

Project Scoping Worksheet

This worksheet is designed to scope actionable analytics and data science projects for organizations (businesses, government agencies, non-profits, social enterprises, and others).

1. Project Name:

2. Organization Name:

3. Names and Titles of Main Organizational Contacts:

4. Problem Description:

- **What is the problem you are facing?**
- **Who/what is affected by this problem?**
- **How many people/organizations/places/etc. and how much are they affected?**
- **Why is solving this problem a priority for your organization?**

5. Goals: What business goals will be accomplished by solving this problem and what constraints do you have? (In order of priority)

- The technical solution that will be built (e.g., predictive model or dashboard or map) is not the business goal - that is the tool that will achieve your goal
- The goal should be specific and measurable
- Achieving the goal should help solve the problem you're tackling
- Typical goals include improving/maximizing/increasing or decreasing/mitigating/reducing some outcome or metric.
- Typical constraints include budget, lack of human capital, legal restrictions, etc.
- Consider tradeoffs between conflicting goals.
- What is the value of a perfect solution? Is it worth the effort required?

	Goal	Constraints
1		
2		
3		

6. Data

- The data has to connect to the goals.
- Typical data science projects use administrative data as the primary data source, and enhance it with publicly available data sources (Census, other open data).

- What data sources do you have internally?

	Data Source 1	Data Source 2	Data Source 3
Name e.g., Hospital Admissions database			
What does it contain? Describe the attributes included in the data source. <i>e.g., admission and discharge records, patient sociodemographic data, insurance, physician, etc.</i>			
What level of granularity? <i>e.g., transaction, person, organization, location</i>			
How frequently is it collected/updated after it's captured? <i>e.g., real time, daily, weekly, monthly, yearly, one off</i>			
Does it have reliable and unique identifiers that can be linked to other data sources? <i>e.g., SSN, National identifier, patient identifier, insurance number, etc.</i>			
Who's the internal owner of the data? <i>e.g., IT department, individual, etc.?</i>			
How is it stored? <i>e.g., in a database, pdfs, excel, SPSS</i>			
Additional comments			

- What data can you get from external, private or public sources?

	Data Source 1	Data Source 2	Data Source 3
Name <i>e.g., Air Quality database</i>			
What does it contain? Describe the attributes included in the data source. <i>e.g., distinct pollution's particle concentration</i>			
What level of granularity? <i>e.g., geolocalized hourly sensor data</i>			
How frequently is it collected/updated after it's captured? <i>e.g., daily</i>			
Does it have unique identifiers that can be linked to other data sources? <i>e.g., sensor identifier</i>			
Who's the internal owner of the data? <i>e.g., NOAA</i>			
How is it stored? <i>e.g., API endpoint from an open data portal</i>			
Additional comments			

- In an ideal world, is there additional data would you want to get/gather that would be relevant to his problem? (surveys, phone records, etc.)

7. Analysis

- Typical data science projects include a combination of analysis, typically including description, detection, prediction, optimization, and/or behavior change.
- Again, the analysis is not the goal of the project - the **analysis** helps you use the **data** you have to inform the **actions** you have access to in order to achieve your **goals**.
- Choose the right set of analyses for each problem
- You must validate the analysis and use a validation process that matches how your analysis will be used in practice

	Analysis 1:	Analysis 2:	Analysis 3:
Analysis type <i>e.g., Description, Prediction, Detection, Behavior Change</i>			
Purpose of the analysis <i>e.g., understand historical behavior of individuals, probability of purchase, identify at risk for failure, etc.</i>			
How will you validate this analysis using existing data? What methodology and what metrics will you use? <i>e.g., Validate on future data using RMSE</i>			

8. Ethical considerations

Privacy Are you working with personal and/or sensitive data that is individually identifiable? Mention them.	
Transparency Which stakeholders should know about which parts of the project? <i>Stakeholders typically include leaders, managers, frontline workers, people who will be affected by the actions, etc</i>	
Discrimination/Equity Are there any specific groups for whom you want to ensure equity of outcomes? <i>e.g., groups of interest defined by gender, age, localization, social class, educational level, urban/rural, ethnicity</i>	
Social License If the country's entire population finds out about your project, will they be ok with it? Why?	
Accountability Who are the people responsible for all the things above?	
Any other considerations such as consent, legal, etc.	

9. Who are the external organizations and internal departments that will need to be involved?

(Typically, data science projects need involvement from data owners, IT infrastructure owners, problem owner, analytics people)

Organization/Department	Description of desired involvement	Name/role of counterpart
<i>IT department</i>	<i>Provide Data infrastructure</i>	<i>Head of IT department</i>
<i>Statistics agency</i>	<i>Provide population data</i>	<i>Head of Department of Statistics</i>

10. Can you develop an experiment to validate the project in real life?

- Define how you will measure the success of the project in real life.
- Describe how the model will be tested.
- Define the duration of the experiment.

This worksheet was modified from [one](#) provided by the Center for Data Science and Public Policy at Carnegie Mellon University. This version of the worksheet has been modified for use by students at Miami University's Farmer School of Business.