# Lecture Notes for
# **Machine Learning in Python**

## Professor Eric Larson
## **Numpy, Pandas, Document Features**

# Class Logistics and Agenda

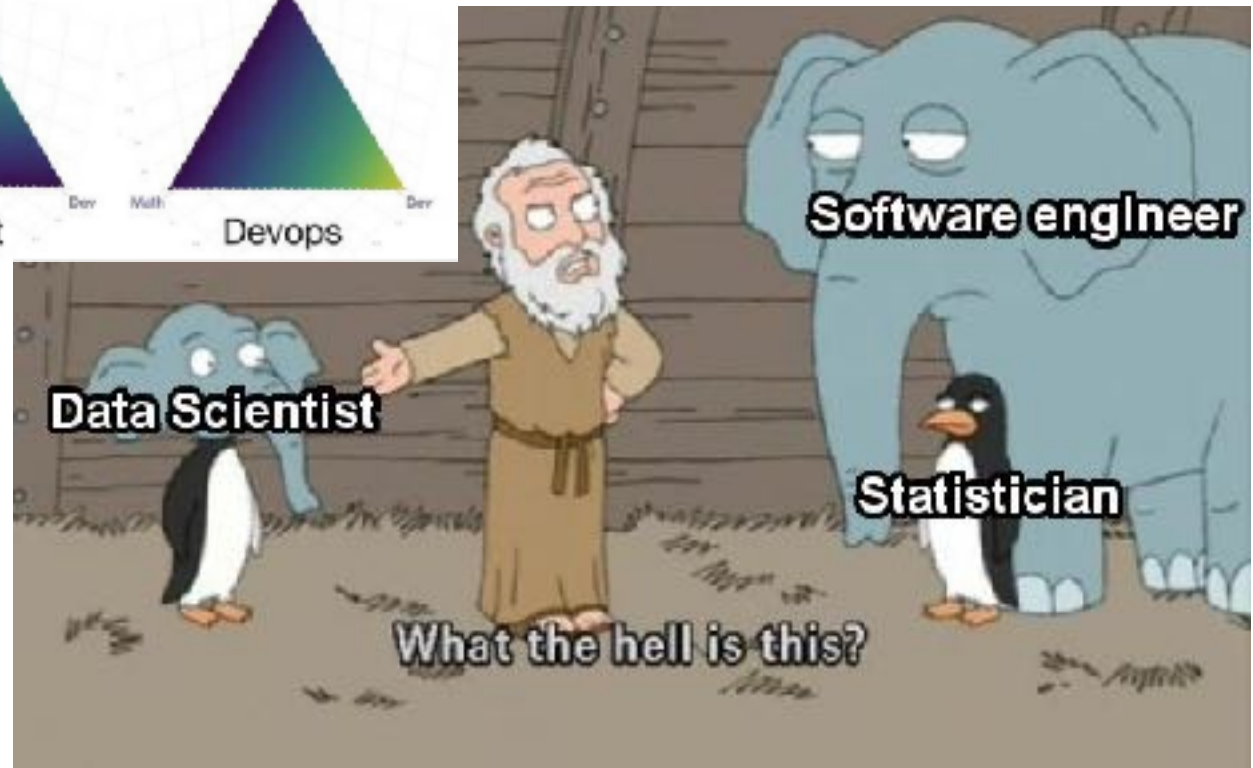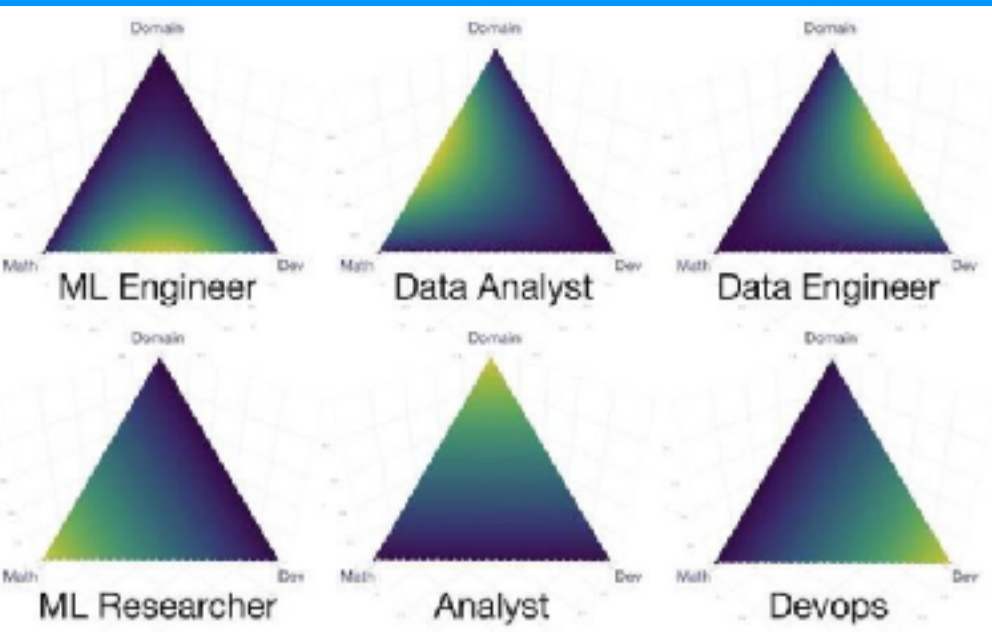- Canvas? Anaconda Installs?
- Red/Blue/Distance
- Agenda:
    - Finish Numpy
    - Data Quality
    - Attributes Representation
        - documents
    - The Pandas eco-system
        - loading and manipulating attributes
- Needing some more help?
    - **fast.ai** has great, free resources

"Finish"
Jupyter Notebooks
and Numpy

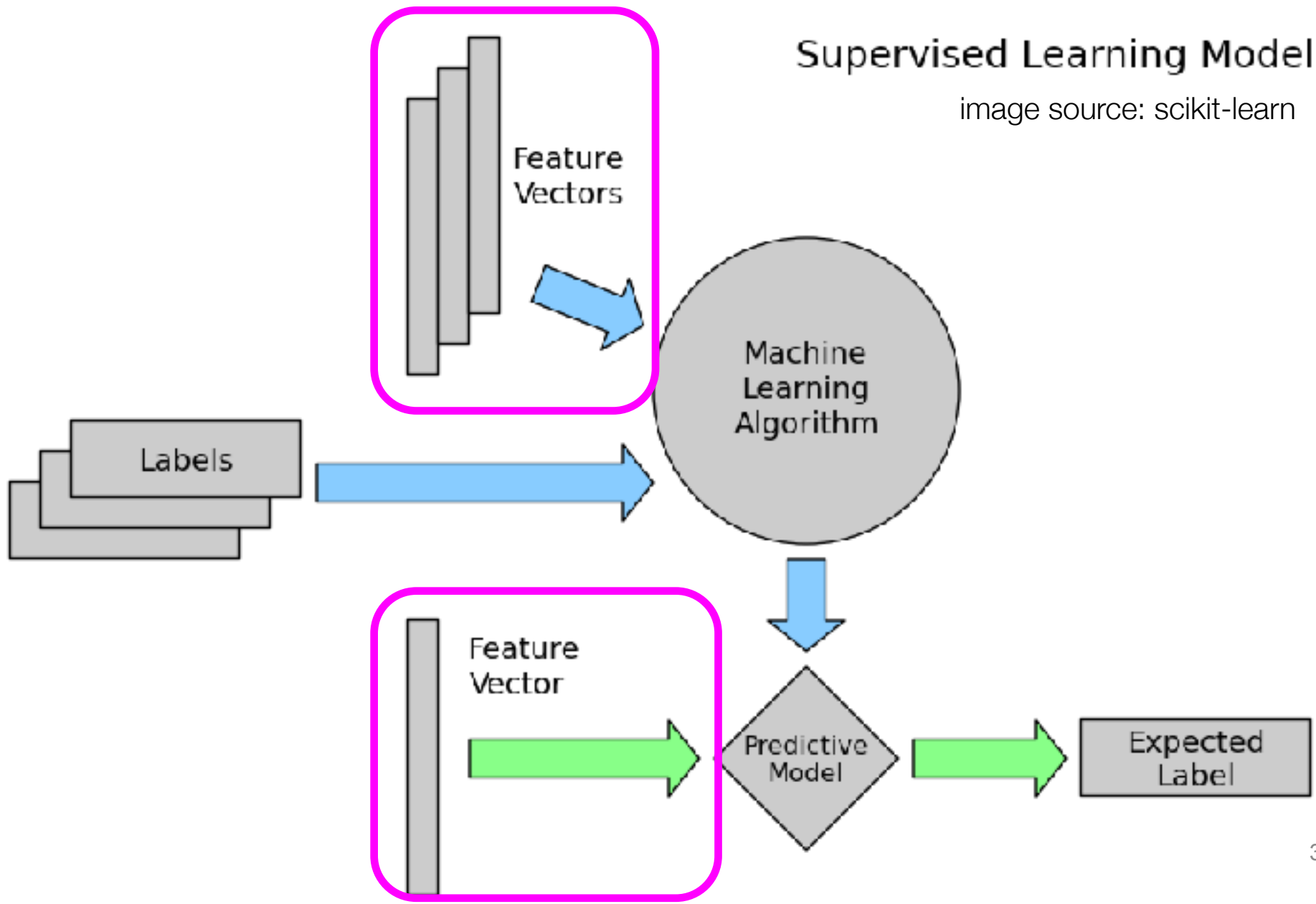`01_Numpy and Pandas Intro.ipynb`

29

# Data Quality

Supervised Learning Model

image source: scikit-learn

# Data Quality Problems

- Missing
  - Easy to find, NaNs
- Duplicated
  - Easy to find, hard to verify
- Noise or Outlier
  - Hard to define
  - Hard to catch

Information is not collected (e.g., people decline to give their age and weight)

Features **not applicable** (e.g., annual income for children)

**UCI ML Repository**: 90% of repositories have missing data

| TID | Hair Color | Height | Age | Arrested |
|-----|------------|--------|-----|----------|
| 1 | Brown | 5'2" | 23 | no |
| 2 | Hazel | 1.5m | 12 | no |
| 3 | Bl | 5 | 999 | no |
| 4 | Brown | 5'2" | 23 | no |

# Handling Issues with Data Quality

☑ **Eliminate** Instance or Feature

☑ **Ignore** the Missing Value During Analysis Replace with all possible values (talk about later)

☑ **Impute** Missing Values  How?
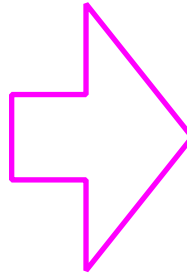
Stats?
mean
median
mode

# Imputation

- When is it probably fine to impute missing data:
  - (A) When there is not much missing data
  - (B) When the missing feature is mostly predictable from another feature
  - (C) When there is not much missing data for each subgroup of the data
  - (D) When it is the class you want to predict

# Split-Impute-Combine

| TID | Pregnant | BMI | Age | Diabetes |
|-----|----------|------|-------|----------|
| 1 | Y | 33.6 | 41-50 | positive |
| 2 | N | 26.6 | 31-40 | negative |
| 3 | Y | 23.3 | ? | positive |
| 4 | N | 28.1 | 21-30 | negative |
| 5 | N | 43.1 | 31-40 | positive |
| 6 | Y | 25.6 | 21-30 | negative |
| 7 | Y | 31.0 | 21-30 | positive |
| 8 | Y | 35.3 | ? | negative |
| 9 | N | 30.5 | 51-60 | positive |
| 10 | Y | 37.6 | 51-60 | positive |

split: pregnant
split: BMI > 32

| TID | Pregnant | BMI | Age | Diabetes |
|-----|----------|------|-------|----------|
| 1 | Y | >32 | 41-50 | positive |
| 8 | Y | >32 | ? | negative |
| 10 | Y | >32 | 51-60 | positive |

Mode: none, can't impute

| TID | Pregnant | BMI | Age | Diabetes |
|-----|----------|------|-------|----------|
| 3 | Y | <32 | ? | positive |
| 6 | Y | <32 | 21-30 | negative |
| 7 | Y | <32 | 21-30 | positive |

Mode: 21-30

# K-Nearest Neighbors Imputation

For K=3, find 3 closest neighbors

| TID | Pregnant | BMI | Age | Diabetes | Distance |
|-----|----------|------|-------|----------|----------|
| 3 | Y | 23.3 | ? | positive | 0 |
| 6 | Y | 25.6 | 21-30 | negative | (0 + 2.3 + 1)/3 |
| 2 | N | 26.6 | 31-40 | negative | (1 + 3.3 + 1)/3 |
| 4 | ? | 28.1 | 21-30 | negative | (4.8 + 1)/2 |

**Imputed Age:** 21-30

**How to calculate distance?**
- Difference for valid features only
- May need to normalize ranges
- Or weight neighbors differently
- Or have min # of valid features
- Euclidean, city-block, etc.

| TID | Pregnant | BMI | Age | Diabetes |
|-----|----------|------|-------|----------|
| 1 | Y | 33.6 | 41-50 | positive |
| 2 | N | 26.6 | 31-40 | negative |
| 3 | Y | 23.3 | ? | positive |
| 4 | ? | 28.1 | 21-30 | negative |
| 5 | N | 43.1 | 31-40 | positive |
| 6 | Y | 25.6 | 21-30 | negative |
| 7 | Y | 31.0 | 21-30 | positive |
| 8 | Y | 35.3 | ? | negative |
| 9 | N | 30.5 | 51-60 | positive |
| 10 | Y | 37.6 | 51-60 | positive |