

Tools for Data Science

The Data Strategy Canvas

The data strategy canvas is a tool that is used in order to structure the implementation of a data science within a company. You will find this canvas especially useful if you are implementing data science for the first time within your company.

The canvas is covering the most important topics that you have read in this book:

- 1) Challenge—Here you have to specify the exact problem you are trying to solve. Don't forget to pose the problem as a question. This forces you to make the problem more concrete and will also help the data scientist better understand what it is that you are trying to do.
- 2) Data sources—Where are you getting your data from?
- 3) Appropriateness, nature, time, and cost—These are the four data considerations that were discussed in Chapter 2 about data management.

- 4) **Method**—While a data scientist would be able to better advise you on this, it is still a useful exercise to try to think how you could solve this problem. It will at least help you make potentially better hiring decisions. However, do not get too stuck on your answer to this question. For example, it is likely that something that might look like a deep learning problem to you can be solved more easily with different techniques. Discuss with the data scientist the best approach, and let them have the final say.
- 5) **Success criteria**—When will you know that the project has finished? The wrong expectations can cause all sorts of problems. If you can't define clear success criteria, then at least think whether you can break down a project into smaller milestones that are more manageable.
- 6) **People**—Use the things you learned from the “Hiring and Managing Data Scientists” section in order to figure out which type of data scientist would be better suited to this problem. Also, think whether it is best to work with a contractor or a full-time employee. Many projects are one-off, so contractors can be a good choice.
- 7) **Culture**—Will a data scientist feel comfortable working with you, or will they feel like second-class citizens? You might spend a long time finding the right person, and then they might leave if the company does not provide the right cultural fit.

The data strategy canvas

Challenge
Data sources
Appropriateness (can the data be used to solve the problem?)
Nature (is the data noisy or other issues?)
Time (how fast can you acquire new data?)
Cost (how expensive is the data?)
Method (how will you solve the problem, machine learning, deep learning, statistical modeling, etc.?)
Success criteria
People (whom should you hire?)
Culture (do you have the right culture for data scientists to come and work for you?)

The Data Science Project Assessment Questionnaire

The purpose of this questionnaire is to help you clarify your objectives and factors that can affect the success of a data science project.

It is important to understand that a data science project is often an exercise in risk management. The problematic part in data science is the “science” bit. Quite often, it is impossible to know how well something will work in advance, until you try things out. The only cases where you can be confident of results before you embark on a project are

- 1) You have worked on the same case with very similar data in the past.
- 2) There is an extensive body of literature on this type of problem.

In any other case, the results cannot be guaranteed. Hence, the right project structure should focus on organizing an “attack plan,” where each step should produce one of either two outcomes:

- 1) Successful completion of the project.
- 2) Learn something about the problem (e.g., maybe a family of methods is not appropriate, which can lead to improved exploration)

In order to better manage the project, expectations, and risk, the following questionnaire can help you list all the important requirements. This questionnaire is both for you and the data scientist, and it will help you build a common understanding of the requirements, the challenges, and how to mitigate risks. It is recommended that you do at least one round of back-and-forth between you and the data scientist in order to make sure that there is appropriate understanding on both sides.

The Data Science Project Assessment Questionnaire

Success criteria

What would you consider successful implementation of the project? Share your own thoughts.

Are there any benchmarks in performance? Please try providing a numerical answer (e.g., “anything above 60% accuracy is good, based on a benchmark X”).

Risk factors

Is the data of appropriate quality? If no, what are the issues?

Is the data of appropriate size? If no, how much do you think this will affect performance? How will you mitigate that risk?

If the data is of high quality and size, then what could possibly prevent a good model from being built (e.g., maybe the domain is particularly difficult)?

Timelines

How time critical is the project?

If the project is difficult and no approach can reach the desired accuracy, which of the following plans seem the most attractive to you and why?

- Keep on trying more advanced approaches (e.g., ensemble modeling).
- Fix data issues and try again (e.g., collect more data).
- Simply use a model that is “good enough,” even if it doesn’t reach the desired performance.

Best/worst possible case

What would be the best possible outcome for you?

If the goals cannot be achieved (e.g., the model doesn’t work to the desired performance or fails completely), what would you do?

Interview Questions for Data Scientists

Interviewing data scientists is not easy for many reasons. As the book explained in Chapter 10, there is not a single data science curriculum that someone should abide by. You get tribes of data scientists that can have different mentalities and use different tools. This makes interviewing data scientists very challenging. This becomes even more challenging when someone is your first hire.

These questions are designed to help you understand how good someone is.

What are the different types of machine learning? Can you explain them to me?

This is a very simple question. Someone who has spent time reading the subject will know that the main types are supervised learning, unsupervised learning, active learning, semi-supervised learning, and reinforcement learning.

Active learning and semi-supervised learning are not very popular, so maybe not everyone knows about them, but you can give bonus points to someone who does.

What is the difference between descriptive and inferential statistics?

Descriptive statistics refers to the statistics that most people do in high school, like summary metrics and graphs. Inferential statistics deals with the problem of inferring something about a population by only getting a sample out of it. Statistical modeling and hypothesis testing, they are all part of inferential statistics.

Which algorithm would be better in a given problem: Random forest or Naïve Bayes?

This is a trick question. Random forest is a much more successful and popular algorithm than Naïve Bayes. However, the no-free lunch theorem tells us that there is the possibility of simple algorithms being better than more complex ones in a given problem. This relates to the inductive bias that an algorithm might have.

While a full technical exposition of this argument is beyond the scope of this book, you will essentially expect two kinds of answers. Data scientists that lack the proper theoretical background will say “random forest.” When asked why, they might not be able to properly justify their choice or say something like “it’s a good algorithm.” Data scientists that fully understand the theory behind machine learning might say something like “it depends” and will give a more thorough answer, which will demonstrate that they understand that a complex algorithm is not always the right choice.

You have a dataset of 1000 rows and 1500 variables. Which algorithms would be suited to this problem?

This is a particular type of challenge where the variables are more than the rows. These problems can be solved through algorithms that are good in handling a large number of variables. Some choices include

- 1) Elastic net
- 2) Random forests
- 3) Gradient boosted trees
- 4) Heavily regularized neural networks

If someone misses this question, it is not very important, but answering it probably demonstrates solid understanding of some important concepts.

Is more data always good?

This is a trick question. While more data in terms of rows is good, more variables is not always a good thing. There is a concept called “the curse of dimensionality.” According to this, adding more variables can make sometimes a problem more difficult to solve, especially when the variables do not contain much information. Hence, you expect that someone who has the right background to answer that more rows of data are good, but more variables might not be necessarily good and whether they are good depends on the problem and the quality of the data.

Other things to ask about

GitHub repository, or samples of code. Samples of code are actually better indicators of someone’s skill and interest in the area than abstract coding exercises, like the ones you see in most interviews.

Kaggle. While many brilliant data scientists might not really have the interest or time to participate in Kaggle competitions, having a Kaggle account is always a good thing.

A degree from a top university.

Side projects that might relate to your particular challenges.

The New Solution Adoption Questionnaire

It is often easy to get carried away by new technologies. Every 1–2 years, there is some new big name in town, and all the companies are racing to adopt the new tools. However, quite often, those solutions are expensive to implement, might not carry any benefits to you, and might be quite immature and untested. There have been countless businesses that rushed to use Hadoop, NoSQL, blockchain, and other technologies without really needing them. This simple questionnaire should help you understand whether you should adopt a new solution or not.

Step 1

Goal: Understand your objective and what you are trying to achieve. For example, what if you want to use a new kind of database, what are the issues with the current one? Long read times, scaling up, or something else?

Step 2

Enumerate each solution, context it was created in, and pros and cons. Read the white paper if there is one available. For example, I’ve filled in the responses for the Ethereum blockchain in the top row of the following table.

Solution	Goal of the solution	Creator/maintainer	Pros	Cons
Ethereum blockchain	Decentralization, immutability, smart contracts	Ethereum Foundation	Very well tested solution	Blockchain solutions face speed issues (whether this is relevant depends on your particular challenge)

Step 3

Analyze costs and risks. Every migration or adoption of a new technology contains unseen risks. Answer the following questions to map out risks and solutions:

Can your current engineers do this or you need to hire more people?

If you need to hire someone new, how easy is it to find someone, and how much would it cost?

What could go wrong? How can you mitigate the risks?

How much money would you save after the solution would be implemented?

How would your service or product improve by adopting the new technology?

Index

A

- Anomaly detection, 85
- Artificial intelligence (AI), 4
 - automated planning, 7
 - general, 21
 - history, 5
 - research, 10, 11
 - vision, 6
 - winters, 9, 10
- Automated planning, 7–8

B

- Bayesian vs. frequentist statistics, 74, 75
- Big five personality traits, 86
- Booking.com case study, 140, 141

C

- Causal inference, 102–103
- Cognitive science, 5, 6, 11
- Computational intelligence (CI), 17, 18
- Computer scientists, 116, 117

D

- Dark data, 138
- Data acquisition
 - choosing data, 25
 - cost, 26

- nature of data, 25
- problems, 26, 27
- time requirement, 26

- Data-Centric, organizational level
 - data collection, 130
 - data, not at core, 131, 132
 - data privacy, 131
 - poor communication, 132

- Data collection
 - B2C apps, 32, 33
 - finance, 34
 - retail, 33
 - sales, 33
 - social media, 34
 - sports, 34
 - types, 24

- Data management, 34
 - bad practices
 - case study, 38, 39
 - connection, 36
 - documentation/data standard, 36
 - lack of objective, 36, 37
 - objective, 36–38
 - definitions, 23
 - good practices, 39
 - awareness, 35
 - data standard, 35
 - establish, goal, 35
 - issue, 40
 - setting goals, 41
 - sources of data, 24

- Data science, 2, 52
 - B2C app, 41
 - core fields, 4, 5
 - definitions, 2
 - entertainment, 42, 43
 - multiple disciplines, 3, 4
 - scientists, 20
 - skills, 19, 20
 - subfields, 17
 - Data science culture
 - analytics and statistics, 126
 - data-centric (see Data-centric, organizational level)
 - data driven, 128
 - data informed, 128
 - data scientists, 129, 130
 - decision-making, 128, 129
 - employee level, 127
 - levels, 127
 - management level, 127
 - organizational level, 127
 - Data science culture, build
 - baby steps approach, 137
 - dark data, 138
 - data driven, 139
 - embedded culture, 140
 - employees benefits, 136
 - friendly environment, 136
 - good behavior; reward, 139
 - resistance to change
 - cultural, 135
 - intellectual, 136
 - personal, 135
 - types, 134, 135
 - scenarios, businesses, 134
 - use case approach, 137
 - Data science process
 - data collection and management, 53
 - example (see Problem solving, data science)
 - life cycle, 52
 - problem, defined, 53
 - solve problem, 54
 - steps, 52
 - value, create, 54
 - Data science project,
 - questionnaire, 147, 148
 - Data science tribes
 - computer scientists, 116, 117
 - deep learning, 122
 - domain specialist, 121
 - experience, 119
 - major, 115
 - quantitative specialists, 117–119
 - self-taught people, 120
 - smaller, 116
 - software platform user, 120
 - statisticians, 117
 - Data scientist, 57, 90
 - code hacking skills, 106
 - disengage, 108, 109
 - domain knowledge, 107
 - evaluation, 122, 123
 - mathematical/statistics knowledge, 107
 - motivation, 108
 - skills and knowledge, 105, 106
 - traditional research, 107
 - Venn diagram, 129
 - Data scientist, job
 - academia relationship, 110, 111
 - problems, 110
 - recruitment drive freeze, 114, 115
 - startup vs. bigger company, 113
 - team, 110
 - technology stack, 110
 - toolbox, 112
 - traditional limitations, avoid, 111
 - young talent, 113
 - Data strategy canvas, 145, 146
 - Deep learning, 13, 14
 - Descriptive statistics, 60–61, 149
 - Dialogflow, 92
 - Dirty data, 45
 - case study, 48, 49
 - causes, 46
 - goals, 46, 47
 - solution, 47, 48
 - Domain expert, 53, 54
 - Domain specialist, 116, 121
- ## E, F, G
- Experimental data collection, 24

H

Hypothesis testing, 63, 91

I, J

Inferential statistics, 61–62, 74

Interpretable machine learning, 87, 88

Interview questions, data scientists, 148–150

K

Knowledge Discovery in
Databases (KDD), 18

L

Lighthill's report, 10

Lisp machines, 10

Local interpretable model-agnostic
explanations (LIME), 87

M

Machine learning, 4
 advantage, 79
 buzzwords and fields, 18
 definitions, 11, 78
 example, 78
 pattern recognition, 18
 problem, 12, 13
 types, 80

Market-basket analysis, 92

Missionaries and cannibals problem, 7, 8

N

Neural networks, 14

O

Observational data collection, 24

OkCupid, 41, 42

P

Pitfalls, data science process
 bad collaboration
 data-related problems, 100
 data scientist, 100

 significance test, 99

 statistical model, 98, 99

 data strategy, 101

 education, 101

 problem, 100

 solution, 101

Problem solving

 data scientist, 90

 example, 95

 heuristics, 90–92

 heuristics, fail

 automate data collection, 93

 data scientists, 94

 right data, lack, 93

 vague project plan, 93

Problem solving, data science

 actionable insights, 58

 build models, 58

 data

 collection, 56, 57

 relevant definitions and
 protocol, 55

 data scientists, 57

 define problem, 55

 identify individuals/culprits, 58

Q

Qualitative research, 28, 29

Quantitative research, 27, 28

Quantitative specialists, 117–119

R

Recommender systems, 19, 92, 128

Regular vs. deep neural network, 13

S

Sampling biases, 73

 area bias, 72

 question bias, 72

 selection bias, 71

 self-selection bias, 72

 social desirability bias, 72

Shapley value explanation, 87

Solution adoption questionnaire,
 146, 150, 151

Speech Understanding Research program, 10
Statistical modeling, 64, 65, 117, 149
Statistical modeling problem, 90
Statistical tests, 64
Statisticians, 117
Statistics, 5, 15, 16
 branches, 59, 60
 chart, 65, 67, 68
 descriptive, 69
 distributions, 70
 examples, 63
 vs. machine learning, 16
 mislead, 65
 skew perception, 68
 tests, 64
 use, 62

Supervised learning, 80, 91
 classification model, 81
 LEGO model, 82
 regression, 82

T

Tinder, 41, 42
Travelling salesman problem, 8

U, V, W, X, Y, Z

Unsupervised learning, 80, 81, 92
 anomaly detection, 85
 clustering, 83, 84
 customer segmentation, 83
 dimensionality reduction, 86
 personality theory, 86