CS 307

F ALL 2023

D ALPIAZ

Week 12

# Some Practical Considerations

# Data Ethics

Just because you can,
Does not mean you should!

"IT'S TOUGH TO MAKE PREDICTIONS,

ESPECIALLY ABOUT THE FUTURE"

— YOGI BERRA

# PREDICT TIME VERSUS TRAIN TIME

- DO YOU EXPECT 'NEW' DATA TO COME FROM THE SAME DISTRIBUTION AS THE TRAIN DATA?

  FORECASTING

- WHAT FEATURES WILL BE AVAILABLE IN 'NEW' DATA WHEN YOU NEED TO MAKE A PREDICTION?

# Missing Data

- Why is it missing?

  - Survey non-response

  - Formatting

  - Measurement malfunction

  - Informative

# Missing Data

- Remove rows with missingness
- Remove columns with high missingness
- Replace with a constant
- Impute with mean/median/mode of column
- Impute with consideration for other columns
- Using missingness as information

# "Look at the data"

- WHAT IS THE SOURCE? HOW WAS COLLECTED / GENERATED?
  ↳ WHAT FEATURES ARE DOMAIN RELEVANT.

- $n$? $p$? SIZE?

- PRINT DATA FRAME
    - head( )         "glimpse"
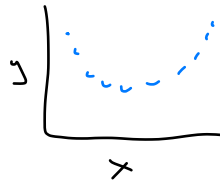    - tail( )

- VARIABLE DATA TYPES                                    PAIRS PLOT

- SUMMARY STATISTICS / VISUALIZATIONS

- HOW MANY MISSING? WHERE MISSING?

CONDITIONAL

# No FREE LUNCH

## Just use RF!

- LINEAR MODEL
- LINEAR MODEL w/ REGULARIZATION

- KAGGLE $\implies$ BOOSTING

# One-Hot vs Dummy vs Ordinal

| X | | ORD | | One-Hot | | | Dummy | |
|---|---|-----|---|---------|---|---|-------|---|
| | | | | $X_A$ | $X_B$ | $X_C$ | $X_B$ | $X_C$ |
| A | | 0 | | 1 | 0 | 0 | 0 | 0 |
| B | | 1 | | 0 | 1 | 0 | 1 | 0 |
| C | | 0 | | 0 | 0 | 1 | 0 | 1 |
| A | | 1 | | 1 | 0 | 0 | 0 | 0 |
| B | | 1 | | 0 | 1 | 0 | 1 | 0 |
| B | | 2 | | 0 | 0 | 0 | 1 | 1 |
| C | | 2 | | 0 | 0 | 1 | 0 | 1 |
| A | | 0 | | 1 | 0 | 0 | 0 | 0 |

# Feature Engineering

- ONE-HOT
- SCALING
- Count Vectorization
- Manual