

Architecture of ML Systems

09 Data Acquisition and Preparation

Matthias Boehm

Graz University of Technology, Austria
Computer Science and Biomedical Engineering
Institute of Interactive Systems and Data Science
BMVIT endowed chair for Data Management

Last update: May 29, 2020

Announcements/Org

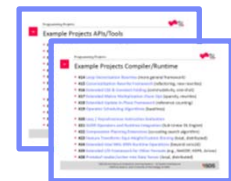
■ #1 Video Recording

- Link in **TeachCenter** & **TUbe** (lectures will be public)
- **Live streaming through TUbe**, starting May 08
- Questions: <https://tugraz.webex.com/meet/m.boehm>



■ #2 AMLS Programming Projects

- **Status:** all project discussions w/ **15 students** (~**8 PRs**)
- Awesome mix of projects (algorithms, compiler, runtime)
- Soft deadline: **June 30**



■ #3 TU Delft DESOSO 2020

- Delft Students on Software Architecture (incl ML systems)

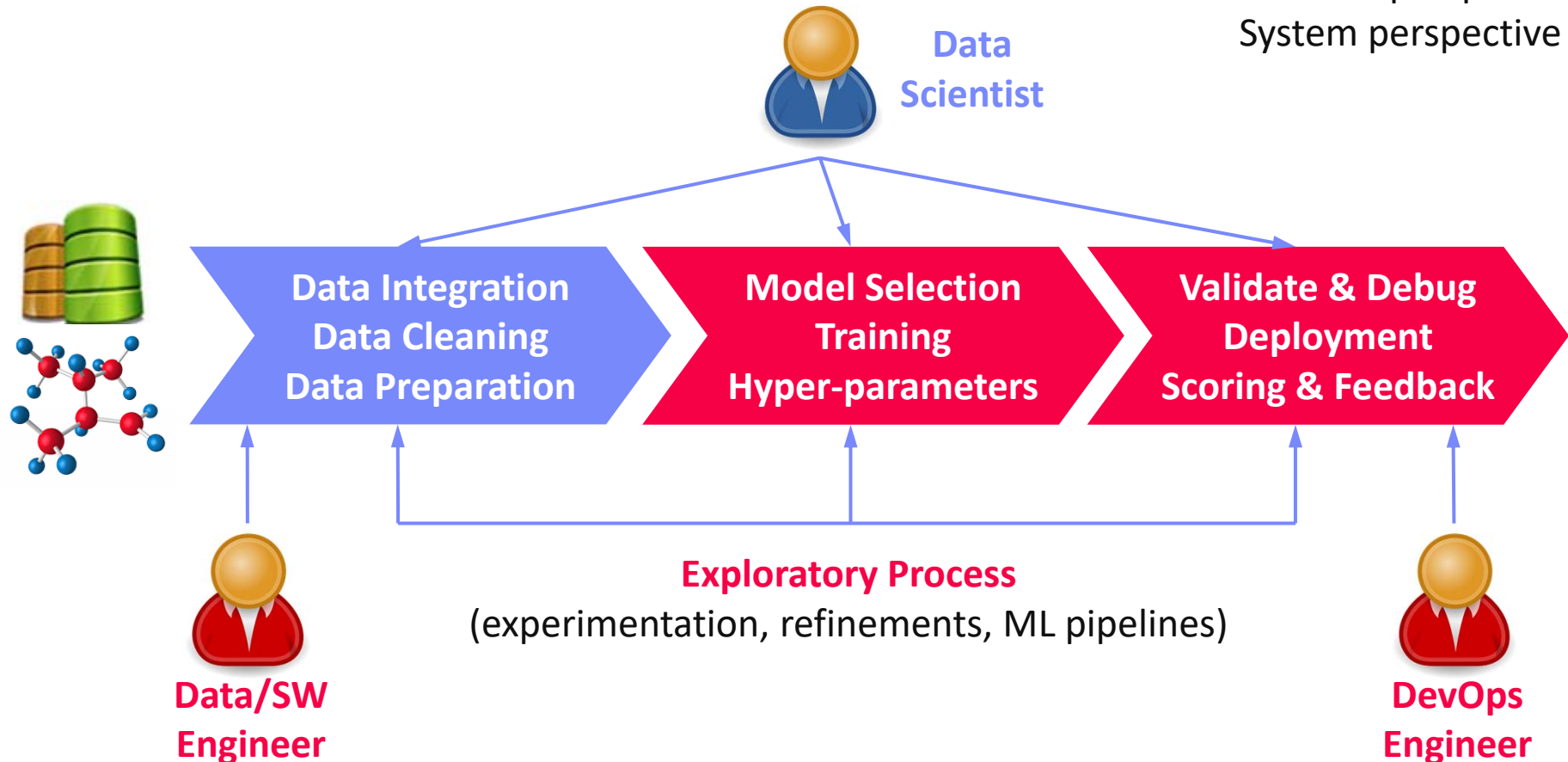
<https://desosa.nl>



Recap: The Data Science Lifecycle

Data-centric View:

Application perspective
Workload perspective
System perspective



The 80% Argument

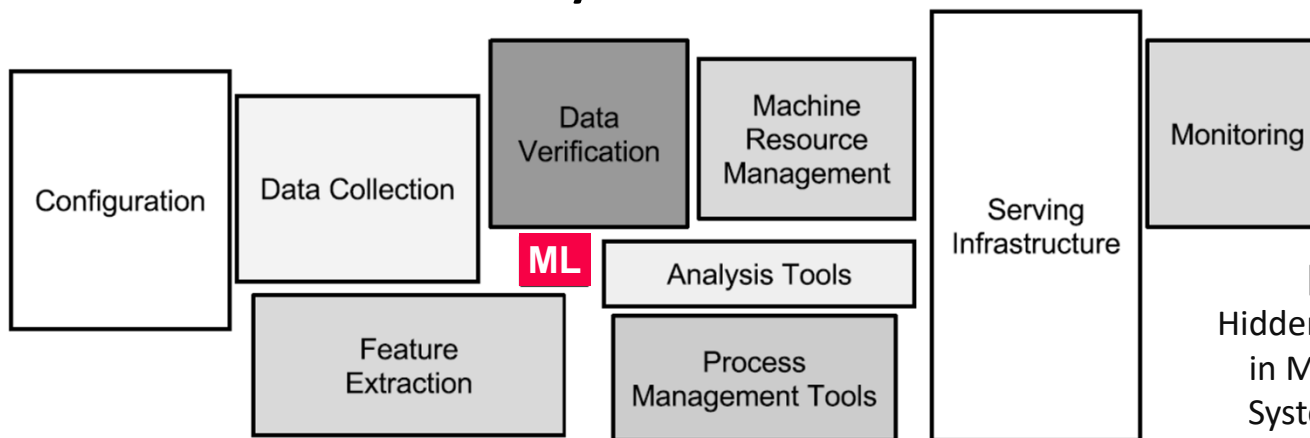
■ Data Sourcing Effort

- Data scientists spend **80-90% time** on finding, integrating, cleaning datasets

[Michael Stonebraker, Ihab F. Ilyas:
Data Integration: The Current
Status and the Way Forward.
IEEE Data Eng. Bull. 41(2) (2018)]



■ Technical Debts in ML Systems



[D. Sculley et al.:
Hidden Technical Debt
in Machine Learning
Systems. NIPS 2015]



- Glue code, pipeline jungles, dead code paths
- Plain-old-data types (arrays), multiple languages, prototypes
- Abstraction and configuration debts
- Data testing, reproducibility, process management, and cultural debts

Agenda

- Data Acquisition and Integration
- Data Preparation and Feature Engineering
- Data Transformation and Cleaning
- Data Augmentation (next week)



“least enjoyable
tasks in data
science lifecycle”

Data Acquisition and Integration

Data Integration for ML and
ML for Data Integration



**Data Integration and
Large-Scale Analysis (DIA)**
(bachelor/master)

Data Sources and Heterogeneity

■ Terminology

- **Integration** (Latin integer = whole): consolidation of data objects / sources
- **Homogeneity** (Greek homo/homoios = same): similarity
- **Heterogeneity**: dissimilarity, different representation / meaning

■ Heterogeneous IT Infrastructure

- Common enterprise IT infrastructure contains >100s of **heterogeneous and distributed systems and applications**
- E.g., health care data management: 20 - 120 systems



■ Multi-Modal Data (example health care)

- **Structured patient data**, patient records incl. prescribed drugs
- **Knowledge base** drug APIs (active pharmaceutical ingredients) + interactions
- **Doctor notes** (text), diagnostic codes, outcomes
- **Radiology images** (e.g., MRI scans), **patient videos**
- **Time series** (e.g., EEG, ECoG, heart rate, blood pressure)

Types of Data Formats

■ General-Purpose Formats

- **CSV** (comma separated values), **JSON** (javascript object notation), **XML**, **Protobuf**
- CLI/API access to DBs, KV-stores, doc-stores, time series DBs, etc

■ Sparse Matrix Formats

- **Matrix market**: text IJV (row, col, value)
- **Libsvm**: text compressed sparse rows
- Scientific formats: **NetCDF**, **HDF5**

```
%%MatrixMarket matrix coordinate real general
% -----
% 0 or more comment lines
% -----
5 5 8
1 1 1.000e+00
2 2 1.050e+01
3 3 1.500e-02
1 4 6.000e+00
4 2 2.505e+02
4 4 -2.800e+02
4 5 3.332e+01
5 5 1.200e+01
```

■ Large-Scale Data Formats

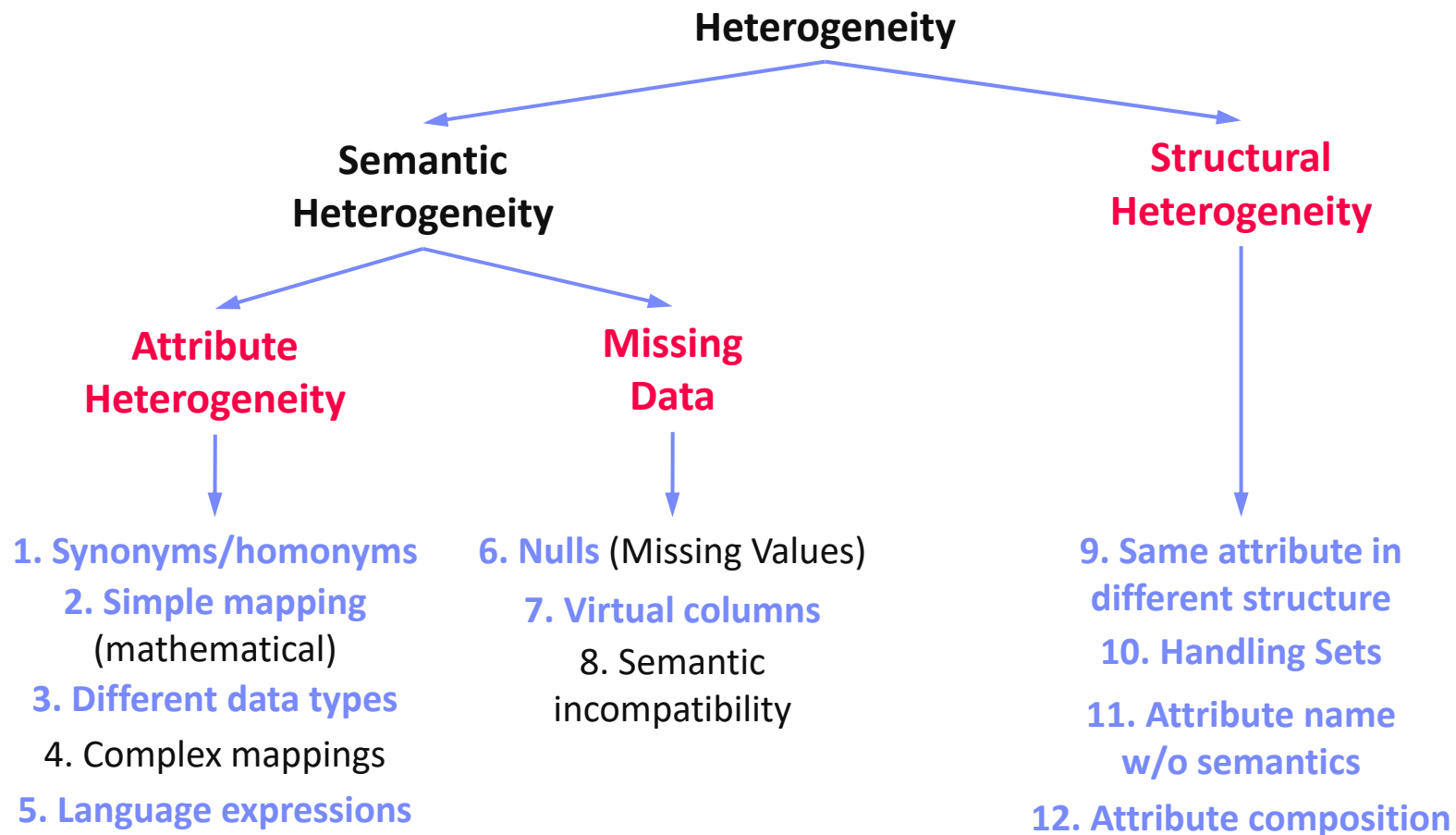
- **Parquet** (columnar file format)
- **Arrow** (cross-platform columnar in-memory data)

■ Domain-Specific Formats

- Health care: **DICOM** images, **HL7** messages (health-level seven, XML)
- Automotive: **MDF** (measurements), **CDF** (calibrations), **ADF** (auto-lead XML)
- Smart production: **OPC** (open platform communications)

Types of Heterogeneity

[J. Hammer, M. Stonebraker, and O. Topsakal:
THALIA: Test Harness for the Assessment of
Legacy Information Integration Approaches. U
Florida, TR05-001, 2005]

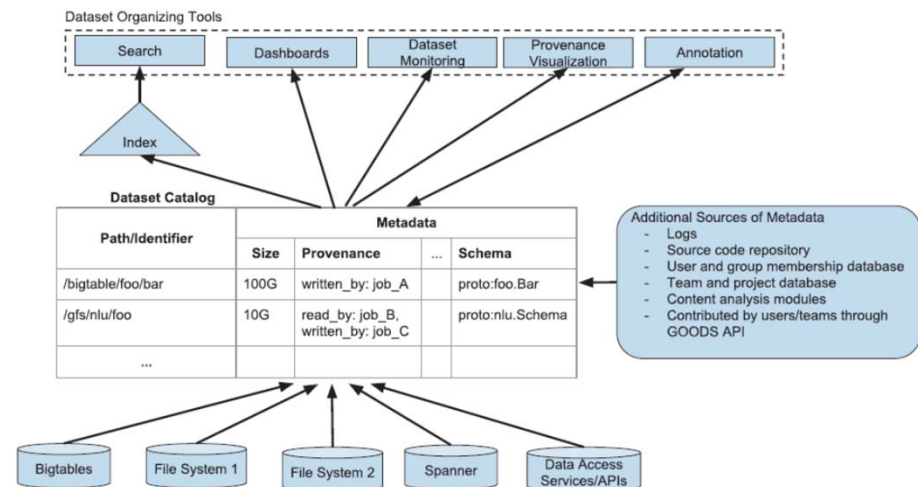


- Data curation in repositories for finding relevant datasets in **data lakes**
- Augment data with open and linked data sources

[Alon Y. Halevy et al: Goods: Organizing Google's Datasets. **SIGMOD 2016**]



Google Data Search



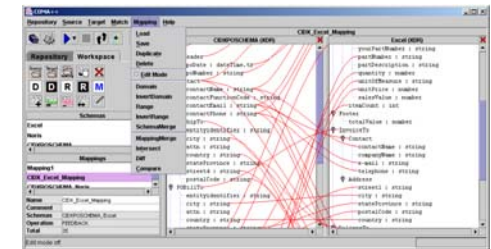
■ Schema Detection

- Sample of the input dataset → infer **syntactic schema** (e.g., data types)
- **Semantic schema** detection (e.g., location, date, rank, name)

- **Schema Matching**

- Semi-automatic mapping of schema S1 to schema S2
- **Approaches:** Schema- vs instance-based; element- vs structure-based; linguistic vs rules
- Hybrid and composite matchers
- Global schema matching (one-to-one): stable marriage problem

[Credit: Erhard Rahm]



- **Schema Mapping**

- Given two schemas and correspondences, generate transformation program
- **Challenges:** complex mappings (1:N cardinality), new values, PK-FK relations and nesting, creation of duplicates, different data types, semantic preserving

Corrupted Data

■ Heterogeneity of Data Sources

- Update anomalies on denormalized data / eventual consistency
- Changes of app/preprocessing over time (US vs us) → inconsistencies

■ Human Error

- Errors in semi-manual data collection, laziness (see default values), bias
- Errors in data labeling (especially if large-scale: crowd workers / users)

■ Measurement/Processing Errors

- Unreliable HW/SW and measurement equipment (e.g., batteries)
- Harsh environments (temperature, movement) → aging

Uniqueness & duplicates

Contradictions & wrong values

Missing Values

Ref. Integrity

[Credit: Felix Naumann]

ID	Name	BDay	Age	Sex	Phone	Zip	Zip	City
3	Smith, Jane	05/06/1975	44	F	999-9999	98120	98120	San Jose
3	John Smith	38/12/1963	55	M	867-4511	11111	90001	Lost Angeles
7	Jane Smith	05/06/1975	24	F	567-3211	98120		

Typos

Examples (aka errors are everywhere)

DM SS'19 (Soccer World Cups)

Commits on Apr 21, 2019	Commits on Apr 19, 2019
[MINOR] Fix 2002 match final scores, squad club mboehm7 committed on Apr 21	Fixed squads issues (resolved null clubs, non-unique clubs, player name) mboehm7 committed on Apr 19
[MINOR] Fixed mapping hansa rostock, and cons mboehm7 committed on Apr 21	Commits on Apr 18, 2019
[MINOR] Fix null in match type (due to input file) mboehm7 committed on Apr 21	[MINOR] Fix squad club-country mapping, unique player names mboehm7 committed on Apr 18
	[MINOR] Fix squad club-country mapping, and spurious spaces mboehm7 committed on Apr 18

DM WS'19/20 (Airports and Airlines)

Commits on Oct 7, 2019	Commits on Oct 30, 2019
New airports and flights datasets (cleaned) ... OlgaOvcharenko authored and mboehm7 committed	Fix data issues: redundant plane types in routes mboehm7 committed 14 days ago
	Fix data issues: referential integrity country names mboehm7 committed 14 days ago
	Fix data issue: spelling united kingdom mboehm7 committed 14 days ago
	<div>- US,DFW,LIT,ER4;M83;M83</div> <div>+ US,DFW,LIT,ER4;M83</div> <div>- Oyo Ollombo Airport,Oyo,Congo (Brazzaville),O</div> <div>- Beni Airport,Beni,Congo (Kinshasa),BNC,FZNP,0.575,?</div> <div>+ Beni Airport,Beni,Democratic Republic of Congo,BNC,</div> <div>- RAF St Athan,4Q,STN,United Kingdom,N</div> <div>+ RAF St Athan,4Q,STN,United Kingdom,N</div>

DM SS'20 (DBLP Publications)

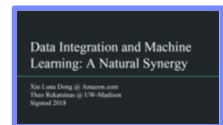
Commits on Mar 13, 2020	Commits on Mar 14, 2020	Commits on Apr 6, 2020	Commits on Apr 5, 2020
Fix conf.csv header meta data (inconsistent number of ...) mboehm7 committed on Mar 14	Extract and clean city/country f mboehm7 committed on Mar 14	Updated dblp publications rea mboehm7 committed on Apr 6	Initial deduplication of person affiliations and thesis schools mboehm7 committed on Apr 5
Fix csv quoting (escaped quotes within fields) mboehm7 committed on Mar 14	Fix various columns by expecte mboehm7 committed on Mar 14	Revert too aggressive matchin mboehm7 committed on Apr 6	Additional country cleaning (for person affiliations) mboehm7 committed on Apr 5
Fix publication titles (punctuation) and csv delimiters mboehm7 committed on Mar 14	Fix person/theses affiliation coi mboehm7 committed on Mar 14	Additional cleaning of instituti mboehm7 committed on Apr 6	Fix country name consistency (UK, Tunisia, The Netherlands, Autralia) mboehm7 committed on Apr 5
Updated dblp publications datasets (DB pubs only, clea mboehm7 committed on Mar 13	Fix conference title normalizati mboehm7 committed on Mar 14	Fix conference venues (consisti mboehm7 committed on Apr 6	Simplify dataset encoding (no quoting, no escaped quoates, etc) mboehm7 committed on Apr 5
	Fix normalization of conference mboehm7 committed on Mar 14	Fix incorrect year in journal vol mboehm7 committed on Apr 6	Fix head Commits on Apr 22, 2020
	Fix affiliation countries via robu mboehm7 committed on Mar 14	Fix handling of special characters beyon mboehm7 committed on Apr 6	Fix special character in french thesis mboehm7 committed on Apr 22

Data Integration for ML and ML for DI

■ #1 Data Extraction

- Extracting structured data from un/semi-structured data
- Rule- and ML-based extractors, combination w/ CNN

[Xin Luna Dong, Theodoros Rekatsinas:
Data Integration and Machine Learning:
A Natural Synergy. **SIGMOD 2018**]



■ #2 Schema Alignment

- Schema matching for consolidating data from heterogeneous systems
- Spatial and Temporal alignment via provenance and query processing (e.g., sensor readings for object along a production pipeline)

■ #3 Entity Linking

- Linking records to entities (deduplication)
- Blocking, pairwise matching, clustering, ML, Deep ML (via entity embedding)

■ #4 Data Fusion

- Resolve conflicts, necessary in presence of erroneous data
- Rule- and ML-based, probabilistic GM, Deep ML (RBMs, graph embeddings)

Data Validation

Sanity checks on **expected** shape
before training first model

[Neoklis Polyzotis, Sudip Roy, Steven
Euijong Whang, Martin Zinkevich: Data
Management Challenges in Production
Machine Learning. Tutorial, **SIGMOD 2017**]



(**Google
Research**)

- **Check a feature's min, max, and most common value**
 - Ex: Latitude values must be within the range $[-90, 90]$ or $[-\pi/2, \pi/2]$
- **The histograms of continuous or categorical values are as expected**
 - Ex: There are similar numbers of positive and negative labels
- **Whether a feature is present in enough examples**
 - Ex: Country code must be in at least 70% of the examples
- **Whether a feature has the right number of values (i.e., cardinality)**
 - Ex: There cannot be more than one age of a person

Data Validation, cont.

[Sebastian Schelter, Dustin Lange, Philipp Schmidt, Meltem Celikel, Felix Bießmann, Andreas Grafberger: Automating Large-Scale Data Quality Verification. **PVLDB 2018**]



Constraints and Metrics for quality check UDFs

constraint	arguments
dimension <i>completeness</i>	
isComplete	column
hasCompleteness	column, udf
dimension <i>consistency</i>	
isUnique	column
hasUniqueness	column, udf
hasDistinctness	column, udf
isInRange	column, value range
hasConsistentType	column
isNonNegative	column
isLessThan	column pair
satisfies	predicate
satisfiesIf	predicate pair
hasPredictability	column, column(s), udf
statistics (can be used to verify dimension <i>consistency</i>)	
hasSize	udf
hasTypeConsistency	column, udf
hasCountDistinct	column
hasApproxCountDistinct	column, udf
hasMin	column, udf
hasMax	column, udf
hasMean	column, udf
hasStandardDeviation	column, udf
hasApproxQuantile	column, quantile, udf
hasEntropy	column, udf
hasMutualInformation	column pair, udf
hasHistogramValues	column, udf
hasCorrelation	column pair, udf
time	
hasNoAnomalies	metric, detector

metric
dimension <i>completeness</i>
Completeness
dimension <i>consistency</i>
Size
Compliance
Uniqueness
Distinctness
ValueRange
DataType
Predictability
statistics (can be used to
Minimum
Maximum
Mean
StandardDeviation
CountDistinct
ApproxCountDistinct
ApproxQuantile
Correlation
Entropy
Histogram
MutualInformation

(Amazon Research)

Organizational Lesson:
benefit of shared vocabulary/procedures

Technical Lesson:
fast/scalable; reduce manual and ad-hoc analysis

Approach

- #1 Quality checks on basic metrics, computed in **Apache Spark**
- #2 **Incremental maintenance** of metrics and quality checks

Data Preparation and Feature Engineering

Overview Feature Engineering

■ Terminology

- Matrix X of m observations (rows) and n features (columns)
- **Continuous features:** numerical values (aka scale features)
- **Categorical features:** non-numerical values, represent groups
- **Ordinal features:** non-numerical values, associated ranking
- Feature space: multi-dimensional space of features → curse of dimensionality

■ Feature Engineering

- Bring multi-modal data and features into numeric representation
- Use domain expertise to expose predictive features to ML model training

■ Excursus: Representation Learning

- Neural networks can be viewed as combined representation learning and model training (pros and cons: learned, repeatable)
- Mostly homogeneous inputs (e.g., image), research on multi-modal learning

➔ **Principle:** If same accuracy, prefer simple model (cheap, robust, explainable)

Recoding

■ Summary

- Numerical encoding of categorical features (arbitrary strings)
- Map distinct values to integer domain (potentially combined w/ one-hot)

City	State
San Jose	CA
New York	NY
San Francisco	CA
Seattle	WA
New York	NY
Boston	MA
San Francisco	CA
Los Angeles	CA
Seattle	WA



Dictionaries

{San Jose : 1,
New York : 2,
San Francisco : 3,
Seattle : 4,
Boston : 5,
Los Angeles : 6}

{CA : 1,
NY : 2,
WA : 3,
MA : 4}

City	State
1	1
2	2
3	1
4	3
2	2
5	4
3	1
6	1
4	3

Feature Hashing

Summary

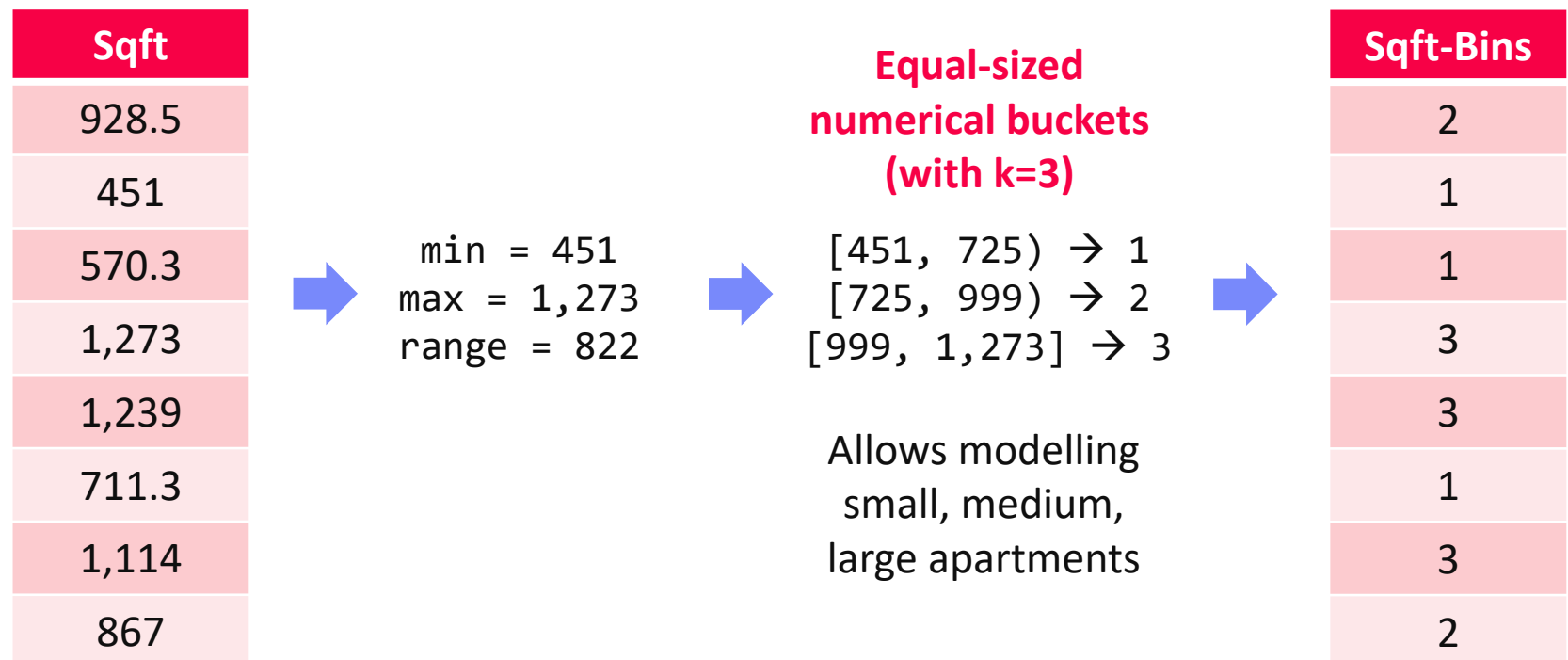
- Numerical encoding of categorical features (arbitrary strings)
- Hash input to k buckets via $\text{hash}(\text{value}) \% k$ (often combined w/ one-hot)

City			City
San Jose	for $k = 5$: Efficient, but collisions	1993955031 % 5 → 1	1
New York		1382994575 % 5 → 0	0
San Francisco		1540367136 % 5 → 1	1
Seattle		-661909336 % 5 → 1	1
New York		1993955031 % 5 → 1	1
Boston		1995575789 % 5 → 4	4
San Francisco		1540367136 % 5 → 1	1
Los Angeles		-425347233 % 5 → 3	3
Seattle		-661909336 % 5 → 1	1

Binning (see also Quantization, Binarization)

Summary

- Encode of numerical features to integer domain (often combined w/ one-hot)
- Equi-width:** split (max-min)-range into k equal-sized buckets
- Equi-height:** compute data-driven ranges for k balanced buckets



One-hot Encoding

■ Summary

- Encode integer feature of cardinality d into sparse 0/1 vector of length d
- Feature vectors of input features concatenated in sequence

City	State		C1	C2	C3	C4	C5	C6	S1	S2	S3	S4
1	1		1	0	0	0	0	0	1	0	0	0
2	2		0	1	0	0	0	0	0	1	0	0
3	1		0	0	1	0	0	0	1	0	0	0
4	3		0	0	0	1	0	0	0	0	1	0
2	2	→	0	1	0	0	0	0	0	1	0	0
5	4		0	0	0	0	1	0	0	0	0	1
3	1		0	0	1	0	0	0	1	0	0	0
6	1		0	0	0	0	0	1	1	0	0	0
4	3		0	0	0	1	0	0	0	0	1	0

Derived Features

■ Intercept Computation

- Add a column of ones to X for computing the intercept as a weight
- Applies to regression and classification

```
X = cbind(X,  
          matrix(1, nrow(X), 1));
```

■ Non-Linear Relationships

- Can be explicitly materialized as feature combinations
- Example: Assumptions of underlying physical system
- Arbitrary complex feature interactions: e.g., $X_1^2 * X_2$

```
// y ~ b1*X1 + b2*X1^2  
X = cbind(X, X^2);
```

NLP Features

Basic NLP Feature Extraction

- **Sentence/word tokenization:** split into sentences/words (e.g., via stop words)
- **Part of Speech (PoS) tagging:** label words verb, noun, adjectives (syntactic)
- **Semantic role labeling:** label entities with their roles in actions (semantic)

Who did **what** to **whom** at **where**?

Bag of Words (BOW) and N-Grams

- Represent sentences as **bag** (multisets)

A B C A B E.
A D E D E D.



A	B	C	D	E
2	2	1	0	1
1	0	0	3	2

- **Bi-grams:** bag-of-words for 2-sequences of words (order preserving)
- **N-grams:** generalization of bi-grams to arbitrary-length sequences

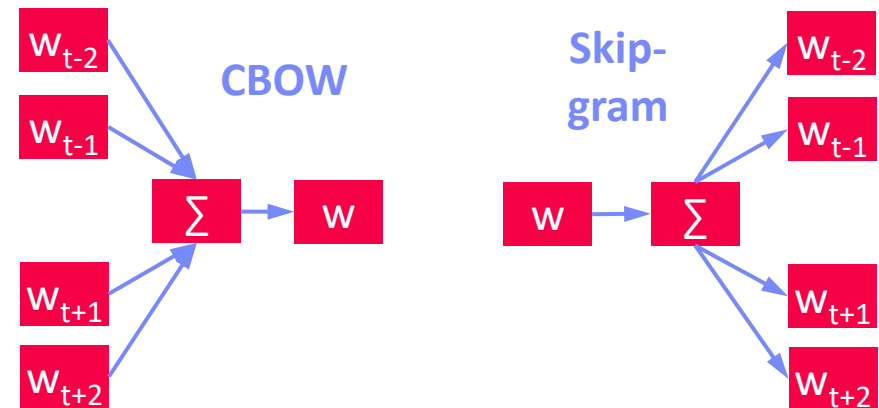
NLP Features, cont.

[Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean:
Efficient Estimation of Word Representations in Vector
github.com/dav/word2vec Space. ICLR (Workshop) 2013]



Word Embeddings

- Trained (word \rightarrow vector) mappings (~ 50 -300 dims)
- Word2vec**: continuous bag-of-words (CBOW) or continuous skip-gram
- Subsampling frequent words
- Semantic preserving arithmetic operations**
(+ \sim * of context distributions)



$$\text{vec}(\text{Paris}) \approx \text{vec}(\text{Berlin}) - \text{vec}(\text{Germany}) + \text{vec}(\text{France})$$

Follow-up Work

- Often pre-trained word embeddings; fine-tuning if necessary for task/domain
- Various extensions/advancements: **Sentence2Vec**, **Doc2Vec**, **Node2Vec**
- BERT**, **RoBERTa**, **ALBERT**, **StructBERT**

[Jacob Devlin et al. : **BERT**: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT (1) 2019]



Example Spark ML



■ API Design

- **Transformers:** Feature transformations and learned models
- **Estimators:** Algorithm that can be fit to produce a transformer
- Compose ML pipelines from chains of transformers and estimators

■ Example Pipeline

```
// define pipeline stages
```

```
tokenizer = Tokenizer(inputCol="text", outputCol="words")
```

```
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(),  
                      outputCol="features")
```

```
lr = LogisticRegression(maxIter=10, regParam=0.001)
```

```
// create pipeline transformer via fit
```

```
pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])
```

```
model = pipeline.fit(training)
```

```
// use of resulting ML pipeline
```

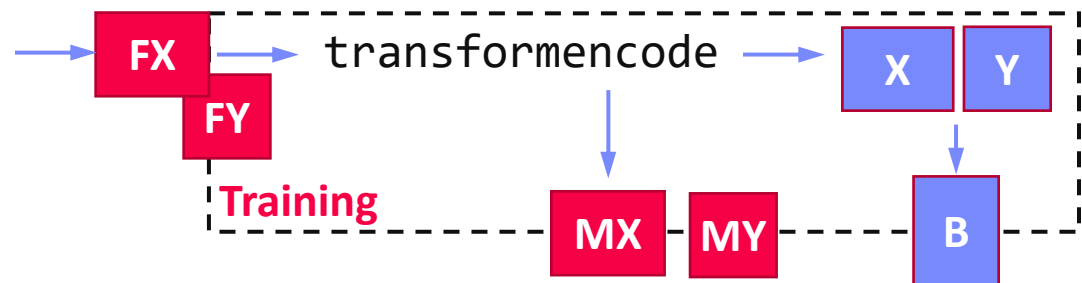
```
prediction = model.transform(test)
```

[<https://spark.apache.org/docs/2.4.3/ml-pipeline.html>]

Example SystemML/SystemDS



Feature Transformation during Training



read tokenized words

```
FX = read("./input/FX", data_type=FRAME); # sentence id, word, count
```

```
FY = read("./input/FY", data_type=FRAME); # sentence id, labels
```

encode and one-hot encoding

```
[X0, MX] = transformencode(target=FX, spec="{recode:[2]}");
```

```
[Y0, MY] = transformencode(target=FY, spec="{recode:[2]}");
```

```
X = table(X0[:,1], X0[:,2], X0[:,3]); # bag of words
```

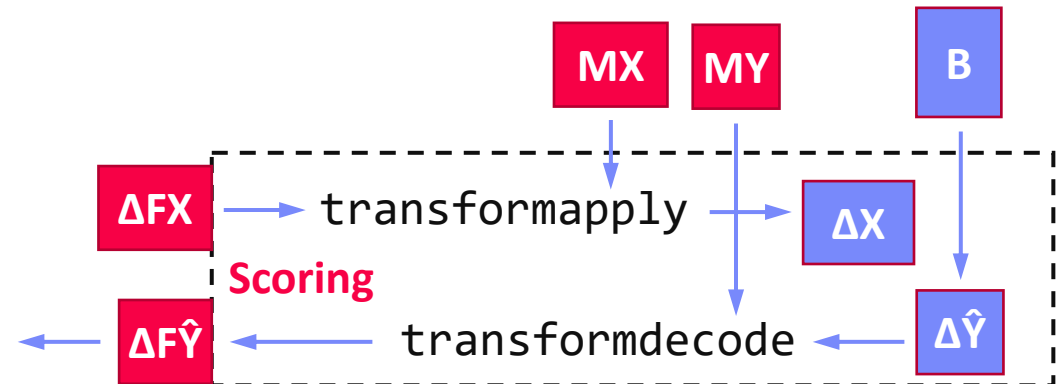
```
Y = table(Y0[:,1], Y0[:,2]); # bag of words
```

model training via multi-label, multi-nominal logical regression

```
B = mlogreg(X, Y);
```

Example SystemML/SystemDS, cont.

■ Feature Transformation during Scoring



```
# read tokenized words of test sentences
```

```
dFX = read("./input/dFX", data_type=FRAME); # sentence id, word, count
```

```
# encode and one-hot encoding
```

```
dX0 = transformapply(target=dFX, spec="{recode:[2]}", meta=MX);
```

```
dX = table(dX0[,1], dX0[,2], dX0[,3]); # bag of words
```

```
# model scoring and postprocessing (reshape, attach sentence ID, etc)
```

```
dYhat = (X %%% B) >= theta; ...;
```

```
# decode output labels: sentence id, label word
```

```
dFYhat = transformdecode(target=dYhat, spec="{recode:[2]}", meta=MY);
```

Data Transformation and Cleaning

Standardization/Normalization

■ #1 Standardization

- Centering and scaling to mean 0 and variance 1

```
X = X - colMeans(X);  
X = X / sqrt(colVars(X));
```

- Ensures well-behaved training

- Densifying operation

```
X = replace(X, pattern=NaN,  
replacement=0); #robustness
```

- Awareness of NaNs

- Batch normalization in DNN: standardization of activations

■ #2 (Min-Max) Normalization

- Rescale values into common range [0,1]

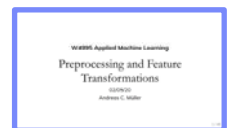
```
X = (X - colMins(X))  
/ (colMaxs(X) - colMins(X));
```

- Avoid bias to large-scale features

- Does not handle outliers

Recommended Reading

[Andreas C. Mueller: Preprocessing and Feature Transformations, **Applied ML Lecture 2020**,
<https://www.youtube.com/watch?v=XpOBSaktb6s>]



Standardization/Normalization, cont.

#3 Deferred Standardization

- Avoid densifying dataset upfront by pushing standardization into inner loop iterations
- Let **matrix-multiplication chain optimization** + rewrites do the rest

[Credit:
Alexandre (Sasha)
V. Evfimievski]

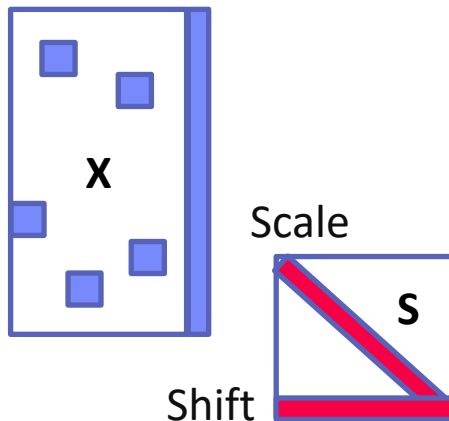


Example GLM/ImCG

operation w/ early standardized X

```
q = t(X) %*% diag(w) %*% X %*% B;
```

Input w/ column of
ones (intercept)



**Substitute X with
X %*% S**

operation w/ deferred standardization

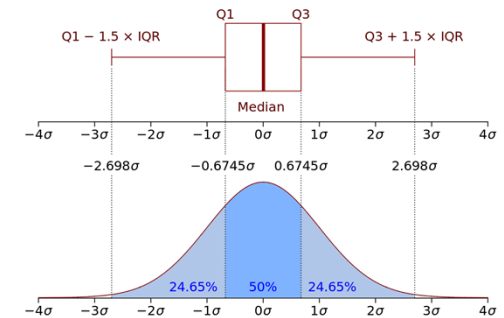
```
q = t(S) %*% t(X) %*% diag(w)
   %*% X %*% S %*% B;
```

```
q = t(S) %*% (t(X) %*% (diag(w)
   %*% (X %*% (S %*% B))));
```

Winsorizing and Trimming

Recap: Quantiles

- Quantile Q_p w/ $p \in (0,1)$ defined as $P[X \leq x] = p$



[Credit: <https://en.wikipedia.org>]

Winsorizing

- Replace** tails of data distribution at user-specified threshold
- Quantiles / std-dev
- ➔ Reduce skew

compute quantiles for lower and upper

```
ql = quantile(X, 0.05);
qu = quantile(X, 0.95);
```

replace values outside [ql,qu] w/ ql and qu

```
Y = ifelse(X < ql, ql, X);
Y = ifelse(Y > qu, qu, Y);
```

SystemDS:
winsorize()
outlier()

Truncation/Trimming

- Remove** tails of data distribution at user-specified threshold

remove values outside [ql,qu]

```
I = X < qu | X > ql;
Y = removeEmpty(X, "rows", select = I);
```

Largest Difference from Mean

determine largest diff from mean

```
I = (colMaxs(X) - colMeans(X))
  > (colMeans(X) - colMins(X));
Y = ifelse(xor(I, op), colMaxs(X), colMins(X));
```

Outliers and Outlier Detection

■ Types of Outliers

- **Point outliers:** single data points far from the data distribution
- **Contextual outliers:** noise or other systematic anomalies in data
- **Sequence (contextual) outliers:** sequence of values w/ abnormal shape/agg
- Univariate vs multivariate analysis
- Beware of underlying assumptions (distributions)

[Varun Chandola, Arindam Banerjee, Vipin Kumar: Anomaly detection: A survey. **ACM Comput. Surv.** 2009]



■ Types of Outlier Detection

- **Type 1 Unsupervised:** No prior knowledge of data, similar to unsupervised **clustering**
→ **expectations:** distance, # errors
- **Type 2 Supervised:** Labeled normal and abnormal data, similar to supervised **classification**
- **Type 3 Normal Model:** Represent normal behavior, similar to **pattern recognition** → **expectations:** rules/constraints

[Victoria J. Hodge, Jim Austin: A Survey of Outlier Detection Methodologies. **Artif. Intell. Rev.** 2004]



Missing Value Imputation

■ Missing Value

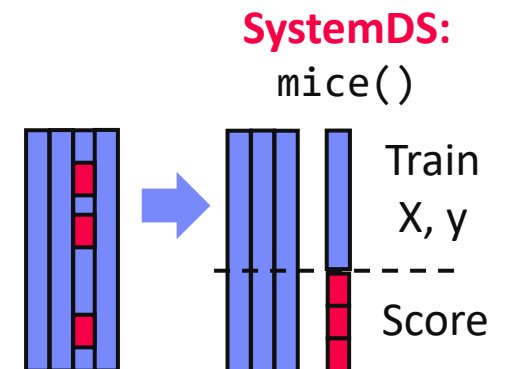
- Application context defines if 0 is missing value or not
- If differences between 0 and missing values, use NA or NaN

■ Basic Value Imputation

- General-purpose: replace by user-specified **constant**
- **Continuous variables**: replace by **mean**
- **Categorical variables**: replace by **median** or **mode**

■ Iterative Algorithms (**chained-equation imputation**)

- Train ML model to predict missing information (feature $k \rightarrow$ label, split data into observed/missing)
- Noise reduction: feature subsets + averaging



■ Dynamic Imputation

- Data exploration w/ on-the-fly imputation
- Optimal placement of imputation operations

[Jose Cambrero, John K. Feser,
Micah Smith, Samuel Madden:
Query Optimization for Dynamic
Imputation. **PVLDB 2017**]



Excursus: Time Series Recovery

■ Motivating Use Case

- Given overlapping weekly aggregates y (daily moving average)
- Reconstruct the original time series X

■ Problem Formulation

- Aggregates y
- Original time series X (unknown)
- Mapping O of subsets of X to y

→ Least squares regression problem

$$\underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}}_O \times \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}}_X = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}}_y$$

■ Advanced Method

- Discrete Cosine Transform (DCT) (sparsest spectral representation)
- Non-negativity and smoothness constraints

[Faisal M. Almutairi et al: HomeRun: Scalable Sparse-Spectrum Reconstruction of Aggregated Historical Data. **PVLDB 2018**]



→ **Use case:** high-precision sensor fusion w/ different data granularity

Selected Research Prototypes

■ ActiveClean (SampleClean)

- Suggest sample of data for manual cleaning (rule/ML-based detectors, **Simpson's paradox**)
- Update dirty model with gradients of cleaned data (weighted gradients of previous clean data and newly cleaned data)

[Sanjay Krishnan et al:
ActiveClean: Interactive Data
Cleaning For Statistical
Modeling. **PVLDB 2016**]



■ HoloClean

- Clean and enrich based on quality rules, value correlations, and reference data
- Probabilistic models for capturing data generation
- HoloDetect
 - **Learn data representations** of errors
 - **Data augmentation** w/ erroneous data from sample of clean data

[Alireza Heidari, Joshua McGrath,
Ihab F. Ilyas, Theodoros Rekatsinas:
HoloDetect: Few-Shot Learning for
Error Detection, **SIGMOD 2019**]

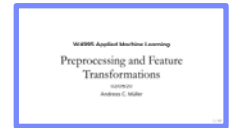


■ Other Systems

- **AlphaClean** (generate data cleaning pipelines) [preprint]
- **BoostClean** (generate repairs for domain value violations) [preprint]

Summary and Q&A

[Andreas C. Mueller: Preprocessing and Feature Transformations, Applied ML Lecture 2020]



- Data Acquisition and Integration
- Data Preparation and Feature Engineering
- Data Transformation and Cleaning

“Coming up with features is difficult, time-consuming, requires expert knowledge. "Applied machine learning" is basically feature engineering”
– Andrew Ng

■ Next Lectures

- [10 Model Selection and Management](#) [Jun 05]
 - Incl Data Augmentation
- [11 Model Debugging Techniques](#) [Jun 12]
- [12 Model Serving Systems and Techniques](#) [Jun 19]