- Target Variable • The y, outcome, or dependent variable
  - The "unknown"
  - what we are trying to predict

**Data Visualization.**
1) Explores data + understanding

2) Communicates to others

**Data Governance**
- Data Policy (security policy)
- Storing, managing & Processing data in a distributed enviornment

### Questions DS CAN ANSWER

**Regression** - How many or how much?

**Classification** - Is this an observation A, or B, or C etc?
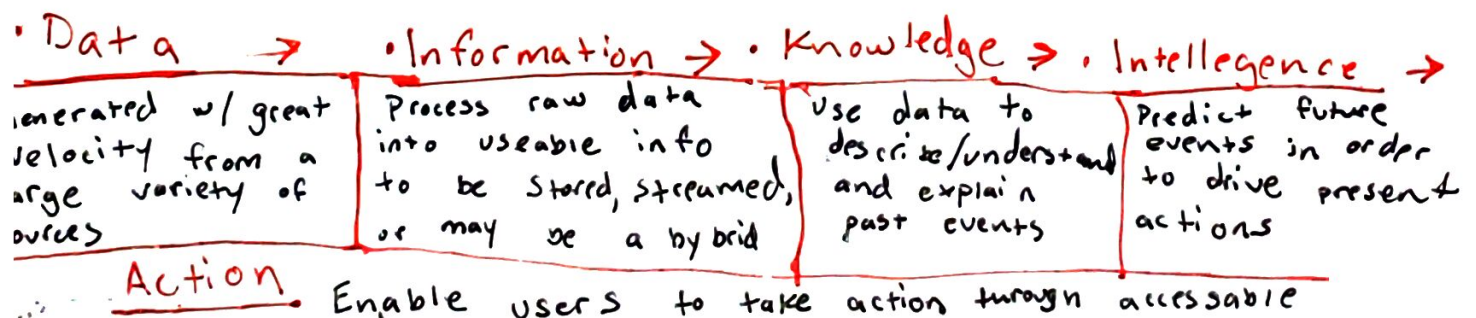
**Clustering** - What grouping s/ relationships exist in the data already?

**Time Series Analysis** - Whats our next likely outcome?

**Anomoly Detection** : Is this weird?

**Data Science Pipeline**
- End to end
- Left to right

| • Data → | • Information → | • Knowledge → | • Intellegence → |
|---|---|---|---|
| Generated w/ great Velocity from a large variety of sources | Process raw data into useable info to be Stored, streamed, or may be a hybrid | Use data to describe/understand and explain past events | Predict future events in order to drive present actions |

**Action** - Enable users to take action through accessable

## Data Planning

**Goal:**
- Clearly define your goal (write it down)
  - Measures of success, and plans on how to achieve that

**Deliverable:**
- Documentation of your goal
- If you haven't defined success, you won't know when you have achieved it

**How You Get There:** By answering questions about the final product & formulating or identifying any initial hypotheses.

---

## Acquisistion

- AKA: Data gathering
  Data Import
  Data Wrangling
  (Acquisistion + Prep)

**Goal:** • Create a path from original data sources to the enviornment in which you will work w/ the data

**Deliverable:** A file, acquire.py, that contains the function(s) needed to reproduce the acquisition of data

**How to get There** : SQL: Clean-up, integration, aggregation or other manipulation of data in the SQL Enviornment

Pylib: pandas

• May use Spark and/or Hive when acquiring data from a distributed enviornment such as HDFS.
• Examples of source types:
  • RDBMS • HDFS  • Static local flat files (csv, txt, xls