# DSND: Introduction To Data Science

→ Lesson 1: The Data Sci. process

How to approach DS problems

→ CRISP-DM

cross industry standard process 4 data mining.

① Recap business understanding
  ↳ understanding problems & pain points

② Data understanding
  ↳ what data will you need to solve your problem?

  ↳ note that ① & ② can happen in either order

③ Prepare data
  ↳ Collect, access data
  ↳ wrangling & feature engineering where possible/needed

④ Data Modeling
    ↳ Perform statistical analysis
    ↳ Answer questions
⑤ Evaluate results
⑥ Deploy / Iterate

what is required
- Curiousity
- the right data
- the right tools

## Modeling

→ instantiate, fit , predict, score
    · Quantitative vs Categorical vars.
Looking at variables
    ↳ correlation matrices, heatmaps, ⎤ Choose
       matrix plots etc ⎦ features
Objective: try to capture as much signal
       in your data as possible

→ Handling missing values

① Drop rows
    ↳ problem, increases bias

② Impute values
    ↳ removes variability & predictive
    power of that column

→ Categorical variables

    ↳ ML models need numerical values
    → one hot encoding
        ↳ each choice in a categorical
        variable becomes a dummy variable
    ↳ Always remember to drop one column

Problems: does not scale well w/ many
        categorical variables.

→ Deploying Model

    ↳ Automate tasks
    ↳ Communicate for others to take
    action
            2) Dashboards
            ↳ etc.

→ Lesson 2: Communication w/ Stakeholders
  ↳ value of a project can be cut short

→ Github
  ↳ Repository for code

  - Covered by README:
        · Motivation
        · versions & installs
        · file desc.
        · Liscensing
→ Medium
  ↳ communicate project w/ other
     people ♂

→ communication
  ↳ Who is the audience?
        ↳ technical vs accessible
        ↳ consider prior knowledge
  ① "Pull in" reader
  ② keep engaged w/ story
  ③ end w/ summary & call to action

① "Pull in"
   ↳ compelling image boosts engagement

  · relevance
     · to others
     · to current events

② keep engaged
  · use paragraphs!
  · pictures & other whitespace
    for pacing

③ Call to action
  · reiterate main points
  · call to action
     ↳ Makes clear how reader
       should act afterward

→ Pick a dataset

→ Pose 3 questions you want answered

→ Analyze data using python
    ↳ Munge
    ↳ Visualize
    ↳ Model
    ↳ etc.

→ Communicate insights

  Deliverable: — Github repo
                   +
             — Medium Blog Post.

Notes:
    — Stick to CRISP DM