

REPRESENTATION

Data comes in in various forms - need to extract feature vectors. This is feature engineering and takes ~70% of your time.

Examples:

vocabulary { "a", "b", "c" }
 ↓ ↓ ↓
 0 1 2
binary vector < $\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \rangle$ \downarrow
for cat data

(too inefficient? use sparse rep)

What are good features?

- Clear and obvious meaning
- Mostly non-zero in the dataset
- No "magic values" (e.g. -1 \rightarrow bool + int)
- Feature values shouldn't change over time
- Should not have extreme outliers (e.g. upstream instability)
 \rightarrow Cap or transform or scale (e.g. $\log(-)$ $\frac{off}{off} = 1000$)

Binning trick: map continuous \rightarrow categorical

e.g. latitude not linearly related to home price, but mapping to city, nbhd, etc. might

★ KNOW YOUR DATA!

- visualize (e.g. histograms, etc.)
- debug (dups? missing vals? outliers?)
- monitor over time (test/training/val sets?)

(Multi-hot: can have multiple present)

One-hot encoding: map each unique string to an integer (a "coefficient").

This is sparse categorical data.

(BUT REPRESENT w/ BOOL VEC!)

Scrubbing incoming (possibly bad) data by writing a program.

- Remove bad examples.
 - Omitted values
 - Duplicates
 - Bad labels or values (tricky to detect)
- Detect aggregate problems w/stats + viz.
 - Histograms
 - Max/min, stdev, etc.
 - List of most common discrete vals

Rules

1. Know how you expect your data to look
2. Verify that the data meets these expectations
3. Check training data agrees w/other sources (e.g. dashboards)

~~*~~ GOOD ML RELIES ON GOOD DATA!