# WELCOME BACK!

Catalit LLC

Humans

Weight (lbs)

Height (in)

Catalit LLC

Internet Service Providers

# COMBINED

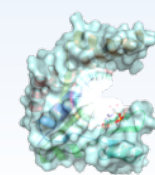| | CONTINUOUS | CATEGORICAL |
|---|---|---|
| SUPERVISED | REGRESSION | CLASSIFICATION |
| UNSUPERVISED | DIMENSION REDUCTION | CLUSTERING |

Exploration
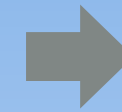
# ML STEPS

| 1. Collection | → | 2. Processing | → | 3. Model Building | → | 4. Evaluation |

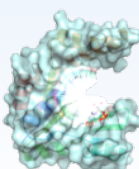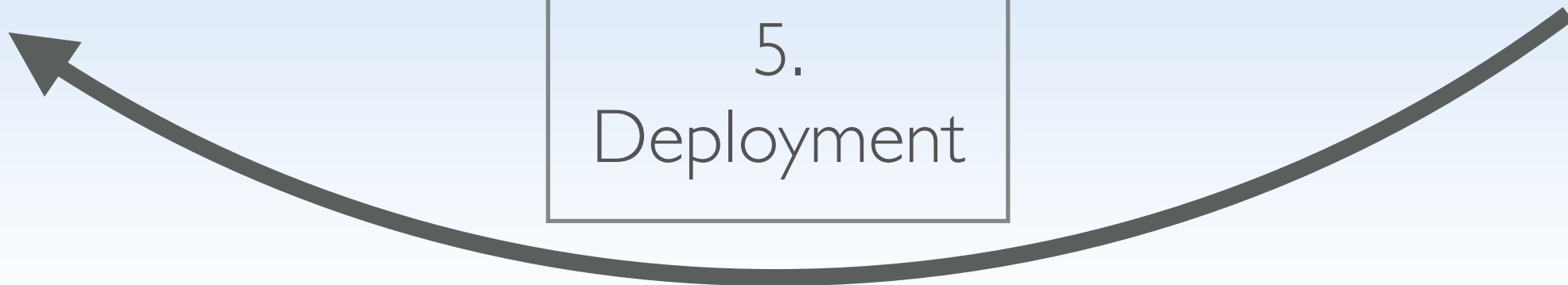| Text Image/Video Transactions User info Revenues … | Clean Transform Impute Features …. | Prediction ✔ ✘ | Score Train/Test Cross Val |

5. Deployment

Catalit LLC

Catalit LLC

Define Hypothesis
Define Cost

Minimize Cost

Catalit LLC

# OVERFITTING



Regularization

# HISTORY OF SUPERVISED LEARNING



Logistic Regression (large scale)

Support Vector Machine

Decision Tree

Neural network (perceptron)

Neural Network

1960    1970    1980    1990    2000    2010    Year

http://www.aboutdm.com/2013/04/history-of-machine-learning.html

Catalit LLC

# TRAIN - TEST SPLIT

# PRECISION - RECALL & ACCURACY

- **Precision:** When test is positive, how often is prediction correct?

  - TP / test yes

- **Recall:** When actual value is positive, how often is prediction correct?

  - TP / actual yes

- **Accuracy:** Overall, how often is it correct?

  - (TP + TN) / total

|  | Condition Positive | Condition Negative |
|---|---|---|
| **Test Positive** | **TRUE POSITIVE** | **FALSE POSITIVE** |
| **Test Negative** | **FALSE NEGATIVE** | **TRUE NEGATIVE** |

Catalit LLC

# DISTANCE & SIMILARITY

| | Age | Gender | Annual Salary | Months in residence | Months in job | Current Debt |
|---|---|---|---|---|---|---|
| Client 1 | 23 | M | $30,000 | 36 | 12 | $5,000 |
| Client 2 | 30 | F | $45,000 | 12 | 12 | $1,000 |
| Client 3 | 19 | M | $15,000 | 3 | 1 | $10,000 |

# CLUSTER VALIDATION



(a) Cohesion.  (b) Separation.

Catalit LLC

# DATA CLEANING

# DATA CLEANING



Other examples include:

Remove inconsistencies
Data type harmonization
Typos correction
Formatting (eg. timestamps)
Sorting

Catalit LLC

# TIME SPENT

<u>O</u>btain Data

<u>S</u>crub Data         80%

<u>E</u>xplore

<u>M</u>odel Algorithms

i<u>N</u>terpret Results   20%

Catalit LLC

# IMPORT DATA

## IO Tools (Text, CSV, HDF5, ...)

The pandas I/O API is a set of top level `reader` functions accessed like `pd.read_csv()` that generally return a `pandas` object.

- *read_csv*
- *read_excel*
- *read_hdf*
- *read_sql*
- *read_json*
- *read_msgpack* (experimental)
- *read_html*
- *read_gbq* (experimental)
- *read_stata*
- *read_sas*
- *read_clipboard*
- *read_pickle*

Catalit LLC

# JSON

**JSON** (JavaScript Object Notation) is:

a lightweight data-interchange format

a string



```
{ "empinfo" :
    {
        "employees" : [
        {
            "name" : "Scott Philip",
            "salary" : £44k,
            "age" : 27,
        },

        {
            "name" : "Tim Henn",
            "salary" : £40k,
            "age" : 27,
        },

        {
            "name" : "Long Yong",
            "salary" : £40k,
            "age" : 28,
        }
        ]
    }
}
```

Catalit LLC

# API



Collection

GET https://api.instagram.com/v1/users/10

Operation

http://www.pythonapi.com/

# CONCATENATE DATA

# MERGE DATA
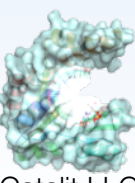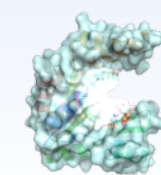
**left**

| | A | B | key |
|---|---|---|---|
| 0 | A0 | B0 | K0 |
| 1 | A1 | B1 | K1 |
| 2 | A2 | B2 | K2 |
| 3 | A3 | B3 | K3 |

**right**

| | C | D | key |
|---|---|---|---|
| 0 | C0 | D0 | K0 |
| 1 | C1 | D1 | K1 |
| 2 | C2 | D2 | K2 |
| 3 | C3 | D3 | K3 |

**Result**

| | A | B | key | C | D |
|---|---|---|---|---|---|
| 0 | A0 | B0 | K0 | C0 | D0 |
| 1 | A1 | B1 | K1 | C1 | D1 |
| 2 | A2 | B2 | K2 | C2 | D2 |
| 3 | A3 | B3 | K3 | C3 | D3 |

**left**

| | A | B | key |
|---|---|---|---|
| 0 | A0 | B0 | K0 |
| 1 | A1 | B1 | K1 |
| 2 | A2 | B2 | K0 |
| 3 | A3 | B3 | K1 |

**right**

| | C | D |
|---|---|---|
| K0 | C0 | D0 |
| K1 | C1 | D1 |

**Result**

| | A | B | key | C | D |
|---|---|---|---|---|---|
| 0 | A0 | B0 | K0 | C0 | D0 |
| 1 | A1 | B1 | K1 | C1 | D1 |
| 2 | A2 | B2 | K0 | C0 | D0 |
| 3 | A3 | B3 | K1 | C1 | D1 |

http://pandas.pydata.org/pandas-docs/stable/merging.html

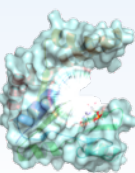# REBUILD MISSING



- Missing at Random?

# REBUILD MISSING

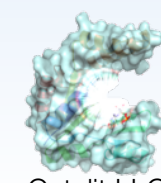| MCAR<br>Missing completely at random | MAR<br>Missing at random | MNAR<br>Missing not at random |
| --- | --- | --- |
| Missing value (y) neither depends on x nor y<br><br>e.g.: some survey questions asked to fewer people | Missing value (y) depends on x, but not y<br><br>e.g. Respondents in service occupations less likely to report income | The probability of a missing value depends on the variable that is missing<br><br>e.g.: Respondents with high income less likely to report income |

# TECHNIQUES

- Imputation, Partial imputation

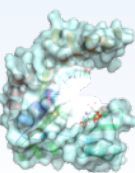- Deletion, Partial deletion
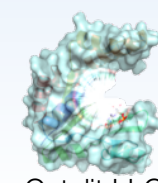
- Analysis

- Interpolation

# STANDARDIZATION

# STANDARDIZATION

- Sep 12th, 2015    9/12/15    12-Sep-15

- USA   United States of America   EU  U.s.a.

- Mr   Mr.   Mister

- etc. etc.

Catalit LLC

# NORMALIZATION

# NORMALIZATION

- STANDARD

    - subtract mean

    - divide by std

- MINMAX

    - subtract min

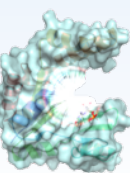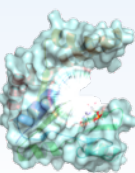    - divide by (max-min)

# DE-DUPLICATE

# FEATURES FROM TEXT

- Bag of Words approach:

  - Term Frequency (TF)

  - Inverse Document Frequency (IDF)

- NLP Approach

  - Stemming

  - Parts of Speech tagging

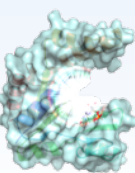  - Named Entity Detection

  - Parsing

Catalit LLC

# TFIDF

- Term frequency

  - Nterm/Nterms in document

- Document frequency

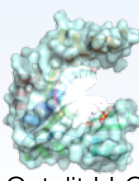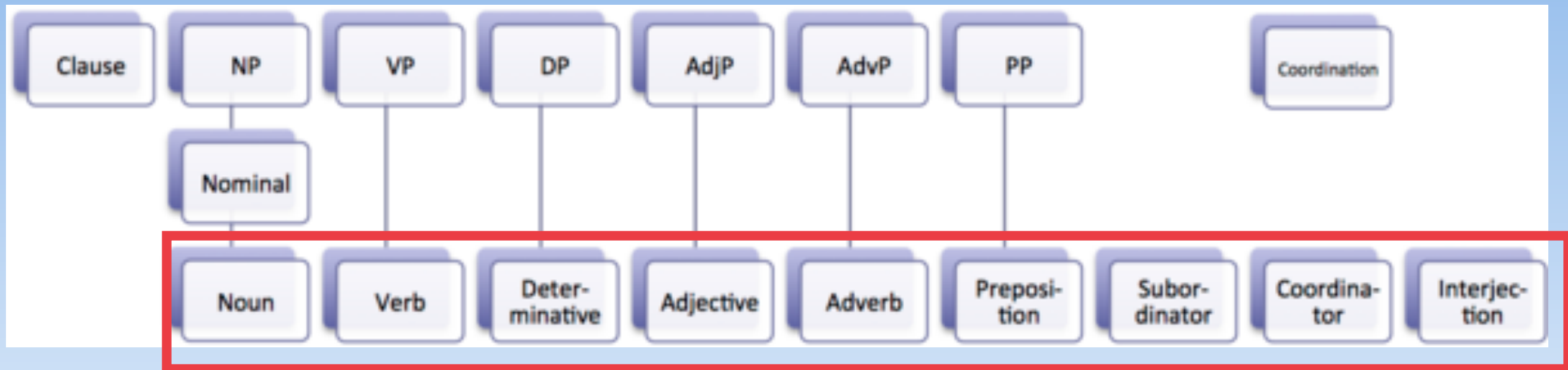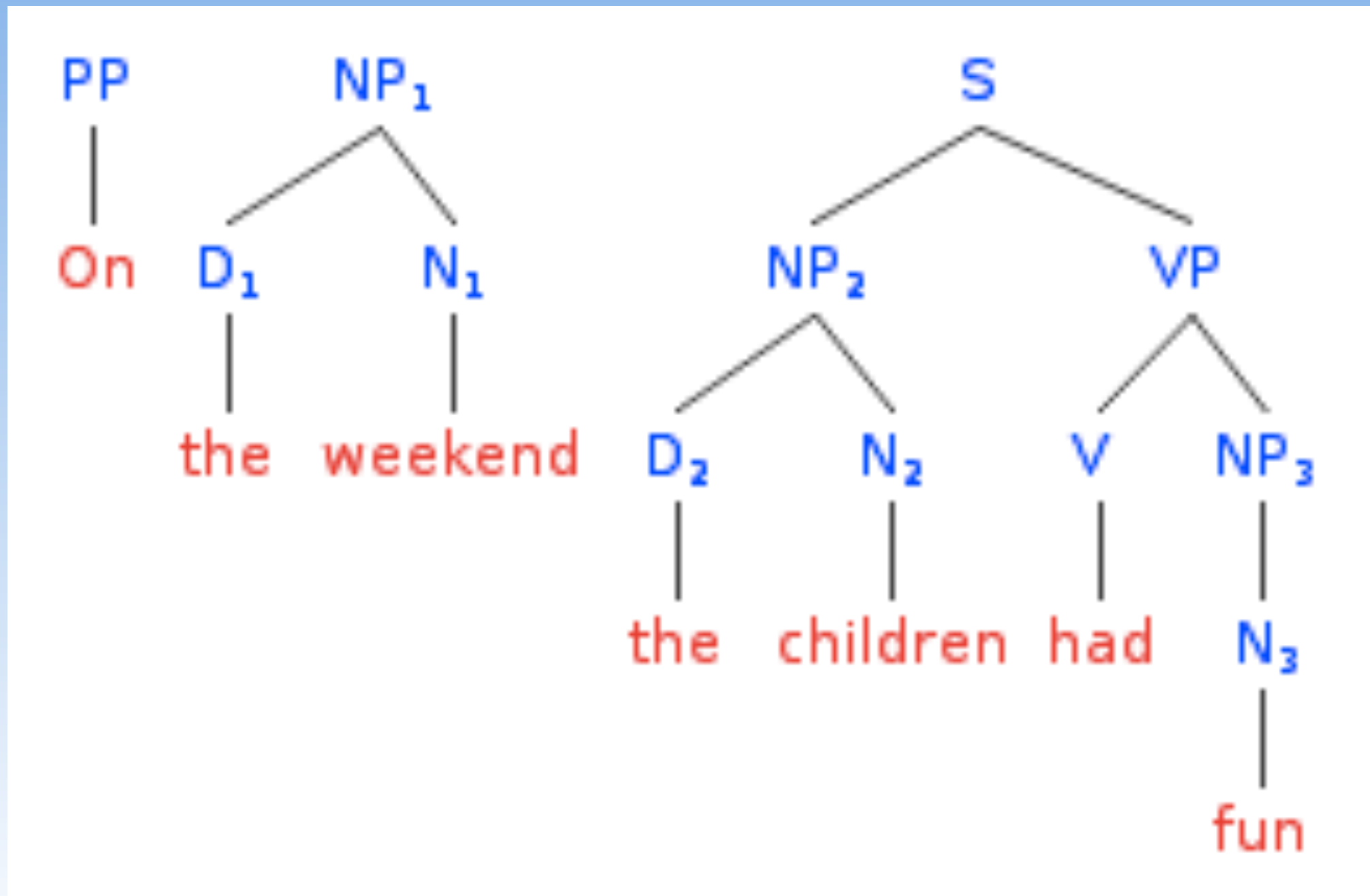  - Ndocuments containing term/Ndocuments

# STEMMING

- science, scientist => scien

- swim, swimming, swimmer => swim

- Porter stemmer

- Very useful to reduce feature set size
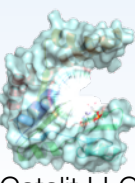
# PART OF SPEECH

# PARSING

# STANDARDIZATION

- you

  - u, ya, yo, yaaaa, yew, you, yoiu, youy, yoooooo, youz, yooouuuu

- Together

  - 2gether, tegetha, tgthr, togather, 2getha, 2gthr, togeter, togehter, t0gether, togeda, 2getter

# LAB CLEANING  + TEXT