

Scoping an Actionable ML Project

Rayid Ghani and Kit Rodolfa

Carnegie Mellon University



Scope

- Goals, Actions, Data, Analysis, Ethics



Data

- Get Data
- Store Data
- Link Data



Exploration

- Entities
- temporal
- Spatial
- ...



Modeling

- Rows
- Labels
- Features
- Models



Model Selection

- Train-Test Splits
- Performance Metrics



Model Interpretation



Dealing with Bias and Fairness



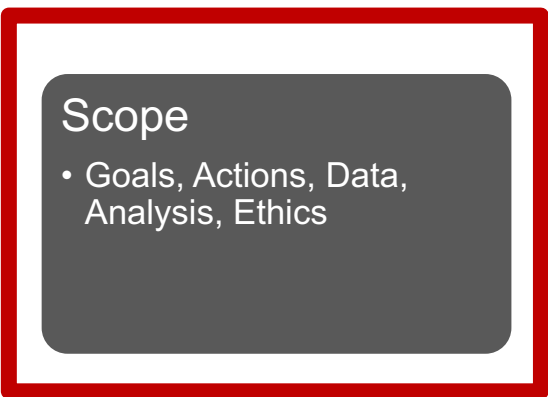
Field Trial Design



Deployment



Monitoring



Why Scoping is Critical

- (Unfortunately) projects do not come as well-defined problems
- A well-scoped project increases the likelihood of your work being used and having an impact
- Shifts focus from “I have some data, what can I do with it” to starting with the problem, informing actions, and improving outcomes

~~The goal of this project is to build a model...~~

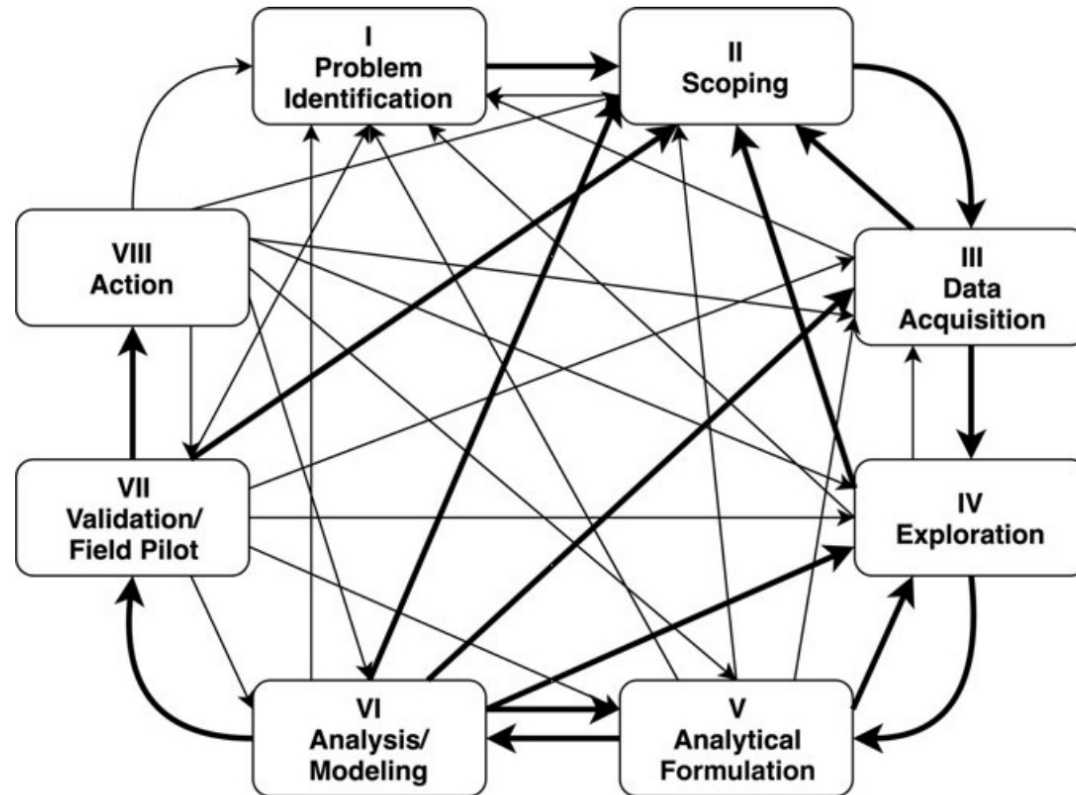
- Doing some type of analysis may be the bulk of what you do during a project but **the analysis is rarely the goal of the project.**
- A well-formulated/scoped project will have the analysis **inform a follow-up action** and **help achieve intended policy/business/research goals**

Before Scoping: Initial Screening Criteria

- Real and significant problem (with clear social impact)
- Priority for the organization
- Commitments in place:
 - Access to the data and to people who understand the data
 - Access to the people who understand the problem
 - Commitment to validate and take actions informed by your work

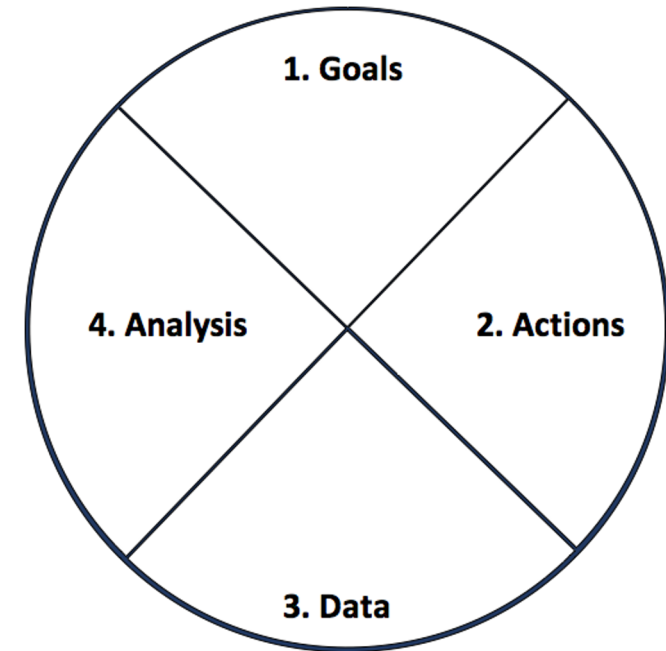
Scoping is an iterative and ongoing process

- Involves different types of people
- Gets refined over time
- Gets modified as a result of later project phases



Actionable and Goal-Driven Project Scope

- 1. Goals:** Define the goal(s) of the project
- 2. Actions:** What actions/interventions will you inform?
- 3. Data:** What data do you have internally? What data do you need? What can you augment from external and public sources?
- 4. Analysis:** What analysis needs to be done? How will it be validated?



Step 1: Determine Goals

If you are successful doing this project, how would the world change?

- Will the process you're trying to improve become more **efficient**?
- Will it be more **effective**?
- Will it be more **equitable/fair**?
- A combination of the above?

Step 1: Determine Goals

- Goals need to be measurable and concrete
- Goal is NOT to build a model or a map or dashboard, make a prediction, etc.
- What are the relative priorities and tradeoffs for each goal?
- What constraints do you face in achieving these goals?
- Need to get different “stakeholders” involved up front

Step 2: Identify Actions to achieve the goal

- What interventions/actions do I have access to?
- What would someone do differently if they had more information or knew where their actions were most likely to be effective?
- Informing these actions:
 - Who? (to target for each action)
 - What? (to say to them)
 - How? (to use different communication channels)

Step 2: Identify Actions to achieve the goal

- Focus on concrete actions
- Existing vs new actions
- Consider the granularity of the actions
 - e.g. students who need help generally vs specific program
- How frequently are interventions taken/planned?
- How far out does planning occur?

Step 3: Data Sources

- What relevant data sources do you have?
- What data do you need?
 - Important to match the granularity, frequency, and time horizon of the actions to the data
- What external data can you augment this with?

Step 3: Data Sources

Types of Data

- Program Level
- Transactional
- Spatial
- Text
- Images/Audio/Video

- Nobody knows what data the entire organization has
- Don't get intimidated by legal acronyms thrown at you
- Data is never perfect – is it useful enough to improve over status quo?

Step 3: Data Sources

- How reliable is the data?
- How current is it?
- How much of it is computer-readable?
- How much of it is stored as notes, audio, photos, videos?
- What resources and authority do you have to collect more?

Step 4: Analysis

- What analysis needs to be done?
- What type of methods should be used?
- How will the analysis be validated?

Types of Analysis Capabilities

- Description (Understand the past)
- Detection (Anomalies, Events, Patterns)
- Prediction (Predict the Future)
- Optimization
- Behavior Change (Causal Inference)

Validation and Implementation Plan

- Go back to the metrics and goals defined at the beginning of the project
- Run a Pilot/Field Trial
- Deploy
- Set up Infrastructure and allocate resources to monitor “lift”

Creating a more equitable society

- How do we **define** equity?
- How do we **detect** inequity?
- How do we **increase** equity?

Data and AI Ethics Issues

Privacy

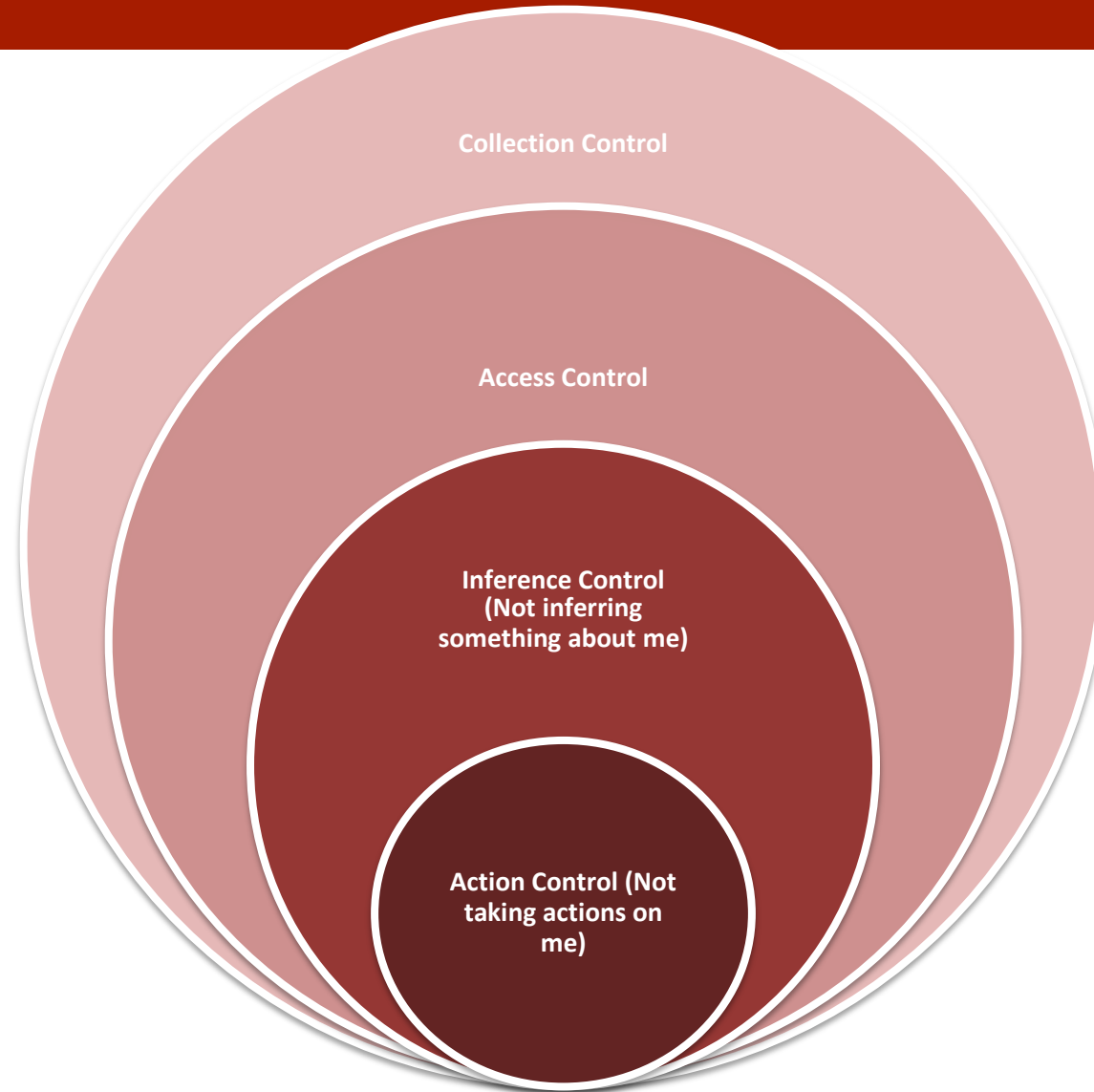
Data Ownership

Bias, Equity, & Fairness

Transparency

Trustworthiness and
Accountability

Levels of control



Data Ethics Questions

- Are you using data for purposes it's intended for?
- How are you protecting the data?
- Do the people who “own” the data know you're using it?
- Do you have their permission? How was it obtained?
- What actions are you taking on individuals based on this data?
- Do the people you're targeting know why and if they're being targeted?
- What recourse do they have?
- Would it make the front page of the national newspaper if they found out what you're doing?

A Few Things to Remember

- Don't be afraid to ask naïve questions
- Spend time discussing goals and metrics – don't forget equity as a goal
- Understand what the current process/solution is
- Communication is critical – before, during, and after
- We need to make sure that we tackle these problems responsibly and ethically
- Data and ML does not solve problems, people do. Is what you're doing helping solve the problem?

Project Scoping Worksheet

<http://bit.ly/dataprojectscoping>

Reminders

- Due Today
 - Team Assignments
 - Make sure you've joined the class slack and filled out the survey
- Next Tuesday
 - Submit weekly review (assignment) before class
 - All the readings for Tuesday are optional