# PyData Ecosystem Webinar:
# Introduction and Best Practices

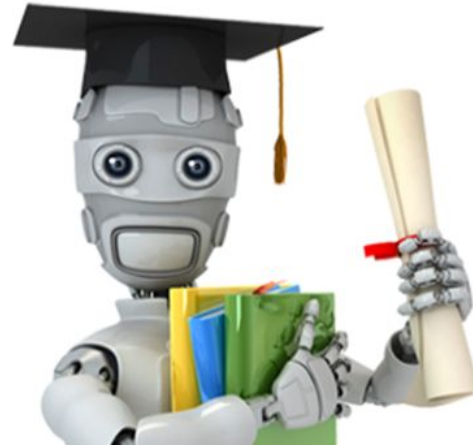**Tim Babych**
Sphere Software Senior Engineer

The field of study that gives computers the ability to learn without being explicitly programmed
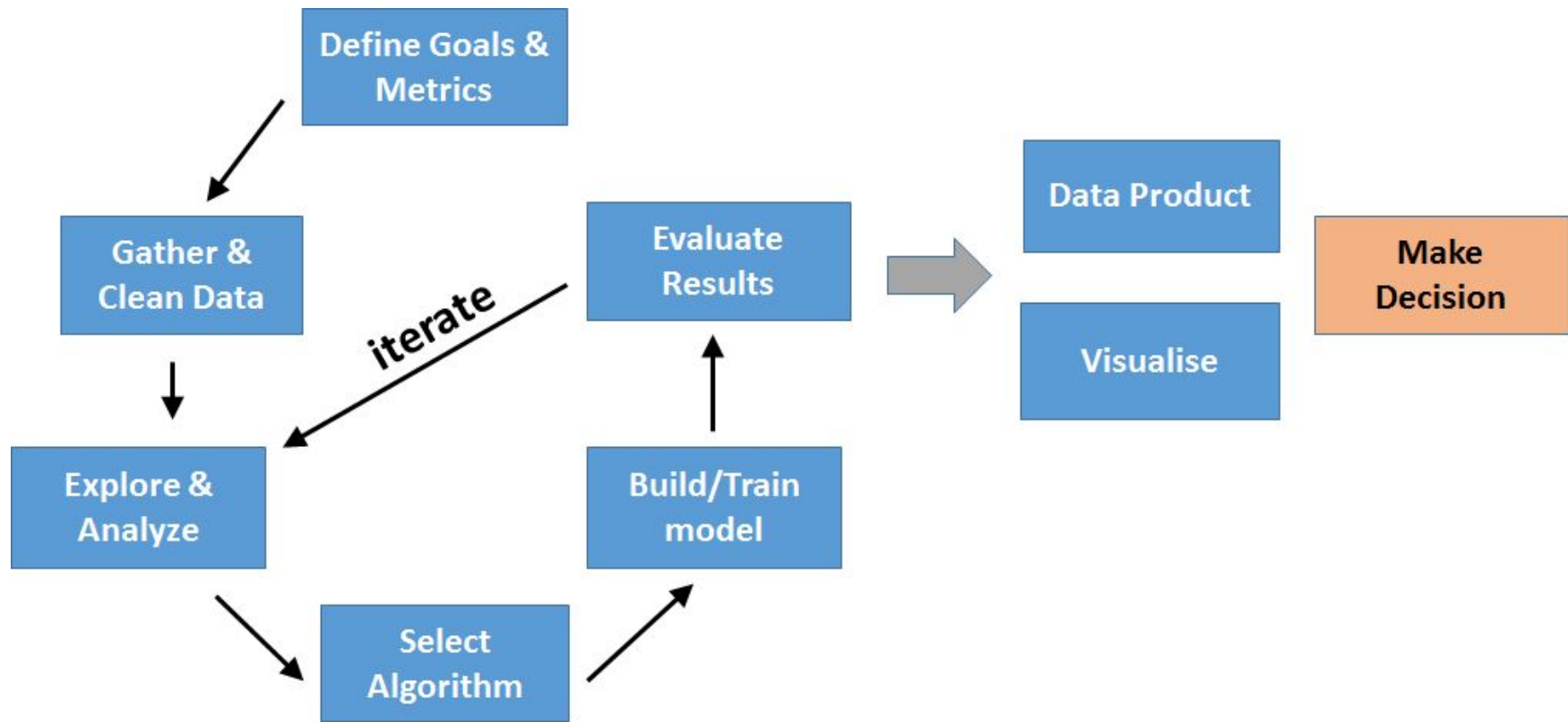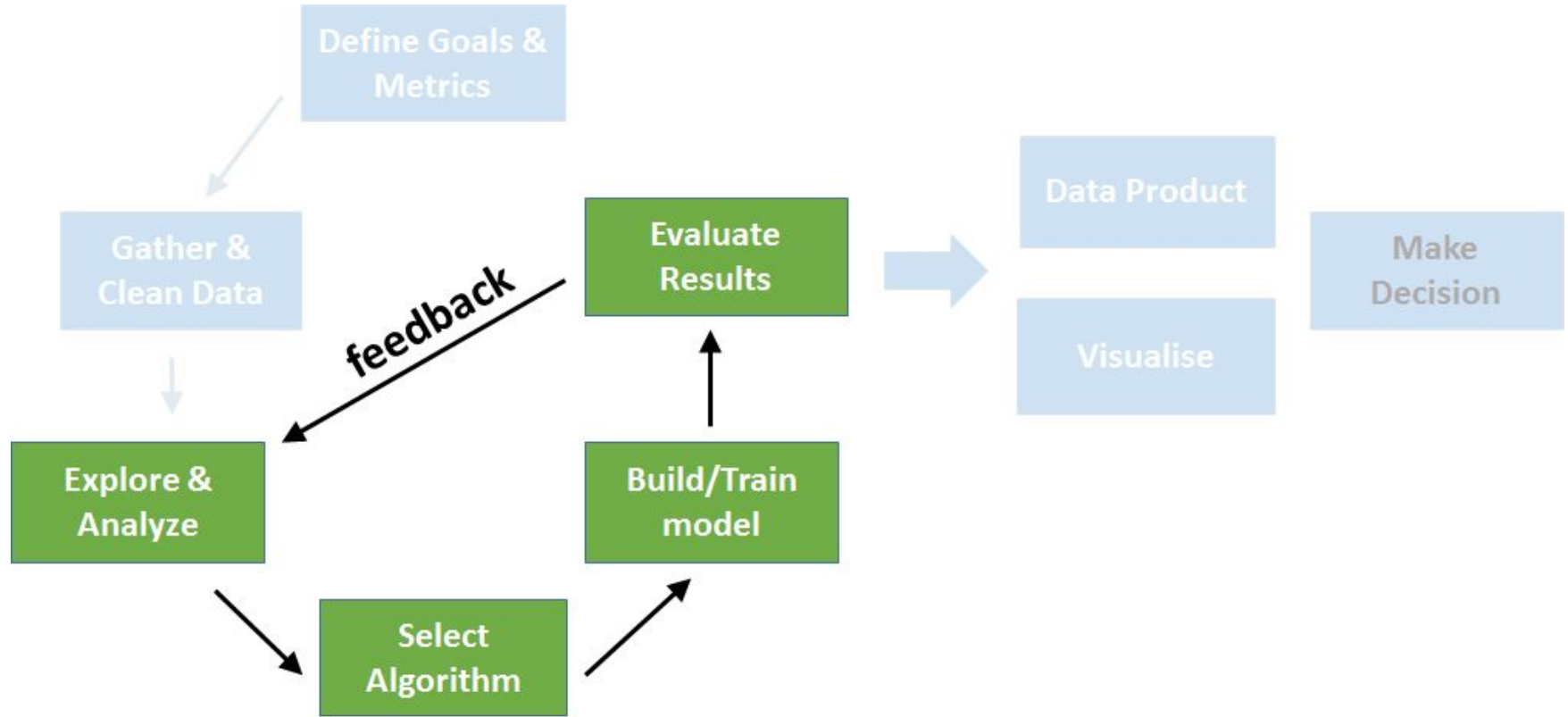
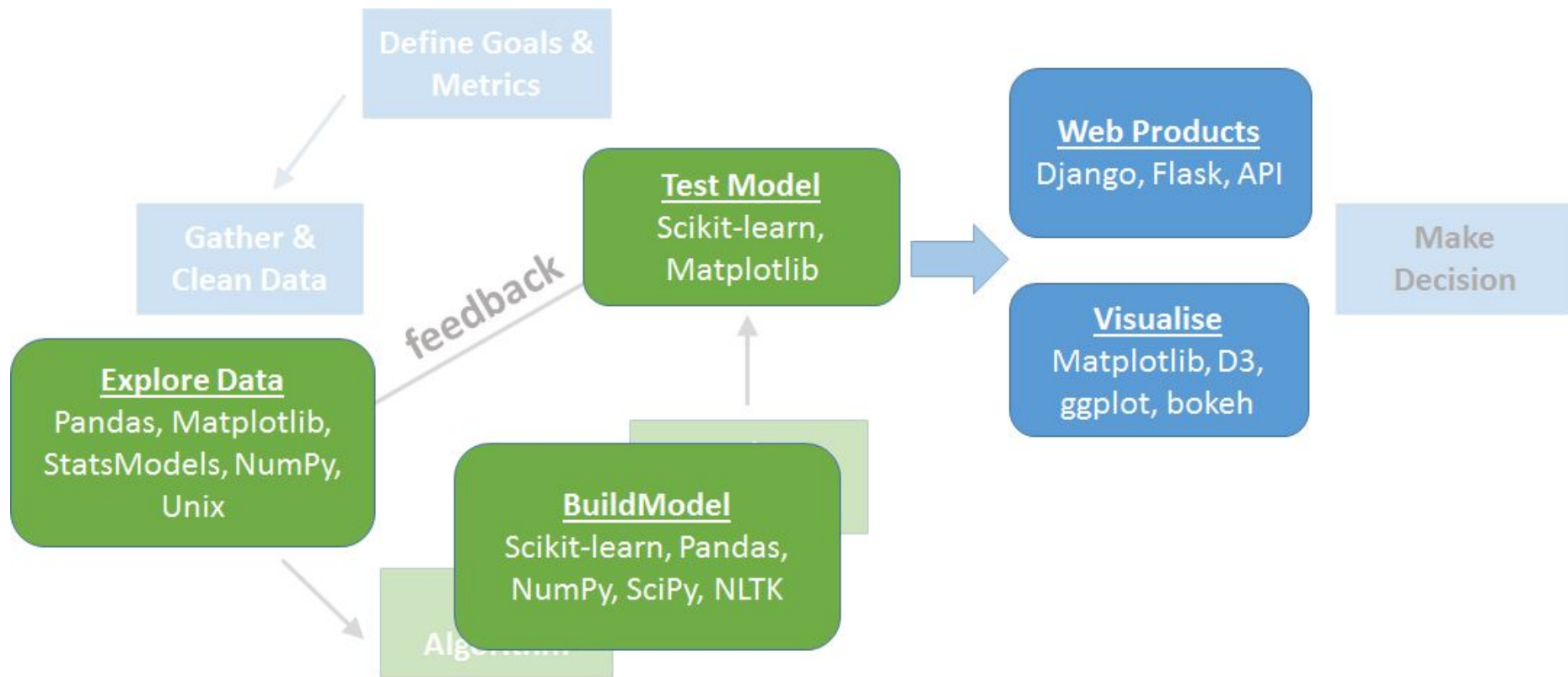--Arthur Samuel, 1959

## "Right answers" do exist

- Spam detectors
- Weather prediction
- Game outcomes
- Medical diagnosis
- Insurance
- Object detection
- Speech recognition

**There are no right answers!  Much harder.**

- Find some structure in the given data

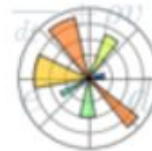- Cluster the data into groups

- Playing games

NumPy

Base
N-dimensional
array package

SciPy library

Fundamental
library for scientific
computing

Matplotlib

Comprehensive 2D
Plotting

IP[y]:
IPython

IPython

Enhanced
Interactive Console

Sympy

Symbolic
mathematics

pandas

Data structures &
analysis

EXAMPLE

**survival**   (0 = No; 1 = Yes)

**pclass**   Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)

**name**

**sex**

**age**

**sibsp**   Number of Siblings/Spouses Aboard

**parch**   Number of Parents/Children Aboard

**ticket**   Ticket Number

**fare**

**cabin**

**embarked**   Port of Embarkation

  (C = Cherbourg; Q = Queenstown; S = Southampton)

pip install numpy scikit-learn pandas matplotlib

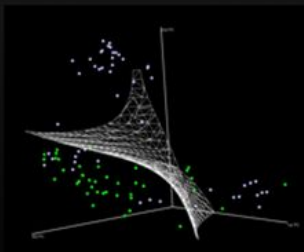pip install "ipython[notebook]"

OR

Use Anaconda distribution

Don't Overdo it:

| Title | Author | About |
|-------|--------|-------|
| Learning From Data | Yaser S. Abu-Mostafa | Small, good for beginners and has an online course |
| Machine Learning: A Probabilistic Perspective | Kevin P. Murphy | Larger, current, and very popular |
| The Elements of Statistical Learning: Data Mining, Inference, and Prediction | Trevor Hastie, Robert Tibshirani, Jerome Friedman | A lot of theory, and has a free PDF edition |

| Title | Author | Website |
| --- | --- | --- |
| Machine Learning | Andy Ng | Coursera.org |
| Intro to Machine Learning | Sebastian Thrun | Udacity.com |
| DataQuest courses | | DataQuest.io |
| Kaggle Competitions and Tutorials | | Kaggle.com |