

Data Science Project Checklist

Created by Jake Strasler

Last updated: January 10, 2024

This checklist is a print-friendly, interactable translation of [DataCamp's Data Science Project Checklist](#). DataCamp describes the checklist as “summarizing data science project management best practices...” and “...combines CRISP-DM project management principles with those of the Agile and Scrum software development frameworks.”

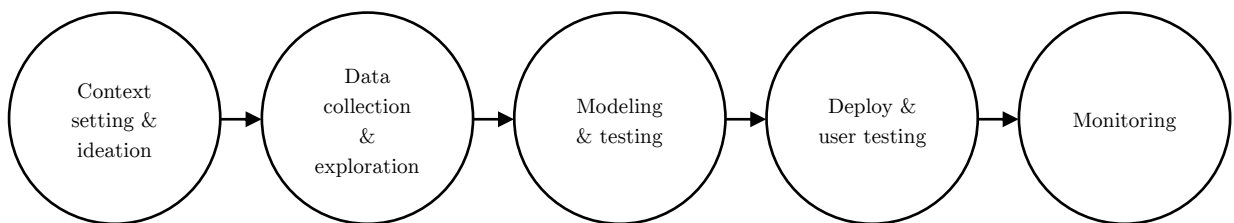
The following characteristics should underpin any data science project:

- **Measurable:** Is the success of a project and its impact on the business quantifiable?
- **Reliable:** What proportion of projects achieved their goals?
- **Scalable:** Can the throughput of projects be increase without significantly degrading reliability?

Additionally, the following principles should be abided by:

- **Iterative**
- **Reusable and recyclable**
- **Reproducible**

The framework followed will be:



Project Name: _____

Project Start Date: _____

Notes: _____

Context Setting & Ideation

Task	Notes
<input type="checkbox"/> Identify the business problem being solved. <i>Clarify why the project is being undertaken as unambiguously as possible.</i>	
<input type="checkbox"/> Identify stakeholders. <i>Roles may include project manager, data scientist, account manager, etc.</i>	
<input type="checkbox"/> Review prior work. <i>Review existing projects that covered similar ground. What were key outcomes? Can anything be reused? What mistakes should be avoided?</i>	
<input type="checkbox"/> Determine key performance indicators to measure success. <i>Metrics should use SMART criteria.</i>	
<input type="checkbox"/> Determine the scope. <i>What are the deliverables? What are the reqs. for the deliverables? What will not be included?</i>	
<input type="checkbox"/> Write a project plan. <i>Set milestones for intermediate steps in the project. Decide on timeline for each milestone. Write a short description of each step.</i>	<p style="text-align: center;">Do external to this document.</p>

<input type="checkbox"/>	<p>Estimate the impact of the project.</p> <p><i>Quantify the benefit to the organization if the goals are reached. List qualitative benefits.</i></p>	
<input type="checkbox"/>	<p>Estimate the effort of the project.</p> <p><i>How much will the project cost? How much time will the project take? What resources will it need?</i></p>	
<input type="checkbox"/>	<p>Estimate the project risks.</p> <p><i>List all the risks to the project. Calculate the risk impact the probability of it occurring times severity of it occurring.</i></p>	
<input type="checkbox"/>	<p>Decide whether or not to proceed with the project.</p> <p><i>Decide to proceed now, put on hold, or cancel.</i></p>	
<input type="checkbox"/>	<p>Determine the responsibility of each stakeholder.</p> <p><i>Use the RACI model for each task.</i></p>	
<input type="checkbox"/>	<p>Determine a communication strategy.</p> <p><i>How will you keep in touch? Meeting cadence?</i></p>	
<input type="checkbox"/>	<p>Identify data sources.</p> <p><i>Do you have access? Where is it stored? What form is it in? Etc.</i></p>	
<input type="checkbox"/>	<p>Anticipate regulatory needs.</p>	

<i>Will anything be audited?</i> <i>Can data sources be legally used?</i>	
<input type="checkbox"/> Decide on a technology stack. <i>Agree on tools for storing, processing, and modeling the data.</i>	
<input type="checkbox"/> Write a project charter. <i>Summarize what was decided in a short document.</i>	Do external to this document.

Data Collection & Exploration

Task	Notes
<input type="checkbox"/> Give data scientists access to all datasets. <i>Organize appropriate permissions. Purchase any commercial datasets or use synthetic data with similar properties.</i>	
<input type="checkbox"/> Ingest the data. <i>For each data source, move it to the analytics environment.</i>	
<input type="checkbox"/> Explore the data. <i>Visualize distribution of each variable with histogram or bar plot. Quantify missing values for each variable. Visualize relationship between features and target variable with scatter, histogram, box, or heatmap plots.</i>	
<input type="checkbox"/> Determine key performance indicators to measure success. <i>For each dataset: provide a summary, describe high-level data quality issues, describe the quality of target variable, describe quality of each feature, describe relationship between each feature and the target variable.</i>	

<input type="checkbox"/>	<p>Decide whether or not to proceed with the project. <i>Continue, pause, or cancel.</i></p>	
<input type="checkbox"/>	<p>Build a data pipeline. <i>Pipeline should automate the ingestion and cleaning process. Should run as batch or stream job.</i></p>	
<input type="checkbox"/>	<p>Document the data pipeline. <i>Draw diagram of steps in the data pipeline and their dependences. Describe what happens in each step.</i></p>	<p>Do external to this document.</p>

Modeling & Testing

Task	Notes
<i>Modeling</i>	
<input type="checkbox"/> <p>Generate a hypothesis. <i>Does the hypothesis make sense in the business domain? Can you measure the outcome? Do you have sufficient data to be significant? Watch for biases.</i></p>	
<input type="checkbox"/> <p>Split your data into training and testing sets. <i>Do this before engineering features to ensure no data leakage.</i></p>	
<input type="checkbox"/> <p>Engineer features. <i>Center/scale numeric values, create categorical variables (incl. binning), apply Box-Cox or Yeo-Johnston transformations to numeric variables to give normal dist., combine rare/related categories of categorical variables, extract or combine parts of datetimes, create new variables from summary statistics, extract quantitative metrics from text and other unstructured data.</i></p>	
<input type="checkbox"/> <p>Fit the model/run an experiment. <i>Fit simplest first then increase complexity.</i></p>	

<input type="checkbox"/>	<p>Report on the results. <i>Regularly provide feedback to stakeholders, adjust language accordingly, report failures as well as successes.</i></p>	
<p><i>Testing</i></p>		
<input type="checkbox"/>	<p>Create a test suite. <i>Unit tests or back test that run automatically to check your model or experiment's performance.</i></p>	
<input type="checkbox"/>	<p>Validate the business impact. <i>Discuss the impact with business stakeholders.</i></p>	
<input type="checkbox"/>	<p>Validate the technical approach. <i>Check assumptions of your model to ensure they're valid. See if results are sensitive to sampled data. Are hyperparameters suitable? Is the model reproducible?</i></p>	
<input type="checkbox"/>	<p>Validate the deployability. <i>Can all input values or use cases be handled? Are all required data sources available in production? Can model fail gracefully? Can predictions be made fast enough?</i></p>	
<input type="checkbox"/>	<p>Preserve null results. <i>Record anything that won't make it to production.</i></p>	

Deployment & User Testing

Task	Notes
<i>Deployment</i>	
<input type="checkbox"/> Develop a data pipeline. <i>Set up DAG for all the data sources to production environment.</i>	
<input type="checkbox"/> Develop a model pipeline. <i>Provide an API to your model that can be accessed by dashboard/website/etc.</i>	
<input type="checkbox"/> Design a monitoring plan. <i>Determine metrics to be tracked, including performance and safety metrics. Determine limits for acceptable ranges for those metrics. Decide how you wish to be alerted if metrics go out of range.</i>	
<input type="checkbox"/> Roll out via A/B test. <i>Provide the new feature/model to random sample of users. Monitor chosen metrics. Declare winner when there is statistical significance.</i>	
<input type="checkbox"/> Analyze and report on A/B test results. <i>Compare the metrics chosen to track for each group. Report results even if test was not a success.</i>	
<input type="checkbox"/> Roll out to most or all users. <i>If successful, rollout to all but small holdout group.</i>	

<i>User Testing</i>	
<input type="checkbox"/> <p>Write an exit report. <i>Summarize status of project and what you learned. Provide overview of the project. Summarize business problem you tried to solve. Describe the data sources and how they were processed. Describe modeling techniques used and their validation. Summarize solution architecture. Outline benefits from the project for company and customer. Describe any learnings around project execution, data science, business domain, and the product. Outline next steps.</i></p>	<p>Do external to this document.</p>
<input type="checkbox"/> <p>Get customer feedback. <i>Conduct surveys and user interviews. Monitor reviews, ratings, and social media.</i></p>	

Monitoring

Task	Notes
<input type="checkbox"/> Develop monitoring pipeline. <i>Set up pipeline to automatically track the performance and safety metrics defined in the monitoring plan.</i>	
<input type="checkbox"/> Create dashboards. <i>Create dashboards to track the changes of metrics over time.</i>	
<input type="checkbox"/> Set up alerts. <i>Set up notifications via email, Slack, etc. when metrics fall out of the acceptable range.</i>	

Project End Date: _____

Time Elapsed: _____