

# A GENERAL FRAMEWORK FOR ADDRESSING “ANY” MACHINE LEARNING PROBLEM

Darren Homrighausen

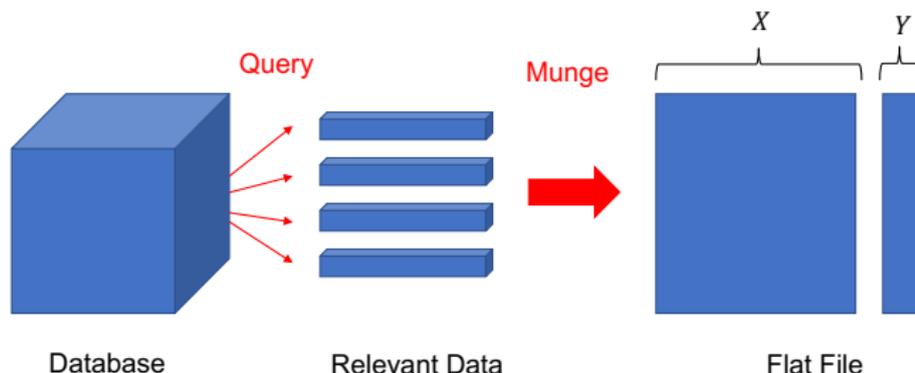
Visiting Assistant Professor, Dept. of Statistical Science, Southern Methodist University

September 21, 2017

# THE SETUP

A machine learning problem can be broken up into two parts:

1. Querying and then cleaning and/or manipulating the data into a format suitable for analysis  
(Sometimes referred to as **munging**)



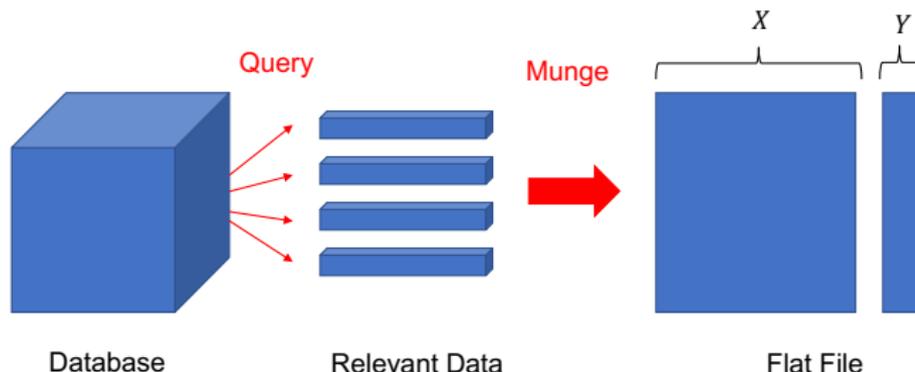
2. Applying machine learning methods to the data

# THE SETUP

A machine learning problem can be broken up into two parts:

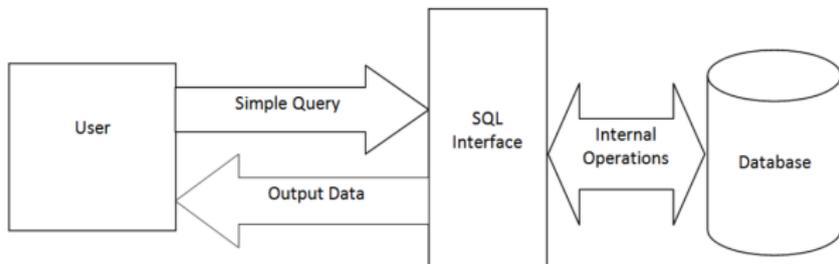
1. Querying and then cleaning and/or manipulating the data into a format suitable for analysis

(Sometimes referred to as **munging**)



2. Applying machine learning methods to the data

# QUERYING



```
INSERT INTO interestingData
SELECT id, trans, city, date
FROM cust_table
WHERE date > 1/1/2015
ORDER BY city, date;
```

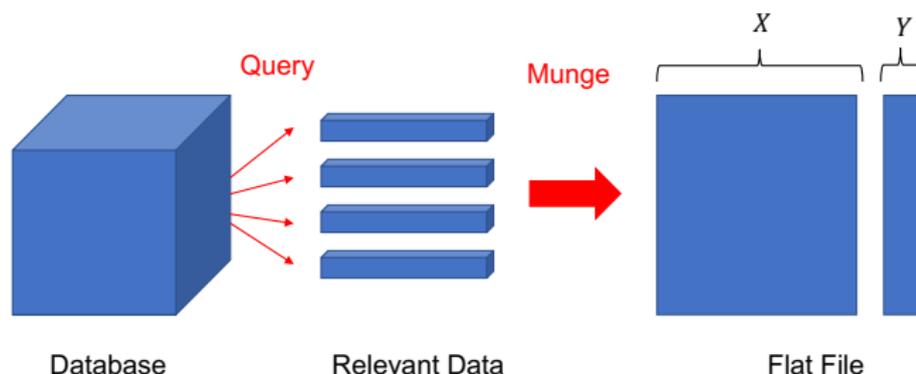
- .
- .
- .

# THE SETUP

A machine learning problem can be broken up into two parts:

1. Querying and then cleaning and/or manipulating the data into a format suitable for analysis

(Sometimes referred to as **munging**)



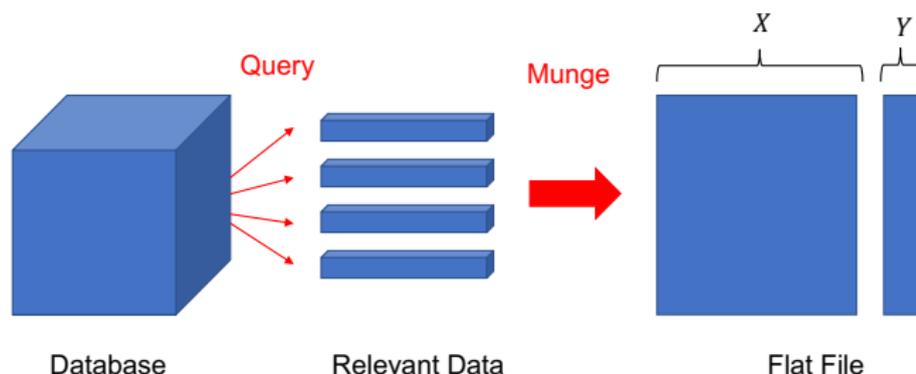
2. Applying machine learning methods to the data

# THE SETUP

A machine learning problem can be broken up into two parts:

1. Querying and then cleaning and/or manipulating the data into a format suitable for analysis

(Sometimes referred to as **munging**)



2. Applying machine learning methods to the data

# THE FEATURES

We need to determine the...

... appropriate processing of  $X$   
(Known as the **features** or **inputs**)

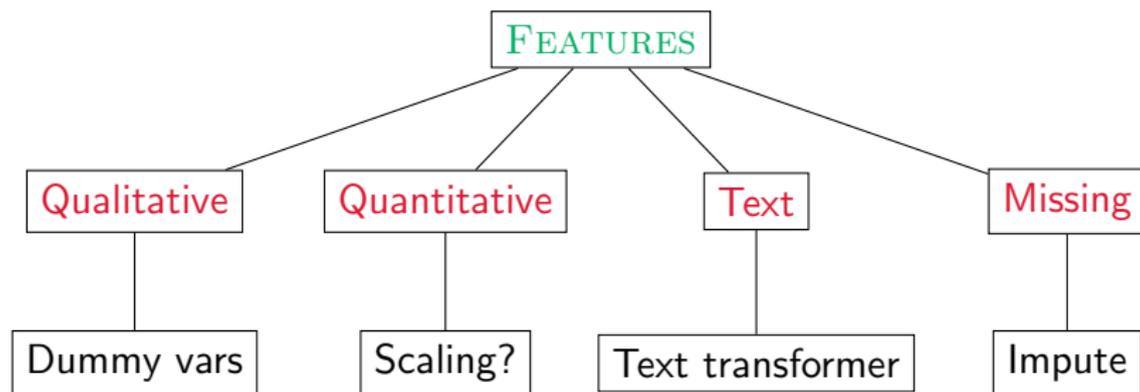
**EXAMPLE:**

$X$ : vectorized sky survey



Semi-supervised learning for photometric supernova classification\*

Joseph W. Richards,<sup>1,2†</sup> Darren Homrighausen,<sup>3</sup> Peter E. Freeman,<sup>3</sup> Chad M. Schafer<sup>3</sup>  
and Dovi Poznanski<sup>1,4</sup>



# THE FEATURES: QUALITATIVE

	x1	x2
1	-0.6264538	no
2	0.1836433	yes
3	-0.8356286	yes
4	1.5952808	no

Gets transformed to...

	x1	x2no	x2yes
1	-0.6264538	1	0
2	0.1836433	0	1
3	-0.8356286	0	1
4	1.5952808	1	0

# THE FEATURES: QUANTITATIVE

Many methods are not invariant to **scale**

The usual way of addressing this is...

... **Do standardize** all features for which scale is meaningful:

$$X \leftarrow \frac{(X - \text{mean}(X))}{\text{sd}(X)}$$

... **Don't standardize** any scale-free nor sparse features

(Care must be taken if normalizing sparse data)

 **NIH Public Access**  
**Author Manuscript**  
Author: Author not certified; available in PMC 2013 September 21.  
Published in final edited form as:  
Analyst. 2012 September 21; 137(18): 4280-4286. doi:10.1039/c2an35578g.

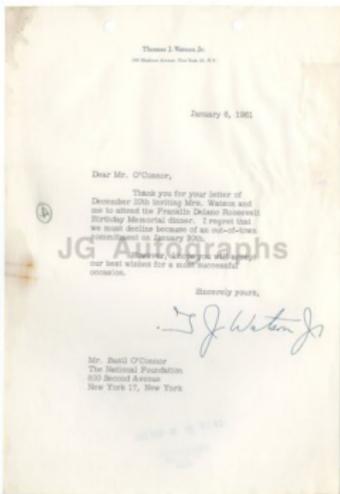
**Surface Enhanced Raman Spectroscopy (SERS) for the discrimination of *Arthrobacter* strains based on variations in cell surface composition**

Kate E. Stephen<sup>a</sup>, Darren Homrighausen<sup>b</sup>, Glen DePalma<sup>c</sup>, Cindy H. Nakatsu<sup>a</sup>, and Joseph Irudayaraj<sup>a</sup>

1	1	1	1	0	0	0	1
0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	1
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1
0	1	0	0	0	0	0	0

}  $X$

# THE FEATURES: TEXT



## Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS)

Item condition: --

Ended: May 27, 2014 16:59:11 PDT

Winning bid: **US\$11.61** [ 6 bids ]

Shipping: **\$3.99** Standard Shipping | [See details](#)

Item location: United States

Ships to: Worldwide

Delivery: Estimated within 3-6 business days Ⓜ

Payments: **PayPal** | [See details](#)

Returns: 14 days money back, buyer pays return shipping | [See details](#)

Guarantee: **ebay** MONEY BACK GUARANTEE | [See details](#)

Get the item you ordered or get your money back.  
Covers your purchase price and original shipping.

### Seller information

**jgautographs** (54927) ★

100% Positive feedback

[Follow this seller](#)

[See other items](#)

Visit store: [JG Autographs](#)

# THE FEATURES: TEXT

## BUYER:

 Always a pleasure! Smooth & pleasant transaction!	r**a ( 3618 ★ )	Jun-10-14 13:52
Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS) (#390846670600)	US \$11.61	<a href="#">View Item</a>

## SELLER:

 Great communication. A pleasure to do business with.	Buyer: r**a ( 3618 ★ )	Jun-05-14 18:59
Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS) (#390846670600)	—	<a href="#">View Item</a>

The  $X$  matrix can then be written as  $X = \begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \end{bmatrix}$  where...

$$X_1^T = \begin{bmatrix} 1 & 2 & 1 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$X_2^T = \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$

---

A text analysis of Ebay auctions

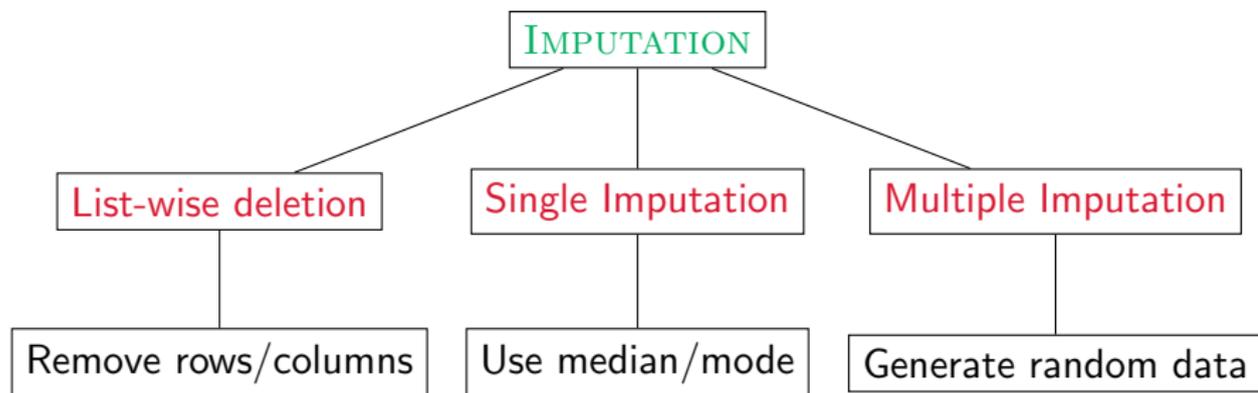
---

## THE FEATURES: MISSING

Corrupted, unrecorded, or unreliable data is commonly referred to as **missing data**

In statistics, correcting for missing data is known as **imputation**

There are many, many techniques available:



# THE FEATURES: MISSING

- Data size/complexity  
(Does it fit in RAM?)
- Business purpose  
(Is data precious? Development time?)
- Are any observations/features missing a large fraction of values?
- Type of features  
(Any sparsity? Is multivariate normality appropriate?)
- Any atypical missing value indicators?  
(e.g. using -1000 for income to indicate a missing value)

# THE SUPERVISOR

We need to determine the...

... nature of  $Y$

(Known as the **supervisor(s)** or **output(s)**)

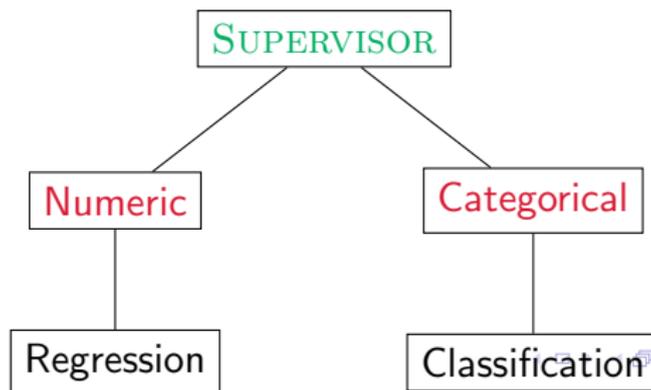
**EXAMPLE:**

$$Y = \begin{cases} 1 & \text{If type 1a supernova} \\ -1 & \text{If not} \end{cases}$$



**Semi-supervised learning for photometric supernova classification\***

Joseph W. Richards,<sup>1,2,†</sup> Darren Hornighausen,<sup>3</sup> Peter E. Freeman,<sup>3</sup> Chad M. Schafer<sup>3</sup> and Dovi Poznanski<sup>1,4</sup>



# EVALUATION METRICS

How to judge success?

Often, this is just **mean square error** or **miss-classification rate**

There can be many others:

**EXAMPLE:** When classifying supernovae, it is bad to incorrectly label a Type Ia supernova

→ Evaluation metric:

$$\left( \frac{1}{\text{Total \#}} \right) \frac{(\# \text{ Correctly labeled})^2}{\# \text{ Correctly labeled} + 3(\# \text{ Incorrectly labeled})}$$

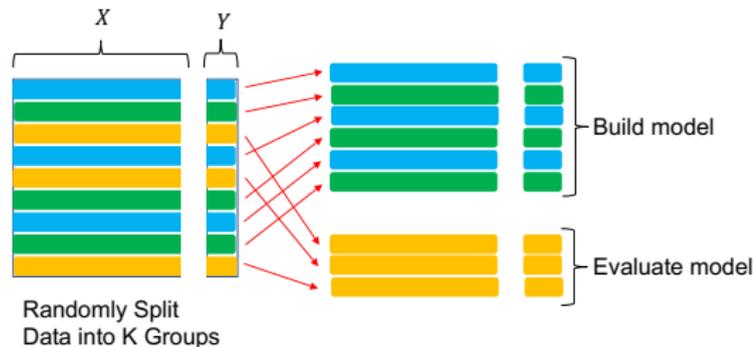
## VALIDATION SET

We need a realistic measure of the **evaluation metric**

If at all possible, set aside a (random) validation set

(Say, 10% of the data)

Alternatively (or additionally)  
a common approach is  
**K-Fold Cross-Validation (CV):**



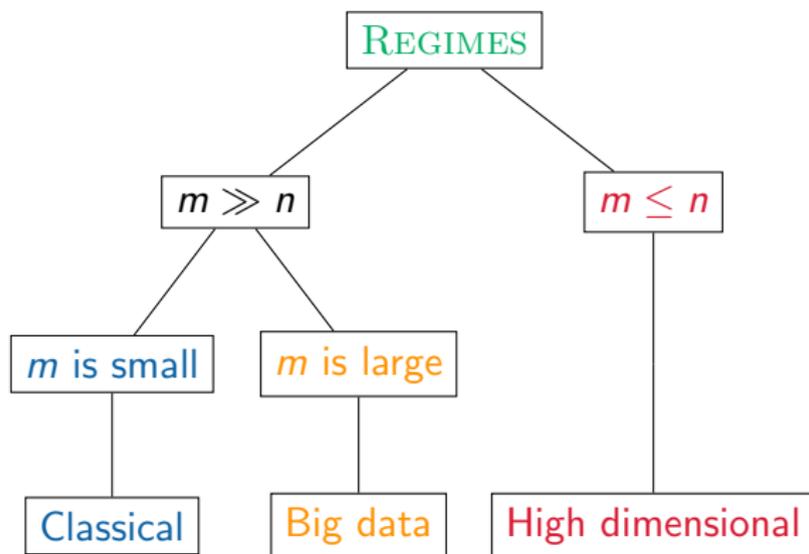
**VERY IMPORTANT:** Make sure to use stratified sampling over

- Any rare, qualitative features
- The supervisor  $Y$   
(If doing classification)

# The analysis

# TURNING THESE IDEAS INTO PROCEDURES

There are roughly **three** regimes of interest, assuming  $X \in \mathbb{R}^{m \times n}$



**ADDITIONALLY:** Is the data sparse?

(i.e. Does it have a lot of zeros?)

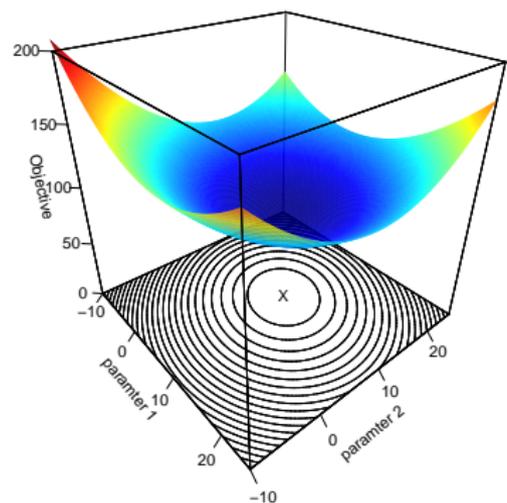
# BIG DATA

**BIG DATA** is usually characterized by 4 “V’s”

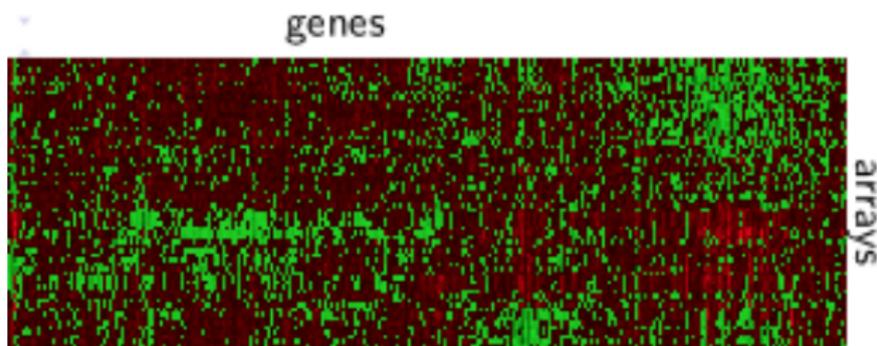
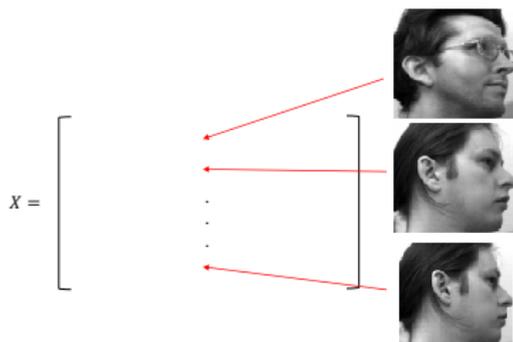
(Volume, Variety, Velocity, Veracity)

Depending on the data and the desired method, we could:

- Combine randomized projections together with in-memory procedures
- Use stochastic gradient descent (or related methods)
- Leverage an iterative implementation for exact computation (e.g. the QR decomposition for least squares)
- Break the computations down into small bits and distribute these to different cores/processors/nodes (e.g. using the MapReduce paradigm)

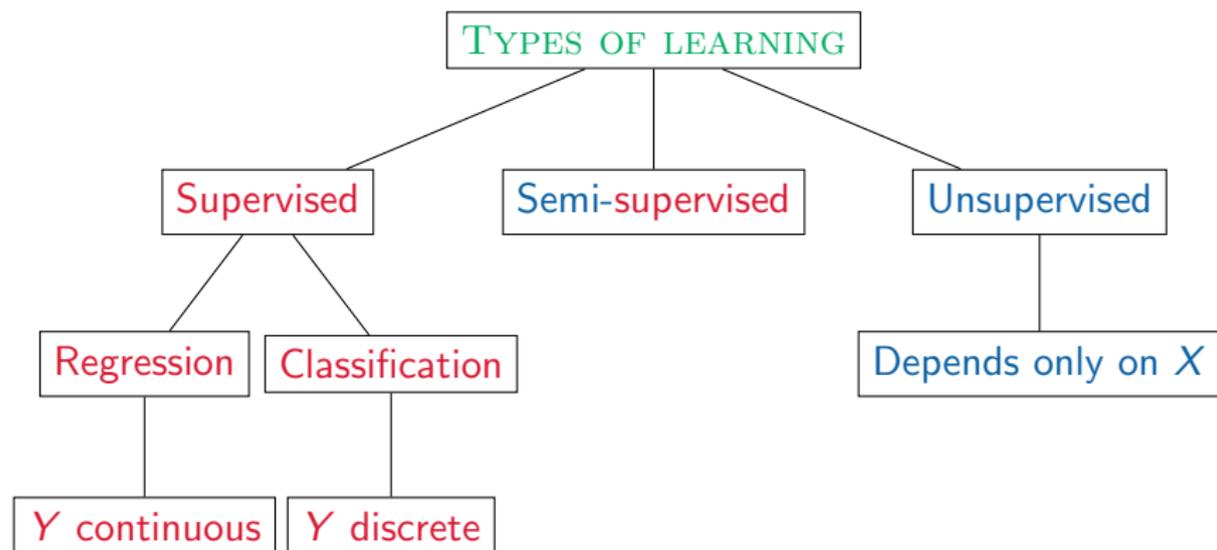


# HIGH DIMENSIONAL REGIME: EXAMPLES



# Methods

# METHODS: TYPES OF LEARNING



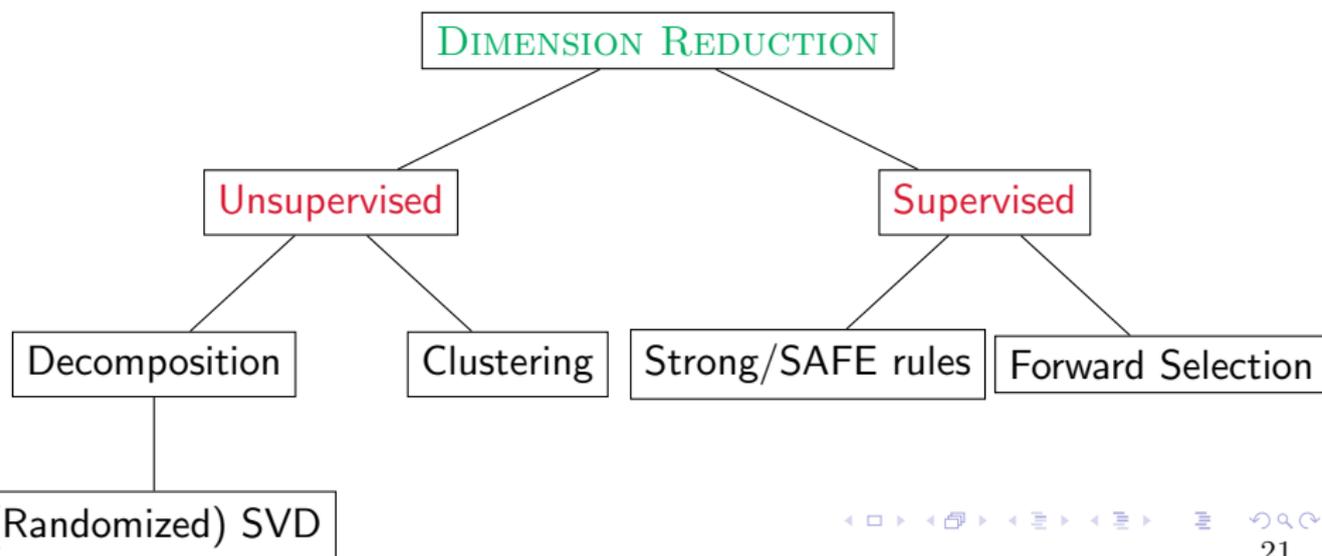
# METHODS: DIMENSION REDUCTION

Dimension reduction can help by ...

... reducing the computational load at later steps

... improving prediction performance

Some examples..



# METHODS: SUPERVISED

- CLASSIFICATION:

- ▶ (Sparse) Logistic Regression  
(I include Linear Discriminant Analysis (LDA) here)
- ▶ Naive Bayes
- ▶ Support Vector Machines (SVM)
- ▶ k-Nearest Neighbors (KNN)

- REGRESSION:

- ▶ (Sparse) Linear Regression  
(I include Elastic Net here)
- ▶ Support Vector Regression

- BOTH:

- ▶ Random Forest
- ▶ Gradient Boosting Machines (GBM)
- ▶ Neural Networks

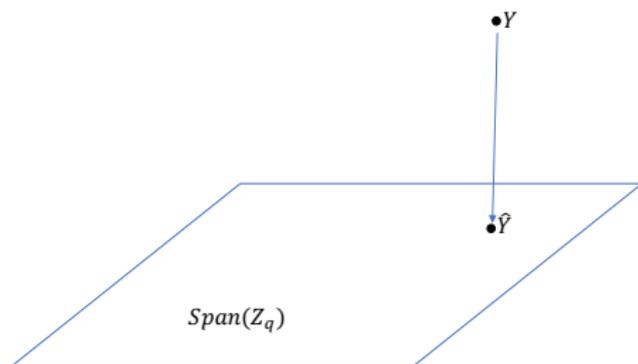
# LEARNING

## EXAMPLE: Semi-supervised learning

1. Form  $X = UDV^T$

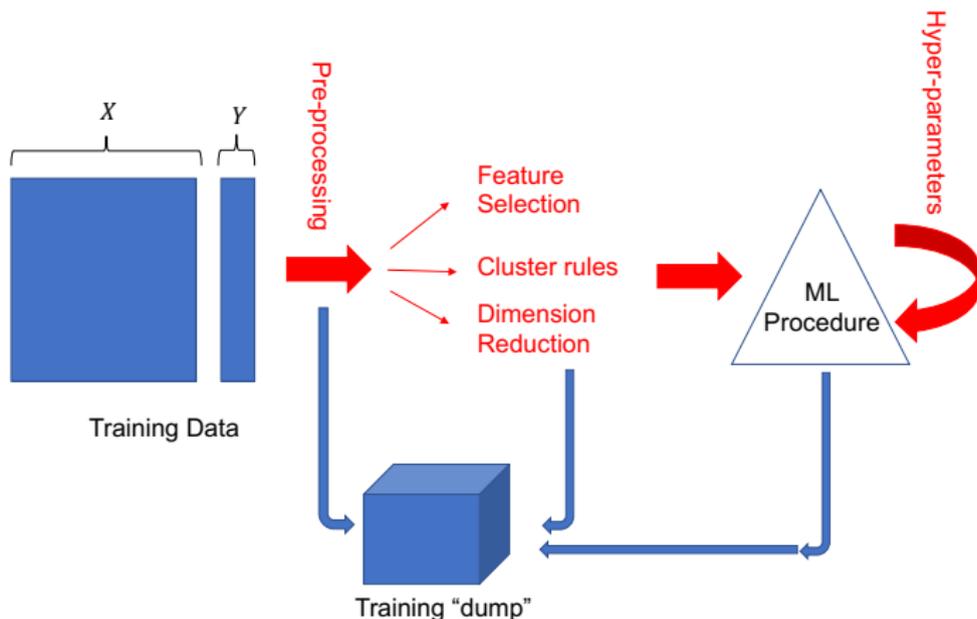
(That is, the singular value decomposition of the (scaled) matrix  $X$ )

2. Project  $Y$  onto the column space spanned by the first  $q$  columns of  $UD$  (call this object  $Z_q$ )



(This is commonly referred to as “principal components regression”)

# ANALYSIS FLOWCHART



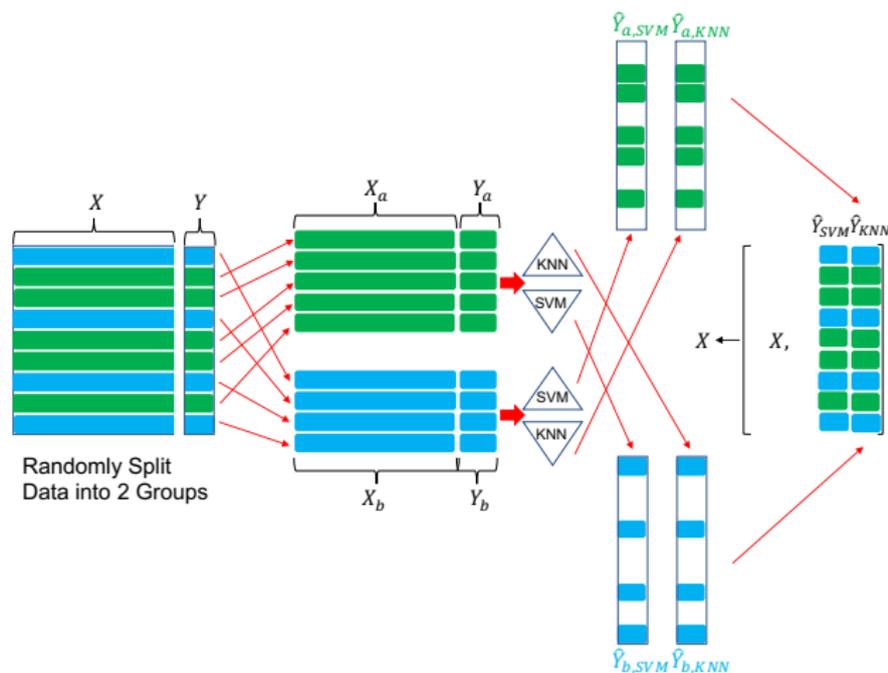
- It is very important to **save** any preprocessing/transformations. These must be applied to the test features
- Choose hyper-parameters via CV or other risk estimator (e.g.  $q$  from principal components regression)

# ENSEMBLE METHODS

Combining supervised methods can result in improved performance

These **ensembles** or **stacks** can be formed in several ways

A **feature** based approach is as follows:



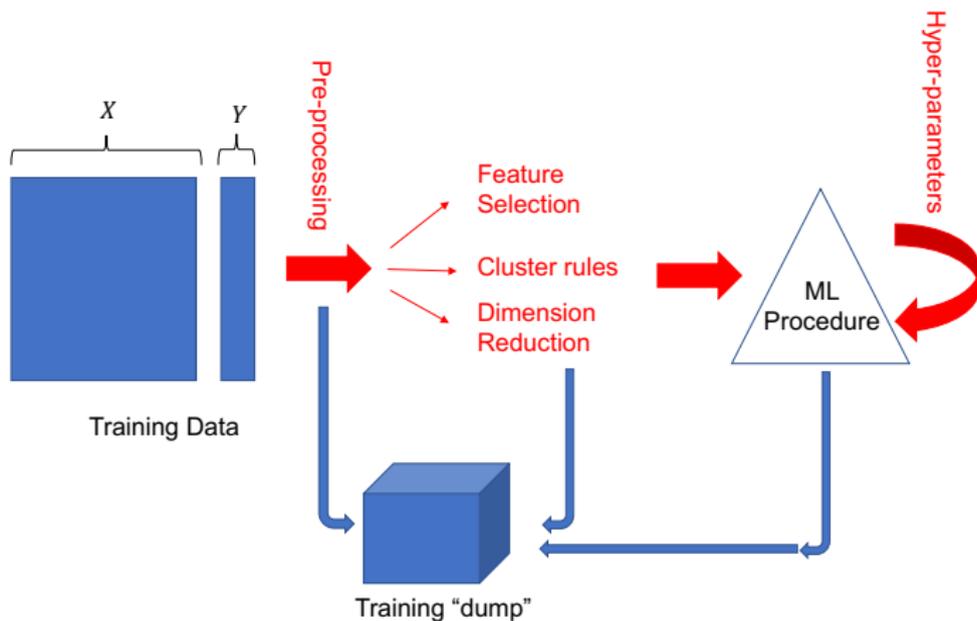


# GOAL

Let's classify documents about renaissance artists and the corresponding **TMNT**

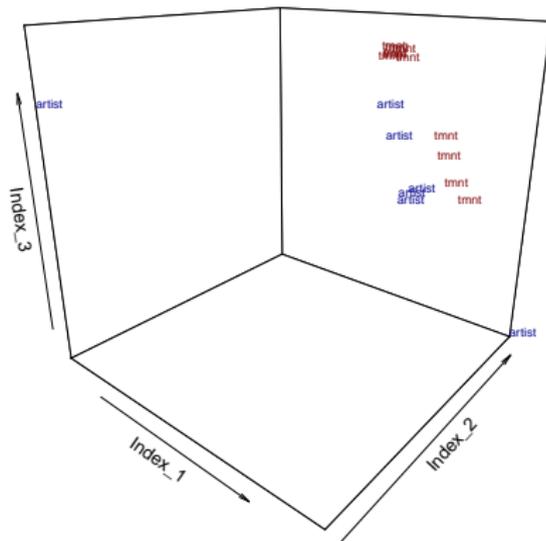


# ANALYSIS FLOWCHART: REMINDER

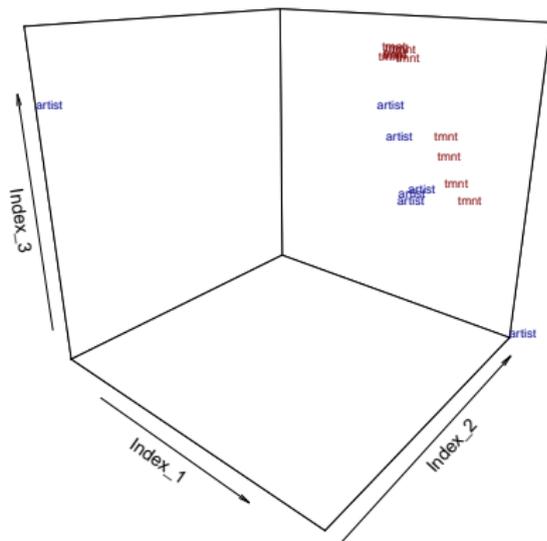




# DIMENSION REDUCTION



# DIMENSION REDUCTION



" ...Category View page ratings Rate  
this page What's this? Trustworthy Objective  
Complete Well-written I am highly knowledgeable  
about this topic (optional) Submit ratings ...

# ML PROCEDURE

- CLASSIFICATION:
  - ▶ (Sparse) Logistic Regression
  - ▶ Naive Bayes
  - ▶ Support Vector Machines (SVM)
  - ▶ k-Nearest Neighbors (KNN)
- REGRESSION:
  - ▶ (Sparse) Linear Regression
  - ▶ Support Vector Regression
- BOTH:
  - ▶ Random Forest
  - ▶ Gradient Boosting Machines (GBM)
  - ▶ Neural Networks

# ML PROCEDURE

- CLASSIFICATION:
  - ▶ (Sparse) Logistic Regression
  - ▶ Naive Bayes
  - ▶ Support Vector Machines (SVM)
  - ▶ k-Nearest Neighbors (KNN)
- REGRESSION:
  - ▶ (Sparse) Linear Regression
  - ▶ Support Vector Regression
- BOTH:
  - ▶ Random Forest
  - ▶ Gradient Boosting Machines (GBM)
  - ▶ Neural Networks

# ML PROCEDURE: SPARSE LOGISTIC REGRESSION

$$\pi(X) = \text{Prob}(Y = \text{artist} \mid X = [\text{accademia}, \text{accent}, \text{accept}, \text{accid}, \dots]^\top)$$

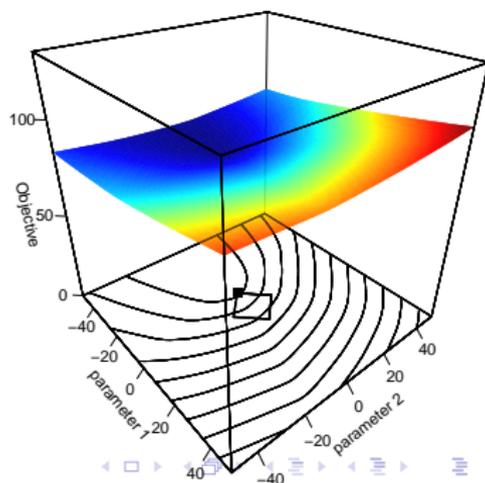
$$\log\left(\frac{\pi(X)}{1 - \pi(X)}\right) = \beta^\top X$$

$$\rightarrow \ell(\beta) = -\sum_{i=1}^m \left[ Y_i X_i^\top \beta - \log\left(1 + e^{X_i^\top \beta}\right) \right] \text{ is the objective}$$

Constrained minimization  $\longleftrightarrow$

$\rightarrow$  Use projected gradient descent

$$(1) \quad \hat{\beta} \stackrel{\text{update}}{=} \hat{\beta} - \eta \nabla \ell|_{\hat{\beta}}$$
$$(2) \quad \hat{\beta} \stackrel{\text{set}}{=} \underset{\text{feasible } \beta}{\text{argmin}} \|\hat{\beta} - \beta\|$$



# ML PROCEDURE: SPARSE LOGISTIC REGRESSION

## INFERENCE:

The magnitude/sign of  $\hat{\beta}$  indicates which words affect the estimated probabilities the most

$$\pi(X) = \text{Prob}(Y = \text{artist} \mid X = [\text{accademia}, \text{accent}, \text{accept}, \text{accid}, \dots]^T)$$



$\hat{\beta}$ : positive, negative

# WRAPPING UP

This approach has performed well on many of the problems I have worked on

Of course, nothing is perfect

It is important to keep on improving on what we have learned...

... just like in machine learning

THANK YOU!