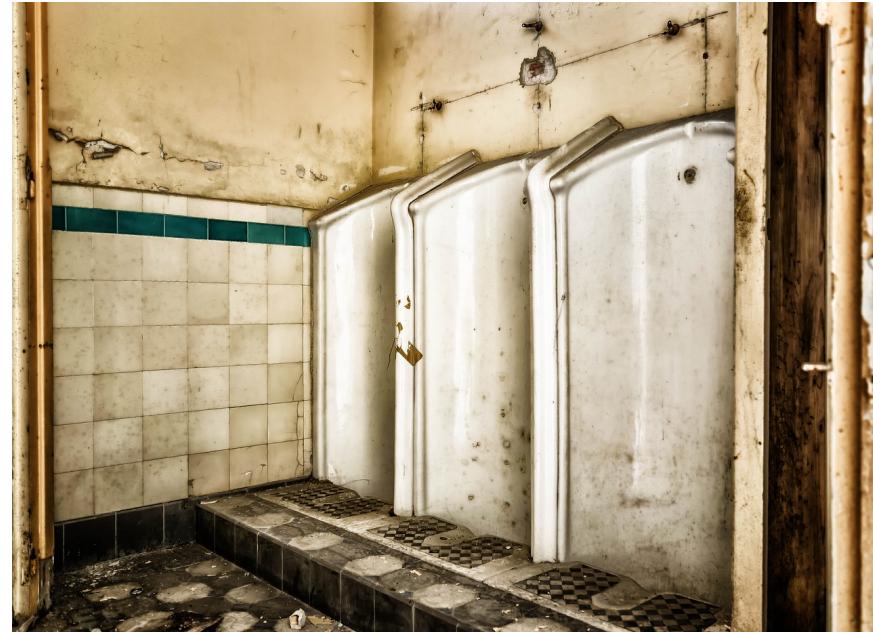


<https://pixabay.com/en/scaffolding-workers-construction-1617969/>

can we do any better? (3/3)

data quality is a continuous process

- data steward
- ... the sexiest job of the XXII century?



<https://pixabay.com/en/lost-places-toilet-urinal-pforphoto-1610652/>

data cleaning as a process

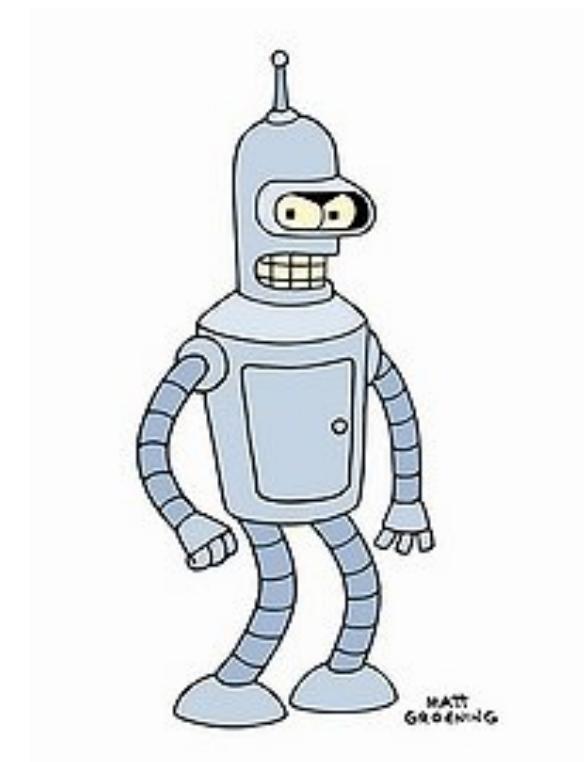
- discrepancy detection
 - validate with metadata (e.g., domain, range, dependency, distribution)
 - check field overloading
 - check uniqueness rule, consecutive rule and null rule
 - commercial tools
 - scrubbing: use simple domain knowledge to detect errors and make corrections
 - e.g. postal code, spell-check
 - auditing: discover rules and relationship to detect violators
 - e.g. correlation and clustering to find outliers
- migration and integration
 - data migration tools
 - allow transformations to be specified
 - ETL (Extraction/Transformation>Loading) tools
 - allow users to specify transformations through a graphical user interface

Summary

- data
- ... quality
 - accuracy, completeness, consistency, timeliness, believability, interpretability
- ... cleaning
 - e.g. missing/noisy values, outliers
- integration from multiple sources
 - entity identification problem is challenging
- reduction
 - curse of dimensionality and dimensionality reduction
 - numerosity reduction
- transformation and discretization

automation

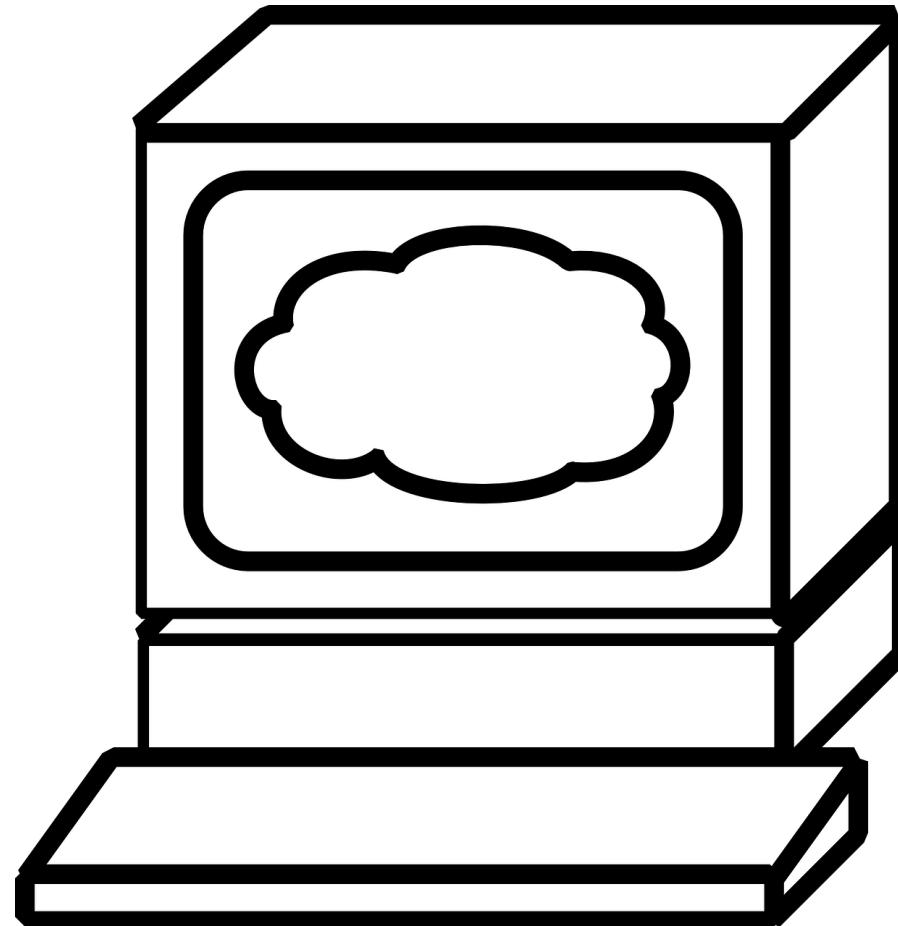
- automl & metalearning
 - some progress on algorithm selection
 - ... early work on workflow
 - including data preparation
 - ... not really data cleaning



[https://pt.wikipedia.org/wiki/Ficheiro:Bender_\(personagem\).jpg](https://pt.wikipedia.org/wiki/Ficheiro:Bender_(personagem).jpg)

DQaaS?

- if automation is possible
 - DQ can become a commodity?
- perhaps there is hope?
- ... many issues
 - confidentiality
 - computational costs



<https://pixabay.com/en/computer-server-internet-network-294036/>