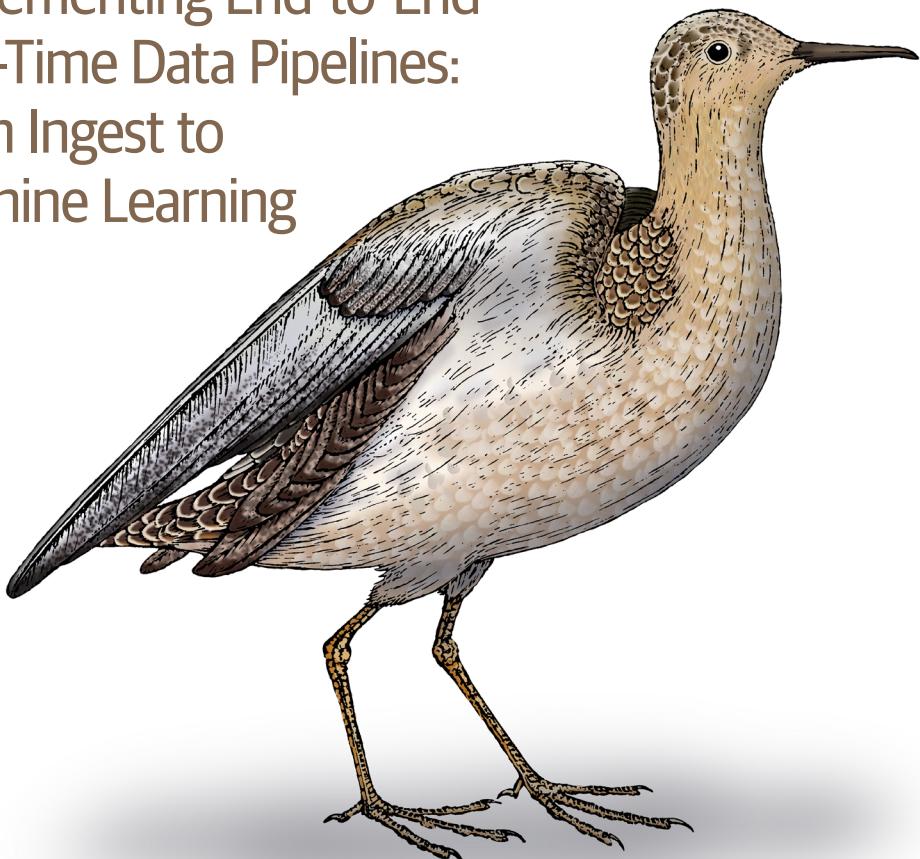


O'REILLY®

2nd Edition

Data Science on the Google Cloud Platform

Implementing End-to-End
Real-Time Data Pipelines:
From Ingest to
Machine Learning



Valliappa Lakshmanan

Data Science on the Google Cloud Platform

Learn how easy it is to apply sophisticated statistical and machine learning methods to real-world problems when you build using Google Cloud Platform (GCP). This hands-on guide shows data engineers and data scientists how to implement an end-to-end data pipeline with cloud native tools on GCP.

Throughout this updated second edition, you'll work through a sample business decision by employing a variety of data science approaches. Follow along by building a data pipeline in your own project on GCP, and discover how to solve data science problems in a transformative and more collaborative way.

You'll learn how to:

- Employ best practices to build highly scalable data and ML pipelines on Google Cloud
- Automate and schedule data ingest using Cloud Run
- Create and populate a dashboard in Data Studio
- Build a real-time analytics pipeline using Pub/Sub, Dataflow, and BigQuery
- Conduct interactive data exploration with BigQuery
- Create a Bayesian model with Spark on Dataproc
- Forecast time series and do anomaly detection with BigQuery ML
- Aggregate within time windows with Dataflow
- Train and deploy explainable machine learning models with TensorFlow on Vertex AI

"Highly recommend this book to anyone who is interested in learning how to apply data science to real-world problems with the latest GCP tools. You will find this book interesting whether you are a data scientist, data engineer, data analyst, or ML engineer."

—Margaret Maynard-Reid
ML Engineer, Tiny Peppers

Valliappa (Lak) Lakshmanan is the director of analytics and AI solutions at Google Cloud, where he leads a team building cross-industry solutions to business problems. His mission is to democratize machine learning so that it can be done by anyone anywhere. Lak is the author or coauthor of *Practical Machine Learning for Computer Vision*, *Machine Learning Design Patterns*, and *Data Governance: The Definitive Guide* (O'Reilly).

DATA SCIENCE

US \$69.99 CAN \$87.99
ISBN: 978-1-098-11895-2



5 6 9 9 9
9 781098 118952

Twitter: @oreillymedia
[linkedin.com/company/oreilly-media](https://www.linkedin.com/company/oreilly-media)
[youtube.com/oreillymedia](https://www.youtube.com/oreillymedia)

SECOND EDITION

Data Science on the Google Cloud Platform

*Implementing End-to-End Real-Time
Data Pipelines: From Ingest to Machine Learning*

Valliappa Lakshmanan

Beijing • Boston • Farnham • Sebastopol • Tokyo

O'REILLY®

Data Science on the Google Cloud Platform, Second Edition

by Valliappa Lakshmanan

Copyright © 2022 Google LLC. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or corporate@oreilly.com.

Acquisition Editor: Jessica Haberman

Indexer: WordCo Indexing Services, Inc.

Development Editor: Michele Cronin

Interior Designer: David Futato

Production Editor: Katherine Tozer

Cover Designer: Karen Montgomery

Copyeditor: Tom Sullivan

Illustrator: Kate Dullea

Proofreader: Piper Editorial Consulting, LLC

January 2018: First Edition

April 2022: Second Edition

Revision History for the Second Edition

2022-03-29: First Release

2022-04-22: Second Release

See <http://oreilly.com/catalog/errata.csp?isbn=9781098118952> for release details.

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *Data Science on the Google Cloud Platform*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the author and do not represent the publisher's views. While the publisher and the author have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the author disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

978-1-098-11895-2

LSI

Table of Contents

Preface.....	xı
1. Making Better Decisions Based on Data.....	1
Many Similar Decisions	4
The Role of Data Scientists	5
Scrappy Environment	7
Full Stack Cloud Data Scientists	8
Collaboration	9
Best Practices	10
Simple to Complex Solutions	10
Cloud Computing	11
Serverless	12
A Probabilistic Decision	13
Probabilistic Approach	14
Probability Density Function	15
Cumulative Distribution Function	16
Choices Made	18
Choosing Cloud	19
Not a Reference Book	19
Getting Started with the Code	20
Agile Architecture for Data Science on Google Cloud	22
What Is Agile Architecture?	23
No-Code, Low-Code	23
Use Managed Services	24
Summary	25
Suggested Resources	26

2. Ingesting Data into the Cloud.....	29
Airline On-Time Performance Data	29
Knowability	31
Causality	31
Training–Serving Skew	32
Downloading Data	33
Hub-and-Spoke Architecture	34
Dataset Fields	35
Separation of Compute and Storage	37
Scaling Up	39
Scaling Out with Sharded Data	41
Scaling Out with Data-in-Place	43
Ingesting Data	46
Reverse Engineering a Web Form	46
Dataset Download	48
Exploration and Cleanup	50
Uploading Data to Google Cloud Storage	51
Loading Data into Google BigQuery	55
Advantages of a Serverless Columnar Database	55
Staging on Cloud Storage	57
Access Control	57
Ingesting CSV Files	61
Partitioning	62
Scheduling Monthly Downloads	63
Ingesting in Python	65
Cloud Run	71
Securing Cloud Run	72
Deploying and Invoking Cloud Run	74
Scheduling Cloud Run	75
Summary	76
Code Break	77
Suggested Resources	78
3. Creating Compelling Dashboards.....	81
Explain Your Model with Dashboards	83
Why Build a Dashboard First?	84
Accuracy, Honesty, and Good Design	86
Loading Data into Cloud SQL	88
Create a Google Cloud SQL Instance	89
Create Table of Data	91
Interacting with the Database	95
Querying Using BigQuery	96

Schema Exploration	96
Using Preview	97
Using Table Explorer	99
Creating BigQuery View	100
Building Our First Model	101
Contingency Table	101
Threshold Optimization	103
Building a Dashboard	106
Getting Started with Data Studio	107
Creating Charts	109
Adding End-User Controls	110
Showing Proportions with a Pie Chart	112
Explaining a Contingency Table	117
Modern Business Intelligence	119
Digitization	119
Natural Language Queries	120
Connected Sheets	122
Summary	123
Suggested Resources	123
4. Streaming Data: Publication and Ingest with Pub/Sub and Dataflow.....	125
Designing the Event Feed	126
Transformations Needed	127
Architecture	128
Getting Airport Information	129
Sharing Data	132
Time Correction	133
Apache Beam/Cloud Dataflow	135
Parsing Airports Data	136
Adding Time Zone Information	139
Converting Times to UTC	141
Correcting Dates	144
Creating Events	146
Reading and Writing to the Cloud	148
Running the Pipeline in the Cloud	150
Publishing an Event Stream to Cloud Pub/Sub	153
Speed-Up Factor	154
Get Records to Publish	155
How Many Topics?	156
Iterating Through Records	157
Building a Batch of Events	158
Publishing a Batch of Events	159

Real-Time Stream Processing	160
Streaming in Dataflow	160
Windowing a Pipeline	162
Streaming Aggregation	162
Using Event Timestamps	165
Executing the Stream Processing	166
Analyzing Streaming Data in BigQuery	168
Real-Time Dashboard	169
Summary	170
Suggested Resources	171
5. Interactive Data Exploration with Vertex AI Workbench.....	173
Exploratory Data Analysis	174
Exploration with SQL	177
Reading a Query Explanation	179
Exploratory Data Analysis in Vertex AI Workbench	184
Jupyter Notebooks	185
Creating a Notebook	186
Jupyter Commands	188
Installing Packages	188
Jupyter Magic for Google Cloud	189
Exploring Arrival Delays	190
Basic Statistics	191
Plotting Distributions	191
Quality Control	194
Arrival Delay Conditioned on Departure Delay	199
Evaluating the Model	204
Random Shuffling	204
Splitting by Date	205
Training and Testing	206
Summary	210
Suggested Resources	210
6. Bayesian Classifier with Apache Spark on Cloud Dataproc.....	211
MapReduce and the Hadoop Ecosystem	211
How MapReduce Works	212
Apache Hadoop	214
Google Cloud Dataproc	214
Need for Higher-Level Tools	216
Jobs, Not Clusters	217
Preinstalling Software	219
Quantization Using Spark SQL	221

JupyterLab on Cloud Dataproc	222
Independence Check Using BigQuery	223
Spark SQL in JupyterLab	225
Histogram Equalization	227
Bayesian Classification	231
Bayes in Each Bin	231
Evaluating the Model	232
Dynamically Resizing Clusters	233
Comparing to Single Threshold Model	235
Orchestration	237
Submitting a Spark Job	238
Workflow Template	238
Cloud Composer	239
Autoscaling	239
Serverless Spark	240
Summary	242
Suggested Resources	243
7. Logistic Regression Using Spark ML.....	245
Logistic Regression	246
How Logistic Regression Works	246
Spark ML Library	249
Getting Started with Spark Machine Learning	250
Spark Logistic Regression	251
Creating a Training Dataset	252
Training the Model	256
Predicting Using the Model	259
Evaluating a Model	260
Feature Engineering	263
Experimental Framework	263
Feature Selection	267
Feature Transformations	271
Feature Creation	274
Categorical Variables	278
Repeatable, Real Time	280
Summary	281
Suggested Resources	282
8. Machine Learning with BigQuery ML.....	283
Logistic Regression	283
Presplit Data	285
Interrogating the Model	286

Evaluating the Model	287
Scale and Simplicity	289
Nonlinear Machine Learning	290
XGBoost	290
Hyperparameter Tuning	292
Vertex AI AutoML Tables	294
Time Window Features	296
Taxi-Out Time	296
Compounding Delays	298
Causality	299
Time Features	300
Departure Hour	300
Transform Clause	302
Categorical Variable	303
Feature Cross	303
Summary	305
Suggested Resources	306
9. Machine Learning with TensorFlow in Vertex AI.....	309
Toward More Complex Models	310
Preparing BigQuery Data for TensorFlow	314
Reading Data into TensorFlow	315
Training and Evaluation in Keras	317
Model Function	317
Features	318
Inputs	320
Training the Keras Model	320
Saving and Exporting	322
Deep Neural Network	322
Wide-and-Deep Model in Keras	323
Representing Air Traffic Corridors	323
Bucketing	324
Feature Crossing	325
Wide-and-Deep Classifier	326
Deploying a Trained TensorFlow Model to Vertex AI	327
Concepts	328
Uploading Model	328
Creating Endpoint	330
Deploying Model to Endpoint	330
Invoking the Deployed Model	331
Summary	332
Suggested Resources	333

10. Getting Ready for MLOps with Vertex AI.....	335
Developing and Deploying Using Python	336
Writing <code>model.py</code>	337
Writing the Training Pipeline	338
Predefined Split	340
AutoML	341
Hyperparameter Tuning	343
Parameterize Model	344
Shorten Training Run	345
Metrics During Training	347
Hyperparameter Tuning Pipeline	347
Best Trial to Completion	349
Explaining the Model	350
Configuring Explanations Metadata	350
Creating and Deploying Model	352
Obtaining Explanations	352
Summary	354
Suggested Resources	355
11. Time-Windows Features for Real-Time Machine Learning.....	357
Time Averages	357
Apache Beam and Cloud Dataflow	358
Reading and Writing	360
Time Windowing	362
Machine Learning Training	367
Machine Learning Dataset	367
Training the Model	373
Streaming Predictions	376
Reuse Transforms	377
Input and Output	379
Invoking Model	380
Reusing Endpoint	381
Batching Predictions	384
Streaming Pipeline	385
Writing to BigQuery	385
Executing Streaming Pipeline	386
Late and Out-of-Order Records	387
Possible Streaming Sinks	393
Summary	400
Suggested Resources	401

12. The Full Dataset.....	403
Four Years of Data	403
Creating Dataset	404
Training Model	409
Evaluation	411
Summary	417
Suggested Resources	417
Conclusion.....	419
Considerations for Sensitive Data Within Machine Learning Datasets.....	423
Index.....	431

Preface

In my current role at Google, I get to work alongside data scientists and data engineers in a variety of industries as they move their data processing and analysis methods to the public cloud. Some try to do the same things they do on premises, the same way they do them, just on rented computing resources. The visionary users, though, rethink their systems, transform how they work with data, and thereby are able to innovate faster.

As early as 2011, an [article in *Harvard Business Review*](#) recognized that some of cloud computing’s greatest successes come from allowing groups and communities to work together in ways that were not previously possible. This is now much more widely recognized. An [MIT survey in 2017](#) found that more respondents (45%) cited increased agility rather than cost savings (34%) as the reason to move to the public cloud. However, it is still not widely achieved. [McKinsey estimated in 2021](#) that companies are leaving behind nearly \$1 trillion of value by not looking at the public cloud as a source of transformative value. Therefore, being able to work on a data science project in the cloud is a skill well worth investing in.

In this book, we walk through an example of a cloud-native, transformative, collaborative way of doing data science. You will learn how to implement an end-to-end data pipeline—we will begin with ingesting the data in a serverless way and work our way through data exploration, dashboards, relational databases, and streaming data all the way to training and making an operational machine learning model. I cover all these aspects of data-based services because data engineers will be involved in designing the services, developing the statistical and machine learning models, and implementing them in large-scale production and in real time.

Who This Book Is For

If you use computers to work with data, this book is for you. You might go by the title of data analyst, database administrator, data engineer, data scientist, or systems programmer today. Although your role might be narrower today (perhaps you do

only data analysis, or only model building, or only DevOps), you want to stretch your wings a bit—you want to learn how to create data science models as well as how to implement them at scale in production systems.

Google Cloud Platform is designed to make you forget about infrastructure. The marquee data services—Google BigQuery, Cloud Dataflow, Cloud Pub/Sub, and Vertex AI—are all serverless and autoscaling. When you submit a query to BigQuery, it is run on thousands of nodes, and you get your result back; you don't spin up a cluster or install any software. Similarly, in Cloud Dataflow, when you submit a data pipeline, and in Vertex AI, when you submit a machine learning job, you can process data at scale and train models at scale without worrying about cluster management or failure recovery. Cloud Pub/Sub is a global messaging service that autoscales to the throughput and number of subscribers and publishers without any work on your part. Even when you're running open source software like Apache Spark that's designed to operate on a cluster, Google Cloud Platform makes it easy with job-specific clusters and serverless Spark. Because of this job-specific infrastructure, there's no need to fear overprovisioning hardware or running out of capacity to run a job when you need it. Plus, data is encrypted, both at rest and in transit, and kept secure. As a data scientist, not having to manage infrastructure is incredibly liberating.

These autoscaled, fully managed services make it easier to implement data science models at scale—which is why data scientists no longer need to hand off their models to data engineers. Instead, they can write a data science workload, submit it to the cloud, and have that workload executed automatically in an autoscaled manner. At the same time, data science packages are becoming simpler and simpler. So, it has become extremely easy for an engineer to slurp in data and use a canned model to get an initial (and often very good) model up and running. With well-designed packages and easy-to-consume APIs, you don't need to know the esoteric details of data science algorithms—only what each algorithm does and how to link algorithms together to solve realistic problems. This convergence between data science and data engineering is why you can stretch your wings beyond your current role.

Rather than simply read this book cover-to-cover, I strongly encourage you to follow along with me by trying out the code. The full source code for the end-to-end pipeline I build in this book is on [GitHub](#). Create a [Google Cloud Platform project](#), and after reading each chapter, try to repeat what I did by referring to the code and to the *README.md* file in each folder of the GitHub repository.



Follow the instructions in the *README.md* files in GitHub to try out the code. The code snippets in the book are often incomplete—for example, I may omit some arguments to cloud commands for clarity or conciseness.

Note that this is not a reference book—the best reference to Google Cloud is its documentation, and there is very little value to be had by simply reproducing that in a book. Instead, this book shows you how to use a variety of tools together to solve a problem. My goal here is to teach you how to think about a problem in order to solve it using Google Cloud, not to comprehensively cover any particular product.

If you find yourself fascinated by a topic in this book and want to dive deeper, you can find a few selected resources at the end of every chapter that provide a deeper dive into topics covered in the chapter. Don’t feel obligated to watch every video or read every article.

Conventions Used in This Book

The following typographical conventions are used in this book:

Italic

Indicates new terms, URLs, email addresses, filenames, and file extensions.

Constant width

Used for program listings, as well as within paragraphs to refer to program elements such as variable or function names, databases, data types, environment variables, statements, and keywords.

Constant width bold

Shows commands or other text that should be typed literally by the user.

Constant width italic

Shows text that should be replaced with user-supplied values or by values determined by context.



This element signifies a tip or suggestion.



This element signifies a general note.



This element indicates a warning or caution.

Using Code Examples

Supplemental material (code examples, exercises, etc.) is available for download at <https://github.com/GoogleCloudPlatform/data-science-on-gcp>.

If you have a technical question or a problem using the code examples, please email bookquestions@oreilly.com.

This book is here to help you get your job done. In general, if example code is offered with this book, you may use it in your programs and documentation. You do not need to contact us for permission unless you're reproducing a significant portion of the code. For example, writing a program that uses several chunks of code from this book does not require permission. Selling or distributing examples from O'Reilly books does require permission. Answering a question by citing this book and quoting example code does not require permission. Incorporating a significant amount of example code from this book into your product's documentation does require permission.

We appreciate, but do not require, attribution. An attribution usually includes the title, author, publisher, and ISBN. For example: “*Data Science on the Google Cloud Platform*” by Valliappa Lakshmanan (O'Reilly). Copyright 2022 Google LLC, 978-1-098-11895-2.”

If you feel your use of code examples falls outside fair use or the permission given above, feel free to contact us at permissions@oreilly.com.

O'Reilly Online Learning



For more than 40 years, *O'Reilly Media* has provided technology and business training, knowledge, and insight to help companies succeed.

Our unique network of experts and innovators share their knowledge and expertise through books, articles, and our online learning platform. O'Reilly's online learning platform gives you on-demand access to live training courses, in-depth learning paths, interactive coding environments, and a vast collection of text and video from O'Reilly and 200+ other publishers. For more information, visit <https://oreilly.com>.

How to Contact Us

Please address comments and questions concerning this book to the publisher:

O'Reilly Media, Inc.
1005 Gravenstein Highway North

Sebastopol, CA 95472
800-998-9938 (in the United States or Canada)
707-829-0515 (international or local)
707-829-0104 (fax)

We have a web page for this book, where we list errata, examples, and any additional information. You can access this page at <https://oreil.ly/data-science-on-gcp>.

Email bookquestions@oreilly.com to comment or ask technical questions about this book.

For news and information about our books and courses, visit <https://oreilly.com>.

Find us on Facebook: <https://facebook.com/oreilly>.

Follow us on Twitter: <https://twitter.com/oreillymedia>.

Watch us on YouTube: <https://www.youtube.com/oreillymedia>.

Acknowledgments

When I took the job at Google in 2014, I had used the public cloud simply as a way to rent infrastructure—so I was spinning up virtual machines, installing the software I needed on those machines, and then running my data processing jobs using my usual workflow. Fortunately, I realized that Google’s big data stack was different, and so I set out to learn how to take full advantage of all the data and machine learning tools on Google Cloud Platform.

The way I learn best is to write code, and so that’s what I did. When a Python meetup group asked me to talk about Google Cloud Platform, I did a show-and-tell of the code that I had written. It turned out that a walk-through of the code to build an end-to-end system while contrasting different approaches to a data science problem was quite educational for the attendees. I wrote up the essence of my talk as a book proposal and sent it to O’Reilly Media.

A book, of course, needs to have a lot more depth than a 60-minute code walk-through. Imagine that you come to work one day to find an email from a new employee at your company, someone who’s been at the company less than six months. Somehow, he’s decided he’s going to write a book on the pretty sophisticated platform that you’ve had a hand in building and is asking for your help. He is not part of your team, helping him is not part of your job, and he is not even located in the same office as you. What is your response? Would you volunteer?

What makes Google such a great place to work is the people who work here. It is a testament to the company’s culture that so many people—engineers, technical leads, product managers, solutions architects, data scientists, legal counsel, directors—across so many different teams happily gave of their expertise to someone they had

never met (in fact, I still haven’t met many of these people in person). This book, thus, is immeasurably better because of (in alphabetical order of last names) William Brockman, Mike Dahlin, Tony DiLoreto, Bob Evans, Roland Hess, Brett Hesterberg, Dennis Huo, Chad Jennings, Puneith Kaul, Dinesh Kulkarni, Manish Kurse, Reuven Lax, Jonathan Liu, James Malone, Dave Oleson, Mosha Pasumansky, Kevin Peterson, Olivia Puerta, Reza Rokni, Karn Seth, Sergei Sokolenko, and Amy Unruh. In particular, thanks to Mike Dahlin, Manish Kurse, and Olivia Puerta for reviewing every single chapter. When the first edition of the book was in early access, I received valuable error reports from Anthonios Partheniou and David Schwantner. Needless to say, I am responsible for any errors that remain.

A few times during the writing of the book, I found myself completely stuck. Sometimes, the problems were technical. Thanks to (in alphabetical order) Ahmet Altay, Eli Bixby, Ben Chambers, Slava Chernyak, Marián Dvorský, Robbie Haertel, Felipe Hoffa, Amir Hormati, Qiming (Bradley) Jiang, Kenneth Knowles, Nikhil Kothari, and Chris Meyers for showing me the way forward. At other times, the problems were related to figuring out company policy or getting access to the right team, document, or statistic. This book would have been a lot poorer had these colleagues not unblocked me at critical points (again alphabetically): Louise Byrne, Apurva Desai, Rochana Golani, Fausto Ibarra, Jason Martin, Neal Mueller, Philippe Poutonnet, Brad Svee, Jordan Tigani, William Vampenebe, and Miles Ward. Thank you all for your help and encouragement.

Five years on, I continue to be humbled by the incredible talent and collaboration of my colleagues. Sagar Baliyara, Filipe Gracio, Polong Lin, and Krishnan Saidapet (in alphabetical order of last names) brought a close eye to the second edition and made many great suggestions.

Thanks also to the O’Reilly team—Marie Beaugureau, Kristen Brown, Ben Lorica, Tim McGovern, Rachel Roumeliotis, and Heather Scherer for believing in me and making the process of moving from draft to the first edition of the book painless. Producing the second edition was greatly streamlined by Katherine Tozer, Michele Cronin, and Tom Sullivan.

The second edition has also greatly benefited from fresh outside perspectives. Colin Dietrich verified much of the code in the book and made numerous pull requests to the GitHub repository. Joy Payton suggested many improvements to make the book more accessible to beginners in data science. Michael Hopkins and Margaret Maynard-Reid scrutinized the book for areas that needed updating. Thanks also to readers of the first edition who left reviews of the book on Amazon, filed issues on GitHub, and reached out to me via email and on Twitter. Your feedback has greatly improved this edition of the book.

Finally, and most important, thanks to Abirami, Sidharth, and Sarada for your understanding and patience even as I became engrossed in writing and coding. You make it all worthwhile.

I am donating 100% of the royalties from this book to [United Way of King County](#), where I live. I strongly encourage you to get involved with a local charity to give, volunteer, and take action to help solve your community's toughest challenges.

Making Better Decisions Based on Data

The primary purpose of data analysis is to make better decisions. There is rarely any need for us to spend time analyzing data if we aren't under pressure to make a decision based on the results of that analysis. When you are purchasing a car, you might ask the seller what year the car was manufactured and the odometer reading. Knowing the age of the car allows you to estimate the potential value of the car. Dividing the odometer reading by the age of the car allows you to discern how hard the car has been driven, and whether it is likely to last the five years you plan to keep it. Had you not cared about purchasing the car, there would have been no need for you to do this data analysis.

In fact, we can go further—the purpose of collecting data is, in many cases, only so that you can later perform data analysis and make decisions based on that analysis (see [Figure 1-1](#)). When you asked the seller the age of the car and its mileage, you were collecting data to carry out your data analysis. But it goes beyond your data collection. The car has an odometer in the first place because many people, not just potential buyers, will need to make decisions based on the mileage of the car. The odometer reading needs to support many decisions—should the manufacturer pay for a failed transmission? Is it time for an oil change? The analysis for each of these decisions is different, but they all rely on the fact that the mileage data has been collected.

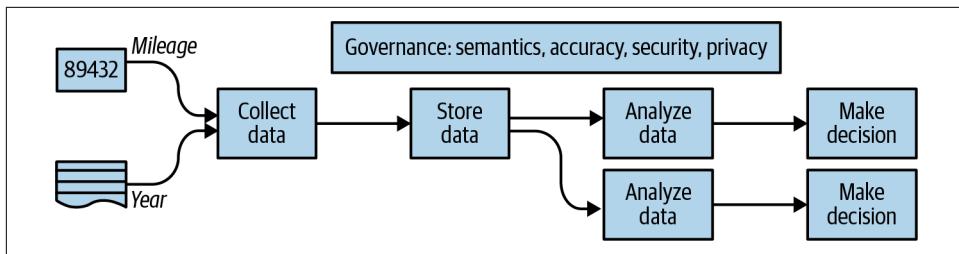


Figure 1-1. The purpose of collecting data is to make decisions with it.

If you are in the business of making a lot of decisions using car mileage, it makes sense to store the data that you have collected so that future decisions are easier to make. Collecting data takes time and effort, whereas storing it is relatively inexpensive. Of course, you have to plan on storing the data in a way that you will know what it *means* when you need it later. This is called capturing the *semantics* of the data and is an important aspect of *data governance*, to ensure that data is useful for decision making.

Collecting data in a form that enables decisions to be made places requirements on the collecting infrastructure and the security of such infrastructure. How does the insurance company that receives an accident claim and needs to pay its customer the car's value know that the odometer reading is accurate? How are odometers calibrated? What kinds of safeguards are in place to ensure that the odometer has not been tampered with? What happens if the tampering is inadvertent, such as installing tires whose size is different from what was used to calibrate the odometer? The *auditability* of data is important whenever there are multiple parties involved, and ownership and use of the data are separate. When data is unverifiable, markets fail, optimal decisions cannot be made, and the parties involved need to resort to signaling and screening.¹

¹ The classic paper on this is George Akerlof's 1970 paper titled "[The Market for Lemons](#)". Akerlof, Michael Spence (who explained signaling), and Joseph Stiglitz (who explained screening) jointly received the 2001 Nobel Prize in Economics for describing this problem. In a transaction that involves asymmetric information, the party with good information signals, whereas the party with poor information screens. For example, the seller of a car (who has good information) might signal that they have a great car by publishing the repair record of the car. The buyer (who has poor information) might screen cars by rejecting any cars from cities that recently experienced a flood.

Not all data is as expensive to collect and secure as the odometer reading of a car.² The cost of sensors has dropped dramatically in recent decades, and many of our daily processes produce so much data that we find ourselves in possession of data that we had no intention of explicitly collecting. Because the hardware to collect, ingest, and store the data has become cheaper, we often default to retaining the data indefinitely, keeping it around for no discernible reason. As the size of data within an organization increases, it becomes more and more essential to organize and catalog it carefully. So, if we're to perform analysis on all of this data that we somehow managed to collect and store, we better have a purpose for it. Labor remains expensive.

Another reason to be purposeful about the data we collect and store is that a lot of it is about people. Knowing the mileage of the car that someone drives gives us a lot of information about them, and this is information that they may not want to share other than for the specific purpose of estimating the market price of their car. Privacy and confidentiality need to be considered even before any data is collected, so that appropriate decisions can be made about what data to collect, how to control access to it, and how long to retain it. This is even more important when the data is provided at significant cost, risk, and/or loss of bodily autonomy, as is the case for much biomedical patient data. Having a data privacy expert examine your schema and data protection practices is an investment that will pay for itself manyfold in terms of regulatory compliance and public relations.

Data analysis is usually triggered because a decision needs to be made. To move into a market or not? To pay a commission or not? How high to bid up the price? How many bags to purchase? Whether to buy now or wait a week? The decisions keep multiplying, and because data is so ubiquitous now, we no longer need to make those decisions based on *heuristics* or simple rules of thumb. We can now make those decisions in a data-driven manner.

Of course, we don't need to create the systems and tools to make every data-driven decision ourselves. The use case of estimating the value of a car that has been driven a certain distance is common enough that there are several companies that provide this as a service—they will verify that an odometer is accurate, confirm that the car hasn't been in an accident, and compare the asking price against the typical selling price of

² The odometer itself might not be all that expensive, but collecting that information and ensuring that it is correct has considerable costs. The last time I sold a car, I had to sign a statement that I had not tampered with the odometer, and that statement had to be notarized by a bank employee with a financial guarantee. This was required by the company that was loaning the purchase amount on the car to the buyer. Every auto mechanic is supposed to report odometer tampering, and there is a state government agency that enforces this rule. All of these costs are significant. Even if you disregard all these external costs, and assume that the hardware and infrastructure exists such that each car has an odometer, there is still a significant cost associated with streaming that data from cars into a central location so that you have real-time odometer readings from all the cars in your fleet. The cost of securing that data to respect the privacy of the drivers can also be quite significant.

cars in your market. The real value, therefore, comes not in making a data-driven decision once, but in being able to do it systematically and provide it as a service.³ This also allows companies to specialize in different business areas and continuously improve the accuracy and value of the decisions that can be made.

Many Similar Decisions

Because of the low costs associated with sensors and storage, there are many industries and use cases that have the potential to support data-driven decision making. If you are working in such an industry, or you want to start a company that will address such a use case, the possibilities for supporting data-driven decision making have just become wider. In some cases, you will need to collect the data. In others, you will have access to data that was already collected, and, in many cases, you will need to supplement the data you have with other datasets that you will need to hunt down. In all these cases, being able to carry out data analysis to support decision making systematically on behalf of users is one of the most important skills to possess.

In this book, I will take a decision that needs to be made and apply different statistical and machine learning methods to gain insight into making that decision. However, we don't want to make that decision just once, even though we might occasionally pose it that way. Instead, we will look at how to make the decision in a systematic manner so that we use the same algorithm to make the decision many, many, many times. Our ultimate goal will be to provide this decision-making capability as a service to our customers—they will tell us the things they reasonably can be expected to know, and we will either know or infer the rest (because we have been systematically collecting data). Based on this data, we will suggest the optimal decision.

Whether or not a decision is a one-off is the primary difference between data analytics and data science. Data analytics is about manually analyzing data to make a single decision or answer a single question. Data science is about developing a technique (called a model or algorithm) so that similar decisions can be made in a systematic way. Often, data science is about automating and optimizing the decision-making process that was first determined through data analysis.⁴

When we are collecting the data, we will need to look at how to make the data secure. This will include how to ensure not only that the data has not been tampered with but also that users' private information is not compromised—for example, if we are systematically collecting odometer mileage and know the precise mileage of the car at

³ Providing it as a service is often the only way to meet the mission of your organization—whether it is to monetize it, support thousands of users, or provide it at low cost to decision makers.

⁴ Contrary to what you may hear, it is not about whether you use SQL or Python! You can do data science in SQL—we will see BigQuery ML later on in the book, and you can use Python for one-off data analysis.

any point in time, this knowledge becomes extremely sensitive information. Given enough other information about the customer (such as the home address and traffic patterns in the city in which the customer lives), the mileage is enough to be able to infer that person's location at all times.⁵ So, the privacy implications of hosting something as seemingly innocuous as the mileage of a car can become enormous. Security implies that we need to control access to the data, and we need to maintain immutable audit logs on who has viewed or changed the data.

It is not enough to simply collect the data or use it as-is. We must understand the data. Just as we needed to know the kinds of problems associated with odometer tampering to understand the factors that go into estimating a vehicle's value based on mileage, our analysis methods will need to consider how the data was collected in real time and the kinds of errors that could be associated with that data. Intimate knowledge of the data and its quirks is invaluable when it comes to doing data science—often the difference between a data-science startup idea that works and one that doesn't is whether the appropriate nuances have all been thoroughly evaluated and taken into account.

When it comes to providing the decision-support capability as a service, it is not enough to simply have a way to do it in some offline system somewhere. Enabling it as a service implies a whole host of other concerns. The first set of concerns is about the quality of the decision itself—how accurate is it typically? What are the typical sources of errors? In what situations should this system not be used? The next set of concerns, however, is about the quality of service. How reliable is it? How many queries per second can it support? What is the latency between some piece of data being available and it being incorporated into the model that is used to provide systematic decision making? In short, we will use this single use case as a way to explore many different facets of practical data science.

The Role of Data Scientists

“Wait a second,” I imagine you saying, “I never signed up for queries-per-second of a web service. We have people who do that kind of stuff. My job is to write SQL queries and create reports. I don’t recognize this thing you are talking about. It’s not what I do at all.” Or perhaps the first part of the discussion was what puzzled you. “Decision making? That’s for the business people. Me? What I do is to design data processing

⁵ In 2014, New York City officials released a [public dataset](#) of New York City taxi trips in response to a Freedom of Information request. However, because it was improperly anonymized, a brute force attack was able to find out the [trips associated with any specific driver](#). It got worse. Privacy researchers were able to [cross-reference paparazzi photos](#) (which revealed the exact location of celebrities at specific times) and figure out which celebrities don’t tip. It got even worse. By looking at people who picked up a taxi cab at the same location every morning, and correlating it with the location from where they got dropped back, they were able to identify New Yorkers who frequented strip clubs.

systems. I can provision infrastructure, tell you what our systems are doing right now, and keep it all secure. Data science sure sounds fancy, but I do engineering. When you said Data Science on the Google Cloud Platform, I was thinking that you were going to talk about how to keep the systems humming and how to offload bursts of activity to the cloud.” A third set of people are wondering, “How is any of this data science? Where’s the discussion of different types of models and of how to make statistical inferences and evaluate them? Where’s the math? Why are you talking to data analysts and engineers? Talk to me, I’ve got a PhD.” These are fair points—I seem to be mixing up the jobs done by different sets of people in your organization.

In other words, you might agree with the following:

- Data analysis is there to support decision making.
- Decision making in a data-driven manner can be superior to heuristics.
- The accuracy of the decision models depends on your choice of the right statistical or machine learning approach.
- Nuances in the data can completely invalidate your modeling, so understanding the data and its quirks is crucial.
- There are large market opportunities in supporting decision making systematically and providing it as a service.
- Such services require ongoing data collection and model updates.
- Ongoing data collection implies robust security and auditing.
- Customers of the service require reliability, accuracy, and latency assurances.

What you might not agree with is whether these aspects are all things that you, personally and professionally, need to be concerned about. Instead, you might think of yourself as a data analyst, a data engineer, or a data scientist and not care about how the other roles do whatever it is that they do.

There are three answers to this objection:

- In any situation where you have small numbers of people doing ambitious things—a scrappy company, an innovative startup, an underfunded nonprofit, or an overextended research lab—you will find yourself playing all these roles, so learn the full lifecycle.
- The public cloud makes it relatively easy to learn all the roles, so why not be a full stack data scientist?
- Even if you work in a large company where these tasks are carried out by different roles, it is helpful to understand the end-to-end process and concerns at each stage. This will help you collaborate with other teams better.

Let's take these answers one by one.

Scrappy Environment

At Google, we look at the role of a data engineer quite expansively. Just as we refer to all our technical staff as engineers, we look at data engineers as an inclusive term for anyone who can “shape business outcomes by performing data analysis.” To perform data analysis, you begin by preparing the data so that you can analyze it at scale. It is not enough to simply count and sum and graph the results using SQL queries and charting software—you must understand the nuances of the data and the statistical framework within which you are interpreting the results. This ability to prepare the data and carry out statistically valid data analysis to solve specific business problems is of paramount importance—the queries, the reports, and the graphs are not the end goal. A verifiably accurate decision is.

Of course, it is not enough to do one-off data analysis. That data analysis needs to scale. In other words, an accurate decision-making process must be repeatable and be capable of being carried out by many users, not just you. The way to scale up one-off data analysis is to make it automated. After a data engineer has devised the algorithm, they should be able to make it systematic and repeatable. Just as it is a lot easier when the folks in charge of systems reliability can make code changes themselves,⁶ it is considerably easier when people who understand statistics and machine learning can code those models themselves. A data engineer, Google believes, should be able to go from building statistical and machine learning models to automating them. They can do this only if they are capable of designing, building, and troubleshooting data processing systems that are secure, responsible, reliable, fault-tolerant, scalable, and efficient.

This desire to have engineers who know data science and data scientists who can code is not Google's alone—it's common at technologically sophisticated organizations and at small companies. When a scrappy company advertises for data engineers or for data scientists, what they are looking for is a person who can do all the three tasks—data preparation, data analysis, and automation—that are needed to make repeatable, scalable decisions on the basis of data.

How realistic is it for companies to expect a Renaissance person, a virtuoso in different fields? Can they reasonably expect to hire data engineers who can do data science? How likely is it that they will find someone who can design a database schema, write SQL queries, train machine learning models, code up a data processing pipeline, and figure out how to scale it all up? Surprisingly, this is a very reasonable

⁶ Google invented the role of Site Reliability Engineers ([SREs](#))—these are folks in charge of keeping systems running. Unlike traditional IT, though, they know the software they are operating and are quite capable of making changes to it.

expectation because the amount of knowledge you need in order to do these jobs has become a lot less than what you needed a few years ago.

Full Stack Cloud Data Scientists

Because of the ongoing movement to the cloud, data scientists can do the job that used to be done by several people with different sets of skills. With the advent of autoscaling, serverless, managed infrastructure that is easy to program, there are more and more people who can build scalable systems. Therefore, it is now reasonable to expect to be able to hire data scientists who are capable of creating holistic data-driven solutions to your thorniest problems. You don't need to be a polymath to be a full stack data scientist—you simply need to learn how to do data science on the cloud.

Saying that the cloud is what makes full stack data scientists possible seems like a very tall claim. This hinges on what I mean by “cloud”—I don't mean simply migrating workloads that run on premises to infrastructure that is owned by a public cloud vendor. I'm talking, instead, about truly autoscaling, managed services that automate a lot of the infrastructure provisioning, monitoring, and management—services such as Google BigQuery, Vertex AI, Cloud Dataflow, and Cloud Run on Google Cloud Platform. When you consider that the scaling and fault-tolerance of many data analysis and processing workloads can be effectively automated, provided the right set of tools is being used, it is clear that the amount of IT support that a data scientist needs dramatically reduces with a migration to the cloud.

At the same time, data science tools are becoming simpler and simpler to use. The wide availability of frameworks like Spark, Pandas, and Keras has made data science and data science tools extremely accessible to the average developer—no longer do you need to be a specialist in data science to create a statistical model or train a random forest. This has opened up the field of data science to people in more traditional IT roles.

Similarly, data analysts and database administrators today can have completely different backgrounds and skill sets because data analysis has usually involved serious SQL wizardry, and database administration has typically involved deep knowledge of database indices and tuning. With the introduction of tools like BigQuery, in which tables are denormalized and the administration overhead is minimal, the role of a database administrator is considerably diminished. The growing availability of turnkey visualization tools like Tableau and Looker that connect to all the data stores within an enterprise makes it possible for a wider range of people to directly interact with enterprise warehouses and pull together compelling reports and insights.

The reason that all these data-related roles are merging together, then, is because the infrastructure problem is becoming less intense and the data analysis and modeling domain is becoming more democratized.

If you think of yourself today as a data scientist, or a data analyst, or a database administrator, or a systems programmer, this is either totally exhilarating or totally unrealistic. It is exhilarating if you can't wait to do all the other tasks that you've considered beyond your ken if the barriers to entry have fallen as low as I claim they have. If you are excited and raring to learn the things you will need to know in this new world of data, welcome!⁷ This book is for you.

If my vision of a blend of roles strikes you as an unlikely dystopian future, hear me out. The vision of autoscaling services that require very little in the form of infrastructure management might be completely alien to your experience if you are in an enterprise environment that is notoriously slow moving—there is no way, you might think, that data roles are going to change as dramatically as all that by the time you retire.

Well, maybe. I don't know where you work or how open to change your organization is. What I believe, though, is that more and more organizations and more and more industries are going to be like digital natives. There will be increasingly more openings for full stack data scientists, and data engineers will be as sought after as data scientists are today. This is because data engineers will be people who can do data science and know enough about infrastructure so as to be able to run their data science workloads on the public cloud. It will be worthwhile for you to learn data science terminology and data science frameworks, and make yourself more valuable for the next decade.

Collaboration

Even if you work in a company with strict separation of responsibilities, it can be helpful to know how the other teams do their work. This is because there are many artifacts that they create that you will use, or that you will create and they will use. Knowing their requirements and constraints will help you be more effective at communicating across organizational boundaries.

The various job roles related to data and machine learning are shown in [Figure 1-2](#). All these roles collaborate in creating a production machine learning model. Between data ingestion and the end-user interface, there are multiple handoffs. Every such

⁷ The words “need to know” are important here. It can sometimes be intimidating to see the breadth and depth of data science and despair of ever understanding everything. Here’s the truth: there is no one who knows the entire field in-and-out. Everyone is, at some level, glossing over some areas. Which areas? Areas that are not important to the problems that they are currently working on. This then gives you a strategy to approach data science—rather than try to learn topics (“I will learn RNNs this month”) or learn how to solve problems (“I will learn how to use AI to complete phrases”). Start with simple approaches, and stop once things become difficult and unintelligible. In most fields of AI, the simple approaches will get you quite far. Also, a deep understanding of the underlying mathematics is usually not required to implement a complex approach using frameworks like Keras.

handoff presents an opportunity for misunderstanding the requirements of the next stage or for an inability to take over what's been created at the previous stage.

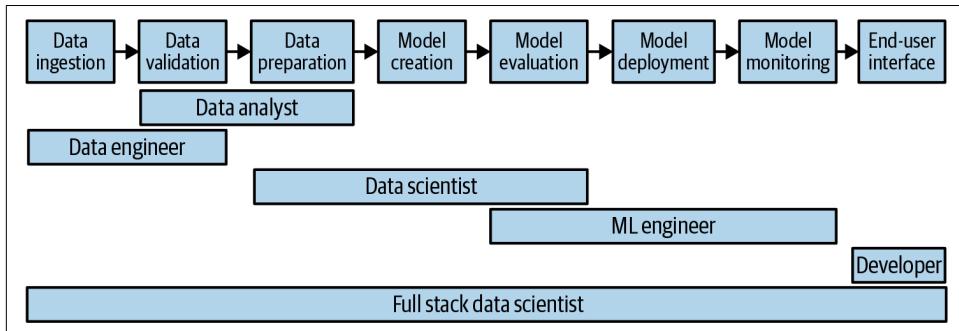


Figure 1-2. There are many job roles that need to collaborate to take a data science solution from idea to production. Every handoff carries a risk of failure.

Understanding the adjacent roles, the tools they work with, and the infrastructure that they use can help you reduce the chances of the baton getting dropped.

That said, it is very difficult to get a clean separation of responsibilities—the best organizations that I know, the ones that have hundreds to thousands of machine learning models in production, employ full stack data scientists that work on a problem from inception to production. They may have specialist data analysts, data engineers, data scientists, or ML engineers, but mostly in a maintenance capacity—the innovation tends to be done by the full stack folks. Even the full stack data scientists have areas in which they are stronger and areas where they collaborate with specialists.

Best Practices

This entire book consists of an extended case study. Solving a real-world, practical problem will help cut through all the hype that surrounds big data, machine learning, cloud computing, and so on. Pulling a case study apart and putting it together in multiple ways illuminates the capabilities and shortcomings of the various big data and machine learning tools that are available to you. A case study can help you identify the kinds of data-driven decisions that you can make in your business and illuminate the considerations behind the data you need to collect and curate, as well as the kinds of statistical and machine learning models you can use. I will attempt, throughout this book, to apply current best practices.

Simple to Complex Solutions

One of the ways that this book mirrors practice is that I use a real-world dataset to solve a realistic scenario and address problems as they come up. So, I will begin with a

decision that needs to be made and apply different statistical and machine learning methods to gain insight into making that decision in a data-driven manner. This will give you the ability to explore other problems and the confidence to solve them from first principles. As with most things, I will begin with simple solutions and work my way to more complex ones. Starting with a complex solution will only obscure details about the problem that are better understood when solving it in simpler ways. Of course, the simpler solutions will have drawbacks, and these will help to motivate the need for additional complexity.

One thing that I do not do, however, is to go back and retrofit earlier solutions based on knowledge that I gain in the process of carrying out more sophisticated approaches. In your practical work, however, I strongly recommend that you maintain the software associated with early attempts at a problem, and that you go back and continuously enhance those early attempts with what you learn along the way. Parallel experimentation is the name of the game. Due to the linear nature of a book, I don't do it, but I heartily recommend that you continue to actively maintain several models. Given the choice of two models with similar accuracy measures, you can then choose the simpler one—it makes no sense to use more complex models if a simpler approach can work with some modifications.⁸ Another reason to have multiple models is that a drop-in replacement is useful to have if you discover that the current production model drops in accuracy or is discovered to have unwanted behaviors.

Cloud Computing

Before I joined Google, I was a research scientist working on machine learning algorithms for weather diagnosis and prediction. The machine learning models involved multiple weather sensors, but were highly dependent on weather radar data. A few years ago, when we undertook a [project to reanalyze historical weather radar data](#) using the latest algorithms, it took us four years to do. However, more recently, my team was able to build rainfall estimates off the same dataset, but were able to traverse the dataset in about two weeks. You can imagine the pace of innovation that results when you take something that used to take four years and make it doable in two weeks.

Four years to two weeks. The reason was that much of the work as recently as five years ago involved moving data around. We'd retrieve data from tape drives, stage it to disk, process it, and move it off to make way for the next set of data. Figuring out

⁸ This goes by the name Principle of Parsimony or Occam's Razor and holds that the simplest explanation, with the fewest assumptions, is best. This is because simpler models are likely to fail less often because they depend on fewer assumptions. In engineering terms, another advantage of simpler models is that they tend to be less costly to implement.

what jobs had failed was time consuming, and retrying failed jobs involved multiple steps including a human in the loop. We were running it on a cluster of machines that had a fixed size. The combination of all these things meant that it took incredibly long periods of time to process the historical archive. After we began doing everything on the public cloud, we found that we could store all of the radar data on cloud storage and, as long as we were accessing it from virtual machines (VMs) in the same region, data transfer speeds were fast enough. We still had to stage the data to disks, carry out the computation, and bring down the VMs, but this was a lot more manageable. Simply lowering the amount of data movement between tape and disk and running the processes on many more machines enabled us to carry out processing much faster; to the credit of *elasticity* (the ability to seamlessly increase the number of resources we can assign to a job in the public cloud).

Was it more expensive to run the jobs on 10 times more machines than we did when we did the processing on premises? No, because the economics are usually in favor of renting on demand rather than buying the processing power outright, especially if you will not be using the machines 24-7. Whether you run 10 machines for 10 hours or 100 machines for 1 hour, the cost remains the same. Why not, then, get your answers in an hour rather than 10 hours?

In this book, we will do all our data science on Google Cloud in order to take advantage of the near-infinite scale that the public cloud offers.

Serverless

When we did our weather data preparation using cloud-based VMs, we were still not taking full advantage of what the cloud has to offer. We should have completely foregone the process of spinning up VMs, installing software on them, and looking for failed jobs—what we should have done was to use an autoscaling data processing framework such as BigQuery or Cloud Dataflow. Had we done that, we would have been able to run our jobs on thousands of machines and might have brought our processing time from two weeks to a few hours. Not having to manage any infrastructure is itself a huge benefit when it comes to trawling through terabytes of data. Having the data processing, analysis, and machine learning autoscale to thousands of machines is a bonus.

The key benefit of performing data engineering in the cloud is the amount of time that it saves you.⁹ You shouldn't need to wait days or months—instead, because many jobs are embarrassingly parallel, you can get your results in minutes to hours by having them run on thousands of machines. You might not be able to afford permanently

⁹ For your organization, any time you save translates to budget savings. You get more accomplished with a smaller budget.

owning so many machines, but it is definitely possible to rent them for minutes at a time. These time savings make autoscaled services on a public cloud the logical choice to carry out data processing.

Running data jobs on thousands of machines for minutes at a time requires fully managed services. Storing the data locally on the virtual machines or persistent disks as with the Apache Hadoop cluster doesn't scale unless you know precisely what jobs are going to be run, when, and where. You will not be able to downsize the cluster of machines if you don't have automatic retries for failed jobs and more importantly, shuffle the data around in remaining data nodes (assuming there is enough free space). The total computation time will be the time taken by the most overloaded worker unless you have dynamic task shifting among the nodes in the cluster. All of these point to the need for autoscaling services that dynamically resize the cluster, split jobs down into tasks, move tasks between compute nodes, and can rely on highly efficient networks to move data to the nodes that are doing the processing.

On Google Cloud Platform, the key autoscaling, fully managed, "serverless" services are BigQuery (for data analytics), Cloud Spanner (for databases), Cloud Dataflow (for data processing pipelines), Cloud Pub/Sub (for message-driven systems), Cloud Bigtable (for high-throughput ingest), Cloud Run or Cloud Functions (for applications, tasks), and Vertex AI (for machine learning).¹⁰ Using autoscaled services like these makes it possible for a data engineer to begin tackling more complex business problems because they have been freed from the world of managing their own machines and software installations whether in the form of bare hardware, virtual machines, or containers. Given the choice between a product that requires you to first configure a container, server, or cluster, and another product that frees you from those considerations, choose the serverless one. You will have more time to solve the problems that actually matter to your business.

A Probabilistic Decision

Imagine that you are about to take a flight and, just before the flight takes off from the runway (and you are asked to switch off your phone), you have the opportunity to send one last text message. It is past the published departure time and you are a bit anxious. [Figure 1-3](#) presents a graphic view of the scenario.

¹⁰ For a word that gets bandied about quite a lot, there is not much agreement on what exactly *serverless* means. In this book, I'll call a service serverless if users of the service have to supply only code and not have to manage the lifecycle of the machines that the code runs on.

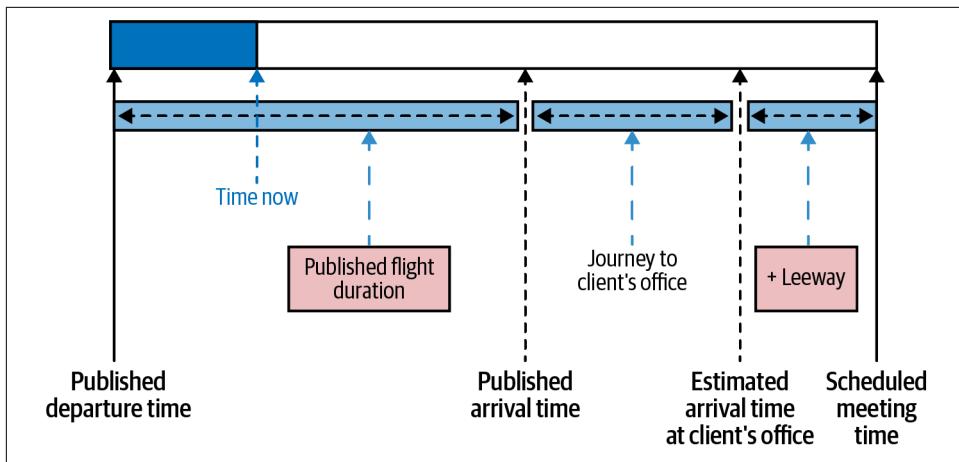


Figure 1-3. A graphic illustration of the case study: if the flight departs late, should the road warrior cancel the meeting?

The reason for your anxiety is that you have scheduled an important meeting with a client at its offices. As befits a rational data scientist,¹¹ you scheduled things rather precisely. You have taken the airline at its word with respect to when the flight would arrive, accounted for the time to hail a taxi, and used an online mapping tool to estimate the time to the client's office. Then, you added some leeway (say 30 minutes) and told the client what time you'd meet them. And now, it turns out that the flight is departing late. So, should you send a text informing your client that you will not be able to make the meeting because your flight will be late or should you not?

This decision could be made in many ways, including by gut instinct and using heuristics. Being very rational people, we (you and I) will make this decision informed by data. Also, we see that this is a decision made by many of the road warriors in our company day in and day out. It would be a good thing if we could do it in a systematic way and have a corporate server send out an alert to travelers about anticipated delays if we see events on their calendar that they are likely to miss. Let's build a data framework to solve this problem.

Probabilistic Approach

If we decide to make the decision in a data-driven way, there are several approaches we can take. Should we cancel the meeting if there is greater than a 30% chance that

¹¹ Perhaps I'm simply rationalizing my own behavior—if I'm getting to the departure gate with more than 15 minutes to spare at least once in about five flights, I decide that I must be getting to the airport too early and adjust accordingly. Fifteen minutes and 20% tend to capture my risk aversion. If you are wondering why my risk aversion threshold is not simply 15 minutes but includes an associated probabilistic threshold, read on.

you will miss it? Or should we assign a cost to postponing the meeting (the client might go with our competition before we get a chance to demonstrate our great product) versus not making it to a scheduled meeting (the client might never take our calls again) and minimize our expected loss in revenue? The probabilistic approach translates to risk, and many practical decisions hinge on risk. In addition, the probabilistic approach is more general because if we know the probability and the monetary loss associated with missing the meeting, it is possible to compute the expected value of any decision that we make. For example, suppose the chance of missing the meeting is 20% and we decide to not cancel the meeting (because 20% is less than our decision threshold of 30%). But there is only a 25% chance that the client will sign the big deal (worth a cool million bucks) for which you are meeting them. Because there is an 80% chance that we will make the meeting, the expected upside value of not canceling the meeting is $0.8 \times 0.25 \times 1$ million, or \$200,000. The downside value of not canceling is that we do miss the meeting. Assuming that the client is 90% likely to blow us off in the future if we miss a meeting with them, the expected downside is $0.2 \times 0.9 \times 0.25 \times 1$ million, or \$45,000. This yields an expected value of \$155,000 in favor of not canceling the meeting. We can adjust these numbers to come up with an appropriate probabilistic decision threshold.

Another advantage of a probabilistic approach is that we can directly take into account human psychology. You might feel frazzled if you arrive at a meeting only two minutes before it starts and, as a result, might not be able to perform at your best. It could be that arriving only two minutes early to a very important meeting doesn't feel like being on time. This obviously varies from person to person, but let's say that this time interval that you need to settle down is 15 minutes. You want to cancel a meeting for which you cannot arrive 15 minutes early. You could also treat this time interval as your personal risk aversion threshold, a final bit of headroom if you will. Thus, you want to arrive at the client's site 15 minutes before the meeting and you want to cancel the meeting if there is a less than 70% chance of doing that. This, then, is our decision criterion:

Cancel the client meeting if the probability of arriving 15 minutes early is 70% or less.

I've explained the 15 minutes, but I haven't explained the 70%. Surely, you can use the aforementioned model diagram ([Figure 1-3](#), in which we modeled our journey from the airport to the client's office), plug in the actual departure delay, and figure out what time you will arrive at the client's offices. If that is less than 15 minutes before the meeting starts, you should cancel! Where does the 70% come from?

Probability Density Function

It is important to realize that the model diagram ([Figure 1-3](#)) of times is not exact. The probabilistic decision framework gives you a way to treat this in a principled way. For example, although the airline company says that the flight duration is 127

minutes and publishes an arrival time, not all flights are exactly 127 minutes long. If the plane happens to take off with the wind, catch a tail wind, and land against the wind, the flight might take only 90 minutes. Flights for which the winds are all precisely wrong might take 127 minutes (i.e., the airline might be publishing worst-case scenarios for the route). Google Maps predicts journey times based on historical data, and the actual journeys by taxi will be centered around those times. Your estimate of how long it takes to walk from the airport gate to the taxi stand might be predicated on landing at a specific gate, and actual times may vary. So, even though the model depicts a certain time between airline departure and your arrival at the client site, this is not an exact number. The actual time between departure and arrival might have a distribution that looks like that shown in [Figure 1-4](#).

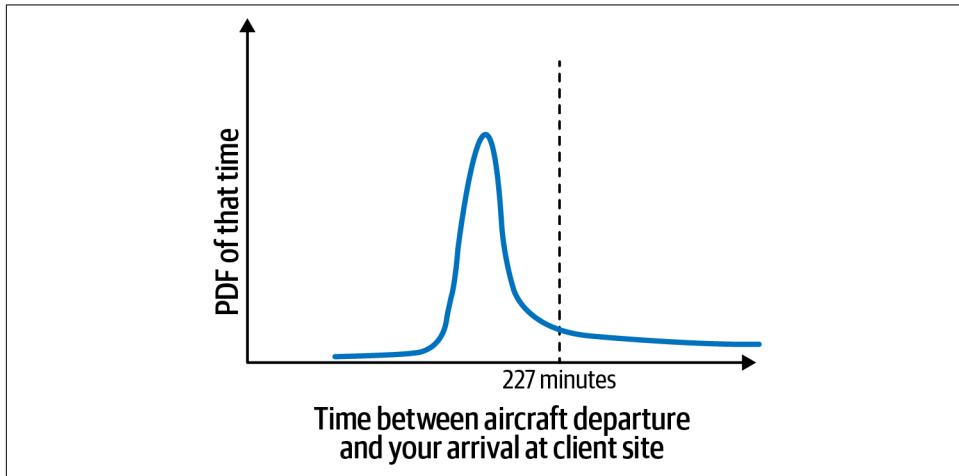


Figure 1-4. There are many possible values for the time differences between aircraft departure and your arrival at a client site, and the distribution of that value is called the probability density function.

The curve in [Figure 1-4](#) is referred to as the probability density function (abbreviated as the PDF). In fact, the PDF can be (and often is) greater than one. In order to get a probability, you will need to integrate the probability density function.¹² A simple way to do this integration is provided by the cumulative distribution function (CDF).

Cumulative Distribution Function

The cumulative probability distribution function of a value x is the probability that the observed value X is less than the threshold x . For example, you can get the

¹² To integrate a function is to compute the area under the curve of that function up to a specific x -value, as shown in [Figure 1-5](#).

cumulative distribution function (CDF) for 227 minutes by finding the fraction of flights for which the time difference is less than 227 minutes, as shown in [Figure 1-5](#).

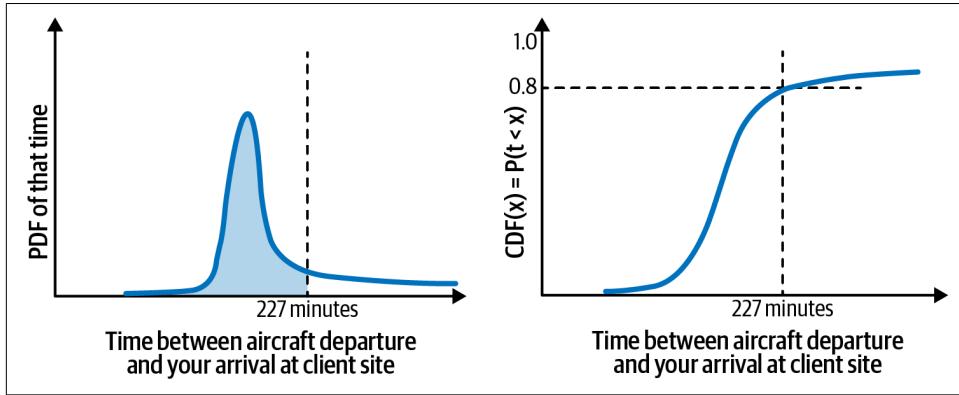


Figure 1-5. The CDF is the area under the curve of the PDF. It is easier to understand and keep track of than the PDF. In particular, it is bounded between 0 and 1, whereas the PDF could be greater than 1.

Let's interpret the graph in [Figure 1-5](#). What does a CDF (227 minutes) = 0.8 mean? It means that 80% of flights will arrive such that we will make it to the client's site in less than 227 minutes—this includes both the situation in which we can make it in 100 minutes and the situation in which it takes us 226 minutes. The CDF, unlike the PDF, is bounded between 0 and 1. The y-axis value is a probability, just not the probability of an exact value. It is, instead, the probability of observing all values less than that value.

Because the time to get from the arrival airport to the client's office is unaffected by the flight's departure delay,¹³ we can ignore it in our modeling. We can similarly ignore the time to walk through the airport, hail the taxi, and get ready for the meeting. So, we need only to find the probability of the arrival delay being more than 15 minutes. If that probability is 0.3 or more, we will need to cancel the meeting. In terms of the CDF, that means that the probability of arrival delays of less than 15 minutes has to be at least 0.7, as presented in [Figure 1-6](#).

Thus, our decision criteria translate to the following:

Cancel the client meeting if the CDF of an arrival delay of 15 minutes is less than 70%.

¹³ This is a simplifying assumption—if the flight was supposed to arrive at 2 p.m., and instead arrives at 4 p.m., the traveler is more likely to hit rush hour traffic.

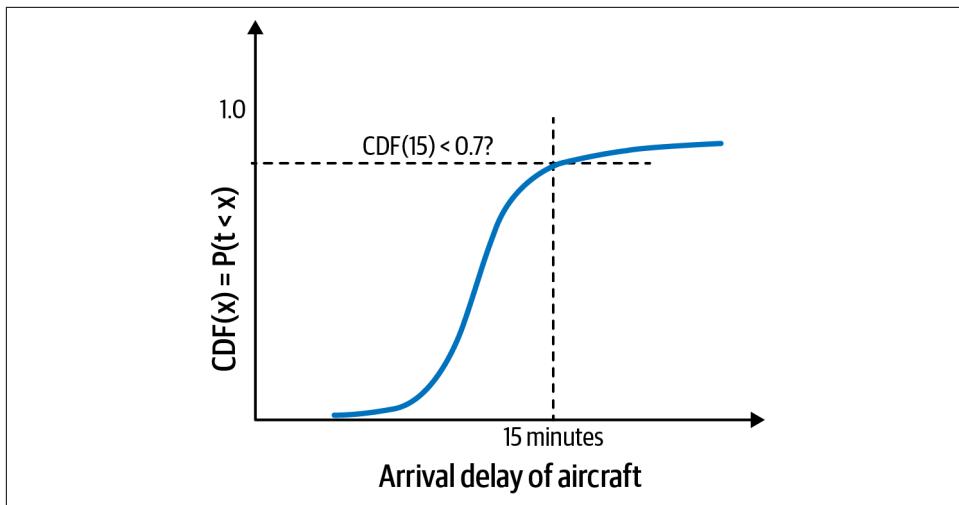


Figure 1-6. Our decision criterion is to cancel the meeting if the CDF of an arrival delay of 15 minutes is less than 70%. Loosely speaking, we want to be 70% sure of the aircraft arriving no more than 15 minutes late.

The rest of this book is going to be about building data pipelines that enable us to compute the CDF of arrival delays using statistical and machine learning models. From the computed CDF of arrival delays, we can look up the CDF of a 15-minute arrival delay and check whether it is less than 70%.

Choices Made

What data do we need to predict the probability of a specific flight delay? What tools shall we use? Should we use Hadoop? BigQuery? Should we do it on my laptop or should we do it in the public cloud? The question about data is easily answered—we will use [historical flight arrival data](#) published by the US Bureau of Transportation Statistics, analyze it, and use it to inform our decision. Often, a data scientist would choose the best tool based on their experience and just use that one tool to help make the decision, but here, I will take you on a tour of several ways that we could carry out the analysis. This will also allow us to model best practice in the sense of picking the simplest tool and analysis that suffices.

Choosing Cloud

On a cursory examination of the data, we discover that there were more than 30.6 million flights in 2015–2019.¹⁴ My laptop, nice as it is, is not going to cut it. We will do the data analysis on the public cloud. Which cloud? We will use the Google Cloud Platform (GCP). Although some of the tools we use in this book (notably Hadoop, Spark, Beam, TensorFlow, etc.) are available on other cloud platforms, the managed services I use (BigQuery, Cloud Dataproc, Cloud Dataflow, Vertex AI, etc.) are specific to GCP. Using GCP will allow me to avoid fiddling around with virtual machines and machine configuration and focus solely on the data analysis. Also, I do work at Google, so this is the platform I know best.

Not a Reference Book

This book is not an exhaustive look at data science—there are other books (often based on university courses) that do that. It is also not a reference book on Google Cloud—the documentation is much more timely and comprehensive. Instead, this book allows you to look over my shoulder as I solve one particular data science problem using a variety of methods and tools. I promise to be quite chatty and tell you what I am thinking and why I am doing what I am doing. Instead of presenting you with fully formed solutions and code, I will show you intermediate steps as I build up to a solution.

This learning material is presented to you in three forms:

- This book that you are reading.
- The code referenced throughout the book [on GitHub](#). Note in particular, the *README.md* file in each folder of the GitHub repository.
- Labs with instructions that allow you to try the code of this book in a sandbox environment, available at <https://qwiklabs.com>.

Rather than simply read this book cover to cover, I strongly encourage you to follow along with me by also taking advantage of the code. After reading each chapter, or major section in each chapter, try to repeat what I did, referring to the code if something's not clear.

¹⁴ Yes, this is the second edition of the book, published in 2022. The first edition of the book used only 2015 data. Here, I'll use 2015–2019. I stopped with 2019 because 2020 was the year of the COVID-19 pandemic, and flights were rather spotty.

Getting Started with the Code

To begin working with the code, follow these steps:

- **Sign up** for the free trial if you haven't already done so. Otherwise, use your existing GCP account.
- Create a new project and give it any name you want. I suggest calling it ds-on-gcp. GCP will assign a unique project ID to your project (see [Figure 1-7](#)). You will need to provide this unique ID whenever you do anything that is billable. Once you are finished working through this book, simply delete the project to stop getting billed.

The screenshot shows the 'Create Project' dialog box. It has two main sections: 'Project name *' containing 'ds-on-gcp' and 'Project ID' containing 'elegant-pipe-343008'. A red arrow points from the text 'Use this unique identifier whenever a script or program asks for a project ID.' to the 'Project ID' field. Below these are 'Location *' (set to 'No organization') and 'Parent organization or folder'. At the bottom are 'CREATE' and 'CANCEL' buttons.

Figure 1-7. When you create a new project, GCP will assign it a unique project identifier. Use this unique identifier whenever a script or program asks for a project ID. You will also be able to get the unique identifier from the dashboard (see [Figure 1-8](#)).

- Open Cloud Shell, your terminal access to GCP. To open Cloud Shell, on the menu bar, click the Activate Cloud Shell icon, as shown in [Figure 1-8](#). Even though the web console is very nice, I typically prefer to script things rather than go through a GUI. To me, web GUIs are great for occasional and/or first-time use, but for repeatable tasks, nothing beats the terminal.

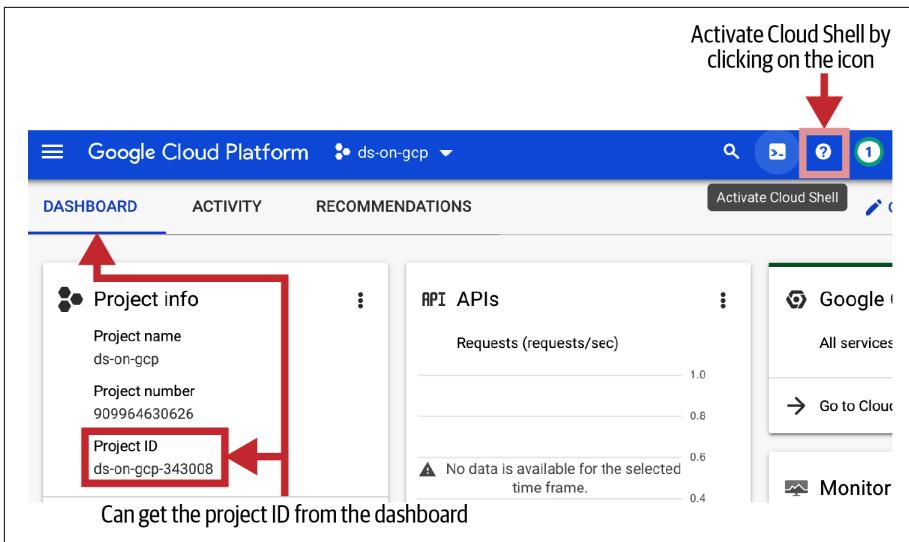


Figure 1-8. Activate Cloud Shell by clicking on the highlighted icon in the top right corner of the GCP web console. Note that the unique project identifier can be obtained at any time from the dashboard.



Cloud Shell is a micro-VM that is alive for the duration of the browser window and gives you terminal access to the micro-VM. Close the browser window, and the micro-VM goes away. The Cloud Shell VM is free and comes loaded with many of the tools that developers on Google Cloud Platform will need. For example, it has Python, Git, the Google Cloud SDK, and Orion (a web-based code editor) installed on it. Although the Cloud Shell VM is ephemeral, it is attached to a persistent disk that is tied to your user account. Files that you store in your home directory are saved across different Cloud Shell sessions.

- In the Cloud Shell window, git clone my repository by typing the following:

```
git clone \
  https://github.com/GoogleCloudPlatform/data-science-on-gcp
  cd data-science-on-gcp
```

Because the Git repository was checked out to the home directory of the Cloud Shell micro-VM, it will be persistent across browser sessions.

- Note that there is a directory corresponding to each chapter of this book (other than Chapters 1 and 12). In each directory, you will find a *README.md* file with directions on how to replicate the steps in that chapter.



Do not copy-paste code snippets from the book. Read the chapters and *then* try out the code by following the steps in each chapter's *README.md* using the *code in the repository*. I recommend that you not copy-paste from electronic versions of this book.

- The book is written for readability, not for completeness. Some flags to cloud tools may be omitted so that we can focus on the key aspect being discussed. The GitHub code will have the full command.
- The GitHub repo will be kept up to date with new versions of cloud tools, Python, etc.
- When following along in the book, it's easy to miss a step.
- Copy-paste of special characters from PDF is problematic.

Developing Locally

If you prefer to do development on your local machine (rather than in Cloud Shell), you will need to install three pieces of software (all three are already present on Cloud Shell, so this is only if you wish to develop on your own laptop):

1. Python version 3.6 or higher
2. The [Google Cloud SDK](#)
3. The version control software Git

That's it. You are now ready to follow along with me. As you do, remember that you need to change my project ID to the ID of your project (you can find this on the dashboard of the Google Cloud web console, as shown in [Figure 1-8](#)) and my bucket-name to your bucket on Cloud Storage (you will create this in [Chapter 2](#); we'll introduce buckets at that time).

Agile Architecture for Data Science on Google Cloud

I will introduce Google Cloud products and technologies as we go along. In this section, I will provide a high-level overview of why I choose what I choose. Do not worry if you don't recognize the names of these technologies (e.g., data warehouse) or products (e.g., BigQuery) since we will cover them in detail as we go along.

What Is Agile Architecture?

One of the [principles of Agile software](#) is that simplicity, by which we mean maximizing the amount of work not done, is essential. Another is that requirements change frequently, and so flexibility is important. An Agile architecture, therefore, is one that gives you:

- Speed of development. You should be able to go from idea to deployment as quickly as possible.
- Flexibility to quickly implement new features. Sometimes speed comes at the expense of flexibility—the architecture might shoehorn you into a very limited set of use cases. You don't want that.
- Low-maintenance. Don't spend your time managing infrastructure.
- Autoscaling and resiliency so that you are not spending your time monitoring infrastructure.

What does such an architecture look like on Google Cloud when it comes to Data Analytics and AI? It will use low-code and no-code services (pre-built connectors, automatic replication, ELT [extract-load-transform], AutoML) so that you get speed of development. For flexibility, the architecture will allow you to drop down to developer-friendly, powerful code (Apache Beam, SQL, TensorFlow) whenever needed. These will run on serverless infrastructure (Pub/Sub, Dataflow, BigQuery, Vertex AI) so that you get low-maintenance, autoscaling, and resiliency.

No-Code, Low-Code

When it comes to architecture, choose no-code over low-code and low-code over writing custom code. Rather than writing ETL (extract-transform-load) pipelines to transform the data you need before you land it into BigQuery, use pre-built connectors to directly land the raw data into BigQuery (see [Figure 1-9](#)). Then, transform the data into the form you need using SQL views directly in the data warehouse. You will be a lot more agile if you choose an ELT approach over an ETL approach.

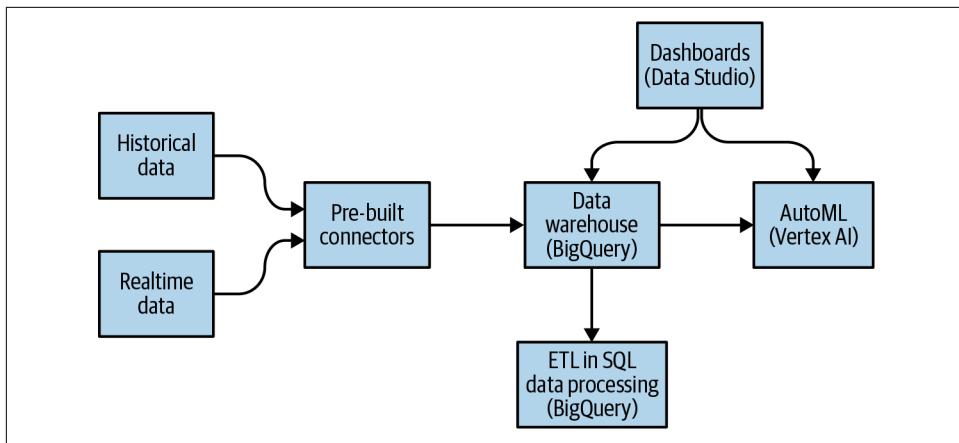


Figure 1-9. Agile architecture for most use cases.

Another place is when you choose your ML modeling framework. Don't start with custom TensorFlow models. Start with AutoML. That's no-code. You can invoke AutoML directly from BigQuery, avoiding the need to build complex data and ML pipelines. If necessary, move on to pre-built models from TensorFlow Hub and pre-built containers on Vertex AI. That's low-code. Build your own custom ML models only as a last resort.

Use Managed Services

You will want to be able to drop down to code if the low-code approach is too restrictive. Fortunately, the no-code architecture described previously is a subset of the full architecture, shown in [Figure 1-10](#), that gives you all the flexibility you need.

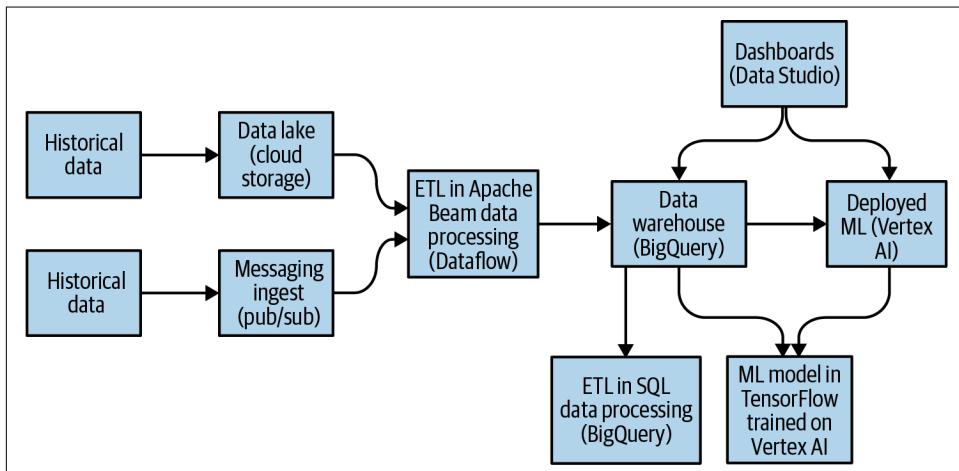


Figure 1-10. Agile architecture for analytics and AI.

When the use case warrants it, you will have the full flexibility of Apache Beam, SQL, and TensorFlow. This is critical—for use cases where the ELT + AutoML approach is too restrictive, you have the ability to drop to a ETL/Dataflow + Keras/Vertex approach.

Best of all, the architecture is unified, so you are not maintaining two stacks. Because the first architecture is a subset of the second, you can accomplish both easy and hard use cases in a unified way.

It is this architecture that we build in this book.

Summary

A key goal of data analysis is to be able to provide data-driven guidance toward making accurate decisions systematically. Ideally, this guidance can be provided as a service, and providing as a service gives rise to questions of service quality—both in terms of the accuracy of the guidance and the reliability, latency, and security of the implementation.

A data engineer needs to be able to go from designing data-based services and building statistical and machine learning models to implementing them as reliable, high-quality services. This has become easier with the advent of cloud services that provide an autoscaling, serverless, managed infrastructure. Also, the wide availability of data science tools has made it so that you don't need to be a specialist in data science to create a statistical or machine learning model. As a result, the ability to work with data has spread throughout an enterprise—no longer is it a restricted skill.

Our case study involves a traveler who needs to decide whether to cancel a meeting depending on whether the flight they are on will arrive late. The decision criterion is that the meeting should be canceled if the probability of arriving within 15 minutes of the scheduled time is less than 70%. To estimate the probability of this arrival delay, we will use historical data from the US Bureau of Transportation Statistics.

To follow along with me throughout the book, create a project on Google Cloud Platform and a clone of the GitHub repository of the source code listings in this book. Alternatively, try the code of this book in a sandbox environment using Qwiklabs. The folder for each of the chapters in GitHub contains a *README.md* file that lists the steps to be able to replicate what I do in the chapters. So, if you get stuck, refer to those README files.



Incidentally, the footnotes in this book are footnotes because they break the flow of the chapter. Some readers of the first edition noted that they realized only toward the middle of the book that many of the footnotes contained useful information. So, this might be a good time to read the footnotes if you have been skipping them.

Suggested Resources

What is data science on Google Cloud? What does the toolkit consist of? The [data science website in Google Cloud](#) contains a set of whitepapers and reference guides that address these topics. Bookmark this page and use it as a starting point for everything data science on GCP.

There are five key autoscaling, fully managed, serverless products for data analytics and AI on Google Cloud. We'll cover them later in the book, but these videos and articles are a great starting point if you want to dive deeper immediately:

- BigQuery is the serverless data warehouse that forms the heart of most data architectures built in Google Cloud. I recommend watching "[Google BigQuery Introduction by Jordan Tigani](#)", one of the founding engineers of BigQuery, even though it is a few years old now.
- Dataflow is the execution service for batch and streaming pipelines written using Apache Beam. Start with "[What Is Dataflow?](#)" by Google Cloud Tech, a 5-minute video that introduces what Dataflow is and how it works. This is part of the [Google Cloud Drawing Board](#) series of videos—they are quick and informative ways to learn about various topics on Google Cloud.
- Pub/Sub is the global messaging service that can be used for use cases ranging from user interaction and real-time event distribution to refreshing distributed caches. Start from the [overview documentation page](#). All Google Cloud products have an overview page that can serve as a launching point to learning not only what a product does but also what it can be used for and how to choose between it and other alternatives.
- Cloud Run provides an autoscaling, serverless platform for containerized applications. You can use it for all kinds of automation and lightweight data transformation. The best way to learn Cloud Run is to try it out, and Qwiklabs provides a [great set of hands-on labs](#) in a sandbox environment. While you are there, check out the Catalog for other quests and labs on the topic of choice.
- Vertex AI is the end-to-end ML development, deployment, and automation platform on Google Cloud. A good place to learn about it is to watch the video that

accompanied its announcement at Google I/O, “Build End-to-End Solutions with Vertex AI”, by the Google Cloud Tech YouTube channel.

There are two key fully managed transaction processing databases on Google Cloud:

- Bigtable is a distributed NoSQL database. You could, of course, learn about it from Google Cloud Tech’s [overview video](#), “[What Is Cloud Bigtable?](#)”, or [the Bigtable documentation](#), but I recommend you read the famous research paper that introduced the idea to the world: Fay Chang et al., “[Bigtable: A Distributed Storage System for Structured Data](#)”, *7th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, USENIX (2006): 205–218.
- Spanner is a distributed SQL database that provides global strong consistency and five nines (99.999%) availability—something that greatly simplifies your architecture if you are in a domain like banking or gaming where you have concurrent users all over the world. “[Why You Should Use Google’s Cloud Spanner for Your Next Game](#)”, a 2019 blog post by Miles Ward, CTO of Google’s partner SADA, is a great starting point for Spanner and Spanner best practices.

For more on probability as applied in information theory and artificial intelligence, read either [Chapter 3 of Deep Learning](#), “Probability and Information Theory,” by Ian Goodfellow et al. (MIT Press) or a [summary of that chapter](#) by William Green on *Medium*. The foundation of information theory was laid by [Claude Shannon in a classic 1948 paper](#). He is also famous for perhaps [the most influential masters thesis in history](#)—showing how to use Boolean algebra to test circuit designs without even building the circuits in the first place.

CHAPTER 2

Ingesting Data into the Cloud

In Chapter 1, we explored the idea of deciding whether to cancel a meeting in a data-driven way. We decided on a probabilistic decision criterion: to cancel the meeting with a client if the probability of the flight arriving within 15 minutes of the scheduled arrival time was less than 70%. To model the arrival delay given a variety of attributes about the flight, we need historical data that covers a large number of flights. Historical data that includes this information from 1987 onward is available from the [US Bureau of Transportation Statistics](#) (BTS). One of the reasons that the government captures this data is to monitor the fraction of flights by a carrier that are on-time (defined as flights that arrive less than 15 minutes late), so as to be able to hold airlines accountable.¹ Because the key use case is to compute on-time performance, the dataset that captures flight delays is called [Airline On-Time Performance Data](#). That's the dataset we will use in this book.



All of the code snippets in this chapter are available in the folder [02_ingest](#) of the book's [GitHub repository](#). See the last section of Chapter 1 for instructions on how to clone the repository, and see the *README.md* file in the *02_ingest* directory for instructions on how to do the steps described in this chapter.

Airline On-Time Performance Data

For nearly 40 years, all major US air carriers have been required to file statistics about each of their domestic flights with the BTS. The data they are required to file includes

¹ See, for example, this [US Senate committee report](#) on the proposed Airline Customer Service Improvement Act. The bill referenced in the report was not enacted into law, but it illustrates Congress's monitoring function based on the statistics collected by the Department of Transportation.

the scheduled departure and arrival times as well as the actual departure and arrival times. From the scheduled and actual arrival times, the arrival delay associated with each flight can be calculated. Therefore, this dataset can give us the true value or “label” for building a model to predict arrival delay.

The actual departure and arrival times are defined rather precisely, based on when the parking brake of the aircraft is released and when it is later reactivated at the destination. The rules even go as far as to define what happens if the pilot forgets to engage the parking brake—in that case, the time that the passenger door is closed or opened is used instead. Because of the precise nature of the rules, and the fact that they are enforced, we can treat arrival and departure times from all carriers uniformly. Had this not been the case, we would have to dig deeper into the quirks of how each carrier defines “departure” and “arrival,” and do the appropriate translations.² Good data science begins with such standardized, repeatable, trustable data collection rules; you should use the BTS’s very well-defined data collection rules as a model when creating standards for your own data collection, whether it is log files, web impressions, or sensor data that you are collecting. The airlines report this particular data monthly, and it is collated by the BTS across all carriers and published as a free dataset on the web.

In addition to the scheduled and actual departure and arrival times, the data includes information such as the origin and destination airports, flight numbers, and nonstop distance between the two airports. It is unclear from the documentation whether this distance is the distance taken by the flight in question or whether it is simply a pre-computed distance—if the flight needs to go around a thunderstorm, is the distance in the dataset the actual distance traveled by the flight or the great-circle distance between the airports?³ This is something that we will need to examine—it should be easy to ascertain whether the distance between a pair of airports remains the same or changes. In addition, a flight is broken into three parts ([Figure 2-1](#))—taxi-out duration, air time, and taxi-in duration—and all three time intervals are reported.

2 For example, weather radar data from before 2000 had timestamps assigned by a radar engineer. Essentially, the engineers would look at their wristwatches and enter a time into the radar products generator. Naturally, this was subject to all kinds of human errors—dates could be many hours off. The underlying problem was fixed by the introduction of network clocks to ensure consistent times between all the radars on the US weather radar network. When using historical weather data, though, time correction is an important preprocessing step.

3 The shortest path between two points on the globe is an arc that passes through the two points and whose focus point is the center of the globe.

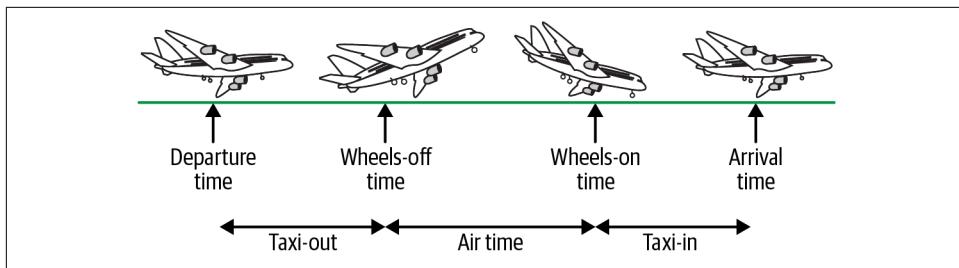


Figure 2-1. A flight is broken into three parts: taxi-out duration, air time, and taxi-in duration.

Knowability

Before we get started with ingesting data, we need to decide what it is that we have to ingest into our model. There are two potential traps—causality and training-serving skew (I'll define them shortly). We should take care to avoid these problems during the ingest phase, in order to save us a lot of heartburn later.

Causality

The causality principle boils down to this key question: what data will we be able to provide to the model at the time that we need to make predictions? If we won't know some piece of information about the flight during prediction, we cannot use that information as an input during training.

Some of the fields in the dataset could form the inputs to our model to help us predict the arrival delay as a function of these variables. Some, but not all. Why? It should be clear that we cannot use taxi-in duration or actual flight distance because at the time the aircraft is taking off, which is when we want to make our decision on whether to cancel the meeting, we will not know either of these things. The in-air flight time between two airports is not known *a priori* given that pilots have the ability to speed up or slow down. Thus, even though we have these fields in our historical dataset, we should not use them in our prediction model. This is called a causality constraint.

The causality constraint is one instance of a more general principle. Before using any field as input to a model, we should consider whether the data will be known at the time we want to make the decision. It is not always a matter of logic as with the taxi-in duration. Sometimes, practical considerations such as security (is the decision maker allowed to know this data?), the latency between the time the data is collected and the time it is available to the model, and the cost of obtaining the information also play a part in making some data unusable. At the same time, it is possible that approximations might be available for fields that we cannot use because of causality—even though, for example, we cannot use the actual flight distance, we should be able to use the great-circle distance between the airports in our model.

Similarly, we might be able to use the data itself to create approximations for fields that are obviated by the causality constraint. Even though we cannot use the actual taxi-in duration, we can use the mean taxi-in duration of this flight at this airport on previous days, or the mean taxi-in duration of all flights at this airport over the past hour, to approximate what the taxi-in duration might be. Over the historical data, this could be a simple batch operation after grouping the data by airport and hour. When predicting in real time, though, this will need to be a moving average on streaming data. Indeed, approximations of unknowable data will be an important part of our models.

Training–Serving Skew

A training–serving skew is the condition in which you use a variable that’s computed differently in your training dataset than in the production model. For example, suppose that you train a model with the distance between cities in miles, but when you predict, the distance that you receive as input is actually in kilometers. That is obviously a bad thing and will result in a bad result from the model because the model will be providing predictions based on the distances being 1.6 times their actual value. Although it is obvious in clear-cut cases such as unit mismatches, the same principle (that the training dataset has to reflect what is done to inputs at prediction time) applies to more subtle scenarios as well.

For example, suppose we determine whether the flight is on a weekday or a weekend and use it as an input to the model. We need to ensure that this calculation is carried out precisely the same way during both training and prediction. For example, if we use the wheels-off time during training but the departure time during prediction, we will suffer from training–serving skew. If the time library we use during training treats the timestamp as being local time, but the library we use during prediction treats the timestamp as being in Coordinated Universal Time (UTC), we will run into training–serving skew.

As our models become increasingly sophisticated—and more and more of a black box—it will become extremely difficult to troubleshoot errors that are caused by a training–serving skew. This is especially true if the code bases for computing inputs for training and during prediction are different and begin to diverge over time. We will always attempt to design our systems in such a way that the possibilities of a training–serving skew are minimized. In particular, we will gravitate toward solutions in which we can use the same code in training (building a model) as in prediction.

The dataset includes codes for the airports (such as ATL for Atlanta) from which and to which the flight is scheduled to depart and land. Planes might land at an airport other than the one they are scheduled to land at if there are in-flight emergencies or if weather conditions cause a deviation. In addition, the flight might be canceled. It is important for us to ascertain how these circumstances are reflected in the dataset—

although they are relatively rare occurrences, our analysis could be adversely affected if we don't deal with them in a reasonable way. The way we deal with these out-of-the-ordinary situations also must be consistent between training and prediction.

The dataset also includes airline codes (such as AA for American Airlines), but it should be noted that airline codes can change over time (for example, United Airlines and Continental Airlines merged and the combined entity began reporting as United Airlines in 2012). If we use airline codes in our prediction, we will need to cope with these changes in a consistent way, too.

Downloading Data

As of August 2021, there were nearly 200 million records in the on-time performance dataset, with records starting in 1987. The last available data was June 2021, indicating that there is more than a month's delay in updating the dataset—this is going to be important when we automate the process of getting new data.

This is not the most helpful way to provide data for download. For one thing, the data can be downloaded only one month at a time. For another, going through a web form is pretty error-prone. Imagine that you want to download all of the data for 2015. In that scenario, you'd painstakingly select the fields you want for January 2015, submit the form, and then have to repeat the process for February 2015. If you forgot to select a field in February, that field would be missing, and you wouldn't know until you began analyzing the data!

Obviously, we can script the download to make it less tiresome and ensure consistency.⁴ However, it is better to download all the raw data, not just a few selected fields. Why? Won't the files be larger if we ask for all the fields? Won't larger files take longer to download?

In this book, our model will use input fields drawn mostly from this dataset, but where feasible and necessary, we will include other datasets such as airport locations and weather. We can download the on-time performance data from the BTS website as comma-separated value (CSV) files. The web interface requires you to filter by geography and period, as illustrated in [Figure 2-2](#). The data itself is offered in two ways: one with all the data in a zipped file and the other containing just the fields that we select in the form.

⁴ Indeed, scripting the field selection and download is what I did in the first edition of the book. If interested, see the [select-and-download code](#) in GitHub in the branch `edition1_tf2`—the key thing is that the web request sends the selected fields inside the POST request and handles the resulting client-side redirect to obtain the ZIP file that is created on demand.

Option 1: all the raw data
But only one month at a time

: Reporting Carrier On-Time Performance (1987-present)

[Data Tables](#) [Table Contents](#)

Download Instructions Filter Geography Filter Year Filter Period

Latest Available Data: November 2021 All 2021 January

Prezipped File % Missing Documentation Terms [Download](#)

Field Name	Description	Support Table
Time Period		
<input type="checkbox"/> Year	Year	
<input type="checkbox"/> Quarter	Quarter (1-4)	Get Lookup Table
<input type="checkbox"/> Month	Month	Get Lookup Table
<input type="checkbox"/> DayofMonth	Day of Month	
<input type="checkbox"/> DayOfWeek	Day of Week	Get Lookup Table
<input type="checkbox"/> FlightDate	Flight Date (yyyyMMdd)	
Airline		
<input type="checkbox"/> Reporting_Airline	Unique Carrier Code. When the same code has been used by multiple carriers, a numeric suffix is used for earlier users, for example, PA, PA(1), PA(2). Use this field for analysis across a range of years.	Get Lookup Table

Option 2: select the fields we want

Figure 2-2. The BTS web interface to download the flights on-time arrival dataset.

Hub-and-Spoke Architecture

Yes, the files will be larger if we download all the fields using the static link. But there is a significant drawback to doing preselection. In order to support the interactive capability of selecting fields, the BTS does server-side processing—it extracts the fields we want, creates a custom ZIP file, and makes the ZIP file available for download. This would make our code reliant on the BTS servers having the necessary uptime and reliability.⁵ Avoiding the server-side processing should help reduce this dependency.⁶

⁵ Over the last 5 years, I have observed that the BTS server that does this ZIP file creation is frequently down. I know, I know. Ideally, they'd use a public cloud to host their website and/or data, but *you* try telling the US government what to do.

⁶ Another thing I am doing to limit the dependence on the BTS website is to host the ZIP files on Google Cloud and have my code hit the Google Cloud server. The code by default will not hit the BTS server anymore. The original BTS URL is still present in the code, just commented out, so change it back if you want to try it out.

An even more salient reason is that best practice in data engineering now is to build ELT (extract-load-transform) pipelines, rather than ETL (extract-transform-load) pipelines. What this means is that we will extract the data from BTS and immediately load the data into a data warehouse rather than rely on the BTS server to do transformation for us before loading it into Google Cloud. This point is important. The recommended modern data architecture is to minimize the preprocessing of data—instead, land all available data as-is into the enterprise data warehouse (EDW) and then carry out whatever transformations are necessary for different use cases (see [Figure 2-3](#)). This is called a hub-and-spoke architecture, with the EDW functioning as the hub.

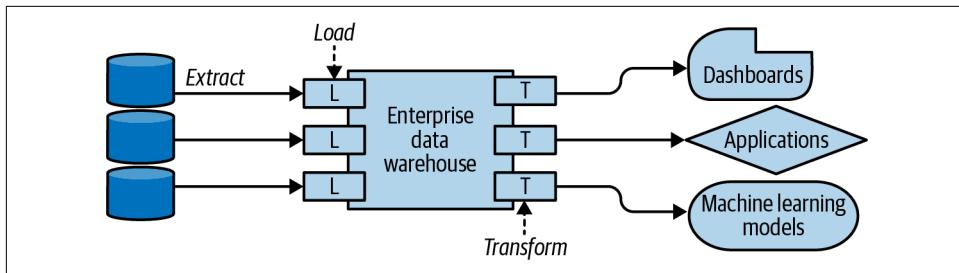


Figure 2-3. The recommended data architecture, whenever you can make it work, is the hub-and-spoke architecture.

Dataset Fields

Even though I'm going to download all the fields, it's worthwhile reading through the [column descriptions provided by BTS](#) to learn more about the dataset and get a preliminary idea about what fields are relevant to our problem and whether there are any caveats. For example, [Table 2-1](#) shows three ways in which the airline is recorded. Which of these should we use?

Table 2-1. The airline operating the flight is recorded in three separate columns

Column name	Description (copied from BTS website)
Reporting_Airline	Unique Carrier Code. When the same code has been used by multiple carriers, a numeric suffix is used for earlier users; for example, PA, PA(1), PA(2). Use this field for analysis across a range of years.
DOT_ID_Report ing_Airline	An identification number assigned by the US Department of Transportation (DOT) to identify a unique airline (carrier). A unique airline (carrier) is defined as one holding and reporting under the same DOT certificate regardless of its Code, Name, or holding company/corporation.
IATA_CODE_Report ing_Airline	Assigned by the International Air Transport Association (IATA) and commonly used to identify a carrier. Because the same code might have been assigned to different carriers over time, the code is not always unique. For analysis, use the Unique Carrier Code.

It's clear that we could use either the `Reporting_Airline` or the `DOT_ID_Report ing_Airline` since they are both unique. Ideally, we'd use whichever one of these corresponds to the common nomenclature (for example, UA or United Airlines). Fortunately, the BTS provides an Analysis link for the columns (see [Figure 2-4](#)), so we don't have to wait until we explore the data to make this decision. It turns out that the `Reporting_Airline` is what we want—the IATA code consists of the number 19977 for United Airlines whereas the `Reporting_Airline` is UA as we would like.

Airline		
Reporting_Airline	Unique Carrier Code. When the same code has been used by multiple carriers, a numeric suffix is used for earlier users, for example, PA, PA(1), PA(2). Use this field for analysis across a range of years.	Analysis
DOT_ID_Report ing_Airline	An identification number assigned by US DOT to identify a unique airline (carrier). A unique airline (carrier) is defined as one holding and reporting under the same DOT certificate regardless of its Code, Name, or holding company/corporation.	Analysis
IATA_CODE_Report ing_Airline	Code assigned by IATA and commonly used to identify a carrier. As the same code may have been assigned to different carriers over time, the code is not always unique. For analysis, use the Unique Carrier Code.	
Tail_Number	Tail Number	
Flight_Number_Report ing_Airline	Flight Number	

Figure 2-4. The BTS provides an Analysis link for some of the columns. These provide a handy way to learn what values a field can take.

The first thing to do in any real-world problem where we are fortunate enough to be provided documentation is to read it!⁷ After reading through the descriptions of the 100-plus fields in the dataset, there are a few fields that appear relevant to the problem of training, predicting, or evaluating flight arrival delay. [Table 2-2](#) presents the fields I shortlisted.

Table 2-2. Selected fields from the airline on-time performance dataset downloaded from the BTS (there is a separate table for each month)

Column name	Description (copied from BTS website)
FlightDate	Flight date (yyyymmdd)
Reporting_Airline	Unique Carrier Code. When the same code has been used by multiple carriers, a numeric suffix is used for earlier users; for example, PA, PA(1), PA(2). Use this field for analysis across a range of years.
Origin	Origin airport

⁷ Normally, you will also have to verify the description since data dictionaries are quite often outdated. The BTS documentation doesn't have this problem—it is correct and corresponds to the version of the data that BTS publishes.

Column name	Description (copied from BTS website)
Dest	Destination airport
CRSDepTime	Computerized reservation system (CRS) departure time (local time: hhmm)
DepTime	Actual departure time (local time: hhmm)
DepDelay	Difference in minutes between scheduled and actual departure time. Early departures show negative numbers.
TaxiOut	Taxi-out duration, in minutes
WheelsOff	Wheels-off time (local time: hhmm)
WheelsOn	Wheels-on time (local time: hhmm)
TaxiIn	Taxi-in duration (minutes)
CRSArrTime	CRS arrival time (local time: hhmm)
ArrTime	Actual arrival time (local time: hhmm)
ArrDelay	Difference in minutes between scheduled and actual arrival time. Early arrivals show negative numbers.
Cancelled	Cancelled flight indicator (1 = Yes)
CancellationCode	Specifies the reason for cancellation
Diverted	Diverted flight indicator (1 = Yes)
Distance	Distance between airports (miles)

Separation of Compute and Storage

There are essentially three options when it comes to processing large datasets (see [Table 2-3](#)), and all three are possible on GCP. Which one you use depends on the problem—in this book, we’ll use the third option because it is the most flexible. However, this option requires a bit of preplanning on our part—we will have to store our data in Google Cloud Storage and in Google BigQuery. To see why we choose this, let’s consider the other two options also.

Table 2-3. How to choose between scaling up, scaling out with data sharding, or scaling out with data in situ

Option	Performance and cost	Required platform capabilities	How to implement on Google Cloud Platform	Example use case
Scaling up	Expensive on both compute and storage; fast, but limited to capability of most powerful machine	Very powerful machines; ability to rent machines by the minute; attachable persistent SSDs	Compute Engine with persistent SSDs Vertex AI Notebooks	Job that requires rereading of data (e.g., training an ML model)

Option	Performance and cost	Required platform capabilities	How to implement on Google Cloud Platform	Example use case
Scaling out with sharding	High storage costs; inexpensive compute; add machines to achieve desired speed, but limited to ability to preshard the data on a cluster of desired size	Data local to the compute nodes; attachable persistent SSDs	Cloud Dataproc (with Spark) and Hadoop Distributed File System (HDFS)	Light compute on a splittable dataset (e.g., creating a search index from thousands of documents). Many data analytics use cases used to be in this segment.
Scaling out with data in situ	Inexpensive storage, compute; add machines to achieve desired speed	Extremely fast networking, cluster-wide filesystem	Cloud Dataproc + Spark on Cloud Storage, BigQuery, Cloud Dataflow, Vertex AI Training, etc.	Any use case where we can configure datasets so that I/O keeps up with computation. Data analytics use cases are increasingly falling into this segment.

Even if you are used to downloading data to your laptop for data analysis and development, you should realize that this is a suboptimal solution. Wouldn't it be great to directly ingest the BTS files into our data analysis programs without having to go through a step of downloading them? Having a single source of truth has many advantages, ranging from security (providing and denying access) to error correction (no need to worry about stale copies of the data lying around). Of course, the reason we don't do this is that we'd have to read the BTS data over the internet, and the public internet typically has speeds of 3 to 10 MBps.⁸ If you are carrying out analysis on your laptop, accessing data via the internet every time you need it will become a serious bottleneck.

Downloading the data has the benefit that subsequent reads happen on the local drive and this is both inexpensive and fast (see [Figure 2-5](#)). For small datasets and short, quick computation, it's perfectly acceptable to download data to your laptop and do the work there. This doesn't scale, though. What if our data analysis is very complex or the data is so large that a single laptop is no longer enough? We have two options: scale up or scale out.

⁸ Public internet as opposed to traffic traveling on private fiber. For example, communication between machines in Google Cloud travels on Google's own cables.

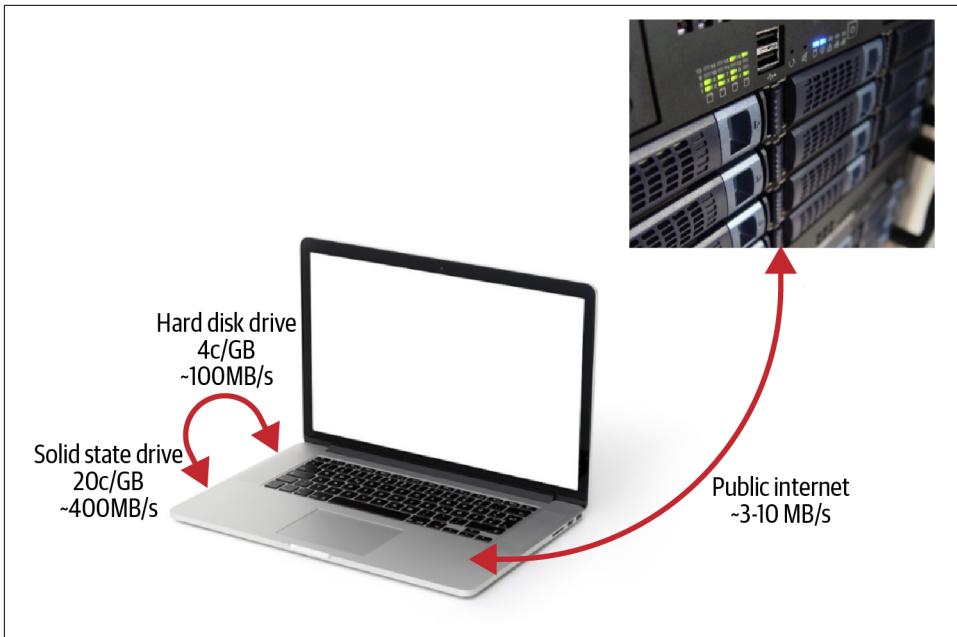


Figure 2-5. Comparison of data access speeds if data is accessed over the public internet versus from a disk drive.

Scaling Up

One option to deal with larger datasets or more difficult computation jobs is to use a larger, more powerful machine with many CPUs/GPUs, lots of RAM, and many terabytes of drive space. This is called *scaling up*, and it is a perfectly valid solution. However, such a computer is likely to be quite expensive. Because we are unlikely to be using it 24 hours a day, we might choose to rent an appropriately large computer from a public cloud provider. In addition, the public cloud offers persistent drives that can be shared between multiple instances and whose data is **geo-replicated to guard against data loss**. In short, then, if you want to do your analysis on one large machine but keep your data permanently in the cloud, a good solution would be to marry a powerful, high-memory Compute Engine instance with a persistent drive, download the data from the external data center (BTS's computer in our case) onto the persistent drive, and start up compute instances on demand, as depicted in [Figure 2-6](#) (cloud prices in [Figure 2-6](#) are estimated monthly charges; actual costs may be higher or lower than the estimate).

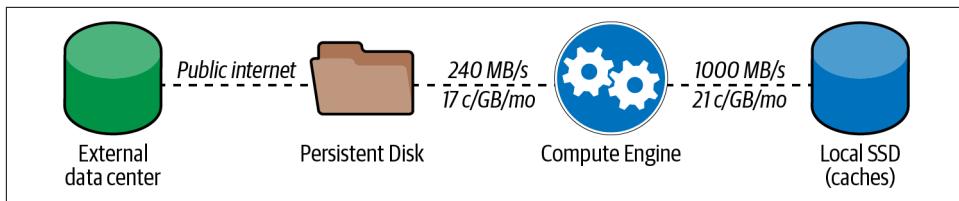


Figure 2-6. One solution to cost-effective and fast data analysis is to store data on a persistent disk that is attached to an ephemeral, high-memory Compute Engine instance.

When you are done with the analysis, you can delete the Compute Engine instance.⁹ Provision the smallest persistent drive that adequately holds your data—temporary storage (or caches) during analysis can be made to an attached SSD that is deleted along with the instance, and persistent drives can always be resized if your initial size proves too small. This gives you all the benefits of doing local analysis but with the ability to use a much more powerful machine at a lower cost. I will note here that this recommendation assumes several things: the ability to rent powerful machines by the minute, to attach resizable persistent drives to compute instances, and to achieve good-enough performance by using solid-state persistent drives. These are true of Google Cloud and other public cloud providers, but are unlikely to be true on premises.

Scaling up is a common approach whenever you have a job that needs to read the data multiple times. This is quite common when training machine learning models, and so scaling up is a common approach in machine learning, especially machine learning on images and video. Indeed, Google Cloud offers special Compute Engine instances, called **Deep Learning VM**, that have accelerators like GPUs and come pre-installed with the libraries that are needed for machine learning. Jupyter Notebook instances are also frequently scaled up as necessary to fit the job. You'd create a Deep Learning VM, and attach to it a network-based persistent disk containing the training data or use local SSD for improved performance.

⁹ You could also just stop (and not delete) the Google Compute Engine instance. Stopping the instance stops the bill associated with the compute machine, but you will continue to pay for storage. In particular, you will continue to pay for the SSD associated with the Compute Engine instance. The key advantage of a stopped instance is that you get to resume exactly where you left off, but this might not be important if you always start from a clean (known) state each time.

Scaling Out with Sharded Data

The solution of using a high-memory Compute Engine instance along with persistent drives and caches might be reasonable for jobs that can be done on a single machine, but it doesn't work for jobs that are bigger than that. Configuring a job into smaller parts so that processing can be carried out on multiple machines is called *scaling out*. One way to scale out a data processing job is to shard the data and store the pieces on the drives attached to multiple compute instances or persistent drives that will be attached to multiple instances.¹⁰ Then, each compute instance can carry out analysis on a small chunk of data at high speeds—these operations are called the *map* operations. The results of the analysis on the small chunks can be combined, after some suitable collation, on a different set of compute nodes—these combination operations are called the *reduce* operations. Together, this model is known as *MapReduce*. This approach also requires an initial download of the data from the external data center to the cloud. In addition, we also need to split the data onto preassigned drives or nodes.

Whenever we need to carry out analysis, we will need to spin up the entire cluster of nodes, reattach the persistent drives, and carry out the computation. Fortunately, we don't need to build the infrastructure to do the sharding or cluster creation ourselves. We could store the data on the Hadoop Distributed File System (HDFS), which will do the sharding for us, spin up a Cloud Dataproc cluster (which has Hadoop, Presto, Spark, etc., preinstalled on a cluster of Compute Engine VMs), and run our analysis job on that cluster. [Figure 2-7](#) presents an overview of this approach.

¹⁰ To shard a large database is to partition it into smaller, more easily managed parts. Whereas normalization of a database table places the columns of a database into different tables, sharding splits the rows of the database and uses different database server instances to handle each part. For more information, go to the [Wikipedia entry for the Shard database architecture](#).

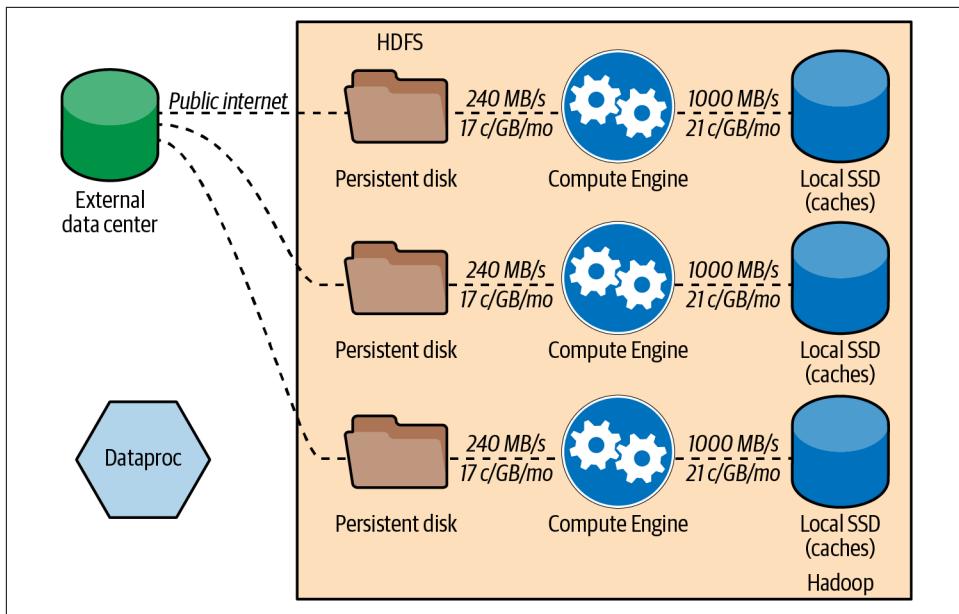


Figure 2-7. For larger datasets, one potential solution is to store data on the HDFS and use an ephemeral Dataproc cluster to carry out the analysis.

A MapReduce framework like the Hadoop ecosystem requires data to be presharded. Because the presharded data must be stored on drives that are attached to compute instances, the scheme can be highly wasteful unless all the data happens to get used all the time by those compute instances. In essence, whenever you need to run a job, the framework ships the code to whichever nodes happen to be storing the data. What the framework should be doing, however, is try to find a machine that has free capacity. Shipping the analysis code to run on storage nodes regardless of their computational load leads to poor efficiency because it is likely that there are long periods during which a node might have nothing to do, and other periods when it is subject to resource contention.

In summary, we have two options to work with large datasets: keep the data as-is and scale up by using a large-enough computer, or scale out by sharding the data and shipping code to the nodes that store the data. Both of these options have some drawbacks. Scaling up is subject to the limitations of whatever the most powerful machine available to you can do. Scaling out is subject to the inefficiencies of resource allocation. Is there a way to keep data-in-place and scale out?

Scaling Out with Data-in-Place

Recall that much of the economics of our case for downloading the data onto nodes on which we can do the compute relied on the slowness of an internet connection as compared to drive speeds—it is because the public internet operates at only 3 to 10 MBps, whereas drives offer two orders of magnitude faster access, that we moved the data to a large Compute Engine instance (scaling up) or sharded it onto persistent drives attached to Compute Engine instances (scaling out).

What if, though, you are operating in an environment in which networking speeds are higher, and files are available to all compute instances at those high speeds? For example, what if you had a job that uses 100,000 servers and those servers could communicate with one another at 1 GBps? This is seriously fast—it is twice the speed of SSDs, 10 times the speed of a local hard drive, and 100 times faster than the public internet. What if, in addition, you have a cluster-level filesystem (not node-by-node) whose metadata is sharded across the data center and replicated on write for durability? Because the total bisection bandwidth of Google's [Andromeda and Jupiter networks](#) in Google's data centers is 125,000 GBps,¹¹ and because Google's next-generation Colossus filesystem operates at the cluster level, this is the scenario that operates if your data is available in Google Cloud Storage and your jobs are running on Compute Engine instances in the same data center as the file. At that point, it becomes worthwhile to treat the entire data center as a single computer. The speed of the network and the design of the storage make both compute and data fungible resources that can be allocated to whichever part of the data center is most free. Scheduling a set of jobs over a single large data center provides much higher utilization than scheduling the same set of jobs over many smaller clusters. This resource allocation can be automatic—there is no need to preshard the data, and if we use an appropriate computation framework (such as BigQuery, Cloud Dataflow, or Vertex AI), we don't even need to instantiate a Compute Engine instance ourselves. [Figure 2-8](#) presents this framework, in which compute and storage are separate.

Therefore, choose between scaling up, scaling out with data sharding, or scaling out with data-in-place depending on the problem that you are solving (see [Table 2-3](#)).

¹¹ The [blog on Google's networking infrastructure](#) is worth a read. One petabit is 1 million gigabits, so the 1 Pbps quoted in the article works out to 125,000 GBps. Networking has only gotten better since 2015, of course.

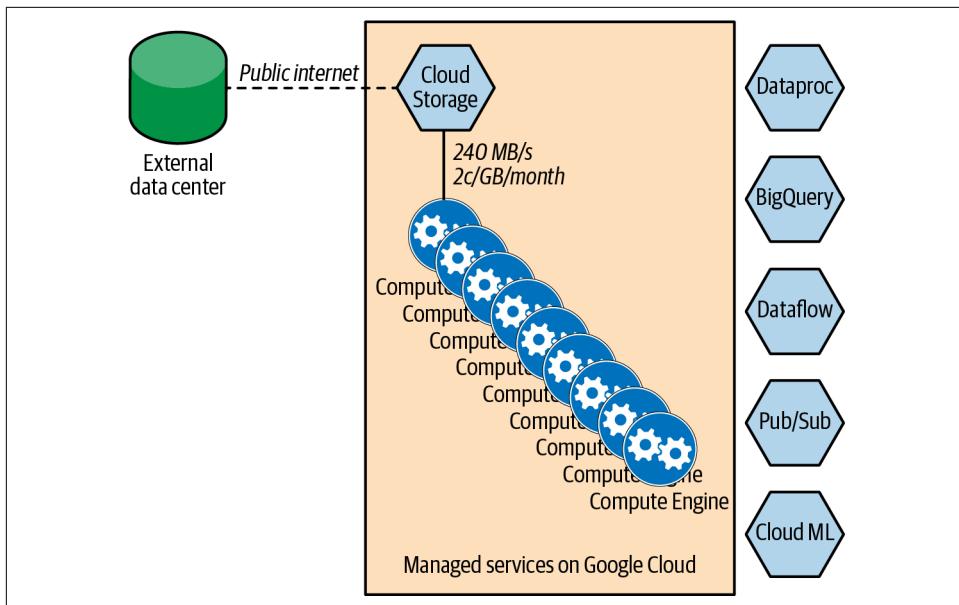


Figure 2-8. On the Google Cloud Platform, the speed of the networking within a data center allows us to store the data persistently and cheaply on Cloud Storage and access it as needed from a variety of ephemeral managed services. This is called separation of compute and storage.

Google Cloud Data Centers Are Different

Google data centers, unlike most other data centers, are optimized for total bisection bandwidth. They are optimized to maximize the network bandwidth between nodes on the backend (“East-West communications” in networking parlance). Most other data centers are optimized to minimize the network time with an outside client sending, for example, a web request (“North-South communications”).

Why would anybody design a data center for East-West communications? Aren’t most applications web applications? You would design a data center for East-West networking only if the amount of network calls you do on the backend in response to a user request is several times the traffic of the request itself. That is true of Google Search (very simple user interface, very complex business logic). Fortunately, this design also comes in extremely useful for data science because it is not necessary to *preshard* the data.

What’s presharding? The Google File System (or GFS, on which the HDFS is based) was built for batch operations and involves storing data in shards close to processing nodes. Colossus (GFS’s successor that is in use in Google data centers) was designed for real-time updates. Although GFS/HDFS suffices for batch processing operations that happen over a few days, Colossus is required to update Google’s search index in

real time—this is why Google Search can now reflect current events. There are several other innovations that were necessary to get to this data processing architecture in which data does not need to be presharded. For example, when performing large fan-out operations, you must be tolerant of latency and design around it. This involves slicing up requests to reduce head-of-line blocking,¹² creating hundreds of partitions per machine to make it easy to move partitions elsewhere, making replicas of heavily used data chunks, using backup requests, and canceling other requests as soon as one is completed, among other strategies. To build a cluster-wide filesystem with high throughput speeds to any compute instance within the data center, it is necessary to minimize the number of network hops within the data center by changing the network definitions through software. Within the Google Cloud Platform, any two machines in the same zone are only one network hop away.

The innovation in networking, compute, and storage at Google and elsewhere is by no means over. Even though the Jupiter network provides bisection bandwidths of 125,000 GBps (the last time Google published this number publicly was in the mid-2010s and it's probably higher now), engineers estimate that 600,000 GBps is what's required to match the performance of local disks. Moreover, jobs are not being sliced finely enough—because I/O devices have response times on the order of microseconds, decisions should be scheduled even more finely than the current milliseconds. Next-generation flash storage is still largely untapped within the data center. Colossus addresses the issue of building a cluster-level filesystem, but there are applications that need global consistency, not just consistency within a single-region cluster. The challenge of building a globally distributed database is being addressed by Cloud Spanner. The ongoing innovations in computational infrastructure promise exciting times ahead.

All of this is in the way of noting (again!) that your mileage will vary if you do your data processing on other infrastructure—there is a reason why the title of this book includes the words “on the Google Cloud Platform.” The hardware optimizations if you implement your data pipelines on premises or in a different cloud provider will typically target different things.¹³ The APIs might look the same, and in many cases, you can run the same software as I do, but the performance characteristics will be different. Google TensorFlow, Apache Beam, and others are open source and portable to on-premises infrastructure and across different cloud providers, but the execution frameworks that make Vertex AI and Cloud Dataflow so powerful may not translate well to infrastructure that is not built the same way as Google Cloud Platform.

Another way to see this is that multicloud software works faster on Google Cloud than on other cloud platforms. BigQuery Omni, although available on Amazon Web

¹² *Head-of-line blocking* is a condition in which network packets need to be delivered in order; thus, a slow packet holds up delivery of later packets.

¹³ Microsoft Azure seems to involve a centralized host layer, for example, while AWS S3 seems to prioritize network latency. You'd design your software for such infrastructure differently.

Services (AWS) and Azure, does not get the performance of BigQuery on GCP. This performance difference is not limited to multicloud software developed by Google. [Databricks notes](#) that startup and certain Spark workloads are faster on GCP than on other clouds. [Actian Avalanche notes](#) that their implementation on GCP is 20% faster than on other cloud platforms.

Ingesting Data

To carry out our data analysis on the on-time performance dataset, we will need to download the monthly data from the BTS website and then upload it to Google Cloud Storage. Doing this manually will be tedious and error-prone, so let's script this operation.

Reverse Engineering a Web Form

How would you script filling out the BTS web form shown in [Figure 2-2](#)? First, verify that the website's terms of use do not bar you from automated downloads! Then, use the Chrome browser's developer tools to find what web calls the form makes. Once you know that, you can repeat those web calls in a script.

The BTS web form is a simple HTML form with no dynamic behavior. This type of form collects all the user selections into a single GET or POST request. If we can create that same request from a script, we will be able to obtain the data without going through the web form.

We can find out the exact HTTP command sent by the browser after we make our selections on the BTS website. You can do this while on the [BTS download website](#) in the Chrome web browser—in the upper-right menu bar of the browser, navigate to the Developer Tools menu, as shown in [Figure 2-9](#).

Now, on the BTS website, select the Prezipped File option, select 2015 and January in the drop-down boxes, and click Download.¹⁴ The Developer tools menu shows us that the browser is now making a GET request for https://transtats.bts.gov/PREZIP/On_Time_Report_Carrier_On_Time_Performance_1987_present_2015_1.zip.

¹⁴ The BTS website is frequently down. Do not be alarmed if you get an error here. You will not need the website to work in order to go through the remaining chapters. If the BTS server is down, you can copy the necessary files from my bucket. Consult the *README.md* file in the GitHub repository for details.

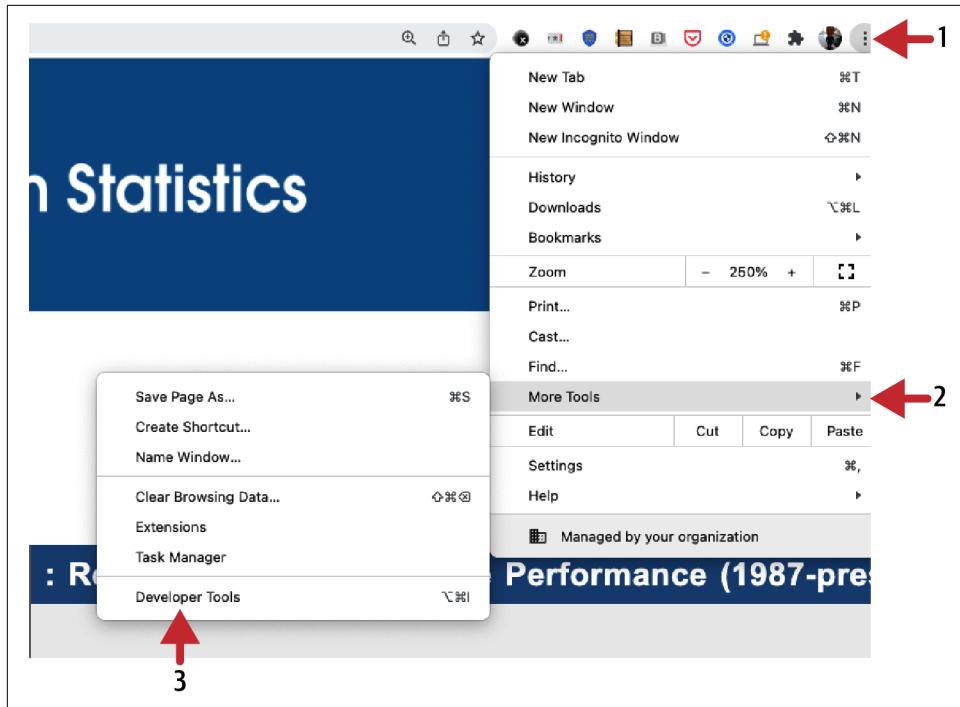


Figure 2-9. Navigating to the Developer Tools menu in Chrome.

It is pretty obvious what the pattern here is. If we issue an HTTP GET for a file with the pattern:

```
 ${BASEURL}_[${YEAR}]_[${MONTH}].zip
```

we should get the data corresponding to a single month. Let's try it from the command line of Cloud Shell:

```
BTS=https://transtats.bts.gov/PREZIP  
BASEURL="${BTS}/On_Time_Reportng_Carrier_On_Time_Performance_1987_present"  
YEAR=2015  
MONTH=3  
curl -k -o temp.zip ${BASEURL}_[${YEAR}]_[${MONTH}].zip
```

We see the data for March 2015 starting to get downloaded. Once the file is downloaded, we can unzip it:

```
unzip temp.zip
```

We then notice that the ZIP file contains a comma-separated values (CSV) file. Peeking at the first few lines of the file using:

```
head -5 *.csv
```

confirms that the file contains flight data from the month of January 2015.

Dataset Download

In the data exploration phase, I'll do most of my processing interactively with Linux command-line tools. I will assume that this is what you are using as well. Adapt the commands as necessary if you are working locally in some other environment (e.g., where I ask you to do a `sudo apt-get install`, you might use the appropriate `install` command for your Linux-like environment). When we have settled on the processing to be done, we'll look at how to make this more automated.

Instead of calling the downloaded file `temp.zip`, let's call it `201501.zip` and place it into a temporary directory. To pad the month 1 to be 01, we can use the `printf` command in bash:¹⁵

```
MONTH2=$(printf "%02d" $MONTH)
```

To create a temporary directory, we can use the Linux command `mktemp`:

```
TMPDIR=$(mktemp -d)
```

Then, to download the file to the temporary directory, we can do:

```
ZIPFILE=${TMPDIR}/${YEAR}_${MONTH2}.zip
curl -o $ZIPFILE ${BASEURL}_${YEAR}_${MONTH}.zip
```

Now, we can unzip the file, extract the CSV file to the current directory (./), and blow out the remaining contents of the ZIP file:

```
unzip -d $TMPDIR $ZIPFILE
mv $TMPDIR/*.csv ./${YEAR}${MONTH2}.csv
rm -rf $TMPDIR
```

I put the preceding commands into a file called `download.sh`, and then in the script `ingest.sh`, I call it from within a `for` loop:

```
for MONTH in `seq 1 12`; do
    bash download.sh $YEAR $MONTH
done
```

On running this, we get a set of CSV files, one for each month in 2015 (see [Figure 2-11](#)).

The [complete download script is on GitHub](#)—if you want to follow along with me, perform these steps:

- Go to <https://console.cloud.google.com>.
- Select the GCP project that you will be working on. This is accomplished by using the drop-down box next to “Example Project” in [Figure 2-10](#).

¹⁵ See the script `02_ingest/download.sh` in the course repository.

- On the top strip, activate Cloud Shell using the button shown in [Figure 2-10](#).



Figure 2-10. The Cloud Shell button on the Google Cloud Platform web console.

- In Cloud Shell, type the following:

```
git clone \
  https://github.com/GoogleCloudPlatform/data-science-on-gcp
```

This downloads the GitHub code to your Cloud Shell home directory.

- Navigate into the flights folder:

```
cd data-science-on-gcp
```

- Make a new directory to hold the data, and then change into that directory:

```
mkdir data
```

```
cd data
```

- Run the code to download the files:

```
for MONTH in `seq 1 12`; do
    bash download.sh 2015 $MONTH
done
```

- When the script completes, run `ls -lrt` to view the downloaded ZIP files, shown in [Figure 2-11](#).

```
vlakshmanan@cloudshell:~/data-science-on-gcp/data (data-science-on-gcp-180606)$ ls -lrt
total 2561352
-rw-r--r-- 1 vlakshmanan vlakshmanan 211633508 Sep 19 2018 201501.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 192791832 Sep 19 2018 201502.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 227017003 Sep 19 2018 201503.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 218600101 Sep 19 2018 201504.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 224003621 Sep 19 2018 201505.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 227418851 Sep 19 2018 201506.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 235038069 Sep 19 2018 201507.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 230183745 Sep 19 2018 201508.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 209229804 Sep 19 2018 201509.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 219158206 Sep 19 2018 201510.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 211122691 Sep 19 2018 201511.csv
-rw-r--r-- 1 vlakshmanan vlakshmanan 216562989 Sep 19 2018 201512.csv
```

Figure 2-11. The `ls -lrt` command shows details of the downloaded files.

This looks quite reasonable—all the files have different sizes and the sizes are robust enough that one would assume they are not just error messages.

Exploration and Cleanup

At this point, I have 12 CSV files. Let's look at the first two lines of one of them to ensure the data matches what we think it ought to be:

head -2 201503.csv

The result is shown in Figure 2-12.

Figure 2-12. The first two lines of the CSV file containing March 2015 data.

There is a header in each CSV file, and the second line looks like data. Some of the fields are enclosed by quotes (perhaps in case the strings themselves have commas), and there are some fields that are missing (there is nothing between successive commas toward the end of the line). There seems to be a pesky extra comma at the end as well.

How many fields are there? Because the second line doesn't have any commas between the quotes, we can check using:

```
head -2 201503.csv | tail -1 | sed 's/,/ /g' | wc -w
```

The number of words is 81, so there are 81 columns (remember there's a comma at the end of the line). Here's how the command works. It first gets the first two lines of the data file (with `head -2`), and the last line of that (with `tail -1`) so that we are looking at the second line of `201503.csv`. Then, we replace all the commas by spaces and count the number of words with `wc -w`.

How much data is there? A quick shell command (`wc` for word count, with an `-l` [lowercase letter L] to display only the line count) informs us that there are between ~43,000 and ~52,000 flights per month:

```
$ wc -l *.csv  
469969 201501.csv  
429192 201502.csv  
504313 201503.csv  
485152 201504.csv  
496994 201505.csv
```

```
503898 201506.csv
520719 201507.csv
510537 201508.csv
464947 201509.csv
486166 201510.csv
467973 201511.csv
479231 201512.csv
5819091 total
```

This adds up to nearly six million flights in 2015! The slowness of this command should tell us that any kind of analysis that involves reading all the data is going to be quite cumbersome. You can repeat this for other years (2016–2019), but let's wait until we have the whole process complete for one year before we add more years.

You might have realized by now that knowing a little Unix Shell scripting can come in very handy at this initial stage of data analysis.¹⁶

Uploading Data to Google Cloud Storage

For durability of this raw dataset, let's upload it to Google Cloud Storage. To do that, you first need to create a *bucket*, essentially a namespace for binary large objects (blobs) stored in Cloud Storage that you typically want to treat similarly from a permissions perspective. You can create a bucket from the [Google Cloud Platform Console](#). For reasons that we will talk about shortly, make the bucket a single-region bucket.

Setting Up a Cost Budget

It is likely that when you create a bucket, you will be prompted to enable billing or connect to a billing account. This is because storage is chargeable. I suggest that you create a brand new Google Cloud project for this book and delete the project once you are done to stop all billing. If you want total peace of mind, you can [set up a cost budget](#) and ask Google Cloud to [cap resource usage](#) once you exceed that budget.

Bucket names must be globally unique (i.e., unique not just within your project or organization, but across Google Cloud Platform). This means that bucket names are globally knowable even if the contents of the bucket are not accessible. This can be problematic. For example, if you created a bucket named `acme_gizmo`, a competitor might later try to create a bucket also named `acme_gizmo`, but fail because the name already exists. This failure can alert your competitor to the possibility that Acme Corp. is developing a new Gizmo. It might seem like it would take Sherlock Holmes-like powers of deduction to arrive at this conclusion, but it's simply best that you

¹⁶ Software Carpentry provides a [good intro to Unix and shell scripting](#).

avoid revealing sensitive information in bucket names. A common pattern to create unique bucket names is to use suffixes on the project ID. Project IDs are globally unique,¹⁷ and thus a bucket name such as `projectid-dsongcp` will also tend to be unique. In my case, my project ID is `cloud-training-demos` and my bucket name is `cloud-training-demos-ml`.

You can create a unique bucket on the command line using:

```
PROJECT=$(gcloud config get-value project)
BUCKET=${PROJECT}-dsongcp
REGION=us-central1 #See https://cloud.google.com/storage/docs/locations
gsutil mb -l $REGION gs://$BUCKET
```

where the first line retrieves the project ID and the second line uses it to form a bucket name that is hopefully unique.

Cloud Storage will also serve as the staging ground to many of the GCP tools and enable collaborative sharing of the data with our colleagues. In my case, to upload the files to Cloud Storage, I type the following in Cloud Shell:

```
gsutil -m cp *.csv gs://cloud-training-demos-ml/flights/raw/
```

This uploads the files to Cloud Storage, specifically to my bucket `cloud-training-demos-ml` in a multithreaded manner (`-m`) and makes me the owner. If you are working locally, another way to upload the files would be to use the [Cloud Platform Console](#).

It is better to keep these as separate files instead of concatenating them into a single large file because Cloud Storage is a [blob store](#), not a regular filesystem. In particular, it is not possible to append to a file on Cloud Storage; you can only replace it. Therefore, although concatenating all 12 files into a single file containing the entire year of data will work for this batch dataset, it won't work as well if we want to later add to the dataset one month at a time, as new data becomes available. Second, because Cloud Storage is blob storage,¹⁸ storing the files separately will permit us to more easily process parts of the entire archive (for example, only the summer months) without having to build slicing into our data processing pipeline. Third, it is generally a good idea to keep ingested data in as raw a form as possible.

¹⁷ You can get your unique project ID from the Cloud Platform Console dashboard; it could be different from the common name that you assigned to your project. By default, Google Cloud Platform tries to give you a project ID that is the same as your project name, but if that name is already taken, you will get an autogenerated, unique project ID. Because of this default, you should be similarly careful about giving projects sensitive names.

¹⁸ In an object store, we don't have random access to data in the middle (technically, this is called a seek), as we would in a filesystem. Rather, we have to download an entire object in order to work with it. Tools such as [GCS FUSE](#) allow you to treat the object store, in certain ways, like a filesystem, but the abstraction is not perfect. There is no seek capability.

It is preferable that this bucket to which we upload the data is a single-region bucket. There are four reasons: first, we will create Compute Engine instances in the same region as the bucket and access it only from this one region. A multiregion bucket would be overkill because we don't need global availability. Second, a single-region bucket is less expensive than a multiregion one. Third, single-region buckets are optimized for low latency and high throughput for data consumers, whereas dual-region and multiregion buckets are optimized for high availability, serving content outside of the Google network (think data analytics versus web traffic). All three of the preceding factors point to using single region buckets for data analytics and machine learning. The fourth reason is helpful here, although it won't come into play for "real-world" uses: at the time of writing, certain US single regions (us-east1, us-west1, and us-central1) offer 5 GB of storage **free**.

Note that both single-region and multiregion buckets in Google Cloud Platform offer strong consistency, so this does not seem like a consideration to choose one over the other. However, the speed differences inherent in being able to offer strong consistency on global buckets points to using single-region buckets if you can. What exactly is strong versus eventual consistency, and why does it matter? Suppose that a worker in a distributed application updates a piece of data, and another worker views that piece of data immediately afterward. Does the second worker always see the updated value? Then, what you have is *strong consistency*. If, on the other hand, there could be a potential lag between an update and availability (i.e., if different viewers can see potentially different values of the data at the same instant in time), what you have is *eventual consistency*. Eventually, all viewers of the data will see the updated value, but that lag will be different for different viewers. Strong consistency is an implicit assumption that is made in a number of programming paradigms. However, to achieve strong consistency, we have to make compromises on scalability and performance (this is called *Brewer's theorem*). For example, we might need to lock readers out of the data while it is being updated so that simultaneous readers always see a consistent and correct value.



Brewer's theorem, also called the CAP theorem, states that no computer system can simultaneously guarantee consistency, availability, and partition resilience. *Consistency* is the guarantee that every reader sees the latest written information. *Availability* is the guarantee that a response is sent to every request (regardless of whether it is the most current information or not). *Partition resilience* is the guarantee that the system continues to operate even if the network connecting readers, writers, and storage drops an arbitrary number of messages. Because network failures are a fact of life in distributed systems, the CAP theorem essentially says that you need to choose between consistency and availability. Neither multiregion buckets nor Cloud Spanner change this: they essentially make choices during partitioning—Cloud Spanner is always consistent and achieves five nines (99.999%, but not perfect) availability despite operating over a wide area. For more details, see this [2017 research paper by Eric Brewer](#).

If you need the performance of a regional bucket, but need to be tolerant to failure (for example, you want to be able to carry out your workload even if a region goes down), there are two options: eventual consistency and dual-region buckets. As an example of eventual consistency consider how DNS servers cache values and have their values replicated across many DNS servers all over the internet. If a DNS value is updated, it takes some time for this modified value to become replicated at every DNS server. Eventually, this does happen, though. Having a centralized DNS server that is locked whenever any DNS value is modified would have led to an extremely brittle system. Because the DNS system is based on eventual consistency, it is highly available and extremely scalable, enabling name lookups for/to millions of internet-capable devices. The other option is to have a dual-region bucket in a multiregion location, so that the metadata remains the same. If, for whatever reason, one region is not available for analytics, computation can be migrated to the other region in a multiregion location (US, EU, Asia). Dual-region buckets are more expensive than either single-region buckets or multiregion buckets, but offer both high performance and reliability.

This being public data, I will ensure that my colleagues can use this data without having to wait on me:

```
gsutil acl ch -R -g google.com:R \
  gs://cloud-training-demos-ml/flights/raw/
```

That changes the access control list (`acl`) recursively (`-R`), applying to the group `google.com` read permission (`:R`) on everything starting from the Cloud Storage URL supplied. You'd of course replace `google.com` with your company's domain in order to share data with your own colleagues. Had there been sensitive information in the dataset, I would have to be more careful. We'll discuss fine-grained security, by

providing views with different columns to different roles in my organization, when we talk about putting the data in BigQuery.

Loading Data into Google BigQuery

On Google Cloud, the best place for structured and semi-structured data is BigQuery, a serverless data warehouse and SQL engine.

Advantages of a Serverless Columnar Database

Most relational database systems, whether commercial or open source, are row oriented in that the data is stored row by row. This makes it easy to append new rows of data to the database and allows for features such as row-level locking when updating the value of a row. For example, if you have an inventory table, you'd lock the row corresponding to the item being purchased so that you can ensure the item is shipped to the person who paid for it.

The drawback is that queries that involve table scans (i.e., any aggregation that requires reading every row in the entire table) can be expensive. For example, if we want to find the number of times an item was purchased by someone living in Belgium, that will involve a table scan. Indexing counteracts this expense by creating a lookup table to map rows to column values, so that `SELECT` queries that involve indexed columns do not have to load unnecessary rows from storage into memory—we might index the country column, for example, if purchases by country is a common query. If you can rely on indexes for fast lookup of your data, a traditional relational database management system (RDBMS) works well. For example, if your queries tend to come from software applications, you will know the queries that will come in (and the columns they are likely to access). So, you can create the appropriate indexes beforehand. This is not an option for use cases like business intelligence for which human users are writing ad hoc queries; therefore, a different architecture is needed.

BigQuery, unlike an RDBMS, is a columnar database—data is stored column by column and each column's data is stored in a **highly efficient compressed format** that enables fast querying. Because of the way data is stored, many common queries can be carried out such that the **query processing time is linear on the size of the relevant data**. For applications such as data warehousing and business intelligence for which the predominant operations are read-only `SELECT` queries requiring full table scans, columnar databases are a better fit. BigQuery, for example, can scan terabytes of data in a matter of seconds. The trade-off is that `INSERT`, `UPDATE`, and `DELETE` statements, although possible in BigQuery, are **significantly more expensive to process** than `SELECT` statements. BigQuery is tuned toward analytics use cases.

BigQuery is serverless, so you don't actually spin up a BigQuery server in your project. You simply submit a SQL query, and it's executed on the cloud. SQL queries that you submit to BigQuery are executed on a large number of compute nodes (called *slots*) in parallel. These slots do not need to be statically allocated beforehand—instead, they are “always on” available on demand, and scale to the size of your job. Because data is kept in-place and not sharded (i.e., not broken into small chunks that are attached to individual compute instances), the total power of the data center can be brought to bear on the problem. Because these resources are elastic and used only for the duration of the query, BigQuery is more powerful and less expensive than a statically preallocated cluster because preallocated clusters will typically be provisioned for the average use case—BigQuery can bring more resources to bear on the above-average computational jobs and utilize fewer resources for below-average ones.

In addition, because you don't need to reserve any compute resources for your data when you are not querying your data, it is extremely cost effective to just keep your data in BigQuery (you'll pay for storage, but storage is inexpensive). Whenever you do need to query the data, the data is immediately available—you can query it without the need to start project-specific compute resources. This on-demand, autoscaling of compute resources is incredibly liberating.

BigQuery Pricing

If an on-demand cost structure (you pay per query) concerns you because costs can fluctuate month over month, you can specify a billing cap for users. For even more cost predictability, it is possible to pay a fixed monthly price for BigQuery—flat-rate pricing means you get a predictable cost regardless of the number of queries run or data processed by those queries. The fixed monthly price essentially buys you access to a specific number of slots.

In short, BigQuery has **two pricing models for analysis**: an on-demand pricing model in which your cost depends on the quantity of data processed, and a flat-rate model in which you pay a fixed amount per month for an unlimited number of queries that will run on a specific set of compute resources. You can augment either with **flex slots** you pay for by the minute. In all these cases, storage is a separate cost and depends on data size.

In summary, BigQuery is a columnar database, making it particularly effective for read-only queries that process all of the data. Because it is serverless, can autoscale to thousands of compute nodes, and doesn't require clusters to be preallocated, it is also very powerful and quite inexpensive.

Staging on Cloud Storage

Although it is possible to ingest files from on-premises hardware directly into BigQuery using the [bq command-line tool](#) that comes with the Google Cloud Software Development Kit (SDK), aka `gcloud`, you should use that capability only for small datasets. To ingest data from outside Google Cloud Platform to BigQuery, it is preferable to first load it into Cloud Storage and use Cloud Storage as the staging ground for BigQuery, as demonstrated in [Figure 2-13](#).

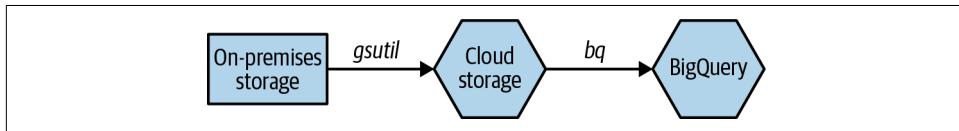


Figure 2-13. Use Cloud Storage as a staging ground to ingest data into BigQuery.

For larger files, it is better to ingest the files into Cloud Storage using `gsutil` first because `gsutil` takes advantage of multithreaded, resumable uploads and is better suited to the public internet. In our case, this is what we did in the previous section when we used `gsutil` to copy the extracted flights CSV files to Cloud Storage. Now that we have the CSV files in Cloud Storage, we can load them into BigQuery.

Cloud Storage or BigQuery?

When should you save your data in Cloud Storage, and when should you store it in BigQuery? First, if the data is not tabular-like (that is: images, videos, and other arbitrary file types), then Google Cloud Storage (GCS) is the right choice. For tabular-like data, the answer boils down to what you want to do with the data and the kinds of analyses you want to perform. If you'll mostly be running custom code that expects to read plain files, or your analysis involves reading the entire dataset, use Cloud Storage. On the other hand, if your desired access pattern is to run interactive SQL queries on the data, store your data in BigQuery. To summarize, if in the pre-cloud world you would use flat files, use Cloud Storage. If you'd put the data in a relational database, put it in BigQuery.

Access Control

The first step to ingest data into BigQuery is to create a BigQuery dataset—a dataset is a container for tables. You can have multiple datasets within a project. Go to the [web console](#) and choose the Create Dataset option. Then, create a dataset called `dsongcp`.

You can also do this from the command line:

```
bq mk dsongcp
```

Datasets in BigQuery are mostly just an organizational convenience—tables are where data resides, and it is the columns of the table that dictate the queries we write. Besides providing a way to organize tables, though, datasets also serve as a convenient *access control point*. You can conveniently provide view or edit access at the dataset level to control access to all the tables in the dataset. Cloud Identity and Access Management (Cloud IAM) on Google Cloud Platform provides a mechanism to control *who* can carry out *what* actions on *which* resource (Figure 2-14).

The “*who*” can be specified in terms of an individual user (identified by their Google account such as a *gmail.com* address, or company email address if the company is a Google Workspace customer), a Google Group (i.e., all current members of the group), or a Google Workspace or **Google Identity domain** (anyone with a Google account in that domain). Google Groups and Google Identity/Workspace domains provide a convenient mechanism for aggregating a number of users and providing similar access to all of them.

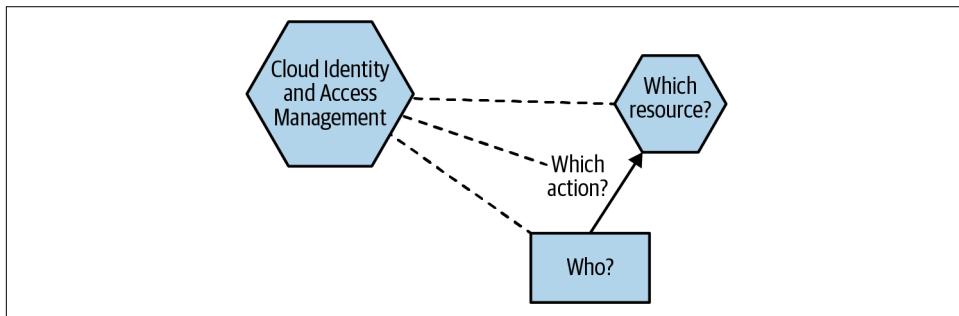


Figure 2-14. Cloud IAM provides a mechanism to control access to resources.

In addition, different logical parts of an application can be assigned separate identities (linked to email addresses) called *service accounts*. Service accounts are a very powerful concept because they allow different parts of a codebase to have permissions that are independent of the access level of the person running that application. For example, you might want an application to be able to query a table but not delete it even if the developer who created the application and the person running the application have that authority.

Careful with Service Accounts

You should use service accounts with care for scenarios in which audit records are mandatory. Providing access at the Google Groups level provides more of an audit trail; because Google Groups don't have login credentials (only individual users do), the user who made a request or action is always recorded, even if their access is provided at the level of a Google Group or Google Workspace domain. However, service accounts are themselves login credentials, and so the audit trail turns cold if you

provide access to service accounts—you will no longer know which user initiated the application request unless that application in turn logs this information.¹⁹ Keep this in mind when granting access to service accounts.

Try to avoid providing service account access to resources that require auditability. If you do provide service account access, you should ensure that the application to which you have provided access itself provides the necessary audit trail by keeping track of the user on behalf of whom it is executing the request. The same considerations apply to service accounts that are part of Google Groups or Google Workspace domains. Because audit trails go cold with service accounts, you should restrict Google Groups and Google Workspace domains to only human users and service accounts that belong to applications that provide any necessary legal auditability guarantees.

Creating single-user projects is another way to ensure that service accounts map cleanly to users, but this can lead to significant administrative overhead associated with shared resources and departing personnel. Essentially, you would create a project that is billed to the same company billing account, but each individual user would have their own project in which they work. You can use the `gcloud` command to script the creation of such single-user projects in which the user in question is an editor (not an owner).

In addition to specific users, groups, domains, and service accounts, there are two wildcard options available. Access can be provided to `allAuthenticatedUsers`, in which case anyone authenticated with either a Google account or a service account is provided access. Because `allAuthenticatedUsers` includes service accounts, it should not be used for resources for which a clear audit trail is required. The other wildcard option is to provide access to `allUsers`, in which case anyone on the internet has access—a common use case for this is to provide highly available static web resources by storing them on Cloud Storage. Be careful about doing this indiscriminately—egress of data from Google Cloud Platform is not free, so you will pay for the bandwidth consumed by the download of your cloud-hosted datasets.

The “what” actions depend on the resource access that is being controlled. The resources themselves fall into a policy hierarchy.

Policies can be specified at an organization level (i.e., to all projects in the organization), at the project level (i.e., to all resources in the project), or at the resource level (i.e., to a Compute Engine instance or a BigQuery dataset). As [Figure 2-15](#) shows,

¹⁹ A service account is tied to a project, but project membership evolves over time. So, even the subset of users who could have invoked the action might not be known unless you have strict governance over who is allowed to be an owner/editor/viewer of a project.

policies specified at higher levels are inherited at lower levels, and the policy in effect is the union of all the permissions granted—there is no way to restrict some access to a dataset to a user who has that access inherited from the project level. Moving a project from one organization to another automatically updates that project’s Cloud IAM policy and ends up affecting all the resources owned by that project.

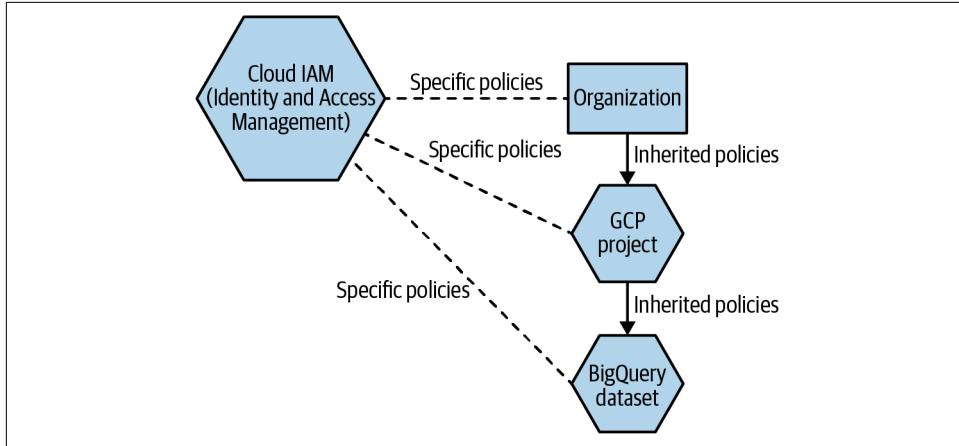


Figure 2-15. Policies specified at higher levels are inherited at lower levels.

What type of actions can be carried out depends on the resource in question. Before Cloud IAM was introduced on the Google Cloud Platform, there were only three roles: owner, editor, and viewer/reader for all resources. Cloud IAM brought with it finer-grained roles, but the original three roles were grandfathered in as primitive roles. [Table 2-4](#) lists some of the roles that are possible for BigQuery datasets. Cloud IAM roles for BigQuery (and by extension for all GCP products) are continuously updated to cater to new use cases. Please refer to [BigQuery access control](#) for the current roles and permissions.

Table 2-4. Some of the available roles in BigQuery

Role	Capabilities	Inherits from
Project Viewer	Execute a query List datasets	
Project Editor	Create a new dataset	Project Viewer
Project Owner	List/delete datasets View jobs run by other project users	Project Editor
bigquery.user	Execute a query List datasets	
bigquery.dataViewer	Read, query, copy, export tables in the dataset	
bigquery.dataEditor	Append, load data into tables in the dataset	Project Editor bigquery.dataViewer

Role	Capabilities	Inherits from
bigrquery.dataOwner	Update, delete on tables in the dataset	Project Owner bigrquery.dataEditor
bigrquery.admin	All	

Ingesting CSV Files

We can load the data directly into BigQuery's native storage using the command-line utility `bq` that comes with the `gcloud` SDK:

```
BUCKET=${PROJECT}-dsongcp
bq load --autodetect --source_format=CSV \
  dsongcp.flights_auto \
  gs://${BUCKET}/flights/raw/201501.csv
```

Here, we are asking BigQuery to autodetect the schema from the CSV file and load the January data into a table named `flights_auto` (if you are following along with me, make sure to change the bucket to reflect the bucket that your files are in).²⁰ Since we are asking BigQuery to autodetect the schema using `--autodetect`, it will read the header line of the CSV and use them as column names. However, the CSV header doesn't specify the column types (string, integer, float, etc.), and so BigQuery will sample a few hundred lines and attempt to guess at the type. This guess will not be perfect, and we might have to fix it.

If you now go to the [BigQuery web console](#) and examine the dataset `dsongcp`, you will see that there is a table named `flights_auto` in it. You can examine the autodetected schema and preview the contents of the table.

We can try querying the data to find the average departure and arrival delays at the busiest airports:

```
SELECT
  ORIGIN,
  AVG(DEP_DELAY) AS dep_delay,
  AVG(ARR_DELAY) AS arr_delay,
  COUNT(ARR_DELAY) AS num_flights
FROM
  dsongcp.flights_auto
GROUP BY
  ORIGIN
ORDER BY num_flights DESC
LIMIT 10
```

The result (see [Table 2-5](#)) starts with Atlanta (ATL), Dallas (DFW), and Chicago (ORD), which is what we would expect.

²⁰ If this is your first time working in BigQuery for the project, you might have to authorize the proper API. That's expected behavior.

Table 2-5. Average arrival and departure delays in January 2015 at the busiest 10 airports

Row	ORIGIN	dep_delay	arr_delay	num_flights
1	ATL	7.265885087329549	1.0802479706819135	29,197
2	DFW	11.761812240572308	9.37162730937924	22,571
3	ORD	19.96205128205128	17.016131923283645	22,316
4	LAX	7.476340878516738	5.542057719380547	17,048
5	DEN	15.506798076352176	11.842324888226543	16,775
6	IAH	9.07378596782721	5.353498597528596	13,191
7	PHX	8.066722908198505	6.197786998616902	13,014
8	SFO	10.328127477406069	9.038424821002382	12,570
9	LAS	8.566096692995435	5.0543525523958595	11,499
10	MCO	9.887440638577354	5.820512820512793	9,867

Autodetection is hit-and-miss, though. This is because the way it works is that BigQuery samples about a hundred rows of data in order to determine what the data type needs to be. If the arrival delay was an integer for all one hundred rows that it saw, but there turns out to be a string (NA) somewhere else in the file, the loading will fail. Autodetection may also fail if many of the fields are empty.

Partitioning

Because of this, autodetection is okay during initial exploration, but we should quickly pivot to actually specifying the schema. At that time, it may be worthwhile to also consider whether this table should be partitioned by date—if most of our queries will be not on the full table, but only a few days, then partitioning will lead to cost savings. If that were the case, we would create the table first, specifying that it should be partitioned by date (don’t do this—we have already created the table and we don’t need daywise partitioning):

```
bq mk --time_partitioning_type=DAY dsongcp.flights_auto
```

When loading the data, we’d need to load each partition separately (partitions are named `flights_auto$20150101`, for example). We can also partition by a column in the data (`FlightDate`, for example).

Currently, we don’t know much about the fields, so we can ask BigQuery to treat all the columns except the `FlightDate` as a string:

```
SCHEMA=Year:STRING,...,FlightDate:DATE,Reporting_Airline:STRING,...
```

Putting all these together, the loading becomes (see `bqlload.sh` in the book’s repo):

```
for MONTH in `seq -w 1 12`; do
  CSVFILE=gs://${BUCKET}/flights/raw/${YEAR}${MONTH}.csv
  bq --project_id $PROJECT \
    load --time_partitioning_field=FlightDate \
```

```
--time_partitioning_type=MONTH \  
--source_format=CSV --ignore_unknown_values \  
--skip_leading_rows=1 --schema=$SCHEMA \  
${PROJECT}:dsongcp.flights_raw\${YEAR}\${MONTH} ${CSVFILE}  
done
```

At this point, we have the CSV files in Cloud Storage and the raw data in BigQuery. We have successfully ingested the 2015 flights data into GCP! If you want, you can repeat this for years 2016–2019 by changing the `for` loop in `ingest.sh` to:²¹

```
for YEAR in `seq 2016 2019`; do
```

Cloud Shell comes with a text editor. You can use it to edit files. Alternatively, use a Unix editor tool such as `nano` or `vim`. However, don't do it just yet—let's develop the code in this book with just 2015 data so that we can move faster. In [Chapter 12](#), we'll expand the analysis and ML models to 2015–2019 data. In [Chapter 3](#), we will start to look at the 2015 data and do useful things with it.

But before we move on, let's digress a little and consider automation.

Scheduling Monthly Downloads

Now that we have some historical flight data in our Cloud Storage bucket, it is natural to wonder how to keep the bucket current. After all, airlines didn't stop flying in 2021, and the BTS continues to refresh its website on a monthly basis. It would be good if we could schedule monthly downloads to keep ourselves synchronized with the BTS.

There are two scenarios to consider here. The BTS could let us know when it has new data, and we could then proceed to ingest the data. The other option is that we periodically monitor the BTS website and ingest new data as it becomes available. The BTS doesn't offer a mechanism by which we can be notified about data updates, so we will need to resort to *polling*. We can, of course, be smart about how we do the polling. For example, if the BTS tends to update its website around the 5th of every month, we could poll at that time.

Where should this ingest program be executed? Realizing that this is a program that will be invoked only once a month (more often if retries are needed if an ingest fails), we realize that this is not a long-running job, but is instead something that should be scheduled to run periodically. The traditional way to do this is to schedule a `cron` job in Unix/Linux.²² To schedule a `cron` job, you add a line to a `crontab` file and then

²¹ Let's ignore 2020 and 2021 because those were the years of the COVID-19 pandemic and historical data would not be very helpful in predicting arrival delays in 2020 or 2021.

²² A shortened form of a misspelling of chronos, Greek for time, `cron` is the name of the Unix daemon process that executes scheduled jobs at specific times.

register it with a Unix daemon that takes care of the scheduling.²³ For example, adding this line:

```
1 2 10 * * /etc/bin/ingest_flights.py
```

to `crontab` will cause the Python program `/etc/bin/ingest_flights.py` (that would carry out the same steps to ingest the flights data that we did on the command line in the previous section) to be run by the system at 02:01 on the 10th of every month.

Although cron jobs are a straightforward solution, there are several problems that all come down to resilience and repeatability:

- The cron job is scheduled on a particular server. If that server happens to be rebooted around 2 a.m. on April 10, the ingest might never take place that month.
- The environment that cron executes in is very restricted. Our task will need to download data from BTS, uncompress it, clean it up, and upload it to the cloud. These impose a variety of requirements in terms of memory, space, and permissions, and it can be difficult to configure cron appropriately. In practice, system administrators configure cron jobs on particular machines and find it difficult to port them to other machines that do not have the same system paths.
- If the ingest job fails (if, for example, the network is down), there is no way to retry it. Retries and other such failure-recovery efforts will have to be explicitly coded in our Python program.
- Remote monitoring and one-time, ad hoc executions are not part of the cron interface. If you need to monitor, troubleshoot, and restart the ingest from a mobile device, good luck.

This litany of drawbacks is not unique to cron. They are implicit in any solution that is tied to specific servers. So, how would you do it on the cloud? What you should not do is to create a Compute Engine VM and schedule a cron job on it—that will be subject to some of the same problems!

For resilience and reliability, we need a serverless way to schedule ingest jobs. Obviously, the ingest job will need to be run on some machine somewhere. However, we shouldn't need to manage that machine at all. This is a job that needs perhaps two minutes of compute resources a month. We should be looking for a way to write the ingest code and let the cloud infrastructure take care of provisioning resources, making retries, and providing for remote monitoring and ad hoc execution.

On Google Cloud Platform, Cloud Scheduler provides a way to schedule periodic jobs in a serverless manner. These jobs can involve hitting an HTTP endpoint (which

²³ Shortened form of cron table.

is what we will do), but can also send a message via Cloud Pub/Sub or trigger a Google Kubernetes Engine or Cloud Dataflow job. [Figure 2-16](#) presents our architecture for the monthly ingest job.

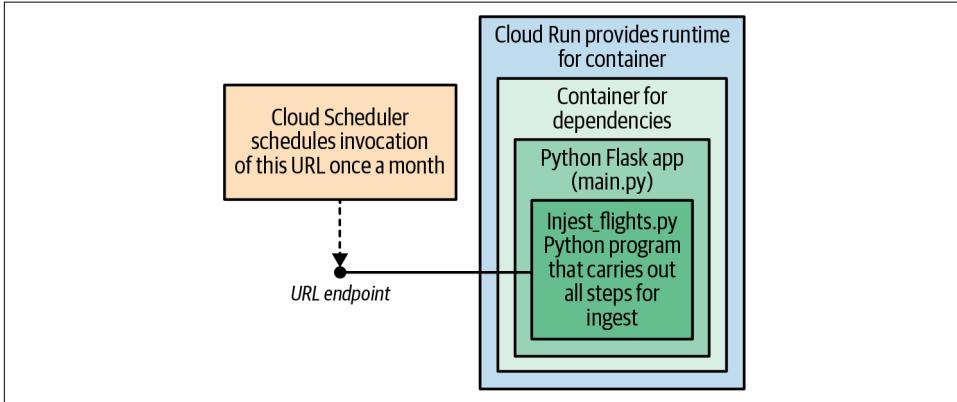


Figure 2-16. The architecture of the monthly ingest job.

First, we will write a standalone `ingest_flights.py` application that is capable of downloading the data for a specific year/month and uploading the data to Cloud Storage. We will invoke the ingest code from a Python Flask application making sure to explicitly capture our dependencies in a `Dockerfile` which describes a Docker container. Cloud Run will run our container.²⁴

The way scheduling works in Cloud Scheduler is that we must specify a URL that will be invoked or a Cloud Pub/Sub topic that must be monitored. Whereas in the previous Linux cron example we specified a script on the server that was running the cron daemon, the Cloud Scheduler endpoint will be a URL that will be visited according to the schedule that we specify (this can be any URL; it doesn't need to be a service that we write). Because our ingest code is a standalone Python program, we will wrap that ingest code into a Python Flask application (`main.py`) so that we can invoke it by using a URL (Flask is a web application framework).

Ingesting in Python

While exploring the data, we carried out the ingest on the command line in Bash. We saved our commands as we went along in the form of Bash scripts. We created our

²⁴ Docker containers are lightweight wrappers around a piece of software (here, the Flask endpoint `main.py`) that contain everything needed to run that software—code (e.g., `ingest_flights.py`), runtime (Python dependencies, etc.), configuration files, and system libraries (here, a specific Linux distribution). Unlike a virtual machine, different containers running on the same machine can share layers of operating system dependencies.

ingest program by simply making a Bash script (*02_ ingest/ingest.sh*) that invokes those intermediate steps:

```
#!/bin/bash
for MONTH in `seq 1 12`; do
    bash download.sh $YEAR $MONTH
done

# upload the raw CSV files to our GCS bucket
bash upload.sh

# load the CSV files into BigQuery as string columns
bash bqload.sh
```

This is the sort of decision that leads to spaghetti-like code that is difficult to unravel and to maintain. There are many assumptions made by this set of Bash scripts in terms of what to download, where the temporary storage resides, and where to upload it. Changing any of these will involve changing multiple scripts. Using Bash to quickly get a handle on the data is a good idea, as is the idea of saving these scripts so as to continue the exploration. But when it comes to making the ingest more systematic and routine, you do not want to use a shell scripting language; a more formal programming language is better.

In this book, we will use Python wherever we can because of its ability to span a wide range of computing tasks, from systems programming to statistics and machine learning. Python is currently the best choice if you need to pick a single language in which to do most of your work. Java is typesafe and performant. Its object-orientation and packaging architectures are suitable for large, multideveloper programs, but it makes the code too verbose. Moreover, the lack of a read–eval–print loop (REPL) interpreter makes Java unwieldy for quick experimentation. C++ is numerically very efficient, but standard libraries for nonnumerical computing are often nonexistent. Scala combines the benefits of Python (easy scriptability, conciseness) with the benefits of Java (type safety, speed), but the tooling for Scala (such as for statistics and visualization) is not as pervasive as it is for Python. Today, therefore, the best choice of programming language is Python. For certain use cases for which speed is important and Python is not performant enough, it might be necessary to use Java.

The ingest program in Python goes through the same four steps as before when we did it manually on the command line:

- Download data from the BTS website to a local file.
- Unzip the downloaded ZIP file and extract the CSV file it contains.
- Upload the CSV file to Google Cloud Storage.
- Load the CSV data into a BigQuery partitioned table.

Whereas our download Bash script got all 12 months of a hardcoded year (2015), our download subroutine in Python will take as input the year and month:

```
def download(YEAR, MONTH, destdir):
    url = os.path.join("https://transtats.bts.gov/PREZIP",
                       "{}_{}.zip".format(YEAR, int(MONTH)))

    filename = os.path.join(destdir, "{}_{}.zip".format(YEAR, MONTH))
    with open(filename, "wb") as fp:
        response = urlopen(url)
        fp.write(response.read())
    return filename
```

Another thing to note is that our Bash script simply downloaded the ZIP file from BTS to the current working directory of the user. However, since our Python script is meant to be executed on demand by the scheduler service, we cannot make assumptions about the directory in which the script will be run. In particular, we don't know whether that directory will be writable and have enough space. Hence, we ask the caller of the function to provide an appropriate destination directory in which to store the downloaded ZIP file.

Here's how to unzip the file and extract the CSV contents:

```
def zip_to_csv(filename, destdir):
    zip_ref = zipfile.ZipFile(filename, 'r')
    cwd = os.getcwd()
    os.chdir(destdir)
    zip_ref.extractall()
    os.chdir(cwd)
    csvfile = os.path.join(destdir, zip_ref.namelist()[0])
    zip_ref.close()
```

Unzipping explodes the size of the file. We can optimize things slightly. Rather than upload the text file, we can gzip it because BigQuery knows how to load gzipped CSV files:

```
gzipped = csvfile + ".gz"
with open(csvfile, 'rb') as ifp:
    with gzip.open(gzipped, 'wb') as ofp:
        shutil.copyfileobj(ifp, ofp)
return gzipped
```

Here's the code to upload the CSV file for a given month to Cloud Storage:

```
def upload(csvfile, bucketname, blobname):
    client = storage.Client()
    bucket = client.get_bucket(bucketname)
    blob = Blob(blobname, bucket)
    blob.upload_from_filename(csvfile)
    gcslocation = 'gs://{}{}'.format(bucketname, blobname)
    print ('Uploaded {} ...'.format(gcslocation))
    return gcslocation
```

The code asks for the `bucketname` (the single-region bucket that was created during our exploration) and a `blobname` (e.g., `flights/201501.csv`) and carries out the upload using the Cloud Storage Python library. Although it can be tempting to simply use the `subprocess` module in Python to invoke `gsutil` operations, it is better not to do so. If you go the `subprocess` route, you will then need to ensure that the Cloud SDK (which `gsutil` comes with) is installed on whichever machine this is going to run on. This won't be a problem in Cloud Run, but might pose problems if you switch the way you provide URL access later (to, say, Google App Engine or Cloud Functions). It is preferable to use pure Python modules when possible and add those modules to `requirements.txt`, as follows:²⁵

```
Flask
google-cloud-storage
google-cloud-bigquery
gunicorn==20.1.0
```

The Flask library will help us handle HTTP requests (covered shortly), and Google Cloud Storage is needed so as to invoke the `get_bucket()` and `upload_from_file_name()` operations. While using the latest version of libraries is okay, it poses the problem that an upgrade to those dependencies might break our code. For production code, it is better to pin the library versions to the ones with which the code has been tested:

```
Flask==2.0.1
google-cloud-storage==1.42.0
google-cloud-bigquery==2.25.1
gunicorn==20.1.0
```

If you do pin libraries, though, you will have to have a process in place to periodically test and upgrade to the latest stable version of your dependencies. Otherwise, your code might go stale or, worse, be insecure because it's using library versions with known vulnerabilities.

We can now write an `ingest()` method that calls the four major steps, plus the verification, in order:

```
def ingest(year, month, bucket):
    """
    ingest flights data from BTS website to Google Cloud Storage
    return cloud-storage-blob-name on success.
    raises DataUnavailable if this data is not on BTS website
    """
    tempdir = tempfile.mkdtemp(prefix='ingest_flights')
    try:
```

²⁵ Hopefully, you know what the required packages are because you installed them using a package manager such as `conda` or `pip`. To find the packages in your development environment, you can use `pip freeze`.

```

        zipfile = download(year, month, tempdir)
        bts_csv = zip_to_csv(zipfile, tempdir)
        gcsloc = f'flights/raw/{year}{month}.csv.gz'
        gcsloc = upload(bts_csv, bucket, gcsloc)
        return bqload(gcsloc, year, month)
    finally:
        print ('Cleaning up by removing {}'.format(tempdir))
        shutil.rmtree(tempdir)

```

The destination directory that we use to stage the downloaded data before uploading to Cloud Storage is obtained using the `tempfile` package in Python. This ensures that if, for whatever reason, there are two instances of this program running at the same time, they will not cause contention issues.

We can try out the code by writing a `main()` that is executed if this program is run on the command line:²⁶

```

if __name__ == '__main__':
    import argparse
    parser = argparse.ArgumentParser(
        description='ingest flights data from BTS website to GCS')
    parser.add_argument('--bucket', help='GCS bucket to upload data to',
                       required=True)
    parser.add_argument('--year', help='Example: 2015.', required=True)
    parser.add_argument('--month', help='01 for Jan.', required=True)
    try:
        args = parser.parse_args()
        gcsfile = ingest(args.year, args.month, args.bucket)
        print (f'Success ... ingested to {gcsfile}')
    except DataUnavailable as e:
        print ('Try again later: {}'.format(e.message))

```

Specifying a valid month ends with a new (or replaced) file on Cloud Storage:

```

$ ./ingest_flights.py --bucket cloud-training-demos-ml \
                     --year 2015 --month 01
...
Success ... ingested to gs://cloud-training-demos-ml/flights/201501.csv

```

Trying to download a month that is not yet available results in an error message:

```

$ ./ingest_flights.py --bucket cloud-training-demos-ml \
                     --year 2029 --month 01
...
HTTP Error 403: Forbidden

```

On Cloud Scheduler, this will result in the call failing and being retried subject to a maximum number of retries. Retries will also happen if the BTS web server cannot be reached.

²⁶ The full program is available as `ingest_flights.py` in the book's GitHub repository—try it out.

At this point, we have the equivalent of our exploratory Bash scripts, but with some additional resilience, repeatability, and fault tolerance built in. Our Python program expects us to provide a year, month, and bucket. However, if we are doing monthly ingestions, we already know which year and month we need to ingest. No, not the current month—recall that there is a time lag between the flight events and the data being reported by the carriers to the BTS. Instead, it is simply the month after whatever files we already have on Cloud Storage!²⁷ So, we can automate this, too:

```
def next_month(bucketname):
    """
        Finds which months are on GCS, and returns next year,month to download
    """
    client = storage.Client()
    bucket = client.get_bucket(bucketname)
    blobs = list(bucket.list_blobs(prefix='flights/raw/'))
    files = [blob.name for blob in blobs if 'csv' in blob.name] # csv files only
    lastfile = os.path.basename(files[-1]) # e.g. 201503.csv
    year = lastfile[:4]
    month = lastfile[4:6]
    dt = datetime.datetime(int(year), int(month), 15) # 15th of month
    dt = dt + datetime.timedelta(30) # will always go to next month
    return '{}'.format(dt.year), '{:02d}'.format(dt.month)
```

To get the next month given that there is a file, say `201503.csv`, on Cloud Storage, we add 30 days to the Ides of March—this gets around the fact that there can be anywhere from 28 days to 31 days in a month, and that `timedelta` requires a precise number of days to add to a date.

By changing the year and month to be optional parameters, we can try out the ingest program’s ability to find the next month and ingest it to Cloud Storage. We simply add:

```
if args.year is None or args.month is None:
    year, month = next_month(args.bucket)
else:
    year = args.year
    month = args.month
    gcsfile = ingest(year, month, args.bucket)
```

Having ingested the data, it is a good idea to verify that the end-to-end pipeline worked as intended. We could count the number of rows in the CSV file and assert that this is appropriately large (at least 10,000, for example) and equal to the number of rows corresponding to the month in BigQuery.

²⁷ Our decision earlier to name the files as `<YYYYMM.csv>` (with zero padding month) was a good data engineering choice. By doing that, we can reliably count the number of months in our archive and check if the data for a specific month has already been downloaded.

Now that we have an ingest program that is capable of updating our Cloud Storage bucket one month at a time, we can move on to building the scaffolding to have it be executed in a serverless way.

Cloud Run

Cloud Run is a serverless framework that provides an autoscaling, resilient runtime for containerized code. The container (see [Figure 2-16](#)) will consist of code that listens for requests or events. Cloud Run abstracts away all the infrastructure management that would otherwise be needed.

Now that we have a Python function that will do the ingest, we will wrap it inside a web application. To write the web application, we will use [Flask](#), which is a light-weight Python web application framework, and as a web server, we will use [Gunicorn](#). Flask provides the ability to invoke Python code in response to an HTTP request, while Gunicorn will listen to HTTP requests and send them to the Flask app. Our container will consist of the Gunicorn server, Flask application, and its dependencies. This is expressed in the form of a Dockerfile:

```
FROM python:3.8-slim

# Copy local code to the container image.
ENV APP_HOME /app
WORKDIR $APP_HOME
COPY . .

# Install production dependencies.
RUN pip install --no-cache-dir -r requirements.txt

# Run the web service on container startup.
# Timeout is set to 0 to disable the timeouts of
# the workers to allow Cloud Run to handle instance scaling.
CMD exec gunicorn --bind :$PORT --workers 1 \
                    --threads 8 --timeout 0 main:app
```

In our `main.py`, we have a function that gets invoked in response to the URL trigger:

```
import logging
from flask import escape
from ingest_flights import *

app = Flask(__name__)

@app.route("/", methods=['POST'])
def ingest_flights(request):
    try:
        json = request.get_json()

        year = escape(json['year']) if 'year' in json else None
```

```

month = escape(json['month']) if 'month' in json else None
bucket = escape(json['bucket']) # required

if year is None or month is None or len(year) == 0 or len(month) == 0:
    year, month = next_month(bucket)
tbleref, numrows = ingest(year, month, bucket)
ok = 'Success ... ingested {} rows to {}'.format(numrows, tbleref)
return ok
except Exception as e:
    logging.exception('Try again later')

```

Essentially, `main.py` has a single function that receives a Flask request object, from which we can extract the JSON (JavaScript Object Notation) payload of the HTTP POST by which the Cloud Run will be triggered. We get the next month by looking to see what months are already in the bucket and then ingest the necessary data using the existing code in the module `ingest_flights`. We can deploy the our codebase as a container to Cloud Run using:

```

NAME=ingest-flights-monthly
REGION=us-central1

gcloud run deploy $NAME --region $REGION --source=$(pwd) \
--platform=managed --timeout 12m

```

But there are a couple of serious security and governance problems if we do this.

Securing Cloud Run

What are the security problems?

- Anyone can invoke the URL and cause our dataset to get updated. We have to disallow unauthenticated users.
- Allowing this code to run with our user account's permissions will pollute any audit logs since we are not actually running the ingest interactively. We need to create a separate account so that the Cloud Run service can run with that identity.
- Allowing this code to run with our user account's permissions is also quite dangerous because our user account will typically have very broad permissions. We'd like to restrict the tasks that this automated service can do: we want it to be able to write only to specific Cloud Storage buckets and BigQuery tables.

The way to address the first point is to disallow unauthenticated users. The way to accomplish the second requirement is to specify that the Cloud Run service will have to run as a service account. A service account is an account whose identity is meant to be taken on by automated services. Like any identity, it can be configured to have specific and limited permissions. Therefore, before we can deploy the Cloud Run service, we will need to create a service account. Service accounts have email addresses of the form `${service-account-name}@${project-id}.iam.gserviceaccount.com`. In our

case, the service name is `svc-monthly-ingest@cloud-training-demos.iam.gserviceaccount.com`.

You can create a service account by going to the IAM area in the web console, but as usual, I prefer to script things:²⁸

```
SVC_ACCT=svc-monthly-ingest
PROJECT_ID=$(gcloud config get-value project)
BUCKET=${PROJECT_ID}-cf-staging
REGION=us-central
SVC_PRINCIPAL=serviceAccount:${SVC_ACCT}@${PROJECT_ID}.iam.gserviceaccount.com

gcloud iam service-accounts create $SVC_ACCT \
    --display-name "flights monthly ingest"
```

Then, we make the service account the admin of the staging GCS bucket so that it can read, write, list, delete, etc., on this bucket (and only this bucket):

```
gsutil mb -l $REGION gs://$BUCKET
gsutil uniformbucketlevelaccess set on gs://$BUCKET
gsutil iam ch ${SVC_PRINCIPAL}:roles/storage.admin gs://$BUCKET
```

We will also allow the service account to create and delete partitions on tables in just the BigQuery dataset `dsongcp` (and no other datasets):

```
bq --project_id=${PROJECT_ID} query --nouse_legacy_sql \
    "GRANT `roles/bigquery.dataOwner` ON SCHEMA dsongcp TO '$SVC_PRINCIPAL' "

gcloud projects add-iam-policy-binding ${PROJECT_ID} \
    --member ${SVC_PRINCIPAL} \
    --role roles/bigquery.jobUser
```

Are these permissions sufficient to carry out all the steps of our data processing pipeline? One way to check is to try to ingest a month of data when running as this service account. To do so, we will have to *impersonate* the service account:²⁹

- Visit the [Service Accounts section](#) of the GCP Console.
- Select the newly created service account `svc-monthly-ingest` and click Manage Keys.
- Add key (Create a new JSON key) and download it to a file named `tempkey.json`. Transfer this key file to your Cloud Shell instance.

²⁸ See `02_ingest/monthlyupdate/01_setup_svc_acct.sh`.

²⁹ See the instructions in `README.md` in `02_ingest`.

- Run:

```
gcloud auth activate-service-account \
--key-file tempkey.json
```

- Try ingesting one month:

```
./ingest_flights.py --bucket $BUCKET \
--year 2015 --month 03 --debug
```

Once you have ensured that the service account has all the necessary permissions, go back to running commands as yourself using `gcloud auth login`.

Deploying and Invoking Cloud Run

Now that we have the code for the Flask application and a service account with the right permissions, we can deploy the code to Cloud Run to run as this service account:³⁰

```
NAME=ingest-flights-monthly
SVC_ACCT=svc-monthly-ingest
PROJECT_ID=$(gcloud config get-value project)
REGION=us-central1
SVC_EMAIL=${SVC_ACCT}@${PROJECT_ID}.iam.gserviceaccount.com

gcloud run deploy $NAME --region $REGION --source=$(pwd) \
--platform=managed --service-account ${SVC_EMAIL} \
--no-allow-unauthenticated --timeout 12m
```

Recall that we started the discussion on securing the Cloud Run instance by saying that we would disallow unauthenticated users and have the Cloud Run service run as a service account. Note how we are turning on both these options when we deploy to Cloud Run.

Once the application has been deployed to Cloud Run, we can try accessing the URL of the service with our authentication details in the header of the web request and a JSON message as its POST:³¹

```
# Feb 2015
echo {"year\":\"2015\", \"month\":\"02\", \"bucket\":\"\${BUCKET}\\"} \
> /tmp/message

curl -k -X POST $URL \
-H "Authorization: Bearer $(gcloud auth print-identity-token)" \
-H "Content-Type:application/json" --data-binary @/tmp/message
```

³⁰ See `02_ ingest/monthlyupdate/02_deploy_cr.sh`.

³¹ See `02_ ingest/monthlyupdate/03_call_cr.sh`.

But what is the URL? Cloud Run generates the URL when we deploy the container, and we can obtain it using:

```
gcloud run services describe ingest-flights-monthly \
    --format 'value(status.url)')
```

Changing the message to provide only the bucket (no year or month) will make the service get the next month:

```
echo {"bucket": "${BUCKET}"} > /tmp/message
curl -k -X POST $URL \
    -H "Authorization: Bearer $(gcloud auth print-identity-token)" \
    -H "Content-Type:application/json" --data-binary @/tmp/message
```

Scheduling Cloud Run

Our intent is to automatically invoke CloudRun once a month. We can do that using Cloud Scheduler, which is also serverless and doesn't require us to manage any infrastructure. We simply specify the schedule and the URL to hit. This URL is what came from the output of the Cloud Run deployment command in the previous section:

```
echo {"bucket": "${BUCKET}"} > /tmp/message
cat /tmp/message

gcloud scheduler jobs create http monthlyupdate \
    --description "Ingest flights using Cloud Run" \
    --schedule="8 of month 10:00" \
    --time-zone "America/New_York" \
    --uri=$URL --http-method POST \
    --oidc-service-account-email $EMAIL \
    --oidc-token-audience=$URL \
    --max-backoff=7d \
    --max-retry-attempts=5 \
    --max-retry-duration=2d \
    --min-backoff=12h \
    --headers="Content-Type=application/json" \
    --message-body-from-file=/tmp/message
```

The preceding parameters would make the first retry happen after 12 hours. Subsequent retries are increasingly farther apart, up to a maximum of 2 days between attempts. We fail the task permanently if it fails five times within a defined time period and the task is more than 7 days old (both limits must be passed for the task to fail).

To try out the Cloud Scheduler, we could wait for the 8th of the month to roll around. Or we could go to the GCP web console and click on Run Now. Unfortunately, when I tried it, it wouldn't work because Cloud Scheduler wanted to run as the service account while I was logged in as myself. So, I gave myself the ability to impersonate the service account by going to the Service Accounts part of the web console. Once I did that, I was able to get Run Now to work.

The monthly update mechanism works if you have the previous month's data on Cloud Storage. If you start out with only 2015 data, updating it monthly means that you will inevitably be many months behind. So, you will need to run it ad hoc until your data is up-to-date and then let the cron service take care of things after that. Alternatively, you can take advantage of the fact that the ingest task is cheap and non-intrusive when there is no new data. So, you can change the schedule to be every day instead of every month. A better solution is to change the ingest task so that if it is successful in ingesting a new month of data, it immediately tries to ingest the next month. This way, your program will crawl month-by-month to the latest available month and then keep itself always up-to-date.

At this point, it is worth reflecting a bit on what we have accomplished. We are able to ingest data and keep it up-to-date by doing just these steps:

- Write some Python code.
- Deploy that Python code to the Google Cloud Platform.

We did not need to manage any infrastructure in order to do this. We didn't install any OS, manage accounts on those machines, keep them up-to-date with security patches, maintain failover systems, and so on—a serverless solution that consists simply of deploying code to the cloud is incredibly liberating. Not only is our ingest convenient, it is also very inexpensive—everything scales down to zero when it is not being used. All this falls comfortably within the free tier or might cost less than 5¢ a month.

Summary

The US BTS collects, and makes publicly available, a dataset of flight information. It includes nearly a hundred fields, including scheduled and actual departure and arrival times, origin and destination airports, and flight numbers of every domestic flight scheduled by every major carrier. We will use this dataset to estimate the likelihood of an arrival delay of more than 15 minutes of the specific flight whose outcome needs to be known in order for us to decide whether to cancel the meeting.

There are three possible data processing architectures on the cloud for large datasets: scaling up, scaling out with sharded data, and scaling out with data in situ. Scaling up is very efficient, but is limited by the size of the largest machine you can get a hold of. Scaling out is very popular but requires that you preshard your data by splitting it among compute nodes, which leads to maintaining expensive clusters unless you can hit sustained high utilization. Keeping data in situ is possible only if your data center supports petabits per second of bisectional bandwidth so that any file can be moved to any compute node in the data center on demand. Because Google Cloud Platform has this capability, we will upload our data to Google Cloud Storage (a blob storage

that is not presharded) and to BigQuery, which will allow us to carry out interactive exploration on large datasets.

To automate the ingest of the files, we reverse engineered the BTS's web form and obtained the format of the POST request that we need to make. With that request in hand, we were able to write a Bash script to pull down 12 months of data, uncompress the ZIP file, and load the data into BigQuery. It is quite straightforward to change this script to loop through multiple years.

We discussed the difference between strong consistency and eventual consistency and how to make the appropriate trade-off imposed by Brewer's CAP theorem. In this case, we wanted strong consistency and did not need global availability. Hence, we chose to use a single-region bucket. We then uploaded the downloaded, unzipped, and cleaned CSV files to Google Cloud Storage.

To schedule monthly downloads of the BTS dataset, we made our download and cleanup Python program and made it callable from Cloud Run so that it was completely serverless. We used Cloud Scheduler to periodically request the Cloud Run application to download BTS data, unzip it, and upload it to both Cloud Storage and BigQuery.

Code Break

This is the point at which you should put this book off to the side and attempt to repeat all the things I've talked about. All the code snippets in this book have corresponding code in the GitHub repository.

I strongly encourage you to play around with the code in *02_ingest* with the goal of understanding why it is organized that way and being able to write similar code yourself. At minimum, though, you should do the following:

- Open Cloud Shell and git clone the book's code repository as explained in [Chapter 1](#).
- Go to the *02_ingest* folder of the repository.
- Go to the Storage section of the GCP web console and create a new regional bucket—choose the region (such as us-central1) that is closest to where you live.
- Run `./ingest.sh` providing the name of the bucket that you just created. This will populate your bucket and BigQuery dataset with data from 2015. Although you can change the year loop in this file to download all the data corresponding to 2015–2019, I recommend that you hold off on downloading all the data until [Chapter 12](#) because processing 5 years of data will make every subsequent thing in this book take five times longer.

- Because software changes, an up-to-date list of the preceding steps is available in the course repository in *02_ ingest/README.md*. This is true for all the following chapters.

Suggested Resources

A key aspect of cloud-native architectures is reliance on fully managed, autoscaling services. That is Principle 3 from Tom Grey's short, but informative 2019 blog post, "[5 Principles for Cloud-Native Architecture](#)". How do you think the other principles apply to data analytics and AI?

When every department in your organization uses a hub-and-spoke architecture, and provides data access to other departments, the architecture you will have across the entire organization is a *data mesh*. This 2020 article from Medium.com, "[Building a Data Platform to Enable Analytics and AI-Driven Innovation](#)", describes the journey to get there.

Cloud Storage is a blob store. This 2020 article from Netapp, "[Storage Options in Google Cloud: Block, Network File, and Object Storage](#)" by Bruno Almeida, is a great introduction to the various storage options available on Google Cloud.

This [handy flowchart](#) summarizes the key considerations when deciding where to store your data. While you are there, take a look at some of the other flowcharts.

[Google BigQuery](#) is the heart of data analytics and AI architectures on Google Cloud. The O'Reilly Media book *[BigQuery: The Definitive Guide](#)* by Valliappa Lakshmanan and Jordan Tigani is a great place to learn about BigQuery.

[Cloud Run](#) is a serverless execution environment for containerized applications. [Cloud Scheduler](#) provides a serverless way to invoke services on a schedule. Having Cloud Scheduler trigger Cloud Run is a common [pattern that is documented](#). There is a list of [guides for Cloud Run](#) and a list of [guides for Cloud Storage](#). Do you see the pattern? Try to find the list of guides for BigQuery.

A lot of data wrangling and early exploration involves Unix tools. A good introduction to the Unix shell and tools is this [online tutorial](#). How to use command-line tools to do data science is the focus of *[Data Science at the Command Line](#)* by Jeroen Janssens (O'Reilly).

“Scaling up” and “scaling out” seem very easy when we say the words or draw neat pictures. However, they are extremely hard problems to solve. As described in this [2018 CIO article by Stephen Watts](#), scaling up requires being able to design application-specific integrated circuits (ASICs) when running into limits on how many transistors we can pack into a chip. The fundamental algorithm that underlies scaling out is the Paxos consensus protocol, which is fiendishly complicated to get

right (see [Deniz Altınbükən's website](#), Paxos Made Moderately Complex). This is another reason to choose a public cloud.

CHAPTER 3

Creating Compelling Dashboards

In Chapter 2, we ingested on-time performance data from the US Bureau of Transportation Statistics (BTS) so as to be able to model the arrival delay given various attributes of an airline flight—the purpose of the analysis is to cancel a meeting if the probability of the flight arriving within 15 minutes of the scheduled arrival time is less than 70%.

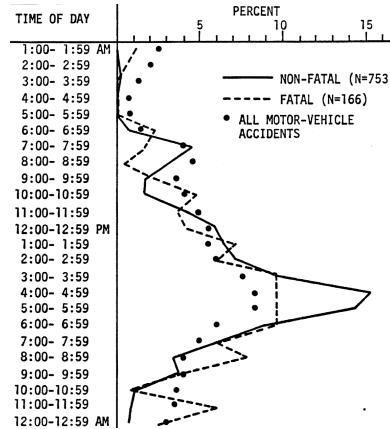
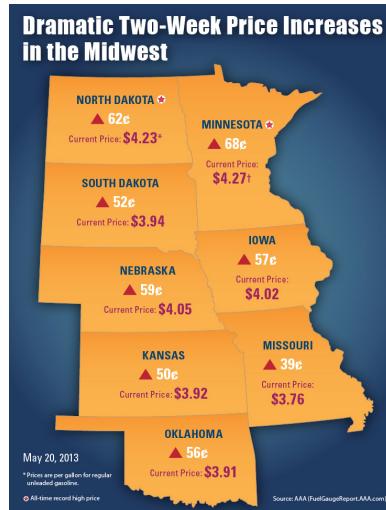
Before we delve into building statistical and machine learning models, it is important to explore the dataset and gain an intuitive understanding of the data—this is called *exploratory data analysis*, and it's covered in more detail in Chapter 5. You should always carry out exploratory data analysis for any dataset that will be used as the basis for decision making. In this chapter, though, I talk about a different aspect of depicting data—of depicting data to end users and decision makers so that they can understand the recommendation that you are making. The audience of these visual representations, called *dashboards*, that we talk about in this chapter is not other data scientists, but is instead the end users. Keep the audience in mind as we go through this chapter, especially if you come from a data science background—the purpose of a dashboard is to explain an existing model, not to develop it. A dashboard is an end-user report that is interactive, tailored to end users, and continually refreshed with new data. See Table 3-1.

Table 3-1. A dashboard is different from exploratory data analysis

	For decision makers	For data scientists
Usage pattern	Dashboards	Exploratory data analysis
Kinds of depictions	Current status, gauges, trendlines	Model fits with error bars, kernel density estimates
What does it explain?	Model recommendations and confidence	Input data, feature importance, model performance, etc.

	For decision makers	For data scientists
Data represented	Subset of dataset, tailored to user's context	Aggregate of historical data
Typical tools	Data Studio, Tableau, Qlik, Looker, plotly, D3, shiny apps, etc.	Jupyter, Python, R Studio, S-plus, matplotlib, seaborn, Matlab, etc.
Mode of interaction	GUI-driven	Code-driven
Update	Real time	Not real time
Covered in	Chapter 3 , Chapter 4	Chapter 5

Example



From AAA safety and educational foundation

From AAA fuel gage report, May 2013

Very often, this step of creating end-user visual depictions goes by the anodyne name of “visualization,” as in visualizing data. However, I have purposefully chosen not to call it by that name because there is more to this than throwing together a few bar graphs and charts. Dashboards are highly visual, interactive reports that have been designed to depict data and explain models. When used in this way, dashboards provide business intelligence (BI).



All of the code snippets in this chapter are available in the folder [03_sqlstudio](#) of the book’s GitHub repository. See the *README.md* file in that directory for instructions on how to do the steps described in this chapter.

Explain Your Model with Dashboards

The purpose of this step in the modeling process is not simply to depict the data but to improve your users' understanding of how the model behaves. Whenever you are designing the display of a dataset, evaluate the design in terms of three aspects:

- Does it accurately and honestly depict the data? This is important when the raw data itself can be a basis for decision making.
- How well does it help envision not just the raw data, but the information content embodied in the data? Will the typical user know whether they need to take action after looking at the graphic? This is crucial for the cases when you are relying on human pattern recognition and interaction to help reveal insights about the environment in which the data was collected.
- Is it constructed in such a way that it explains the model being used to provide recommendations?

You want to build displays that are always accurate and honest. At the same time, the displays need to be interactive so as to provide viewers with the ability to play with the data and gain insights. Insights that users have gained should be part of the display of that information going forward in such a way that those insights can be used to explain the data.

The last point, that of explanatory power, is very important. The idea is to disseminate data understanding throughout your company. A statistical or machine learning model that you build for your users will be considered a black box, and while you might get feedback on when it works well and when it doesn't, you will rarely get pointed suggestions on how to actually improve that model in the field. In many cases, your users will use your model at a much more fine-grained level than you ever will because they will use your model to make a single decision, whereas in both training and evaluation, you would have been looking at model performance as a whole. Explainability is also critical to catch situations where the model is amplifying unfair bias.¹

Although this holistic overview is useful for statistical rigor, you need people taking a close look at individual cases, too. Because users are making decisions one at a time, they are analyzing the data one scenario at a time. If you provide your users not just with your recommendation, but also with an explanation of why you are

¹ When Amazon built a hiring tool to help select resumes using machine learning, they were able to use explainability to recognize what the model was keying off—apparently, the machine learning model **penalized resumes** that included terms more commonly found in women's resumes. Amazon was able to catch this error and not use the ML model to actually evaluate candidates. Had explainability not been part of the workflow, many women would have been unfairly not considered for jobs at Amazon.

recommending it, they will begin to develop insights into your model. However, your users will only be able to develop such insights into the problem and your recommendations if you give them ways to observe the data that went into your model. Give enough users ways to view and interact with your data, and you will have unleashed a never-ending wave of innovation as users suggest improvements and factors the model should be considering.

Your users have other activities that require their attention. Why would they spend their time looking at your data? One of the ways to entice them to do that is by making the depiction of the information compelling. In my experience,² the most compelling displays are displays of real-time information in context. You can show people the average airline delay at JFK on January 12, 2012, and no one will care. But show a traveler in Chicago the average airline delay at ORD *right now* and you will have their interest—the difference is that the data is in context (O’Hare Airport, or ORD, for a traveler in Chicago) and that it is real-time information.

In this chapter, therefore, we will look at building dashboards that combine accurate depictions with explanatory power and interactivity in a compelling package. This seems to be a strange time to be talking about building dashboards—shouldn’t the building of a dashboard wait until after we have built the best possible predictive model?

Why Build a Dashboard First?

Building a dashboard when building a machine learning model is akin to building a form or survey tool to help you build the machine learning model. To build powerful machine learning models, you need to understand the dataset and devise features that help with prediction. By building a dashboard, you get to rope in the eventual users of the predictive model to take a close look at your data. Their fine-grained look at the data (remember that everyone is looking at the data corresponding to their context) will complement your overarching look at it. As they look at the data and keep sending suggestions and insights about the data to your team, you will be able to incorporate them into the machine learning model that you actually build.

² When I worked on developing machine learning algorithms for weather prediction, nearly every one of the suggestions and feature requests that I received emerged when the person in question was looking at the real-time radar feed. There would be a storm, my colleague would watch it go up on radar, observe that the tracking of the storm was patchy, and let me know what aspect of the storm made it difficult to track. Or, someone would wake up, look at the radar image, and discover that birds leaving to forage from their roost had been wrongly tagged as a mesocyclone. It was all about real-time data. No matter how many times I asked, I never ever got anyone to look at how the algorithms performed on historical data. It was also often about Oklahoma (where our office was) because that’s what my colleagues would concentrate on. Forecasters from around the country would derisively refer to algorithms that had been hypertuned to Oklahoma supercells.

In addition, when presented with a dataset, you should be careful that the data is the way you imagine it to be. There is no substitute for exposing and exploring the data to ensure that. Doing such exploratory data analysis with an immediately attainable milestone—building a dashboard from your dataset—is a fine way to do something real with the data and develop awareness of the subtleties of your data. Just as you often understand a concept best when you explain it to someone, building an explanatory display for your data is one of the best ways to develop your understanding of a dataset. The fact that you have to visualize the data in order to do basic pre-processing such as outlier detection makes it clear that building visual representations is work you will be doing anyway. If you are going to be doing it, you might as well do it well, with an eye toward its eventual use in production.

Eventual use in production is the third reason why you should develop the dashboard first instead of leaving it as an afterthought. Building explanatory power should be constantly on your mind as you develop the machine learning model. Giving users just the machine learning model will often go down poorly—they have no insight into why the system is recommending whatever it does. Adding explanations to the recommendations is more likely to succeed. For example, if you accompany your model prediction with five of the most salient features presented in an explanatory way, it will help make the model output more believable and trustworthy.

Even for cases in which the system performs poorly, you will receive feedback along the lines of “the prediction was wrong, but it is because Feature #3 was fishy. I think maybe you should also look at Factor Y.” In other words, shipping your machine learning model along with an explanation of its behavior gets you more satisfied users, and users whose criticism will be a lot more constructive. It can be tempting to ship the machine learning model as soon as it is ready, but if there is a dashboard already available (because you were building it in parallel), it is easier to counsel that product designers consider the machine learning model and its explanatory dashboard as the complete product.

Finally, creating a dashboard is fun. It will help you build a user base quickly and show end users what’s in it for them.



Explanations can be a double-edged sword because humans are not fully rational beings. Explanations can be the result of *apophenia*—the tendency of humans to see meaningful patterns even when there are none. This can lead to *motivated reasoning*—the tendency of humans to create justifications for decisions that are more desirable at an emotional level. The combination of apophenia and motivated reasoning can lead to just-so stories that attempt to justify whatever the state of the world is on the basis of spurious explanations. As data scientists, we should realize that we too are human. We need to be careful to set aside these biases, consider counterfactuals, and be willing to revise our initial judgments. Easier said than done, of course.

Where should these dashboards be implemented? Find out the environment that gets the largest audience of experts and eventual users and build your dashboard to target that environment.

Your users might already have a visualization interface with which they are familiar. Especially when it comes to real-time data, your users might spend their entire work-day facing a visualization program that is targeted toward power users—this is true of weather forecasts, air traffic controllers, and options traders. If that is the case, look for ways to embed your visualizations into that interface. In other cases, your users might prefer that your visualizations be available from the convenience of their web browser. If this is the case, look for a visualization tool that lets you share the report as an interactive, commentable document (not just a static web page). In many cases, you might have to build multiple dashboards for different sets of users (don't shoe-horn everything into the same dashboard).

Accuracy, Honesty, and Good Design

Because the explanatory power of a good dashboard is why we are building visualizations, it is important to ensure that our explanations are not misleading. In this regard, it is best not to do anything too surprising. Although modern-day visualization programs are chock-full of types of graphs and palettes, it is best to pair any graphic with the idiom for which it is appropriate. For example, some types of graphics are better suited to relational data than others, and some graphics are better suited to categorical data than to numerical data.

Broadly, there are four fundamental types of graphics: relational (illustrating the relationship between pairs of variables), time series (illustrating the change of a variable over time), geographical maps (illustrating the variation of a variable by location), and narratives (to support an argument). Narrative graphics are the ones in magazine spreads, which win major design awards. The other three are more worker-like.

You have likely seen enough graphical representations to realize intuitively that the graph is somehow wrong when you violate an accuracy, honesty, or aesthetic principle,³ but this section of the book lists a few of the canonical ones. For example, it is advisable to choose line graphs or scatter plots for relational graphics and to ensure that autoscaling of the axes doesn't portray a misleading story about your data. A good design principle is that your time series graphs should be more horizontal than vertical, and that it is the data lines and not "chart junk" (grid lines, labels, etc.) that ought to dominate your graphics. Maximizing the ratio of data to space and ink is a principle that will stand you in good stead when it comes to geographical data—ensure that the domain is clipped to the region of interest, and go easy on place names and other text labels.

Just as you probably learned to write well by reading good writers, one of the best ways to develop a feel for accurate and compelling graphics is to increase your exposure to good exemplars. *The Economist* newspaper has a [Graphic Detail blog](#)⁴ that is worth following—they publish a chart, map, or infographic every weekday, and these run the gamut of the fundamental graphics types. [Figure 3-1](#) shows a graphic from the blog published on Nov. 25, 2016.

The graphic depicts the increase in the number of coauthors on scientific papers over the past two decades. The graphic itself illustrates several principles of good design. It is a time series, and as you'd expect of this type of graphic, the time is on the horizontal axis and the time-varying quantity (number of authors per article or the number of articles per author) is on the vertical axis. The vertical axis values start out at zero, so that the height of the graphs is an accurate indicator of magnitude. Note how minimal the chart junk is—the axes labels and gridlines are very subtle and the title doesn't draw attention to itself. The data lines, on the other hand, are what pop out. Note also the effective use of repetition—instead of all the different disciplines (Economics, Engineering, etc.) being on the same graph, each discipline is displayed on its own panel. This serves to reduce clutter and makes the graphs easy to interpret. Each panel has two graphs, one for authors per article and the other for articles per author. The colors remain consistent across the panels for easy comparison, and the placement of the panels also encourages such comparisons. We see, for example, that the increase in number of authors per article is not accompanied by an increase in articles per author in any of the disciplines, except for Physics & Astronomy. Perhaps the physicists and astronomers are gaming the system?

³ If you are not familiar with design principles, I recommend *The Visual Display of Quantitative Information* by Edward Tufte (Graphics Press).

⁴ *The Economist* is published weekly as an 80-page booklet stapled in the center, and each page is about the size of a letter paper. However, for historical reasons, the company refers to itself as a newspaper rather than a magazine.

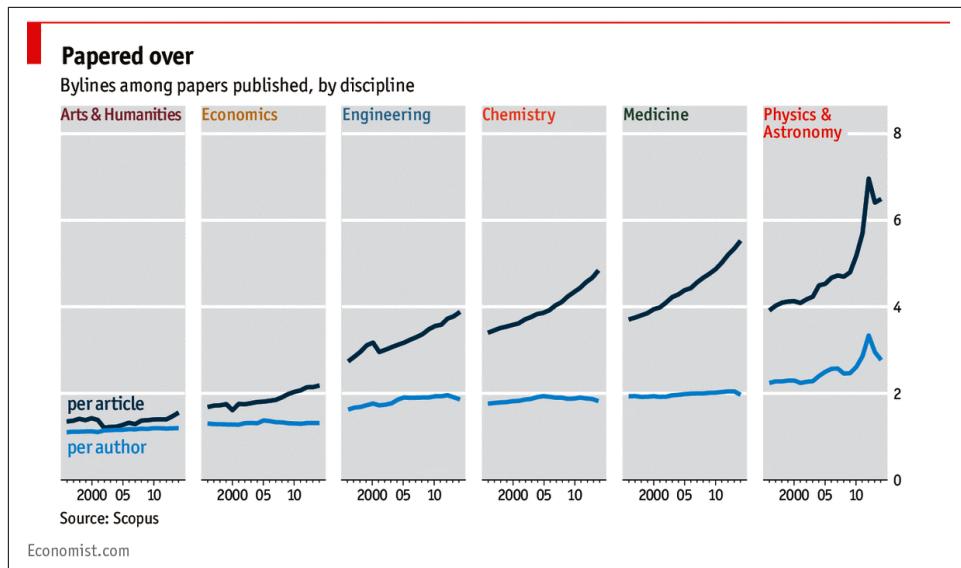


Figure 3-1. This graphic from The Economist shows an increase in the number of authors of papers in various academic disciplines over time.

The graphic does, however, subtly mislead viewers who are in a hurry. Take a moment and try to critique the graphic—figure out how a viewer might have been misled. It has to do with the arrangement of the panels. It appears that the creator of the graphic has arranged the panels to provide a pleasing upward trend between the panels, but this upward trend is misleading because there is no relationship between the number of authors per article in Economics in 2016 and the same quantity in Engineering in 1996. This misdirection is concerning because the graph is supposed to support the narrative of an increasing number of authors, but the increase is not from one author to six authors over two decades—the actual increase is much less dramatic (for example, from four to six in Medicine). However, a viewer who only glances at the data might wrongly believe that the increase in the number of authors is depicted by the whole graph and is therefore much more than it really is.

Loading Data into Cloud SQL

To create dashboards to allow interactive analysis of the data, we will need to store the data in a manner that permits fast random access and aggregations. Because our flight data is tabular, SQL is a natural choice, and if we are going to be using SQL, we should consider whether a relational database meets our needs. Relational databases are a mature technology and remain the tool of choice for many business problems. Relational database technology is well known and comes with a rich ecosystem of

interoperable tools. The problem of standardized access to relational databases from high-level programming languages is pretty much solved.

PostgreSQL is a very popular open source relational database that is used in production at many enterprises. In addition to its high performance, PostgreSQL is easy to program against—it supports ANSI SQL, geographic information system (GIS) functionality, client libraries in a variety of programming languages, and standard connector technologies such as Open Database Connectivity (ODBC) and Java Database Connectivity (JDBC).

Create a Google Cloud SQL Instance

Google **Cloud SQL** offers a managed database service that supports PostgreSQL, MySQL, and SQL Server. Cloud SQL manages backups, patches, updates, and even replication while providing for global availability, automatic failover, and high uptime. For best performance, choose a machine whose RAM is large enough to hold your largest table in memory—as of this writing, available machine types range from a single CPU with less than 4 GB of memory all the way to a 96 CPU machine with 624 GB of memory. Balance this desire for speed with the monthly cost of a machine, of course.

Let's configure a Cloud SQL instance, create a database table in it, and load the table with the data we ingested into Cloud Storage. You can do all these things on the command line using `gcloud`, but let's begin by using the **SQL section of Cloud Platform Console** and select Create Instance. Choose PostgreSQL and then fill out the form as follows (see [Figure 3-2](#)):

- Call the instance “flights.”
- Generate a strong password by clicking on the GENERATE button.
- Choose the default PostgreSQL version.
- Choose the region where your bucket of CSV data exists.
- Choose a single zone instance since we are just trying it out. We won't take this to production.
- Choose a Standard machine type with 2 vCPU.
- Click Create Instance, accepting all the other defaults.

Create a PostgreSQL instance

Instance info

Instance ID * flights

Use lowercase letters, numbers, and hyphens. Start with a letter.

Password * GENERATE

Set a password for the default admin user "postgres". [Learn more](#)

Database version * PostgreSQL 13

Choose region and zonal availability

For better performance, keep your data close to the services that need it. Region is permanent, while zone can be changed any time.

You have the location org policy in effect that is restricting your regions.

Region us-central1 (Iowa)

Zonal availability

Single zone
In case of outage, no failover. Not recommended for production.

Multiple zones (Highly available)
Automatic failover to another zone within your selected region. Recommended for production instances. Increases cost.

SPECIFY ZONES

Customize your instance

You can also customize instance configurations later

SHOW CONFIGURATION OPTIONS

CREATE INSTANCE CANCEL

Summary	
Region	us-central1 (Iowa)
DB Version	PostgreSQL 13
vCPUs	4 vCPU
Memory	26 GB
Storage	100 GB
Network throughput (MB/s) <small>?</small>	1,000 of 2,000
Disk throughput (MB/s) <small>?</small>	Read: 48.0 of 240.0 Write: 48.0 of 240.0
IOPS <small>?</small>	Read: 3,000 of 15,000 Write: 3,000 of 15,000
Connections	Public IP
Backup	Automated
Availability	Single zone
Point-in-time recovery	Enabled

Figure 3-2. Creating a PostgreSQL instance using the web console.

Interacting with Google Cloud Platform

Instead of filling out the dialog box by hand, we could have used the command-line tool `gcloud` from Cloud Shell (or any other machine that has `gcloud` installed); here's how to do that:

```
gcloud sql instances create flights \
--database-version=POSTGRES_13 --cpu=2 --memory=8GiB \
--region=us-central1 --root-password=somestrongpassword
```

In the rest of the book, I will show you just the `gcloud` commands, but you don't need to memorize them. You can use `--help` at any time on `gcloud` to get a list of options. For example:

```
gcloud sql instances create --help
```

will give you all the options available to create Cloud SQL instances (the database-version, its zone, etc.), whereas:

```
gcloud sql instances --help
```

will give you all the ways in which you can work with Cloud SQL instances (`create`, `delete`, `restart`, `export`, etc.).

In general, everything you can do on the command line is doable using the Cloud Platform Console, and vice versa. In fact, both the Cloud Platform Console and the `gcloud` command invoke representational state transfer (REST) API actions. You can invoke the same REST APIs from your programs (the APIs are documented on the Google Cloud Platform website). Here is the REST API call to create an instance from Bash:⁵

```
ACCESS_TOKEN=$(gcloud auth application-default print-access-token)"  
curl --header "Authorization: Bearer ${ACCESS_TOKEN}" \  
--header 'Content-Type: application/json' \  
--data '{"name":"flights", "settings":  
    {"database-version":"POSTGRES_13", ...}}' \  
https://www.googleapis.com/sql/v1beta4/projects/[PROJECT-ID]/instances \  
-X POST
```

Alternatively, you can use the Cloud Client Library (available for a [variety of programming languages](#)) to issue the REST API calls. We saw this in [Chapter 2](#) when we used the `google.cloud.storage` Python package to interact with Cloud Storage.

In summary, there are three ways to interact with Google Cloud Platform:

1. The web console:

The `gcloud` command from the command line in Cloud Shell
or a machine that has the `gcloud` SDK installed

2. Directly invoke the REST API

3. Google Cloud Client Library (available as of this writing for Go, Java, Node.js, Python, Ruby, PHP, and C#)

In this book, I use primarily Option 2 (from the shell) and Option 4 (from Python programs).

Create Table of Data

In order to import data into a Postgres table, we first have to create an empty database and a table with the correct schema.

⁵ I looked up the URL and the format of the message from the [REST API reference documentation](#). All the other methods of interacting with a service (web console, `gcloud` command line, Cloud Client Library) end up invoking the REST API.

In the Cloud (web) console, navigate to the databases section of Cloud SQL, select instance “flights,” and create a new database called bts. This will be where we load our data.

Next, we have to create a file with the following syntax, to create a column for every field in the CSV file:

```
drop table if exists flights;

CREATE TABLE flights (
    "Year" TEXT,
    "Quarter" TEXT,
    "Month" TEXT,
    "DayofMonth" TEXT,
    "DayOfWeek" TEXT,
    "FlightDate" DATE,
    "Reporting_Airline" TEXT,
    "DOT_ID_Reported_Airline" TEXT,
    "IATA_CODE_Reported_Airline" TEXT,
    ...
)
```

For your convenience, the file I created is already in the Git repository, so just go to Cloud Shell, change into the `03_sqlstudio` directory, and do the following steps:

- Stage the file into Google Cloud Storage (changing the bucket to one that you own):

```
gsutil cp create_table.sql \
    gs://cloud-training-demos-ml/flights/ch3/create_table.sql
```

- In the web console, navigate to the flights instance of Cloud SQL and select IMPORT. In the form, specify the location of `create_table.sql` and specify that you want to create a table in the database `bts` (see [Figure 3-3](#)).

[Import data from Cloud Storage](#)

Source

Choose a file to import from. Make sure you have read access first. [Learn more](#)

bucket-name/file-name *
 ai-analytics-solutions-dsogcp/flights/ch3/create_table.sql [BROWSE](#)

Browse for a Cloud Storage file or enter the path to one (bucket/folder/file)

File format

SQL
A plain text file with a sequence of SQL commands, like the output of pg_dump

CSV
If your Cloud Storage file is a CSV file, select CSV. The CSV file should be a plain text file with one line per row and comma-separated fields.

Destination

Choose the database in this instance for your file to import into. [Learn more](#)

Database *
bts

[▼ SHOW USER OPTIONS](#)

When you import, a Cloud SQL service account will be granted read access to the selected file and bucket, which will be reflected in your permissions.

[IMPORT](#) [CANCEL](#)

Figure 3-3. Creating an empty table.

A few seconds later, the empty table will be created.

We can now load the CSV files into this table. Start by loading the January data by browsing to `201501.csv` in your bucket and specifying CSV as the format, `bts` as the database, and `flights` as the table (see [Figure 3-4](#)).

[← Import data from Cloud Storage](#)

Source

Choose a file to import from. Make sure you have read access first. [Learn more](#)

bucket-name/file-name *
 ai-analytics-solutions-dsongcp/flights/raw/201501.csv [BROWSE](#)

Browse for a Cloud Storage file or enter the path to one (bucket/folder/file)

File format

SQL
A plain text file with a sequence of SQL commands, like the output of pg_dump

CSV
If your Cloud Storage file is a CSV file, select CSV. The CSV file should be a plain text file with one line per row and comma-separated fields.

Destination

Choose the database and table in your instance for this file to import into. [Learn more](#)

Database *
bts

Table *
flights

Enter the name of an existing table in the database to house your CSV file

[▼ SHOW USER OPTIONS](#)

i When you import, a Cloud SQL service account will be granted read access to the selected file and bucket, which will be reflected in your permissions.

[IMPORT](#) [CANCEL](#)

Figure 3-4. Populating the table with data from January.



Note that the user interface doesn't provide a way to skip the first line, so the header will also get loaded as a row in the table. Fortunately, our schema calls all the fields as text, so this doesn't pose a problem—after loading the data, we can delete the row corresponding to the header. If we have a more realistic schema, we will have to remove the header line before loading the file.

Interacting with the Database

We can connect to the Cloud SQL instance from Cloud Shell using:⁶

```
gcloud sql connect flights --user=postgres
```

In the prompt that comes up, we connect to the `bts` database:

```
\c bts;
```

Then, we can run a query to obtain the five busiest airports:

```
SELECT "Origin", COUNT(*) AS num_flights
FROM flights GROUP BY "Origin"
ORDER BY num_flights DESC
LIMIT 5;
```

While this is performant because the dataset is relatively small (only January!), as I added more months, the database started to slow down.

Relational databases are particularly well suited to smallish datasets on which we wish to do ad hoc queries that involve searching and that return small subsets of data. For larger datasets, we can tune the performance of a relational database by indexing the columns of interest. Further, because relational databases typically support transactions and guarantee strong consistency, they are an excellent choice for data that will be updated often.

However, a relational database is a poor choice if your data is primarily read-only, if your dataset sizes go into the terabyte range, if you have a need to scan the full table (such as to compute the maximum value of a column), or if your data streams in at high rates. This describes our flight delay use case. So, let's switch from a relational database to an analytics data warehouse—BigQuery. The analytics data warehouse will allow us to use SQL and is much more capable of dealing with large datasets and ad hoc queries (i.e., doesn't need the columns to be indexed).

⁶ If your organization has set up a security policy to allow access only from authorized networks, you might have to [use a SQL proxy](#) to connect to the instance. At the time of writing, this is available only in the beta version, so use the following command: `gcloud beta sql connect flights --user=postgres`.

If you are following along with me, delete the Cloud SQL instance. We won't need it any further in this book.⁷

Querying Using BigQuery

In [Chapter 2](#), we loaded the CSV data into BigQuery into a table named `flights_raw` in the dataset `dsongcp`. Let's explore that dataset a bit—this is not a full exploratory data analysis, which I will do in [Chapter 5](#).

My goal here is to do “just enough” analysis on the data and then quickly pivot to building my first model. Once I have the model, I will be able to build a dashboard to explain that model. The idea is to get a first iteration out in front of users as quickly as possible. Going from ingested data to minimum viable outputs (model, dashboard, etc.) quickly is what agile development in data science looks like.

Teams that wait until they build a fully capable model before incorporating it into decision tools often build the wrong model (i.e., they solve the wrong problem because of misunderstanding how the decision will be used) or choose unviable technology (that is hard to get into production). Avoid these traps by testing your work with real users as quickly as possible!

Schema Exploration

Navigate to the [BigQuery section](#) of the Google Cloud (web) console and select the `flights_raw` table. On the right side of the window, select Schema (see [Figure 3-5](#)). Which fields do you think are relevant to predicting flight arrival delays?

⁷ Considering that we were using BigQuery and we are going to use BigQuery, what was the point of this detour into Cloud SQL? In many cases, your data will originally be in a relational database management system (RDBMS). If that dataset is small, you can power the dashboard off the RDBMS. Even though BigQuery is the better choice for the flights delay dataset, it won't always be the best choice for all the projects you will work on in Google Cloud. Part of my philosophy in this book is to “show,” not “tell.” So, the reason for this section was to show this trade-off between RDBMS and a data warehouse. This will be true later on in the book as well. I'll do a Spark ML model in [Chapter 7](#), but it won't scale once we add categorical columns—in that sense, that entire chapter is a digression! I'll stream into Bigtable in [Chapter 11](#) and show that the resulting performance improvement is overkill. In summary, the point of the Cloud SQL detour was to explore the possibilities, show what the problems are, and help you develop the intuition for the trade-offs involved.

The screenshot shows the BigQuery web interface with the following details:

- Table Name:** flights_raw
- Table Type:** Partitioned table (indicated by a note: "This is a partitioned table. [Learn more](#)")
- Schema Tab:** Selected tab.
- Details Tab:** Unselected tab.
- Preview Tab:** Unselected tab.
- Table Explorer Tab:** Unselected tab.
- Table Schema:**
 - Filter bar: "Filter Enter property name or value"
 - Table header: "Field name", "Type", "Mode", "Policy Tags ?"
 - List of columns:

Field name	Type	Mode	Policy Tags
Year	STRING	NULLABLE	
Quarter	STRING	NULLABLE	
Month	STRING	NULLABLE	
DayofMonth	STRING	NULLABLE	
DayOfWeek	STRING	NULLABLE	
FlightDate	DATE	NULLABLE	
Reporting_Airline	STRING	NULLABLE	
DOT_ID_Reported_Airline	STRING	NULLABLE	
IATA_CODE_Reported_Airline	STRING	NULLABLE	
Tail_Number	STRING	NULLABLE	

Figure 3-5. The schema of the `flights_raw` table that we loaded into BigQuery in Chapter 2.

Just looking at the schema is not enough. For example, do we really need the Year, Month, DayofMonth, and so on? Isn't FlightDate enough? It's best to not have duplicative data—the more columns we have, the more work we have to do to keep analysis consistent.

Similarly, which of the various Airline columns do we need? For the Airline columns, we did read the description on the BTS website in [Chapter 2](#), and will probably follow their recommendation that `Reporting_Airline` be the one that we use. Still, it's worth verifying why that is.

To make decisions like this, we can use two features—the Preview tab and the Table Explorer tab (see [Figure 3-6](#)).

Using Preview

The best way to quickly look at a BigQuery table is to use the Preview functionality. The Preview is free, whereas doing a `SELECT * FROM ... LIMIT 10` will incur a querying cost.

Looking at the preview (see [Figure 3-6](#)), the Year, Month, etc., columns do seem to be redundant. (If you are following along with me, you may see different rows because the Preview just picks whatever is most handy.)

SCHEMA		DETAILS		PREVIEW		TABLE EXPLORER				
Row	Year	Quarter	Month	DayofMonth	DayOfWeek	FlightDate	Reporting_Airline	DOT_ID_Reported_Airline	IATA_CODE_Reported_Airline	
1	2015	1	1	4	7	2015-01-04	DL	19790	DL	
2	2015	1	1	5	1	2015-01-05	DL	19790	DL	
3	2015	1	1	6	2	2015-01-06	DL	19790	DL	
4	2015	1	1	7	3	2015-01-07	DL	19790	DL	
5	2015	1	1	8	4	2015-01-08	DL	19790	DL	

Figure 3-6. Preview of the flights_raw table that we loaded into BigQuery in Chapter 2.

Let's check whether we can resurrect the FlightDate from the other columns and extract the date pieces from the FlightDate. We can do that with SQL:

```
SELECT
  FORMAT("%s-%02d-%02d",
    Year,
    CAST(Month AS INT64),
    CAST(DayofMonth AS INT64)) AS resurrect,
  FlightDate,
  CAST(EXTRACT(YEAR FROM FlightDate) AS INT64) AS ex_year,
  CAST(EXTRACT(MONTH FROM FlightDate) AS INT64) AS ex_month,
  CAST(EXTRACT(DAY FROM FlightDate) AS INT64) AS ex_day,
FROM dsongcp.flights_raw
LIMIT 5
```

The result appears to bear this out:

Row	resurrect	FlightDate	ex_year	ex_month	ex_day
1	2015-02-19	2015-02-19	2015	2	19
2	2015-02-20	2015-02-20	2015	2	20
3	2015-02-22	2015-02-22	2015	2	22
4	2015-02-23	2015-02-23	2015	2	23
5	2015-02-25	2015-02-25	2015	2	25

But we have to be sure. Let's print out rows where the extracted data from FlightDate is *not* identical:

```
WITH data AS (
SELECT
  FORMAT("%s-%02d-%02d",
    Year,
    CAST(Month AS INT64),
```

```

        CAST(DayofMonth AS INT64)) AS resurrect,
FlightDate,
CAST(EXTRACT(YEAR FROM FlightDate) AS INT64) AS ex_year,
CAST(EXTRACT(MONTH FROM FlightDate) AS INT64) AS ex_month,
CAST(EXTRACT(DAY FROM FlightDate) AS INT64) AS ex_day,
FROM dsongcp.flights_raw
)
SELECT * FROM data
WHERE resurrect != CAST(FlightDate AS STRING)

```

This query returns an empty result set, so we are sure that we can safely keep only the FlightDate column.

Using Table Explorer

How about the Airline code? Switch to the Table Explorer tab and select the three airline columns as shown in [Figure 3-7](#).

	Year	STRING
	Quarter	STRING
	Month	STRING
	DayofMonth	STRING
	DayOfWeek	STRING
	FlightDate	DATE
<input checked="" type="checkbox"/>	Reporting_Airline	STRING
<input checked="" type="checkbox"/>	DOT_ID_Reported_Airline	STRING
<input checked="" type="checkbox"/>	IATA_CODE_Reported_Airline	STRING

Figure 3-7. Selecting fields for Table Explorer.

BigQuery analyzes the full dataset and shows the unique values in the table, as shown in [Figure 3-8](#).

Reporting_Airline	Value	Count
AA	725984	
OO	588353	
EV	571977	
UA	515723	
MQ	294632	
B6	267048	
US	198715	
AS	172521	

DOT_ID_Reported_Airline	Value	Count
19393	1261...	
19790	875881	
19805	725984	
20304	588353	
20366	571977	
19977	515723	
20398	294632	
20409	267048	

IATA_CODE_Reported_Airline	Value	Count
AA	725984	
OO	588353	
EV	571977	
UA	515723	
MQ	294632	
B6	267048	
US	198715	
AS	172521	

Figure 3-8. Distinct values for the three Airline fields.

It is clear from the Table Explorer that we want to use either the `Reporting_Airline` or the `IATA_CODE_Reported_Airline`. As before, checking to see if there are rows where these are different indicates that `Reporting_Airline` is sufficient.

Creating BigQuery View

Based on such analysis on the remaining fields, I came up with the following sets of operations I want to do to the raw data to make it more usable. For example, the Departure Delay should be a floating point number and not a string. The Cancellation Code should be a boolean and not 1.00:

```
CREATE OR REPLACE VIEW dsongcp.flights AS

SELECT
    FlightDate AS FL_DATE,
    Reporting_Airline AS UNIQUE_CARRIER,
    OriginAirportSeqID AS ORIGIN_AIRPORT_SEQ_ID,
    Origin AS ORIGIN,
    DestAirportSeqID AS DEST_AIRPORT_SEQ_ID,
    Dest AS DEST,
    CRSDepTime AS CRS_DEP_TIME,
    DepTime AS DEP_TIME,
    CAST(DepDelay AS FLOAT64) AS DEP_DELAY,
    CAST(TaxiOut AS FLOAT64) AS TAXI_OUT,
    WheelsOff AS WHEELS_OFF,
    WheelsOn AS WHEELS_ON,
    CAST(TaxiIn AS FLOAT64) AS TAXI_IN,
    CRSArrTime AS CRS_ARR_TIME,
    ArrTime AS ARR_TIME,
    CAST(ArrDelay AS FLOAT64) AS ARR_DELAY,
    IF(Cancelled = '1.00', True, False) AS CANCELLED,
    IF(Diverted = '1.00', True, False) AS DIVERTED,
    DISTANCE
FROM dsongcp.flights_raw;
```

In order to avoid repeating these casts in all queries from here on out, I am creating a view that consists of the SELECT statement (see the first line in the preceding listing). A view is a virtual table—we can query the view just as if it were a table:

```
SELECT
    ORIGIN,
    COUNT(*) AS num_flights
FROM dsongcp.flights
GROUP BY ORIGIN
ORDER BY num_flights DESC
LIMIT 5
```

Any queries that happen on the view are rewritten by the database engine to happen on the original table—conceptually, a view works as if the SQL corresponding to the view was to be inserted into every query that uses the view.

What if the view includes a WHERE clause so that the number of rows is much less? In such cases, it would be far more efficient to export the results into a table and query that table instead:

```
CREATE OR REPLACE TABLE dsongcp.flights AS  
SELECT
```

But what if you export the results into a table and then the original table has a new month of data added to it? We'd have to rerun the table creation statement to make the extracted table up-to-date. In the case of a view, we wouldn't have to do anything special—all new queries would automatically be querying the entire raw table and thus include the new month of data.

Can we have our cake and eat it too? Can we get the “live” nature of a view, but the query efficiency of a table? Yes. It's called a *materialized view*:

```
CREATE MATERIALIZED VIEW dsongcp.flights AS  
SELECT
```

The view is materialized into a table and kept up-to-date by BigQuery. While views are free, materialized views carry an extra cost because of the extra storage and compute overhead they involve.

In this book, I'll use a regular view during development, since it's easy to come back and add new columns, etc. Later on, once we go to production, it's quite simple to change it over to a materialized view—none of the client code will need to change.

Building Our First Model

Intuitively, we feel that if the flight departure is delayed by 15 minutes, it will also tend to arrive 15 minutes late. So, our model could be that we cancel the meeting if the departure delay of the flight is 15 minutes or more. Of course, there is nothing here about the probability (recall that we wanted to cancel if the probability of an arrival delay of 15 minutes was greater than 30%). Still, it will be a quick start and give us something that we can ship now and iterate upon.

Contingency Table

Suppose that we need to know how often we will be making the right decision if our decision rule is the following:

If $\text{DEP_DELAY} \geq 15$, cancel the meeting; otherwise, go ahead.

There are four possibilities in the *contingency table* or the *confusion matrix*, which you can see in [Table 3-2](#).

Table 3-2. Confusion matrix for the decision to cancel the meeting

	Arrival delay < 15 minutes	Arrival delay \geq 15 minutes
We did not cancel meeting (departure delay < 15 min)	Correct (true positive)	False positive
We canceled meeting (departure delay \geq 15 min)	False negative	Correct (true negative)

If we cancel the meeting and it turns out that the flight arrived on time (let's call that a "positive"), it is clear that we made a wrong decision. It is arbitrary whether we refer to it as a false positive (treating the on-time arrival as a positive event) or as a false negative (treating the late arrival as a negative event). Because the dataset is called "on-time arrivals," let's term on-time arrival the positive event. How do we find out how often the decision rule of thresholding the departure delay at 15 minutes will tend to be correct? We can evaluate the first box in the confusion matrix using BigQuery:

```
SELECT
    COUNT(*) AS true_positives
FROM dsongcp.flights
WHERE dep_delay < 15 AND arr_delay < 15
```

There are 4,430,885 such flights.

To compute all four values in a single statement, move the WHERE clause into the SELECT itself:

```
SELECT
    COUNTIF(dep_delay < 15 AND arr_delay < 15) AS true_positives,
    COUNTIF(dep_delay < 15 AND arr_delay  $\geq$  15) AS false_positives,
    COUNTIF(dep_delay  $\geq$  15 AND arr_delay < 15) AS false_negatives,
    COUNTIF(dep_delay  $\geq$  15 AND arr_delay  $\geq$  15) AS true_negatives,
    COUNT(*) AS total
FROM dsongcp.flights
WHERE arr_delay IS NOT NULL AND dep_delay IS NOT NULL
```

Each of the COUNTIF statements counts the number of rows that match the given criterion, and COUNT(*) counts all rows. This way, we get to scan the table just once, and still manage to collect the four numbers that form the confusion matrix:

Row	true_positives	false_positives	false_negatives	true_negatives	total
1	4430885	232701	219684	830738	5714008

Recall that these numbers assume that we are making a decision by thresholding the departure delay at 15 minutes. But is that the best threshold?

Threshold Optimization

Ideally, we want to try out different values of the threshold and pick the one that provides the best results. To do so, we can declare a variable called THRESH and use it in the query. This way, there is just one number to change when we want to try out a different threshold:

```
DECLARE THRESH INT64;
SET THRESH = 15;

SELECT
    COUNTIF(dep_delay < THRESH AND arr_delay < 15) AS true_positives,
    COUNTIF(dep_delay < THRESH AND arr_delay >= 15) AS false_positives,
    COUNTIF(dep_delay >= THRESH AND arr_delay < 15) AS false_negatives,
    COUNTIF(dep_delay >= THRESH AND arr_delay >= 15) AS true_negatives,
    COUNT(*) AS total
FROM dsongcp.flights
WHERE arr_delay IS NOT NULL AND dep_delay IS NOT NULL
```

Still, I'd rather not run the query several times, once for each threshold. It's not about the drudgery of it—I could avoid the manual work by using a `for` loop in a script. What I'm objecting to is scanning the table four times. The better way to do this in SQL is to declare an array of possible thresholds and then group by them:

```
SELECT
    THRESH,
    COUNTIF(dep_delay < THRESH AND arr_delay < 15) AS true_positives,
    COUNTIF(dep_delay < THRESH AND arr_delay >= 15) AS false_positives,
    COUNTIF(dep_delay >= THRESH AND arr_delay < 15) AS false_negatives,
    COUNTIF(dep_delay >= THRESH AND arr_delay >= 15) AS true_negatives,
    COUNT(*) AS total
FROM dsongcp.flights, UNNEST([5, 10, 11, 12, 13, 15, 20]) AS THRESH
WHERE arr_delay IS NOT NULL AND dep_delay IS NOT NULL
GROUP BY THRESH
```

This way, we get to run a single query, which scans the table just once and still manages to create contingency tables for all the thresholds we want to try. The result consists of the four contingency table values for each of the seven values of the threshold:

Row	THRESH	true_positives	false_positives	false_negatives	true_negatives	total
1	5	3931979	144669	718590	918770	5714008
2	10	4242286	184944	408283	878495	5714008
3	11	4288279	193912	362290	869527	5714008
4	12	4329146	203068	321423	860371	5714008
5	13	4366641	212498	283928	850941	5714008
6	15	4430885	232701	219684	830738	5714008
7	20	4542475	291791	108094	771648	5714008



Learn SQL. You'll thank me later.

This is all well and good, but recall that our goal (see [Chapter 1](#)) is to cancel the client meeting if the probability of arriving 15 minutes late is 30% or more. How close do we get with each of these thresholds?

To know this, we need to compute the fraction of times a decision is wrong. We can do this by calling the preceding result the contingency table, and then computing the necessary ratios:

```
WITH contingency_table AS (
    SELECT
        THRESH,
        COUNTIF(dep_delay < THRESH AND arr_delay < 15) AS true_positives,
        COUNTIF(dep_delay < THRESH AND arr_delay >= 15) AS false_positives,
        COUNTIF(dep_delay >= THRESH AND arr_delay < 15) AS false_negatives,
        COUNTIF(dep_delay >= THRESH AND arr_delay >= 15) AS true_negatives,
        COUNT(*) AS total
    FROM dsongcp.flights, UNNEST([5, 10, 11, 12, 13, 15, 20]) AS THRESH
    WHERE arr_delay IS NOT NULL AND dep_delay IS NOT NULL
    GROUP BY THRESH
)

SELECT
    ROUND((true_positives + true_negatives)/total, 2) AS accuracy,
    ROUND(false_positives/(true_positives+false_positives), 2) AS fpr,
    ROUND(false_negatives/(false_negatives+true_negatives), 2) AS fnr,
    *
FROM contingency_table ORDER BY accuracy ASC
```

The result now includes the accuracy, false positive rate, and false negative rate:

Row	accuracy	fpr	fnr	THRESH	true_positives	false_positives	false_negatives	true_negatives	total
1	0.85	0.04	0.44	5	3931979	144669	718590	918770	5714008
2	0.9	0.04	0.32	10	4242286	184944	408283	878495	5714008
3	0.9	0.04	0.29	11	4288279	193912	362290	869527	5714008
4	0.91	0.04	0.27	12	4329146	203068	321423	860371	5714008
5	0.91	0.05	0.25	13	4366641	212498	283928	850941	5714008
6	0.92	0.05	0.21	15	4430885	232701	219684	830738	5714008
7	0.93	0.06	0.12	20	4542475	291791	108094	771648	5714008

We want to cancel the meeting whenever we think the flight will be late. Our decision will not be perfect. There are times that the decision will be wrong. What is our tolerance for error? It's 30%. This means that:

- Flights should arrive on time (when we cancel the meeting) less than 30% of the time. So, we want the false positive rate to be 0.3 or less.
- Flights should arrive late (when we go ahead with the meeting) less than 30% of the time. So, we want the false negative rate 0.3 or less.

Looking at the preceding contingency table, which of these criteria looks like it could be a problem? That's right—the false negative rate. The false positive rate, at 0.04 or so, is comfortably within our error tolerance.

If we are going to make the decision based on the departure delay, our choice of departure delay threshold will have to be such that the false negative rate is 0.3.

It is clear from the preceding table that if we want our decision to have a false negative rate of 30%, the departure delay threshold needs to be 10 or 11 minutes (in the dataset, departure delay is an integer, so an intermediate threshold like 10.6 minutes does not make sense). We could choose 11 minutes on the grounds that, at 11 minutes, the FNR is less than 0.3. Or we could choose 10 minutes on the grounds that it's a nice, round number and our model is not so sophisticated that we can make decisions at a 1 minute precision.⁸

If we choose a threshold of 10 minutes, we will make the correct decision 96% of the time when we don't cancel the meeting and 68% of the time when we cancel the meeting. Overall, we are correct 90% of the time.

Note that 10 minutes is not the threshold that maximizes the overall accuracy. Had we chosen a threshold of 20 minutes, we'd cancel far fewer meetings (108,000 versus 408,000) and be correct more often overall (93%). However, that would be very conservative. Since it is not our goal to be correct 88% of the times we cancel the meeting—we only want to be correct 70% of the time—10 minutes is the right threshold.

However, we could also consider that if we can increase the threshold to 20 minutes, we would be correct far more often with very little impact on the false negative rate. Until we looked at the data, we didn't know what was achievable, and it is possible that the original target was set in a fog of uncertainty. It might be worthwhile asking

⁸ If this argument surprises you, imagine that we are talking about a threshold of 0.1 versus a threshold of 0.11. It's the same principle—don't use thresholds that have a misleading number of significant digits. This is important because I'm going to show the threshold to end users in a dashboard. If I were not showing the threshold, then I could use arbitrary precision. But when showing numbers to end users, keep the number of significant digits in mind.

our stakeholders whether they are really wedded to the 30% false negative rate and whether we have leeway to change the trade-offs available to users of our application—a dashboard that shows the impact of a threshold is an excellent way to gauge this. If the stakeholders don’t know, it might be worth doing an A/B test with a focus group, and that’s what we are about to do next.

Is This Machine Learning?

What we did here—trying different thresholds—is at the heart of machine learning. Our model is a simple rule that has a single parameter (the departure delay threshold), and we can tune it to maximize some objective measure (here, the desired precision). We can (and will) use more complex models, and we can definitely make our search through the parameter space a lot more systematic, but this process of devising a model and tuning it is the gist of machine learning. We haven’t evaluated the model (we can’t take the 70% we got on the 12 months of 2015 data and claim that to be our model performance—we need an independent dataset to do that), and that is a key step in machine learning. However, we can plausibly claim that we have built a simple machine learning model to provide guidance on whether to cancel a meeting based on historical flight data.

Building a Dashboard

Even this simple model is enough for us to begin getting feedback from end users. Recall that my gut instinct at the beginning of the previous section was that I needed to use a 15-minute threshold on the departure delay. Analysis of the contingency table, however, indicated that the right threshold to use was 10 minutes. I’m satisfied with this model as a first step, but will our end users be? Let’s go about building a dashboard that explains the model recommendations to end users. Doing so will also help clarify what I mean by explaining a model to end users.

There are a large number of business intelligence and visualization tools available, and many of them connect with data sources like BigQuery and Cloud SQL on Google Cloud Platform. In this chapter, we build dashboards using Data Studio, which is free and comes as part of Google Cloud Platform, but you should be able to do similar things with Tableau, QlikView, Looker, and so on.

Looker or Data Studio?

Google Cloud has two business intelligence tools—Looker and Data Studio. Data Studio is free and much more suitable for self-service use. Looker is much more capable and more suitable for enterprise use.

What do I mean by enterprise use? Here are a few examples of things that Looker can do that Data Studio can't:

Consistency

It is possible for one team to define a semantic layer consisting of standard nomenclature for columns and ways of computing key performance metrics. The rest of the organization then builds dashboards starting from the semantic layer rather than from the raw data.

Multicloud

Looker can access data in BigQuery, Amazon Redshift, Azure SQL Data Warehouse, and Snowflake *simultaneously* in the same report.

Embedded analytics

Have you been to a website where you can see charts and graphs of your activity? This is provided by a lot of B2B applications, such as marketplaces allowing sellers to visualize their own data. Looker allows you to embed analytics in another website.

Alerts and updates

Data Studio requires the user to visit the Data Studio web page and refresh the graphics. With Looker, you can push reports on a schedule or whenever an event happens.

That said, in our case, all we want is a self-serve dashboard, and Data Studio fits the bill perfectly.

For those of you with a data science background, I'd like to set expectations here—a dashboard is a way for end users to quickly come to terms with the current state of a system and is not a full-blown, completely customizable, statistical visualization package. Think about the difference between what's rendered in the dashboard of a car versus what would be rendered in an engineering visualization of the aerodynamics of the car in a wind tunnel—that's the difference between what we will do in Data Studio versus what we will use Vertex AI Notebooks for in later chapters. Here, the emphasis is on providing information effectively to end users—thus, the key aspects are interactivity and collaboration. With Data Studio, you can share reports similarly to Google Workspace documents; that is, you can give different colleagues viewing or editing rights, and colleagues with whom you have shared a visualization can refresh the charts to view the most current data.

Getting Started with Data Studio

To work with Data Studio, navigate to [the Data Studio home page](#). There are two key concepts in Data Studio: reports and data sources. A report is a set of charts and commentary that you create and share with others. The charts in the report are built from

data that is retrieved from a data source. The first step, therefore, is to set up a data source. Because our data is in BigQuery, the data source we need to set up is for Data Studio to connect to BigQuery.

On the Data Studio home page, click on the Create button, click the Data Source menu item, and choose the BigQuery button, as illustrated in [Figure 3-9](#).⁹

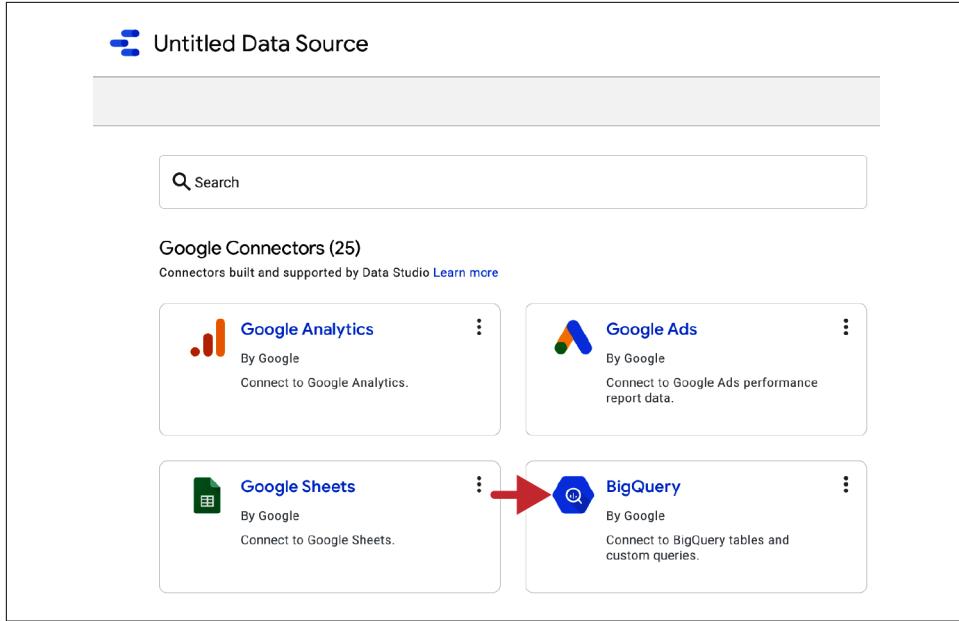


Figure 3-9. Choose BigQuery from the Data Source menu item in Data Studio.

Select your project, the `dsongcp` dataset, and `flights` as the table. Then, click on the Connect button. Recall that `flights` is the view that we have set up with the streamlined set of fields.

A list of fields in the table displays, with Data Studio inferring something about the fields based on a sampling of the data in that table. We'll come back and correct some of these, but for now, just click Create Report, accepting all the prompts.

⁹ Graphical user interfaces are often the fastest-changing parts of any software. So, if the user interface has changed from these screenshots by the time this book gets into your hands, please hunt around a bit. There will be some way to add a new data source.

Creating Charts

On the top ribbon, select the scatter chart icon from the “Add a chart” pulldown (see [Figure 3-10](#)) and draw a rectangle somewhere in the main window; Data Studio will draw a chart. The data that is rendered is pulled from some rather arbitrary columns.

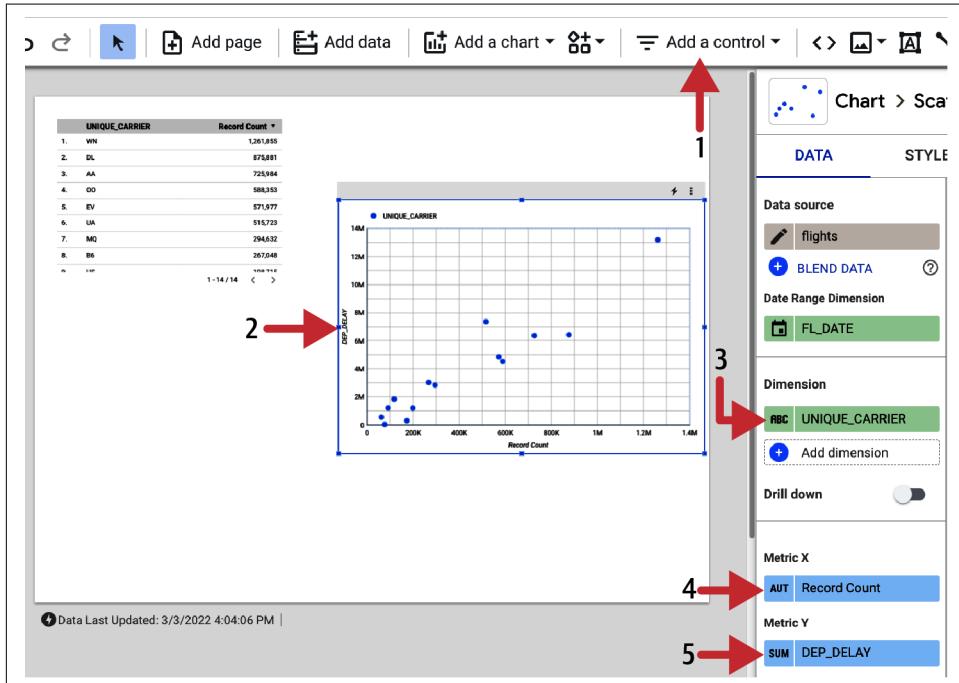


Figure 3-10. Initial chart rendered by Data Studio.

Ignoring the Date Range Dimension for now, there are three columns being used: the Dimension is the quantity being plotted; Metric X is along the x-axis; and Metric Y is along the y-axis. Change (if necessary) Dimension to UNIQUE_CARRIER, Metric X to DEP_DELAY, Metric Y to ARR_DELAY, and change the aggregation metric for both Metric X and Metric Y to Average. Ostensibly, this should give us the average departure delay and arrival delay of different carriers. Click the Style tab and add in a linear trendline and show the data labels. [Figure 3-11](#) depicts the resulting chart.

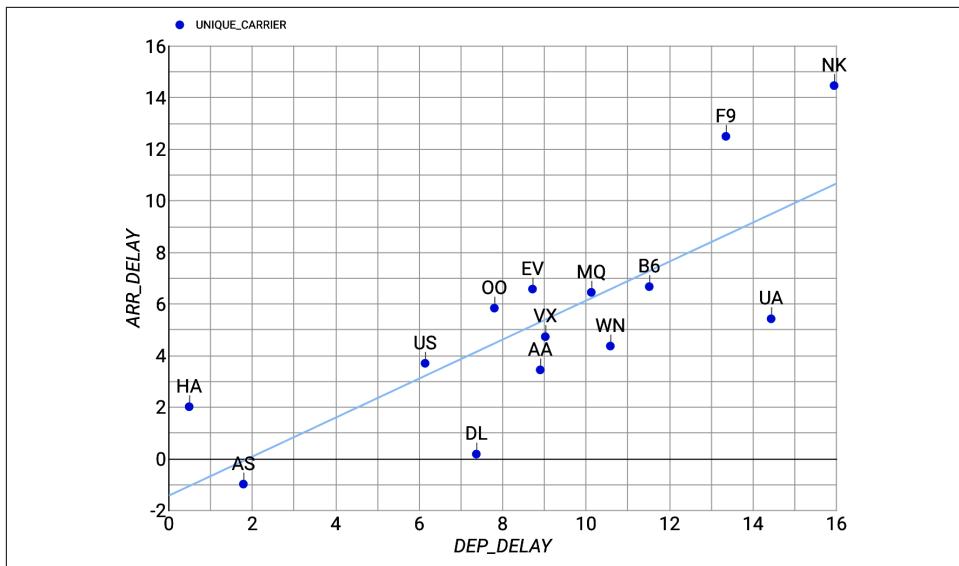


Figure 3-11. Chart after changing Metric, Dimension, and Style.

Adding End-User Controls

So far, our chart is static—there is nothing for end users to interact with. They get to see a pretty picture, but do not get to change anything about our graph. To permit the end user to change something about the graph, we should add *controls* to our graph.

Let's give our end users the ability to set a date range. On the top icon ribbon, click the “Date range control” button, as illustrated in Figure 3-12.

On your chart, place the rectangle where you'd like the control to appear. Change the time window to be Fixed and set the Start Date to Jan. 1, 2015, and end date to Dec. 31, 2019.¹⁰ This is how the report will initially appear to users.

In the upper-right corner, change the toggle to switch to the View mode. This is the mode in which users interact with your report. Change the data range to Jan 1, 2015 to May 31, 2015 (see Figure 3-13) and you should see the chart immediately update.

¹⁰ Or whatever the last month that you downloaded is. Just so that you don't have to wait a long time for the data to be available in your Google Cloud project, the *ingest.sh* script in Chapter 2, by default, downloads only 2015 data. Change the YEAR loop in that script to download 2015 to 2019.

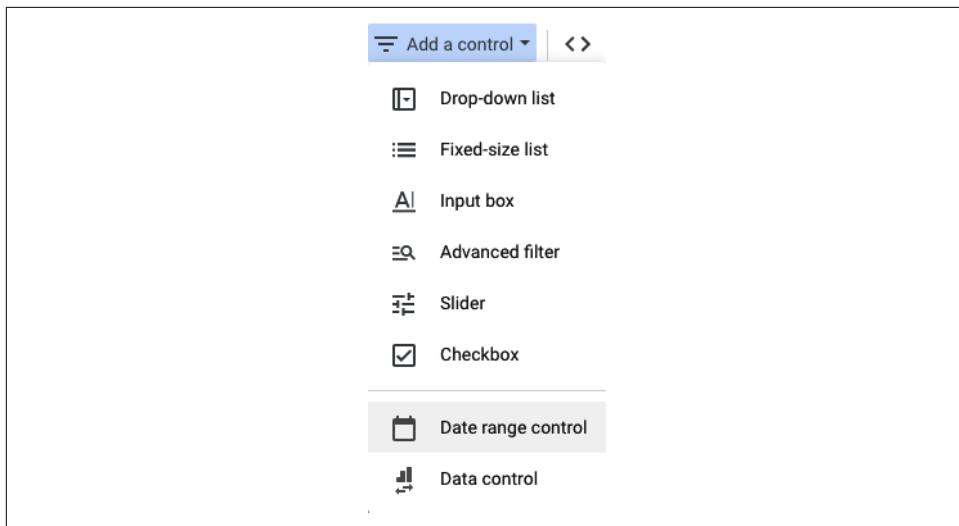


Figure 3-12. The Date range control on the top icon ribbon.

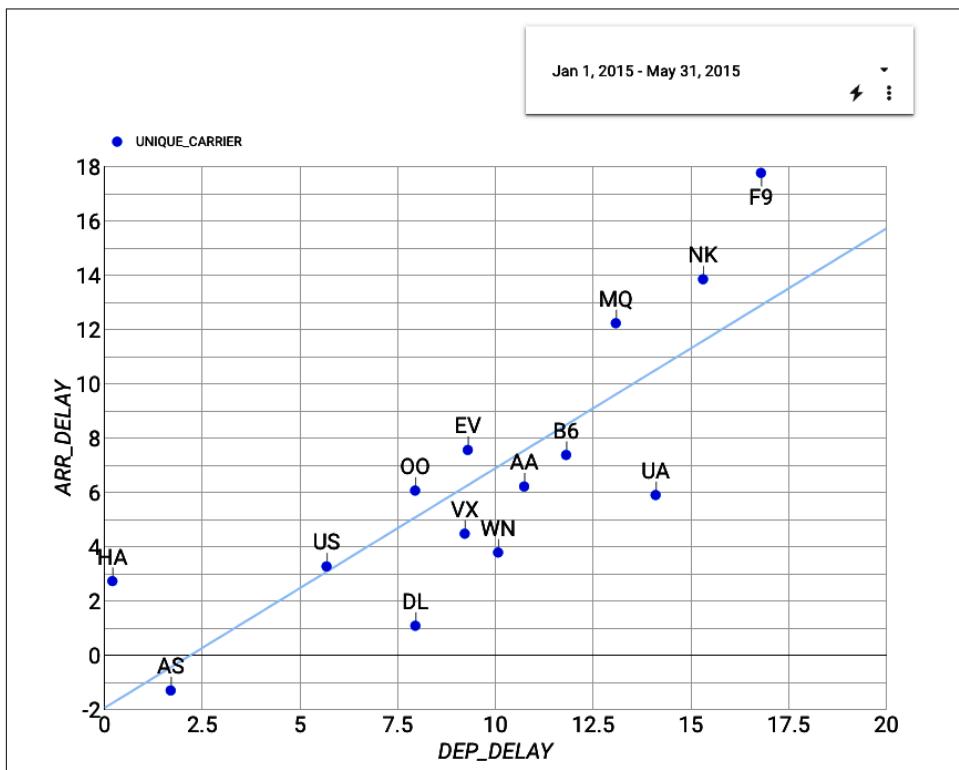


Figure 3-13. The chart in View mode.

Pause a bit here and ask yourself what kind of a model the chart in [Figure 3-13](#) would explain. Because there is a line, it strongly hints at a linear model. If we were to recommend meeting cancellations based on this chart, we'd be suggesting, based on the linear trend of arrival delay with departure delay, that departure delays of more than 20 minutes lead to arrival delays of more than 15 minutes. That, of course, was not our model—we did not do linear regression, and certainly not airline by airline. Instead, we picked a departure threshold based on a contingency table over the entire dataset. So, we should not use the preceding graph in our dashboard—it would be a misleading description of our actual model.

Showing Proportions with a Pie Chart

How would you explain our contingency table-based thresholds to end users in a dashboard? Recall that the choice comes down to the proportion of flights that arrive more than 15 minutes after their scheduled time. That is what our dashboard needs to show.

One of the best ways to show a proportion is to use a pie chart.¹¹ Switch back to the Edit mode, and from the pull-down menu, select the Donut Chart button (this is a type of pie chart), and then, on your report, draw a square where you'd like the donut chart to appear (it is probably best to delete the earlier scatter plot from it). As we did earlier, we need to edit the dimensions and metrics to fit what it is that we want to display. Perhaps things will be clearer if you see what the end product ought to look like. [Figure 3-14](#) gives you a glimpse.

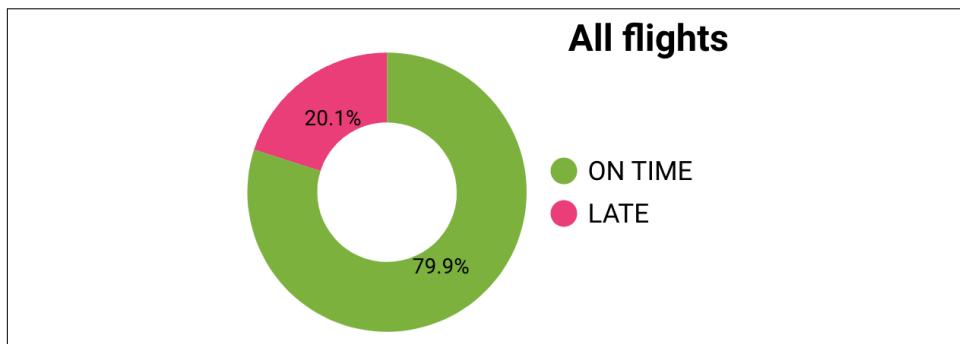


Figure 3-14. Desired end result is a chart that shows the proportion of flights that are late versus on time.

¹¹ An alternative way to show proportions, especially of a time-varying whole, is a stacked column chart.

In this chart, we are displaying the proportion of flights that arrived late versus those that arrived on time. The labeled field ON TIME versus LATE is the Dimension. The number of flights is the metric that will be apportioned between the labels. So, how do you get these values from the BigQuery view?

It is clear that there is no column in the database that indicates the total number of flights. However, Data Studio has a special value Record Count that we can use as the metric, after making sure to change the aggregate from the default Sum to Count.

The “islate” value, though, will have to be computed as a formula. Conceptually, we need to **add a new calculated field** to the data that looks like this:

```
CASE WHEN  
(ARR_DELAY < 15)  
THEN  
"ON TIME"  
ELSE  
"LATE"  
END
```

Click on the current Dimension column and click on “Create Field.” Give the field the name **is_late**, enter the preceding formula, and change the type to Text (see **Figure 3-15**).

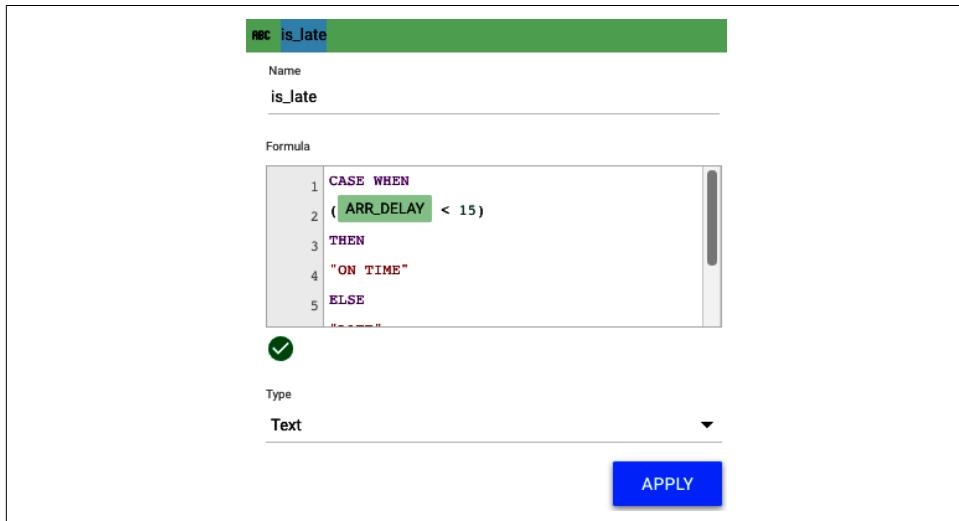


Figure 3-15. How to set up the *is_late* definition.

The pie chart is now complete and reflects the proportion of flights that are late versus those that are on time. You can switch over to the Style tab if you'd like to change the appearance of the pie chart to be similar to mine.

Because the proportion of flights that end up being delayed is the quantity on which we are trying to make decisions, the pie chart translates quite directly to our use case. However, it doesn't tell the user what the typical delay would be. To do that, let's create a bar (column) chart that looks like the one shown in [Figure 3-16](#).

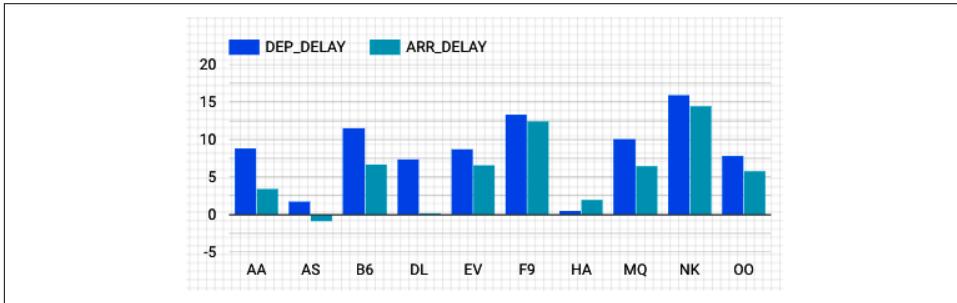


Figure 3-16. Typical delay for each carrier.

Here, the labeled quantity (or Dimension) is the `Carrier`. There are two metrics being displayed: the `DEP_DELAY` and `ARR_DELAY`, both of which are aggregated to their averages over the dataset. [Figure 3-17](#) shows the specifications.

Note the Sort column at the end—it is important to have a reliable sort order in dashboards so that users become accustomed to finding the information they want in a known place. Also, the default is to use different axes for the two variables.

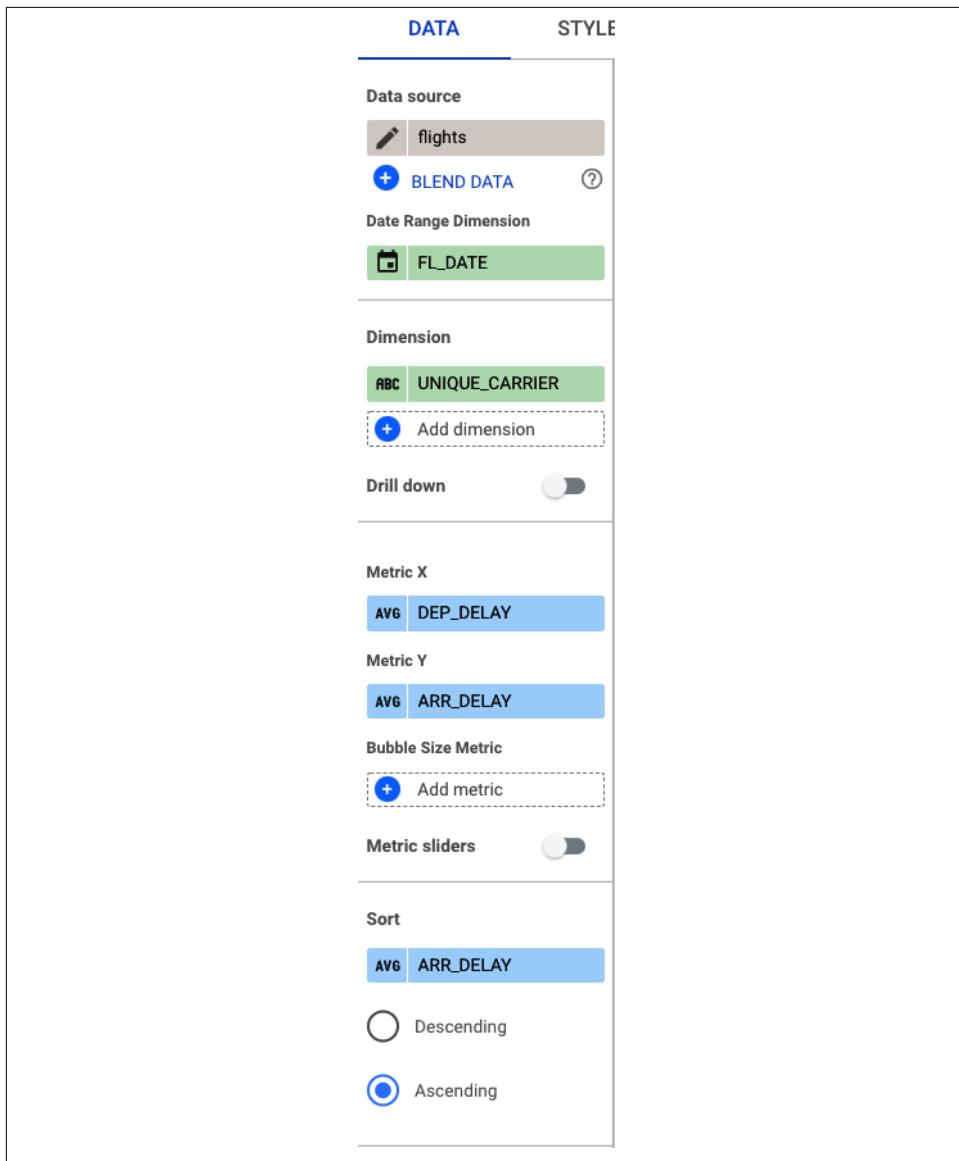


Figure 3-17. How to set the bar chart properties to generate the desired chart.

Switch over to the Style tab and change this to use a single axis. Finally, Data Studio defaults to 10 bars. In the Style tab, change this to reflect that we expect to have up to 20 unique carriers (Figure 3-18).

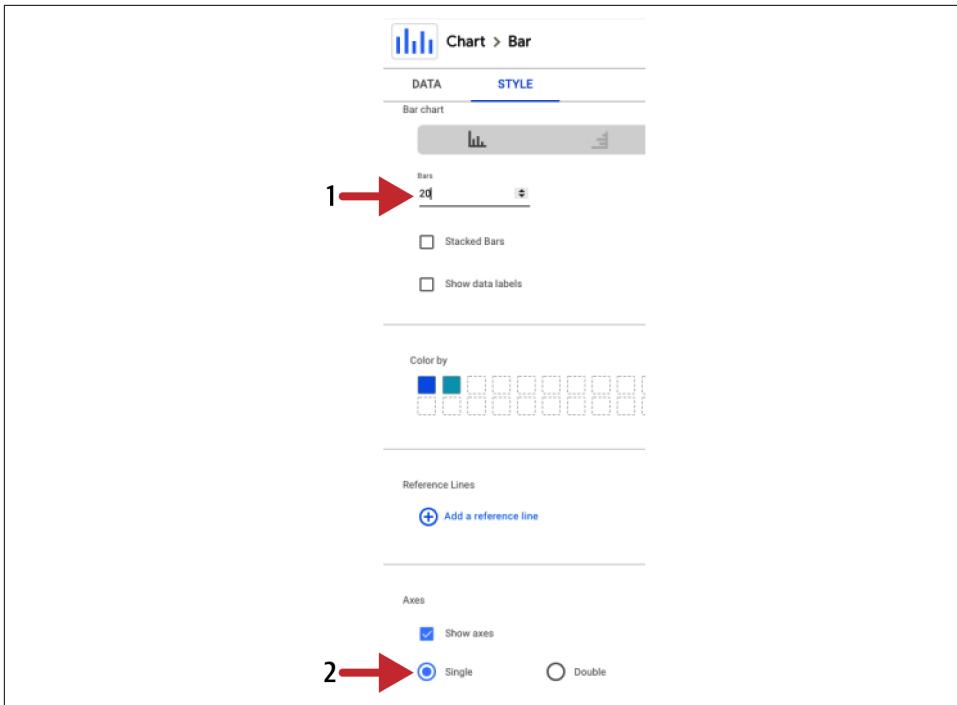


Figure 3-18. How to set the bar chart properties to generate the desired chart.

Of course, we can now add in a date control as we did earlier to end up with the report in [Figure 3-19](#) (“All flights” in the diagram is just a text label that I added).

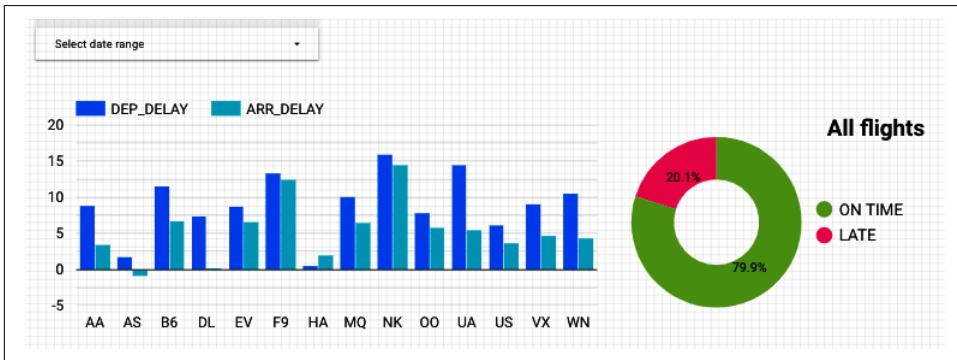


Figure 3-19. Resulting dashboard consisting of a pie chart and bar chart.

It appears that, on average, about 80% of flights are on time and that the typical arrival delay varies between airlines but lies in a range of 0 to 15 minutes.

Explaining a Contingency Table

Even though the dashboard we just created shows users the decision-making criterion (proportion of flights that will be late) and some characteristics of that decision (the typical arrival delay), it doesn't actually show our model. Recall that our model involved a threshold on the departure delay. We need to show that. [Figure 3-20](#) shows what we want the dashboard to look like.

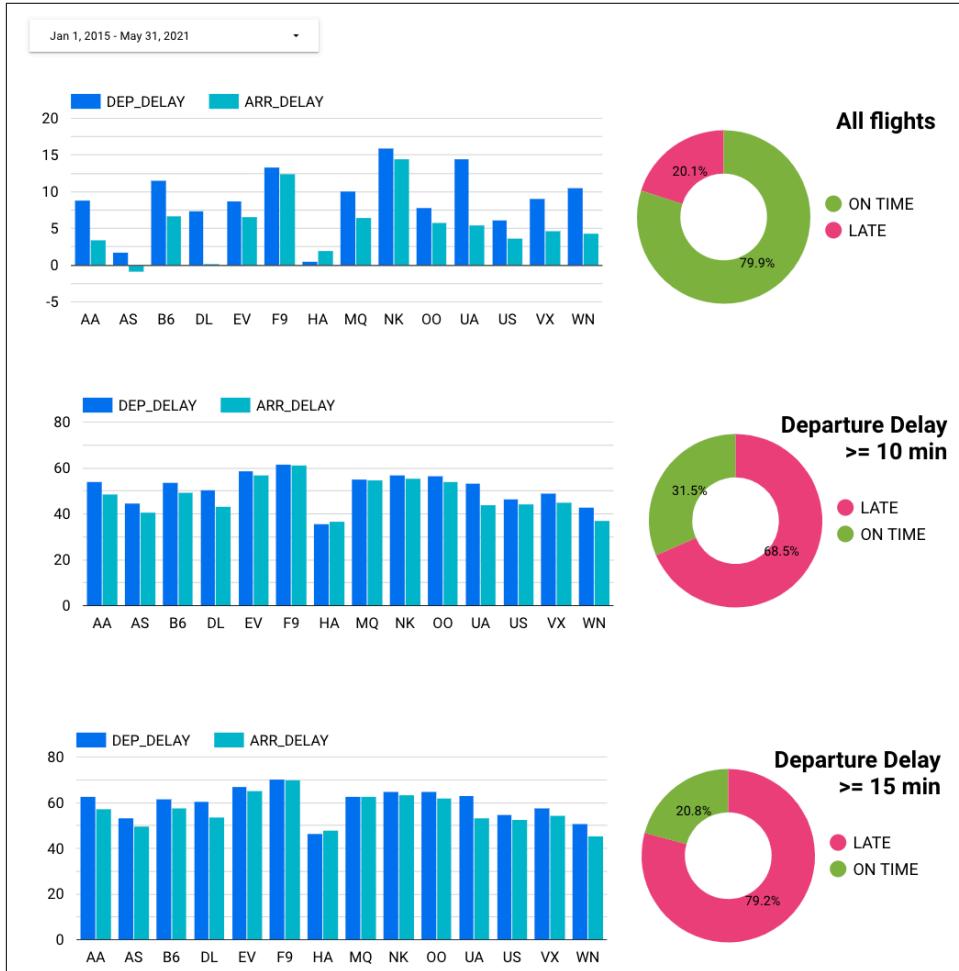


Figure 3-20. Dashboard consisting of three pairs of pie charts and bar charts along with a date control.

In other words, we want to show the same two charts, but for the decision thresholds that we considered—departure delays of 10, 15, and 20 minutes or more.

To get there, we need to change our data source. No longer can we populate the chart from the entire table. Instead, we should populate it from a query that pulls only those flights whose departure delay is greater than the relevant threshold. In BigQuery, we can create the views we need and use those views as data sources.¹² Here's how:

```
CREATE OR REPLACE VIEW dsongcp.delayed_10 AS
SELECT * FROM dsongcp.flights WHERE dep_delay >= 10;

CREATE OR REPLACE VIEW dsongcp.delayed_15 AS
SELECT * FROM dsongcp.flights WHERE dep_delay >= 15;

CREATE OR REPLACE VIEW dsongcp.delayed_20 AS
SELECT * FROM dsongcp.flights WHERE dep_delay >= 20;
```

Alternatively, we can add a filter on the departure delay and allow the end user to try out different thresholds. However, because we want to explain the decision model, it is better to ensure that the 10-minute threshold is explicitly present in the dashboard.

Looking at the resulting pie chart for a 10-minute threshold ([Figure 3-20](#)), we see that it comes quite close to our target of 30% on-time arrivals. The bar chart for the 10-minute delay explains why the threshold is important. Hint: it is not about the exact numeric value of 10 minutes. It is about what the 10-minute delay is indicative of. Can you decipher what is going on? Still stuck? Look at the y-axis of the three bar charts.

Although the typical departure delay of a flight is only about 5 minutes (see the chart corresponding to all flights that we created earlier), flights that are delayed by more than 10 minutes fall into a separate statistical regime. The typical departure delay of an aircraft that departs more than 10 minutes late is around 50 minutes! A likely explanation is that a flight that is delayed by 10 minutes or more typically has a serious issue that will not be resolved quickly. If you are sitting in an airplane and it is more than 10 minutes late in departing, you might as well cancel your meeting—you are going to be sitting at the gate for a while.¹³

At this point, we have created a very simple model and created dashboards to explain the model to our end users. Our end users have a visual, intuitive way to see how often our model is correct and how often it is wrong. The model might be quite simple, but the explanation of why the model works is a satisfying one.

¹² Data Studio does support a BigQuery query as a data source, but it is preferable to read from a view because views are more reusable.

¹³ Road warriors know this well. Ten minutes in, and they whip out their phones to try to get on a different flight.

There is one teeny, tiny thing missing, though. Context. The dashboard that we have built so far is all about historical data, whereas real dashboards need to be timely. Our dashboard shows aggregates of flights all over the country, but our users will probably care only about the airport from which they are departing and the airport to which they are going. We have a wonderfully informative dashboard, but without such time and location context, few users would care. In [Chapter 4](#), we look at how to build real-time, location-aware dashboards—unfortunately, however, there is a problem with our dataset that prevents us from doing so immediately. So, we'll spend the first part of [Chapter 4](#) doing some data wrangling.

Before we move on to [Chapter 4](#), let's expand our discussion of business intelligence beyond dashboards.

Modern Business Intelligence

Modern business intelligence (BI) is more than dashboards. BI is data science for business users. As such, the trends in BI are toward digitization (there is data on more and more things), democratization (more and more people can get insights from data), and integration (access to sophisticated insights from familiar tools).

Let's look at these trends.

Digitization

The [release notes of Data Studio](#) are a great way to stay up-to-date with product updates and new visualization options. All Google Cloud products have release notes that are published online—take a look at the [release notes for BigQuery](#), for example.

Once upon a time, release notes were nothing more than pages of text. But now, they are digitized and queryable. As befits a Data Cloud, the release notes for BigQuery (and all other Google Cloud products) are available as a public dataset in [BigQuery](#). Let's query it:

```
SELECT
    product_name,
    DATE_TRUNC(published_at, MONTH) AS month,
    COUNT(*) AS releases
FROM `bigquery-public-data.google_cloud_release_notes.release_notes`
GROUP BY product_name, month
```

The result looks like this:

Row	product_name	month	releases
1	BigQuery	2012-05-01	8
2	BigQuery	2012-07-01	5
3	BigQuery	2012-08-01	5

Row	product_name	month	releases
4	Dataflow	2017-10-01	6
5	Dataflow	2019-02-01	3

Natural Language Queries

We can explore the data with Data Studio, of course, but let's try using natural language queries. At the time of writing, this was still in alpha and one could only query one's own tables (not a public dataset). So let's export the preceding query to a table in our own project:

```
CREATE OR REPLACE TABLE dsongcp.monthly_releases AS

SELECT
    product_name,
    DATE_TRUNC(published_at, MONTH) AS month,
    COUNT(*) AS releases
FROM `bigquery-public-data.google_cloud_release_notes.release_notes`
GROUP BY product_name, month
```

Now, select the option to “Ask Question.” The UI immediately proposes a few questions (see Figure 3-21).

The screenshot shows the 'Natural language query' interface in 'ALPHA' mode. On the left, there's a dropdown menu labeled 'Select table *' containing 'dsongcp.monthly_releases'. Below it is a blue button labeled 'GENERATE EQUIVALENT SQL'. To the right, there's a text input field labeled 'Question *' with a blue outline. A list of generated questions is shown in a dropdown menu:

- Total releases for product name Anthos
- Top 3 month by sum of releases
- Distribution of month
- Unique count of product_name
- Most frequent month

Figure 3-21. Some of the questions automatically generated based on the dataset.

Let's pick a different question, though: “top three products by total releases.” The UI generates the following SQL query automatically (see Figure 3-22):

```
SELECT
    product_name AS product_name,
    (SUM(releases)) AS SUM_releases
FROM
    `ai-analytics-solutions.dsongcp.monthly_releases`
```

```
GROUP BY product_name  
ORDER BY SUM_releases DESC  
LIMIT 3;
```

Isn't that something? Notice that even though I said "top three products," the engine was smart enough to pick up the product_name column. The result, in case you are curious, is:

Row	product_name	SUM_releases
1	Google Kubernetes Engine	749
2	Dataproc	625
3	App Engine standard environment Java	472

Natural language query ALPHA

Type a question to generate your SQL query (on a single table or view). If you do not see your table here, you need to enable it. Please visit [Data QnA setup](#) to get started.

Select table * Question *

GENERATE EQUIVALENT SQL

RESULT

Interpretation
Top 3 sum of releases by product_name

Equivalent SQL

```
SELECT  
    product_name AS product_name,  
    (SUM(releases)) AS SUM_releases  
FROM  
    `ai-analytics-solutions.dsongcp.monthly_releases`  
GROUP BY product_name  
ORDER BY SUM_releases DESC  
LIMIT 3;
```

OPEN IN QUERY EDITOR

Figure 3-22. Natural language querying in BigQuery.

Connected Sheets

Try this. Open up a new Google Sheet (you can do so by visiting <https://sheets.new>) and select Data > Data Connectors > Connect to BigQuery using the menu.¹⁴

Connect to the monthly_releases table in your project. What we now have is a *Connected Sheet*—all operations we carry out in the sheet are actually executed in BigQuery! This allows us to interact with tables and query results that may have millions of rows!

Select “Discover Data Insights” and follow the prompts to create a chart of Google Cloud releases over time (see Figure 3-23). The machines are getting smarter, and the pace of change is increasing!

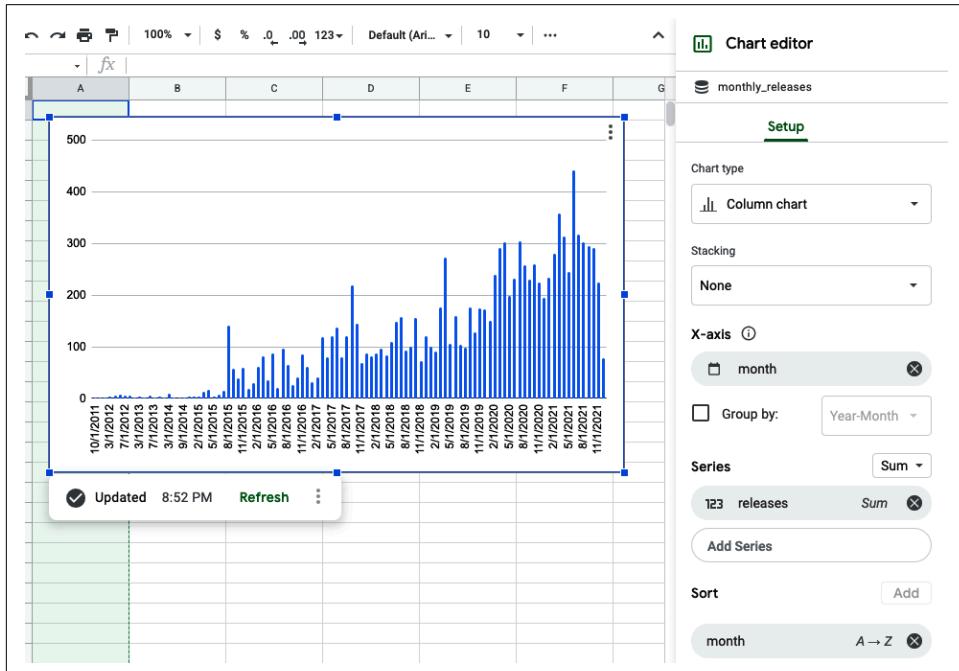


Figure 3-23. Google Cloud releases over time.

¹⁴ Make sure you are using the same Google account as the one you are using for your Google Cloud project. At the time of writing, this capability is only available to Enterprise Google Workspace accounts. If you are trying this section in a personal Google Account, you likely won't have the Data Connectors option in the Data menu.

Summary

In this chapter, we discussed the importance of bringing the insights of our end users into our data modeling efforts as early as possible. Bringing their insights is possible only if you make it a point to explain your models in context from the get-go.

We tried using Cloud SQL, a transactional, relational database whose management is simplified by virtue of it running on the cloud and being managed by Google Cloud Platform. However, it stopped scaling once we got to millions of flights. Transactional databases are not built for queries that involve scanning the entire table. For such queries, we want to use an analytics data warehouse. Hence, we switched to using BigQuery.

Within BigQuery, we previewed the table, selected a subset of columns, and created a view to make downstream analysis simpler.

The first model that we built was to suggest that our road warriors cancel their immediately scheduled meeting if the departure delay of the flight was more than 10 minutes. At this threshold, flights arrive late (when we go ahead with the meeting) less than 30% of the time. This would enable them to make 70% of their meetings with 15 minutes to spare.

We then built a dashboard in Data Studio to explain the contingency table model. Because our choice of threshold was driven by the proportion of flights that arrived late given a particular threshold, we illustrated the proportion using a pie chart for two different thresholds. We also depicted the average arrival delay given some departure delay—this gives users an intuitive understanding of why we recommend a 10-minute threshold.

Finally, we looked at trends in business intelligence.

Suggested Resources

An influential three-part series on business intelligence (BI) by Forrester Analyst [Boris Evelson](#) recommends a “layer-cake” model for modernizing BI. Unfortunately, at the time of writing, it’s behind a \$1,495 paywall.¹⁵ So, perhaps read a [summary of that article in Forbes](#) by Shant Hovsepian. According to Hovsepian, Evelson suggests that organizations modernize their use of BI by:

- Bringing BI to the data, rather than doing BI on extracts of the data (don’t do “[data cubes](#)”)
- Infusing AI such as [natural language](#) into BI

¹⁵ No, not a typo. It’s not \$14.95. It’s \$1495.

- Moving BI to the public cloud to take advantage of elasticity and separation of compute and storage

Data Studio is great as a self-service dashboard, but enterprises are often short of SQL experts. Looker provides a self-service business intelligence workflow by interposing a level of indirection, called LookML. This 2021 article by Clay Porter, “[Using Big-Query & Data Studio? You Should Check Out Looker](#)”, provides an excellent explanation. There is an [online course](#) if you want to learn how to create dashboards in Looker.

Once you create a dashboard, you might want to embed the graphics within a website. This is called embedded analytics—it can be done through iframes or using Looker’s application programming interface. See the 2020 Google blog post “[A Step-by-Step Guide to Building and Delivering Embedded Analytics](#)” by Sharon Zhang for more information.

Streaming Data: Publication and Ingest with Pub/Sub and Dataflow

In Chapter 3, we developed a dashboard to explain a contingency table-based model of suggesting whether to cancel a meeting. However, the dashboard that we built lacked immediacy because it was not tied to users' context. Because users need to be able to view a dashboard and see the information that is relevant to them at that point, we need to build a real-time dashboard with location cues.

How would we add context to our dashboard? We'd have to show maps of delays in real time. To do that, we'll need locations of the airports, and we'll need real-time data. Airport locations can be obtained from the US Bureau of Transportation Statistics (BTS; the same US government agency from which we obtained our historical flight data). Real-time flight data, however, is a commercial product. If we were to build a business out of predicting flight arrivals, we'd purchase that data feed. For the purposes of this book, however, let's just simulate it.

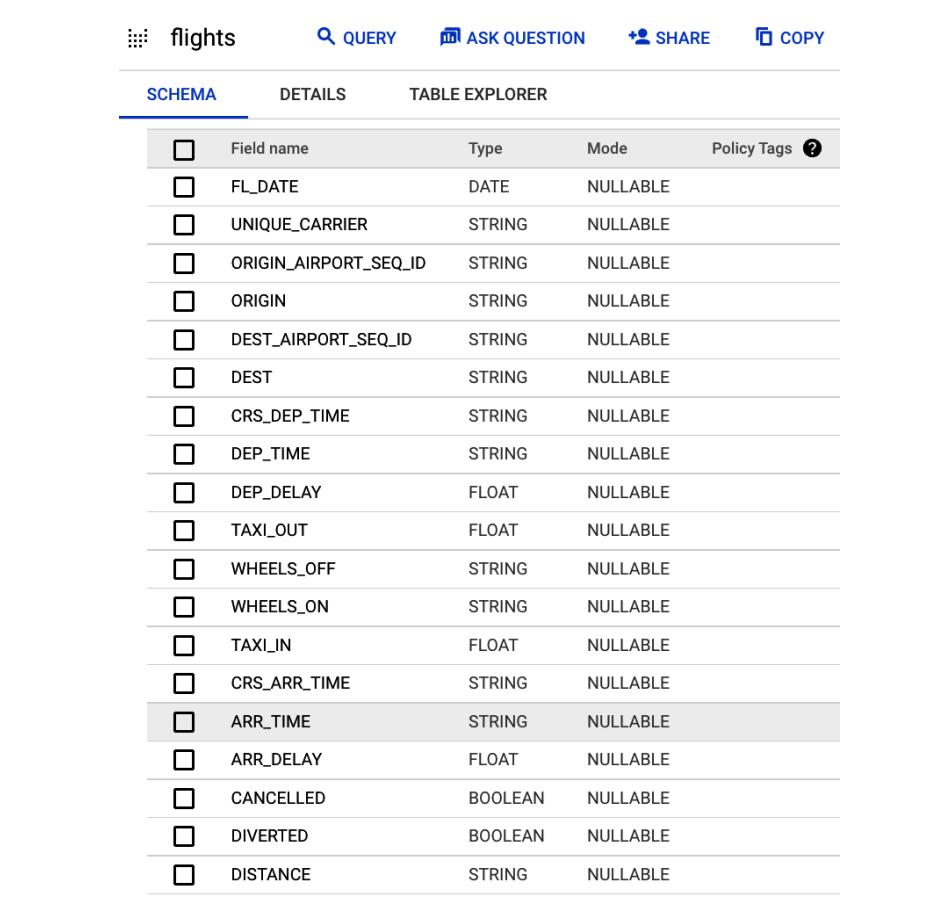
Simulating the creation of a real-time feed from historical data has the advantage of allowing us to see both sides of a streaming pipeline (production as well as consumption). In the following section, we look at how we could stream data into the database if we were to receive it in real time.



All of the code snippets in this chapter are available in the folder [04_streaming](#) of the book's GitHub repository. See the *README.md* file in that directory for instructions on how to do the steps described in this chapter.

Designing the Event Feed

Let's assume that we wish to create an event feed, not with all 100 fields in the raw BTS dataset, but with only the few fields that we selected in [Chapter 3](#) as being relevant to the flight delay prediction problem (see [Figure 4-1](#)).



The screenshot shows the 'flights' view in BigQuery. At the top, there are navigation links: 'flights' (highlighted), 'QUERY', 'ASK QUESTION', 'SHARE', and 'COPY'. Below this is a tab bar with 'SCHEMA' (selected), 'DETAILS', and 'TABLE EXPLORER'. The main area displays a table of fields:

<input type="checkbox"/>	Field name	Type	Mode	Policy Tags
<input type="checkbox"/>	FL_DATE	DATE	NULLABLE	
<input type="checkbox"/>	UNIQUE_CARRIER	STRING	NULLABLE	
<input type="checkbox"/>	ORIGIN_AIRPORT_SEQ_ID	STRING	NULLABLE	
<input type="checkbox"/>	ORIGIN	STRING	NULLABLE	
<input type="checkbox"/>	DEST_AIRPORT_SEQ_ID	STRING	NULLABLE	
<input type="checkbox"/>	DEST	STRING	NULLABLE	
<input type="checkbox"/>	CRS_DEP_TIME	STRING	NULLABLE	
<input type="checkbox"/>	DEP_TIME	STRING	NULLABLE	
<input type="checkbox"/>	DEP_DELAY	FLOAT	NULLABLE	
<input type="checkbox"/>	TAXI_OUT	FLOAT	NULLABLE	
<input type="checkbox"/>	WHEELS_OFF	STRING	NULLABLE	
<input type="checkbox"/>	WHEELS_ON	STRING	NULLABLE	
<input type="checkbox"/>	TAXI_IN	FLOAT	NULLABLE	
<input type="checkbox"/>	CRS_ARR_TIME	STRING	NULLABLE	
<input type="checkbox"/>	ARR_TIME	STRING	NULLABLE	
<input type="checkbox"/>	ARR_DELAY	FLOAT	NULLABLE	
<input type="checkbox"/>	CANCELLED	BOOLEAN	NULLABLE	
<input type="checkbox"/>	DIVERTED	BOOLEAN	NULLABLE	
<input type="checkbox"/>	DISTANCE	STRING	NULLABLE	

Figure 4-1. In [Chapter 3](#), we created a view in BigQuery with the fields relevant to the flight delay prediction problem. In this chapter, we will simulate a real-time stream of this information.

To simulate a real-time stream of the flight information shown in [Figure 4-1](#), we can begin by using the historical data in the `flights` view in BigQuery but will need to transform it further. What kinds of transformations are needed?

Transformations Needed

Note that FL_DATE is a Date while DEP_TIME is a STRING. This is because FL_DATE is of the form 2015-07-03 for July 3, 2015, whereas DEP_DATE is of the form 1406 for 2:06 p.m. local time. This is unfortunate. I'm not worried about the separation of date and time into two columns—we can remedy that. What's unfortunate is that there is no time zone offset associated with the departure time. Thus, in this dataset, a departure time of 1406 in different rows can be different times depending on the time zone of the origin airport.

The time zone offsets (there are two, one for the origin airport and another for the destination) are not present in the data. Because the offset depends on the airport location, we need to find a dataset that contains the time zone offset of each airport and then mash this data with that dataset.¹ To simplify downstream analysis, we will then put all the times in the data in a common time zone—Coordinated Universal Time (UTC) is the traditional choice of common time zone for datasets. We cannot, however, get rid of the local time—we will need the local time in order to carry out analysis, such as the typical delay associated with morning flights versus evening flights. So, although we will convert the local times to UTC, we will also store the time zone offset (e.g., -3,600 minutes) to retrieve the local time if necessary.

Therefore, we are going to carry out two transformations to the original dataset. First, we will convert all the time fields in the raw dataset to UTC. Second, in addition to the fields present in the raw data, we will add three fields to the dataset for the origin airport and the same three fields for the destination airport: the latitude, longitude, and time zone offset. These fields will be named:

```
DEP_AIRPORT_LAT, DEP_AIRPORT_LON, DEP_AIRPORT_TZOFFSET  
ARR_AIRPORT_LAT, ARR_AIRPORT_LON, ARR_AIRPORT_TZOFFSET
```

The third transformation that we will need to carry out is that, for every row in the historical dataset, we will need to publish multiple events. This is because it would be too late if we wait until the aircraft has arrived to send out a single event containing all the row data. If we do this at the time the aircraft departs, our models will be violating causality constraints. Instead, we will need to send out events corresponding to each state the flight is in. Let's choose to send out five events for each flight: when the flight is first scheduled, when the flight departs the gate, when the flight lifts off, when the flight lands, and when the flight arrives. These five events cannot have all the same data associated with them because the knowability of the columns changes during the flight. For example, when sending out an event at the departure time, we will

¹ This is a common situation. It is only as you start to explore a dataset that you discover you need ancillary datasets. Had I known beforehand, I would have ingested both datasets. But you are following my workflow, and as of this point, I knew that I needed a dataset of time zone offsets but hadn't yet searched for it!

not know the arrival time. For simplicity, we can notify the same structure, but we will need to ensure that unknowable data is marked by a `null` and not with the actual data value.

Architecture

Table 4-1 lists when those events can be sent out and the fields that will be included in each event.

Table 4-1. Fields that will be included in each of the five events that will be published.

Compare the order of the fields with those in the schema in [Figure 4-1](#).

Event	Sent at (UTC)	Fields included in event message
Scheduled	CRS_DEP_TIME minus 7 days	FL_DATE, UNIQUE_CARRIER, ORIGIN_AIRPORT_SEQ_ID, ORIGIN, DEST_AIRPORT_SEQ_ID, DEST, CRS_DEP_TIME [nulls], CRS_ARR_TIME [nulls], DISTANCE
Departed	DEP_TIME	All fields available in scheduled message, plus: <ul style="list-style-type: none">• DEP_TIME, DEP_DELAY CANCELLED• CANCELLATION_CODE• DEP_AIRPORT_LAT, DEP_AIRPORT_LON, DEP_AIRPORT_TZOFFSET
Wheelsoff	WHEELS_OFF	All fields available in departed message, plus: TAXI_OUT and WHEELS_OFF
Wheelson	WHEELS_ON	All fields available in wheelsoff message, plus: <ul style="list-style-type: none">• WHEELS_ON• DIVERTED• ARR_AIRPORT_LAT, ARR_AIRPORT_LON, ARR_AIRPORT_TZOFFSET
Arrived	ARR_TIME	All fields available in wheelson message, plus: ARR_TIME and ARR_DELAY

We will carry out the transformations needed and then store the transformed data in a database so that it is ready for the event simulation code to use. [Figure 4-2](#) shows the steps we are about to carry out in our extract-transform-load (ETL) pipeline and the subsequent steps to simulate an event stream from these events, and then create a real-time dashboard from the simulated event stream.

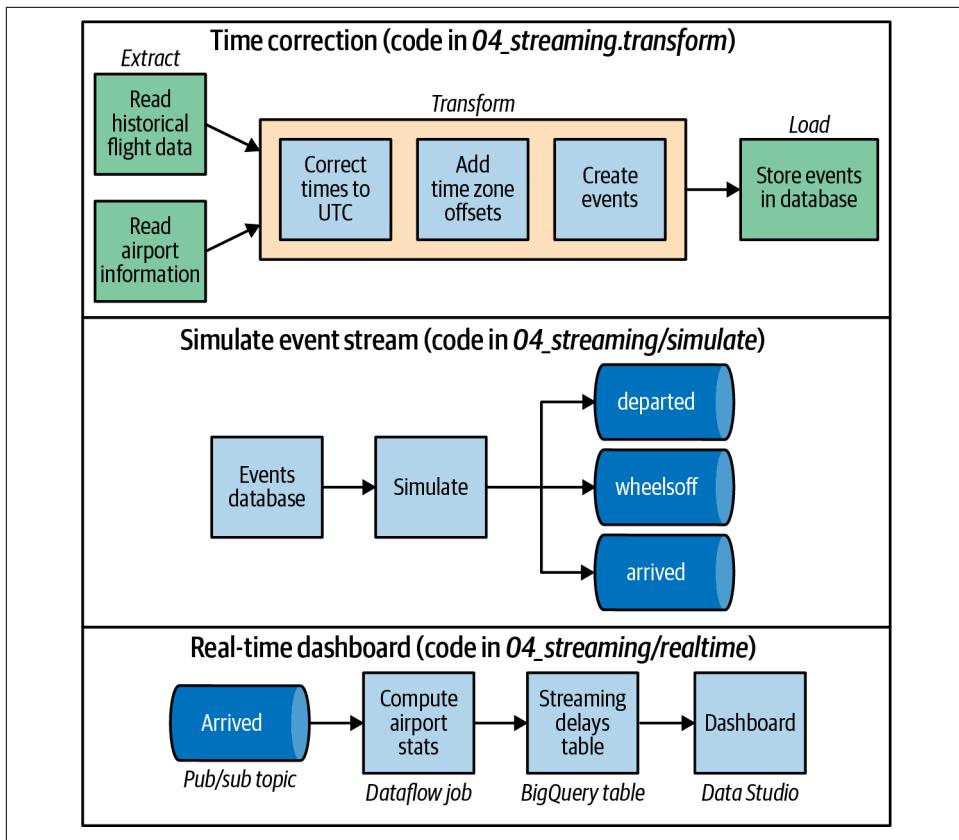


Figure 4-2. Steps in our ETL (extract-transform-load) pipeline to (a) transform the raw data into events, (b) simulate the event stream, and (c) process the event stream to populate a real-time dashboard.

Getting Airport Information

In order to do the time correction, we need to obtain the latitude and longitude of each airport. The BTS has a [dataset](#) that contains this information, which we can use to do the lookup. For convenience, I've downloaded the data and made it publicly available at <gs://data-science-on-gcp/edition2/raw/airports.csv>.

Let's examine the data to determine how to get the latitude and longitude of the airports. In [Chapter 2](#), when I needed to explore the flights data to create the first delays model, I loaded the data into BigQuery.

Do we have to import all the data that is shared with us into our BigQuery dataset in order to do exploration? Of course not. We can query BigQuery datasets in other projects without having to make our own copies of the data. In the `FROM` clause of the

BigQuery query, all that we have to do is to specify the name of the project that the dataset lives in:

```
SELECT
  airline,
  AVG(departure_delay) AS avg_dep_delay
FROM `bigquery-samples.airline_onime_data.flights`
GROUP BY airline
ORDER by avg_dep_delay DESC
```

What if someone shares a comma-separated values (CSV) file with us, though? Do we have to load the data into BigQuery in order to see what's in the file? No.

BigQuery allows us to query data in Cloud Storage through its *federated query* capabilities. This is the ability of BigQuery to query data that is not stored within the data warehouse product, but instead operate on data sources such as Google Sheets (a spreadsheet on Google Drive) or files on Cloud Storage. Thus, we could leave the files as CSV on Cloud Storage, define a table structure on it, and query the CSV files directly. Recall that we suggested using Cloud Storage if your primary analysis pattern involves working with your data at the level of flat files—this is a way of occasionally applying SQL queries to such datasets.

The first step is to get the schema of these files. Let's look at the first line:

```
gsutil cat gs://data-science-on-gcp/edition2/raw/airports.csv | head -1
```

We get:

```
"AIRPORT_SEQ_ID","AIRPORT_ID","AIRPORT","DISPLAY_AIRPORT_NAME",
"DISPLAY_AIRPORT_CITY_NAME_FULL","AIRPORT_WAC_SEQ_ID2","AIRPORT_WAC",
"AIRPORT_COUNTRY_NAME","AIRPORT_COUNTRY_CODE_ISO","AIRPORT_STATE_NAME",
"AIRPORT_STATE_CODE","AIRPORT_STATE_FIPS","CITY_MARKET_SEQ_ID","CITY_MARKET_ID",
"DISPLAY_CITY_MARKET_NAME_FULL","CITY_MARKET_WAC_SEQ_ID2","CITY_MARKET_WAC",
"LAT_DEGREES","LAT_HEMISPHERE","LAT_MINUTES","LAT_SECONDS","LATITUDE",
"LONGITUDE","LON_DEGREES","LON_HEMISPHERE","LON_MINUTES","LON_SECONDS","LONGITUDE",
"UTC_LOCAL_TIME_VARIATION","AIRPORT_START_DATE","AIRPORT_THRU_DATE",
"AIRPORT_IS_CLOSED","AIRPORT_IS_LATEST"
```

Use this header to write a BigQuery schema string of the format (specify STRING for any column you are not sure about, since you can always CAST it to the appropriate format when querying the data):

```
AIRPORT_SEQ_ID:INTEGER,AIRPORT_ID:STRING,AIRPORT:STRING, ...
```

Alternately, if you have a similar dataset lying around, start from its schema and edit it:

```
bq show --format=prettyjson dsongcp.sometable > starter.json
```

Once we have the schema of the GCS files, we can make a table definition for the federated source:²

```
bq mk --external_table_definition= \
./airport_schema.json@CSV(gs://data-science-on-gcp/edition2/raw/airports.csv \
dsongcp.airports_gcs
```

If you visit the BigQuery web console now, you should see a new table listed in the dsongcp dataset (reload the page if necessary). This is a federated data source in that its storage remains the CSV file on Cloud Storage. Yet you can query it just like any other BigQuery table:

```
SELECT
AIRPORT_SEQ_ID, AIRPORT_ID, AIRPORT, DISPLAY_AIRPORT_NAME,
LAT_DEGREES, LAT_HEMISPHERE, LAT_MINUTES, LAT_SECONDS, LATITUDE
FROM dsongcp.airports_gcs
WHERE DISPLAY_AIRPORT_NAME LIKE '%Seattle%'
```

In the preceding query, I am trying to find in the file which airport column and which latitude column I need to use. The result indicates that AIRPORT and LATITUDE are the columns of interest, but that there are several rows corresponding to the airport SEA:

Row	AIRPORT_SEQ_ID	AIRPORT_ID	AIRPORT	DISPLAY_AIRPORT_NAME	LAT_DEGREES	LAT_HEMISPHERE	LAT_MINUTES	LAT_SECONDS	LATITUDE
				PORT_NAME					
1	1247701	12477	JFB	Seattle 1st National.Bank Helipad	47	N	36	25	47.60694444
2	1474701	14747	SEA	Seattle Inter national	47	N	26	50	47.44722222
3	1474702	14747	SEA	Seattle/ Tacoma Inter national	47	N	26	57	47.44916667
4	1474703	14747	SEA	Seattle/ Tacoma Inter national	47	N	27	0	47.45

Fortunately, there is a column that indicates which row is the latest information, so what I need to do is:

```
SELECT
AIRPORT, LATITUDE, LONGITUDE
FROM dsongcp.airports_gcs
WHERE AIRPORT_IS_LATEST = 1 AND AIRPORT = 'DFW'
```

² See *04_streaming/design/mktbl.sh* for the actual syntax; we've made adjustments here for printing purposes.

Don't get carried away by federated queries, though. The most appropriate uses of federated sources involve frequently changing, relatively small datasets that need to be joined with large datasets in BigQuery native tables. Because the columnar storage in BigQuery is so fundamental to its performance, we will load most data into BigQuery's native format.

Sharing Data

Now that we have the *airports.csv* in Cloud Storage and the airports' dataset in BigQuery, it is quite likely that our colleagues will want to use this data too. Let's share it with them—one of the benefits of bringing your data to the cloud (and more specifically into a data warehouse) is to allow the mashing of datasets across organizational boundaries. So, unless you have a clear reason not to do so, like security precautions, try to make your data widely accessible.

Costs of querying are borne by the person submitting the query to the BigQuery engine, so you don't need to worry that you are incurring additional costs for your division by doing this. It is possible to make a GCS bucket "requester-pays" to get the same sort of billing separation for data in Cloud Storage.

Sharing a Cloud Storage dataset

To share some data in Cloud Storage, use `gsutil`:

```
gsutil -m acl ch -r -u abc@xyz.com:R gs://$BUCKET/data
```

In the preceding command, the `-m` indicates multithreaded mode, the `-r` provides access recursively starting with the top-level directory specified, and the `-u` indicates that this is a user being granted read (`:R`) access.

We could provide read access to the entire organization or a Google Group using `-g`:

```
gsutil -m acl ch -r -g xyz.com:R gs://$BUCKET/data
```

Sharing a BigQuery dataset

BigQuery sharing can happen at the granularity of a column, a table, or a dataset. None of our BigQuery tables hold personally identifiable or confidential information. Therefore, there is no compelling access-control reason to control the sharing of flight information at a column or table level. So, we can share the `dsongcp` dataset that was created in [Chapter 2](#), and we can make everyone in the organization working on this project a `bigrquery.user` so that they can carry out queries on this dataset. You can do this from the BigQuery web console from the dataset menu.

In some cases, you might find that your dataset or table contains certain columns that have personally identifying or confidential information. You might need to restrict access to those columns while leaving the remainder of the table accessible to a wider

audience.³ Whenever you need to provide access to a subset of a table in BigQuery (whether it is specific columns or specific rows), you can use views. Put the table itself in a dataset that is accessible to a very small set of users. Then, create a view on this table that will pull out the relevant columns and rows and save this view in a separate dataset that has wider accessibility. Your users will query only this view, and because the personally identifying or confidential information is not even present in the view, the chances of inadvertent leakage are lowered.

Another way to restrict access at the level of a [BigQuery table](#) is to use Cloud IAM. To control access at the level of a column, you'd use [policy tags](#) and Data Catalog.

Dataplex and Analytics Hub

Once you get into the habit of sharing data widely, governance can become problematic. It is better if you can administer data across Cloud Storage in a consistent manner and track lineage, etc. That's what Dataplex is for.

It can be rather cumbersome to share tables and datasets one at a time with one user or one group at a time. To implement sharing at scale and get statistics on how people are using the data you have shared, use Analytics Hub.

Time Correction

Correcting times reported in local time to UTC is not a simple endeavor. There are several steps:

- Local time depends on, well, the location. The flight data that we have records only the name of the airport (e.g., ALB for Albany). We, therefore, need to obtain the latitude and longitude given an airport code. The BTS has a dataset that contains this information, which we can use to do the lookup.
- Given a latitude/longitude pair, we need to look up the time zone from a [map of global time zones](#). For example, given the latitude and longitude of the airport in Albany, we would need to get back `America/New_York`. There are several web services that do this, but the Python package `timezonefinder` is a more efficient option because it works completely offline. The drawback is that this package does not handle oceanic areas and some historical time zone changes,⁴ but that's a trade-off that we can make for now.

³ Or make a copy or view of the table with anonymized column values—we cover safeguarding personally identifiable information in [Chapter 7](#) and in the Appendix.

⁴ For example, the time zone of Sevastopol changed in 2014 from Eastern European Time (UTC+2) to Moscow Time (UTC+4) after the annexation of Crimea by the Russian Federation.

- The time zone offset (from Greenwich Mean Time [GMT/UTC]) at a location changes during the year due to daylight savings corrections. In New York, for example, it is six hours in summer and five hours in winter behind UTC. Given the time zone (`America/New_York`), therefore, we also need the local departure date and time (say Jan. 13, 2015, 2:08 p.m.) in order to find the corresponding time zone offset. The Python package `pytz` provides this capability by using the underlying operating system.

The problem of ambiguous times still remains—every instant between 01:00 and 02:00 local time occurs twice on the day that the clock switches from daylight savings time (summer time) to standard time (winter time). So, if our dataset has a flight arriving at 01:30, we need to make a choice of what time that represents. In a real-world situation, you would look at the typical duration of the flight and choose the one that is more likely. For the purposes of this book, I'll always assume the winter time (i.e., `is_dst` is `False`) on the dubious grounds that it is the standard time zone for that location.

The complexity of these steps should, I hope, convince you to follow best practices when storing time.

Best Practices When Storing Time

Always try to store two columns for every timestamp:

1. The timestamp in UTC so that you can merge data from across the world if necessary.
2. The currently active time zone offset so that you can carry out analysis that requires the local time. For example, is there a spike associated with traffic between 5 p.m. and 6 p.m. local time?

Apache Beam/Cloud Dataflow

The canonical way to build data pipelines on Google Cloud Platform is to use Cloud Dataflow. Cloud Dataflow is an externalization of technologies called [Flume](#) and [Mill-Wheel](#) that have been in widespread use at Google for several years. It employs a programming model that handles both batch and streaming data in a uniform manner, thus providing the ability to use the same codebase both for batch and continuous stream processing. The code itself is written in [Apache Beam](#), either in Java, Python, or Go,⁵ and it is portable in the sense that it can be executed on multiple execution environments, including [Apache Flink](#) and [Apache Spark](#). On GCP, Cloud Dataflow provides a fully managed (serverless) service that is capable of executing Beam pipelines. Resources are allocated on demand, and they autoscale so as to achieve both minimal latency and high resource utilization.

Beam programming involves building a pipeline (a series of data transformations) that is submitted to a runner. The runner will build a graph and then stream data through it. Each input dataset comes from a source and each output dataset is sent to a sink. [Figure 4-3](#) illustrates the Beam pipeline that we are about to build.

Compare the steps in [Figure 4-2](#) with the block diagram of the ETL (extract-transform-load) pipeline in [Figure 4-3](#). Let's build the data pipeline piece by piece.

⁵ The Java API is much more mature and performant, but Python is easier and more concise. In this book, we will use Python.

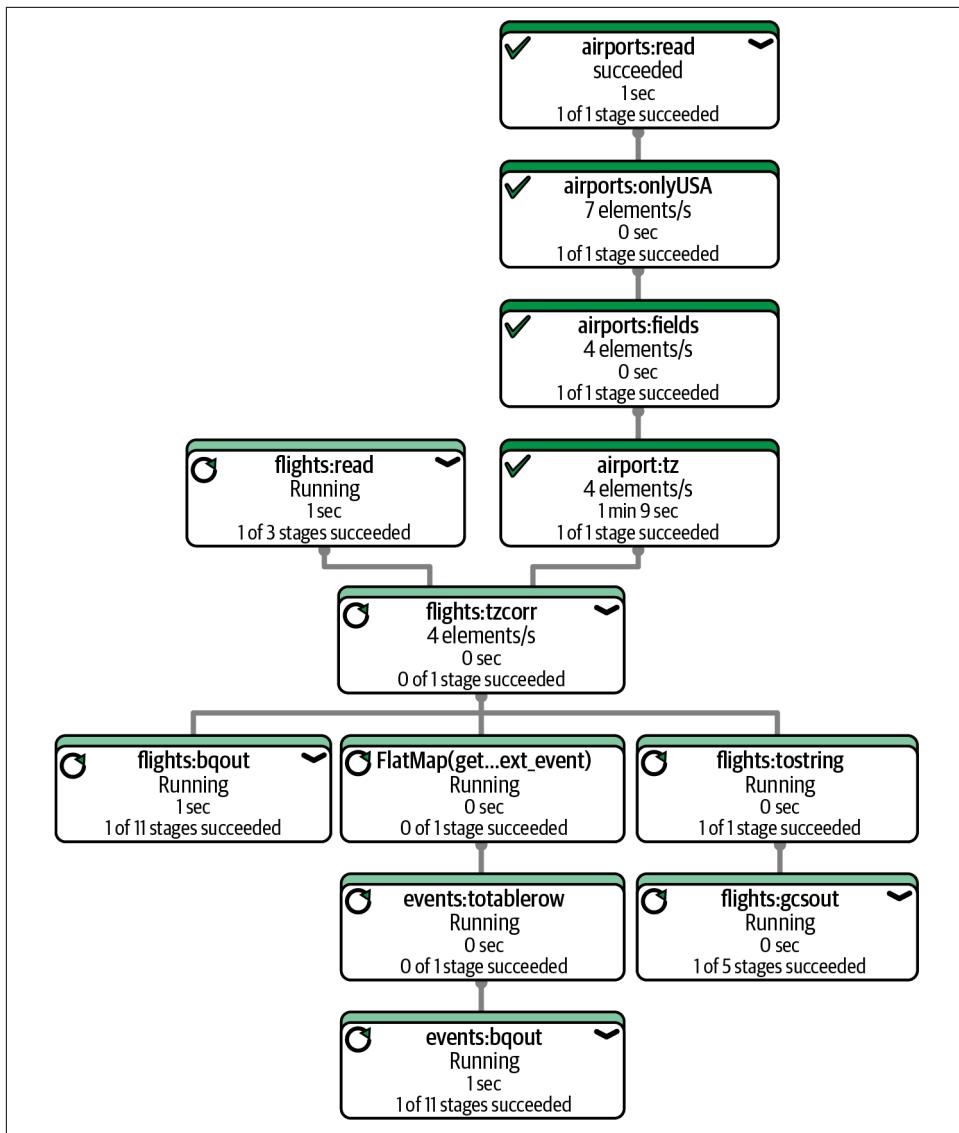


Figure 4-3. The Dataflow pipeline that we are about to build.

Parsing Airports Data

You can [download information about the location of airports](#) from the BTS website. I selected all of the fields, downloaded the CSV file to my local hard drive, extracted it, and compressed it with gzip. The gzipped airports file is available in the [GitHub repository for this book](#).

In order to use Apache Beam from Cloud Shell, we need to [install it](#) into our Python environment. Go ahead and install the time zone packages also at this time:⁶

```
virtualenv ~/beam_env
source ~/beam_env/bin/activate
python3 -m pip install --upgrade \
    pytz \
    apache-beam[gcp]
```

The `Read` transform in the Beam pipeline that follows reads in the airports file line by line:⁷

```
with beam.Pipeline('DirectRunner') as pipeline:
    airports = (pipeline
        | beam.io.ReadFromText('airports.csv.gz')
        | beam.Map(lambda line: next(csv.reader([line])))
        | beam.Map(lambda fields: (fields[0], (fields[21], fields[26]))))
    )
```

Apache Beam Python Syntax

The Apache Beam code may look nothing like any Python you have seen before. Let's break it down.

The first line creates a Beam pipeline that will be executed on the local machine ("DirectRunner"):

```
with beam.Pipeline('DirectRunner') as pipeline:
    # some code here
```

Next, all the lines within the `with` block are executed. However, those lines of code only create the execution graph. Conceptually, you can think of it as the pipeline getting compiled. It is only when the graph corresponding to the full `with` block has been created that the pipeline is executed.

Let's look at the lines of code within the `with` block. These consist of data transformation methods executed one after the other, with the output of one transform being fed in as the input to the next transform. For example, these two lines:

```
beam.io.ReadFromText('airports.csv.gz')
| beam.Map(lambda line: next(csv.reader([line])))
```

mean that the file `airports.csv.gz` should be read. `ReadFromText` reads the input file line-by-line. Each line is then sent to the `Map` transform. The pipe symbol (`|`) is what

⁶ If you are using an ephemeral shell like Cloud Shell, you will have to run the activate line every time you start a new session. This will load up the [virtual environment](#) that you were using earlier. This way, you will not need to reinstall the Python packages every time.

⁷ This code is in `04_streaming/transform/df01.py` of the GitHub repository of this book.

says that the output of `ReadFromText` has to be sent in as the input to `Map`. You might be familiar with this idiomatic use of the pipe symbol in the Linux command line.

The `Map` transform applies some user-specified function to the data. What function are we applying? We could have written:

```
| beam.Map(parse_line)
```

and defined the function `parse_line` as follows:

```
def parse_line(line):
    return next(csv.reader([line]))
```

But doing this would have involved writing a whole bunch of small functions. The lambda syntax in Python allows us to define a function in-line without giving the function a name:

```
| lambda line: some_code_with(line))
```

The last strange thing might be the function body itself:

```
next(csv.reader([line]))
```

We are invoking the `reader()` function in the Python module called `csv` and passing it an array of lines. The array here has only one item (`[line]`). The `reader()` will give us back an iterator to an array of parsed rows. When we call `next()` in Python, we are asking the iterator to get us the next item. By calling `next()`, therefore, we get the first parsed row.

This is all quite terse, but quite idiomatic. Pick up a Python book if any of this syntax was new to you.

For example, suppose that one of the input lines read out of the text file source is the following:

```
1000401,10004,"04A","Lik Mining Camp","Lik, AK",101,1,"United
States","US","Alaska","AK","02",3000401,30004,"Lik,
AK",101,1,68,"N",5,0,68.08333333,163,"W",10,0,-163.16666667,",2007-07-01,,0,1,
```

The first `Map` takes this line and passes it to a CSV reader that parses it (taking into account fields like `Lik, AK` that have commas in them) and pulls out the fields as a list of strings. These fields are then passed to the next transform. The second `Map` takes the fields as input and outputs a tuple of the form (the extracted fields are shown in bold in the previous example):

```
(1000401, (68.08333333, -163.16666667))
```

The first number is the unique airport code (we use this, rather than the airport's three-letter code, because airport locations can change over time), and the next two numbers are the latitude/longitude pair for the airport's location. The variable `airports`, which is the result of these three transformations, is not a simple in-memory

list of these tuples. Instead, it is an immutable collection, termed a `PCollection`, that you can take out-of-memory and distribute.

We can write the contents of the `PCollection` to a text file to verify that the pipeline is behaving correctly:

```
(airports
| beam.Map(lambda airport_data: '{},{}'.format(airport_data[0], ',' \
    .join(airport_data[1])))
| beam.io.WriteToText('extracted_airports')
)
```

Try this out: the code, in *04_streaming/transform/df01.py*, is just a Python program that you can run from the command line. First, install the Apache Beam package if you haven't yet done so and then run the program *df01.py* while you are in the directory containing the GitHub repository of this book:

```
cd 04_streaming/simulate
./install_packages.sh
python3 ./df01.py
```

This runs the code in *df01.py* locally. Later, we will change the pipeline line to:

```
with beam.Pipeline('DataflowRunner') as pipeline:
```

and get to run the pipeline on the Google Cloud Platform using the Cloud Dataflow service. With that change, simply running the Python program launches the data pipeline on multiple workers in the cloud. As with many distributed systems, the output of Cloud Dataflow is potentially sharded to one or more files. You will get a file whose name begins with "extracted_airports" (mine was *extracted_airports-00000-of-00001*), a few of whose lines might look something like this:

```
1000101,58.10944444,-152.90666667
1000301,65.54805556,-161.07166667
```

The columns are `AIRPORT_SEQ_ID`, `LATITUDE`, and `LONGITUDE`—the order of the rows you get depends on which of the parallel workers finished first, so it could be different.

Adding Time Zone Information

Let's now change the code to determine the time zone corresponding to a latitude/longitude pair. In our pipeline, rather than simply emitting the latitude/longitude pair, we emit a list of three items: latitude, longitude, and time zone:

```
airports = (pipeline
| beam.Read(bean.io.ReadFromText('airports.csv.gz'))
| beam.Map(lambda line: next(csv.reader([line])))
| beam.Map(lambda fields: (fields[0], addtimezone(fields[21], fields[26])))
)
```

The `lambda` keyword in Python sets up an anonymous function. In the case of the first use of `lambda` in the preceding snippet, that method takes one parameter (`line`) and returns the stuff following the colon. We can determine the time zone by using the `timezonefinder` package:⁸

```
def addtimezone(lat, lon):
    import timezonefinder
    tf = timezonefinder.TimezoneFinder()
    lat = float(lat)
    lon = float(lon)
    return (lat, lon, tf.timezone_at(lng=lon, lat=lat))
```

The location of the import statement in the preceding example might look strange (most Python imports tend to be at the top of the file), but this import-within-the-function pattern is recommended by Cloud Dataflow so that,⁹ when we submit it to the cloud, pickling of the main session doesn't end up pickling imported packages also.¹⁰

For now, though, we are going to run this (`df02.py`) locally. This will take a while because the time zone computation involves a large number of polygon intersection checks and because we are running locally, not (yet!) distributed in the cloud. So, let's speed it up by adding a filter to reduce the number of airport locations we have to look up:

```
| beam.io.ReadFromText('airports.csv.gz')
| beam.Filter(lambda line: "United States" in line
              and line[-2:] == '1,')
```

The BTS flight delay data is only for US domestic flights, so we don't need the time zones of airports outside the United States. The reason for the second check is that airport locations change over time, but we are interested only in the current location of the airport. For example, here are the airport locations for ORD (or Chicago):

```
1393001,...,"ORD",...,41.97805556,...,-87.90611111,...,1950-01-01,2011-06-30,0,0,
1393002,...,"ORD",...,41.98166667,...,-87.90666667,...,2011-07-01,2013-09-30,0,0,
1393003,...,"ORD",...,41.97944444,...,-87.90750000,...,2013-10-01,2015-09-30,0,0,
1393004,...,"ORD",...,41.97722222,...,-87.90805556,...,2015-10-01,,0,1,
```

The first row captures the location of Chicago's airport between 1950 and June 30, 2011.¹¹ The second row is valid from July 1, 2011, to September 30, 2013. The last

⁸ This code is in `04_streaming/transform/df02.py` of the GitHub repository of this book.

⁹ See the answer to the question “How do I handle `NameErrors`? ” in the [Google documentation](#).

¹⁰ Saving Python objects is called [pickling](#).

¹¹ Chicago's airport didn't pack up and move on June 30. Most likely, a new terminal or runway was opened at that time, and this changed the location of the centroid of the airport's aerial extent. Notice that the change is just 0.0036 in latitude degrees. At Chicago's latitude, this translates to about 400 meters.

row, however, is the current location and this is marked by the last column (the AIRPORT_IS_LATEST field) being 1.

That's not the only line we are interested in, however! Flights before 2015-10-01 will report the ID of the second to last row. We could add a check for this, but this looks rather dicey for a slight bit of optimization. So, I'll remove that last check, so that we have only:

```
| beam.io.ReadFromText('airports.csv.gz')  
| beam.Filter(lambda line: "United States" in line)
```

Once I do this and run *df02.py*, the extracted information for the airports looks like this:

```
1672301,62.03611111,-151.45222222,America/Anchorage  
1672401,43.87722222,-73.41305556,America/New_York  
1672501,40.75722222,-119.21277778,America/Los_Angeles
```

The last column in the extracted information has the time zone, which was determined from the latitude and longitude of each airport.

Converting Times to UTC

Now that we have the time zone for each airport, we are ready to tackle converting the times in the flights data to UTC. At the time that we are developing the program, we'd prefer not to process all the months we have in BigQuery—waiting for the query each time we run the program will be annoying. Instead, we will create a small sample of the flights data in BigQuery against which to develop our code:¹²

```
SELECT *  
FROM dsongcp.flights  
WHERE RAND() < 0.001
```

This returns about 6,000 rows. We can use the BigQuery web UI to save these results as a JavaScript Object Notation (JSON) file. However, I prefer to script things out:¹³

```
bq query --destination_table dsongcp.flights_sample \  
--replace --nouse_legacy_sql \  
'SELECT * FROM dsongcp.flights WHERE RAND() < 0.001'  
  
bq extract --destination_format=NEWLINE_DELIMITED_JSON \  
dsongcp.flights_sample \  

```

¹² Normally, the recommended way to sample a BigQuery table is to do `SELECT * FROM dsongcp.flights WHERE TABLESAMPLE SYSTEM (0.001)` because table sampling isn't cached, so we will get different results each time. However, at the time of writing, table sampling works only on tables and flights is a View. Besides, in our current use case, we don't care whether or not we get different samples each time we run the query. That's why I'm using `RAND()`.

¹³ See the file *04_streaming/transform/bqsample.sh*.

```
gs://${BUCKET}/flights/ch4/flights_sample.json  
gsutil cp gs://${BUCKET}/flights/ch4/flights_sample.json
```

This creates a file named *flight_sample.json*, a row of which looks similar to this:

```
{"FL_DATE":"2015-04-28","UNIQUE_CARRIER":"EV","ORIGIN_AIRPORT_SEQ_ID":1013503,  
"ORIGIN": "ABE", "DEST_AIRPORT_SEQ_ID":1039705, "DEST": "ATL",  
"CRS_DEP_TIME":1600, "DEP_TIME":1555, "DEP_DELAY": -5, "TAXI_OUT": 7,  
"WHEELS_OFF":1602, "WHEELS_ON":1747, "TAXI_IN":4, "CRS_ARR_TIME":1809,  
"ARR_TIME":1751, "ARR_DELAY": -18, "CANCELLED":false, "DIVERTED":false,  
"DISTANCE":692.00}
```

Reading the flights data starts out similar to reading the airports data:¹⁴

```
flights = (pipeline  
| 'flights:read' >> beam.io.ReadFromText('flights_sample.json')  
| 'flights:parse' >> beam.Map(lambda line: json.loads(line))
```

This is the same code as when we read the *airports.csv.gz* file, except that I am also giving a name (*flights:read*) to this transform step and using a JSON parser instead of a CSV parser. Note the syntax here:

```
| 'name-of-step' >> transform_function()
```

The next step, though, is different because it involves two PCollections. We need to join the flights data with the airports data to find the time zone corresponding to each flight. To do that, we make the airports PCollection a “side input.” Side inputs in Beam are like views into the original PCollection, and are either lists or dicts (dictionaries). In this case, we will create a dict that maps airport ID to information about the airports:

```
flights = (pipeline  
| 'flights:read' >> beam.io.ReadFromText('flights_sample.json')  
| 'flights:parse' >> beam.Map(lambda line: json.loads(line))  
| 'flights:tzcorr' >> beam.FlatMap(tz_correct,  
                                beam.pvalue.AsDict(airports))  
)
```



The fact that the PCollection has to be a Python list or a Python dict means that side inputs have to be small enough to fit into memory. If you need to join two large PCollections that will not fit into memory, use a [CoGroupByKey](#).

¹⁴ This code is in *04_streaming/transform/df03.py* of the GitHub repository of this book.

The `FlatMap()` method calls out to a method `tz_correct()`, which takes the parsed content of a line from `flights_sample.json` (containing a single flight's information) and a Python dictionary (containing all the airports' time zone information):

```
def tz_correct(fields, airport_timezones):
    try:
        # convert all times to UTC
        # ORIGIN_AIRPORT_SEQ_ID is the name of JSON attribute
        dep_airport_id = fields["ORIGIN_AIRPORT_SEQ_ID"]
        arr_airport_id = fields["DEST_AIRPORT_SEQ_ID"]
        # airport_id is the key to airport_timezones dict
        # and the value is a tuple (lat, lon, timezone)
        dep_timezone = airport_timezones[dep_airport_id][2]
        arr_timezone = airport_timezones[arr_airport_id][2]

        for f in ["CRS_DEP_TIME", "DEP_TIME", "WHEELS_OFF"]:
            fields[f] = as_utc(fields["FL_DATE"], fields[f], dep_timezone)
        for f in ["WHEELS_ON", "CRS_ARR_TIME", "ARR_TIME"]:
            fields[f] = as_utc(fields["FL_DATE"], fields[f], arr_timezone)

        yield json.dumps(fields)
    except KeyError as e:
        logging.exception(" Ignoring " + line +
                           " because airport is not known")
```

Why `FlatMap()` instead of `Map` to call `tz_correct()`? A `Map` is a 1-to-1 relation between input and output, whereas a `FlatMap()` can return 0–N outputs per input. The way it does this is with a Python generator function (i.e., the `yield` keyword—think of the `yield` as a return that returns one item at a time until there is no more data to return). Using `FlatMap` here allows us to ignore any flight information corresponding to unknown airports—even though this doesn't happen in the historical data we are processing, a little bit of defensive programming doesn't hurt.

The `tz_correct()` code gets the departure airport ID from the flight's data and then looks up the time zone for that airport ID from the airport's data. After it has the time zone, it calls out to the method `as_utc()` to convert each of the datetimes reported in that airport's time zone to UTC:

```
def as_utc(date, hhmm, tzone):
    try:
        if len(hhmm) > 0 and tzone is not None:
            import datetime, pytz
            loc_tz = pytz.timezone(tzone)
            loc_dt = loc_tz.localize(datetime.datetime.strptime(date, '%Y-%m-%d'),
                                    is_dst=False)
            loc_dt += datetime.timedelta(hours=int(hhmm[:2]),
                                         minutes=int(hhmm[2:])))
            utc_dt = loc_dt.astimezone(pytz.utc)
            return utc_dt.strftime('%Y-%m-%d %H:%M:%S')
    else:
```

```

        return '' # empty string corresponds to canceled flights
    except ValueError as e:
        print('{}'.format(date, hhmm, tzone))
        raise e

```

As before, you can run this locally. To do that, run `df03.py`. A line that originally (in the raw data) looked like:

```
{"FL_DATE": "2015-11-05", "UNIQUE_CARRIER": "DL", "ORIGIN_AIRPORT_SEQ_ID": "1013503",
"ORIGIN": "ABE", "DEST_AIRPORT_SEQ_ID": "1039705", "DEST": "ATL",
"CRS_DEP_TIME": "0600", "DEP_TIME": "0556", "DEP_DELAY": -4, "TAXI_OUT": 12,
"WHEELS_OFF": "0608", "WHEELS_ON": "0749", "TAXI_IN": 10, "CRS_ARR_TIME": "0818",
"ARR_TIME": "0759", "ARR_DELAY": -19, "CANCELLED": false,
"DIVERTED": false, "DISTANCE": "692.00"}
```

now becomes:

```
{"FL_DATE": "2015-11-05", "UNIQUE_CARRIER": "DL",
"ORIGIN_AIRPORT_SEQ_ID": "1013503", "ORIGIN": "ABE",
"DEST_AIRPORT_SEQ_ID": "1039705", "DEST": "ATL",
"CRS_DEP_TIME": "2015-11-05 11:00:00", "DEP_TIME": "2015-11-05 10:56:00",
"DEP_DELAY": -4, "TAXI_OUT": 12, "WHEELS_OFF": "2015-11-05 11:08:00",
"WHEELS_ON": "2015-11-05 12:49:00", "TAXI_IN": 10,
"CRS_ARR_TIME": "2015-11-05 13:18:00", "ARR_TIME": "2015-11-05 12:59:00",
"ARR_DELAY": -19, "CANCELLED": false, "DIVERTED": false, "DISTANCE": "692.00"}
```

All the times have been converted to UTC. For example, the 0759 time of arrival in Atlanta has been converted to UTC to become 12:59:00.

Correcting Dates

Look carefully at the following line involving a flight from Honolulu (HNL) to Dallas–Fort Worth (DFW). Do you notice anything odd?

```
{"FL_DATE": "2015-03-06", "UNIQUE_CARRIER": "AA",
"ORIGIN_AIRPORT_SEQ_ID": "1217302", "ORIGIN": "HNL",
"DEST_AIRPORT_SEQ_ID": "1129803", "DEST": "DFW",
"CRS_DEP_TIME": "2015-03-07 05:30:00", "DEP_TIME": "2015-03-07 05:22:00",
"DEP_DELAY": -8, "TAXI_OUT": 40, "WHEELS_OFF": "2015-03-07 06:02:00",
"WHEELS_ON": "2015-03-06 12:32:00", "TAXI_IN": 7,
"CRS_ARR_TIME": "2015-03-06 12:54:00", "ARR_TIME": "2015-03-06 12:39:00",
"ARR_DELAY": -15, "CANCELLED": false, "DIVERTED": false, "DISTANCE": "3784.00"}
```

Examine the departure time in Honolulu and the arrival time in Dallas—the flight is arriving the day before it departed! That's because the flight date (2015-03-06) is the date of departure in local time. Add in a time difference between airports, and it is quite possible that it is not the date of arrival. We'll look for these situations and add 24 hours if necessary. This is, of course, quite a hack (have I already mentioned that times ought to be stored in UTC?!):

```
def add_24h_if_before(arrtime, deptime):
    import datetime
```

```

if len(arrrtime) > 0 and len(deptime) > 0 and arrtime < deptime:
    adt = datetime.datetime.strptime(arrrtime, '%Y-%m-%d %H:%M:%S')
    adt += datetime.timedelta(hours=24)
    return adt.strftime('%Y-%m-%d %H:%M:%S')
else:
    return arrtime

```

The 24-hour hack is called just before the yield in `tz_correct`.¹⁵ Now that we have new data about the airports, it is probably wise to add it to our dataset. Also, as remarked earlier, we want to keep track of the time zone offset from UTC because some types of analysis might require knowledge of the local time. Thus, the new `tz_correct` code becomes the following:

```

def tz_correct(line, airport_timezones):
    fields = json.loads(line)
    try:
        # convert all times to UTC
        dep_airport_id = fields["ORIGIN_AIRPORT_SEQ_ID"]
        arr_airport_id = fields["DEST_AIRPORT_SEQ_ID"]
        dep_timezone = airport_timezones[dep_airport_id][2]
        arr_timezone = airport_timezones[arr_airport_id][2]

        for f in ["CRS_DEP_TIME", "DEP_TIME", "WHEELS_OFF"]:
            fields[f], deptz = as_utc(fields["FL_DATE"], fields[f], dep_timezone)
        for f in ["WHEELS_ON", "CRS_ARR_TIME", "ARR_TIME"]:
            fields[f], arrtz = as_utc(fields["FL_DATE"], fields[f], arr_timezone)

        for f in ["WHEELS_OFF", "WHEELS_ON", "CRS_ARR_TIME", "ARR_TIME"]:
            fields[f] = add_24h_if_before(fields[f], fields["DEP_TIME"])

        fields["DEP_AIRPORT_TZOFFSET"] = deptz
        fields["ARR_AIRPORT_TZOFFSET"] = arrtz
        yield json.dumps(fields)
    except KeyError as e:
        logging.exception(" Ignoring " + line + " because airport is not known")

```

When I run `df04.py`, which has these changes applied to it, the flight from Honolulu to Dallas becomes:

```
{"FL_DATE": "2015-03-06", "UNIQUE_CARRIER": "AA",
"ORIGIN_AIRPORT_SEQ_ID": "1217302", "ORIGIN": "HNL",
"DEST_AIRPORT_SEQ_ID": "1129803", "DEST": "DFW",
"CRS_DEP_TIME": "2015-03-07 05:30:00", "DEP_TIME": "2015-03-07 05:22:00",
"DEP_DELAY": -8, "TAXI_OUT": 40, "WHEELS_OFF": "2015-03-07 06:02:00",
"WHEELS_ON": "2015-03-07 12:32:00", "TAXI_IN": 7,
"CRS_ARR_TIME": "2015-03-07 12:54:00", "ARR_TIME": "2015-03-07 12:39:00",
"ARR_DELAY": -15, "CANCELLED": false, "DIVERTED": false, "DISTANCE": "3784.00",
"DEP_AIRPORT_TZOFFSET": -36000.0, "ARR_AIRPORT_TZOFFSET": -21600.0}
```

¹⁵ This code is in `04_streaming/transform/df04.py` of the GitHub repository of this book.

As you can see, the dates have now been corrected (see the bolded parts).

Creating Events

After we have our time-corrected data, we can move on to creating events to publish into Pub/Sub. For now, we'll limit ourselves to just the departed and arrived messages—we can rerun the pipeline to create the additional events if and when our modeling efforts begin to use other events:

```
def get_next_event(fields):
    if len(fields["DEP_TIME"]) > 0:
        event = dict(fields) # copy
        event["EVENT_TYPE"] = "departed"
        event["EVENT_TIME"] = fields["DEP_TIME"]
        for f in ["TAXI_OUT", "WHEELS_OFF", "WHEELS_ON",
                  "TAXI_IN", "ARR_TIME", "ARR_DELAY", "DISTANCE"]:
            event.pop(f, None) # not knowable at departure time
        yield event
    if len(fields["ARR_TIME"]) > 0:
        event = dict(fields)
        event["EVENT_TYPE"] = "arrived"
        event["EVENT_TIME"] = fields["ARR_TIME"]
        yield event
```

Essentially, we pick up the departure time and create a `departed` event at that time after making sure to remove the fields (such as arrival delay) we cannot know at the departure time. Similarly, we use the arrival time to create an `arrived` event, as shown in [Figure 4-4](#).

In the pipeline, the event creation code is called on the `flights` PCollection after the conversion to UTC has happened:

```
flights = (pipeline
    | 'flights:read' >> beam.io.ReadFromText('flights_sample.json')
    | 'flights:tzcorr' >> beam.FlatMap(tz_correct,
                                         beam.pvalue.AsDict(airports))
)
events = flights | beam.FlatMap(get_next_event)
```

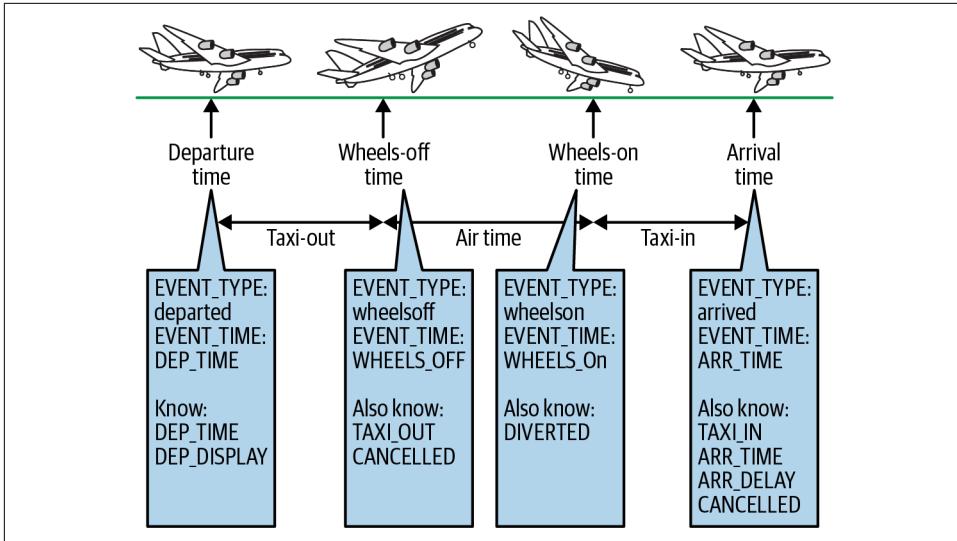


Figure 4-4. Events, when they are published, and some of the fields present in those events.

If we now run the pipeline,¹⁶ we will see two events for each flight:

```
{"FL_DATE": "2015-04-28", "UNIQUE_CARRIER": "EV",
"ORIGIN_AIRPORT_SEQ_ID": "1013503", "ORIGIN": "ABE",
"DEST_AIRPORT_SEQ_ID": "1039705", "DEST": "ATL",
"CRS_DEP_TIME": "2015-04-28 20:00:00", "DEP_TIME": "2015-04-28 19:55:00",
"DEP_DELAY": -5, "CRS_ARR_TIME": "2015-04-28 22:09:00", "CANCELLED": false,
"DIVERTED": false, "DEP_AIRPORT_TZOFFSET": -14400.0,
"ARR_AIRPORT_TZOFFSET": -14400.0, "EVENT_TYPE": "departed",
"EVENT_TIME": "2015-04-28 19:55:00"}
{"FL_DATE": "2015-04-28", "UNIQUE_CARRIER": "EV",
"ORIGIN_AIRPORT_SEQ_ID": "1013503", "ORIGIN": "ABE",
"DEST_AIRPORT_SEQ_ID": "1039705", "DEST": "ATL",
"CRS_DEP_TIME": "2015-04-28 20:00:00", "DEP_TIME": "2015-04-28 19:55:00",
"DEP_DELAY": -5, "TAXI_OUT": 7, "WHEELS_OFF": "2015-04-28 20:02:00",
"WHEELS_ON": "2015-04-28 21:47:00", "TAXI_IN": 4,
"CRS_ARR_TIME": "2015-04-28 22:09:00", "ARR_TIME": "2015-04-28 21:51:00",
"ARR_DELAY": -18, "CANCELLED": false, "DIVERTED": false, "DISTANCE": "692.00",
"DEP_AIRPORT_TZOFFSET": -14400.0, "ARR_AIRPORT_TZOFFSET": -14400.0,
"EVENT_TYPE": "arrived", "EVENT_TIME": "2015-04-28 21:51:00"}
```

The first event is a departed event and is to be published at the departure time, while the second event is an arrived event and is to be published at the arrival time. The

¹⁶ This code is in `04_streaming/transform/df05.py` of the GitHub repository of this book.

departed event has a number of missing fields corresponding to data that is not known at that time.

Once we have this code working, let's add a third event that will be sent when the plane takes off:

```
if len(fields["WHEELS_OFF"]) > 0:  
    event = dict(fields) # copy  
    event["EVENT_TYPE"] = "wheelsoff"  
    event["EVENT_TIME"] = fields["WHEELS_OFF"]  
    for f in ["WHEELS_ON", "TAXI_IN",  
              "ARR_TIME", "ARR_DELAY", "DISTANCE"]:  
        event.pop(f, None) # not knowable at departure time  
    yield event
```

At this point, we haven't created a wheelsdown event yet.

Reading and Writing to the Cloud

So far, we have been reading and writing local files. However, once we start to run our code in production, in a serverless environment, the concept of a local drive no longer makes sense. We have to read and write from Cloud Storage. Also, because this is structured data, it is preferable to read and write to BigQuery—recall that we loaded our full dataset into BigQuery in [Chapter 2](#). Now, we'd like to put the transformed (time-corrected) data there as well.

Fortunately, all this involves is changing the source or the sink. The rest of the pipeline stays the same. For example, in the previous section (see *04_streaming/transform/df05.py*), we read the *airports.csv.gz* as:

```
| 'airports:read' >> beam.io.ReadFromText('airports.csv.gz')
```

Now, in order to read the equivalent file from Cloud Storage, we change the corresponding code in *04_streaming/transform/df06.py* to be:

```
airports_filename = 'gs://{}//flights/airports/airports.csv.gz'.format(  
    bucket)  
...  
| 'airports:read' >> beam.io.ReadFromText(airports_filename)
```

Of course, we'll have to make sure to upload the file to Cloud Storage and make it readable by whoever is going to run this code. Having the data file be available in our GitHub repository was not going to scale anyway—Cloud Storage (or BigQuery) is the right place for data.

In *df05.py*, I had to read a local file that contained the JSON export of a smart part of the dataset and use a JSON parser to obtain a dict:

```
| 'flights:read' >> beam.io.ReadFromText('flights_sample.json')  
| 'flights:parse' >> beam.Map(lambda line: json.loads(line))
```

In *df06.py*, the corresponding code becomes simpler because the BigQuery reader returns a dict where the column names of the result set are the keys:

```
'flights:read' >> beam.io.ReadFromBigQuery(  
    query='SELECT * FROM dsongcp.flights WHERE rand() < 0.001',  
    use_standard_sql=True)
```

Of course when we run it for real, we'll change the query to remove the sampling (`rand() < 0.001`) so that we can process the entire dataset.

Similarly, where before we wrote to a local file using:

```
| 'flights:tostring' >> beam.Map(lambda fields: json.dumps(fields))  
| 'flights:out' >> beam.io.textio.WriteToText('all_flights')
```

we'll change the code to write to Cloud Storage using:

```
flights_output = 'gs://{}//flights/tzcorr/all_flights'.format(bucket)  
...  
| 'flights:tostring' >> beam.Map(lambda fields: json.dumps(fields))  
| 'flights:gcsout' >> beam.io.textio.WriteToText(flights_output)
```

We can write the same data to a BigQuery table also:

```
flights_schema = \  
    'FL_DATE:date,UNIQUE_CARRIER:string,...CANCELLED:boolean'  
...  
| 'flights:bqout' >> beam.io.WriteToBigQuery(  
    'dsongcp.flights_tzcorr', schema=flights_schema,  
    write_disposition=beam.io.BigQueryDisposition.WRITE_TRUNCATE,  
    create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED  
)
```

Note that we need to provide a schema when writing to BigQuery, and specify what to do if the table already exists (we ask for the table to be truncated and contents replaced) and if the table doesn't already exist (we ask for the table to be created).

We can try running this code, but the pipeline will require a few extra parameters. So where we used to have:

```
with beam.Pipeline('DirectRunner') as pipeline:
```

we now need:

```
argv = [  
    '--project={0}'.format(project),  
    '--staging_location=gs://{}//flights/staging/'.format(bucket),  
    '--temp_location=gs://{}//flights/temp/'.format(bucket),  
    '--runner=DirectRunner'  
]  
with beam.Pipeline(argv=argv) as pipeline:
```

The reason is that when we read from BigQuery, we are providing a query:

```
'flights:read' >> beam.io.ReadFromBigQuery(  
    query='SELECT * FROM dsongcp.flights WHERE rand() < 0.001',  
    use_standard_sql=True)
```

So, we need to provide the project that needs to be billed. In addition, and this is an implementation detail, some temporary data needs to be staged and cached in Cloud Storage, and we need to provide the pipeline a place to store this temporary data—we will never be sure which operations will require staging or caching, so it's a good idea to always specify a scratch location in Cloud Storage for this purpose.

We can run *df06.py* and then check that new tables are created in BigQuery. So far, we have been running the code locally, either on your laptop or in Cloud Shell.

Next, let's look at how to run this in Cloud Dataflow, which is the GCP managed service for running Apache Beam pipelines.

Running the Pipeline in the Cloud

That last run took a few minutes on the local virtual machine (VM), and we were processing only a thousand lines! Let's change the code (see *df07.py*) to process all the rows in the BigQuery view:

```
'flights:read' >> beam.io.ReadFromBigQuery(  
    query='SELECT * FROM dsongcp.flights',  
    use_standard_sql=True)
```

Now that we have much more data, we need to distribute the work, and to do that, we will change the runner from `DirectRunner` (which runs locally) to `DataflowRunner` (which lobs the job off to the cloud and scales it out):

```
argv = [  
    '--project={0}'.format(project),  
    '--job_name=ch04timecorr',  
    '--save_main_session',  
    '--staging_location=gs://{0}/flights/staging/'.format(bucket),  
    '--temp_location=gs://{0}/flights/temp/'.format(bucket),  
    '--setup_file=./setup.py',  
    '--max_num_workers=8',  
    '--region={}'.format(region),  
    '--runner=DataflowRunner'  
]  
  
pipeline = beam.Pipeline(argv=argv)
```

Notice that there are a few extra parameters now:

- The job name provides the name by which this job will be listed in the GCP console. This is so that we can troubleshoot the job if necessary.
- We ask the Dataflow submission code to save our main session. This is needed whenever we have global variables in our Python program.

- The file `setup.py` should list the Python packages that we needed to install (`tzzonefinder` and `pytz`) as we went along—Cloud Dataflow will need to install these packages on the Compute Engine instances that it launches behind the scenes:

```
REQUIRED_PACKAGES = [
    'tzzonefinder',
    'pytz'
]
```

- By default, Dataflow autoscales the number of workers based on throughput—the more lines we have in our input data files, the more workers we need. This is called **Horizontal Autoscaling**. To turn off autoscaling, we can specify `--autoscaling_algorithm=None`, and to constrain it somewhat, we can specify the maximum number of workers.
- We specify the region in which the Dataflow pipeline needs to run.
- The runner is no longer `DirectRunner` (which runs locally). It is now `Dataflow Runner`.

Running the Python program submits the job to the cloud. Cloud Dataflow auto-scales each step of the pipeline based on throughput, and streams the events data into BigQuery (see [Figure 4-3](#)). You can monitor the running job on the Cloud Platform Console in the Cloud Dataflow section.

Even as the events data is being written out, we can query it by browsing to the BigQuery console and typing the following:

```
SELECT
  ORIGIN,
  DEP_TIME,
  DEST,
  ARR_TIME,
  ARR_DELAY,
  EVENT_TIME,
  EVENT_TYPE
FROM
  dsongcp.flights_simevents
WHERE
  (DEP_DELAY > 15 and ORIGIN = 'SEA') or
  (ARR_DELAY > 15 and DEST = 'SEA')
ORDER BY EVENT_TIME ASC
LIMIT
  5
```

This returns:

Row	ORI	DEP_TIME	DEST	ARR_TIME	ARR_DELAY	EVENT_TIME	EVENT_TYPE
GIN							
1	SEA	2015-01-01 08:21:00 UTC	IAD	null	null	2015-01-01 08:21:00 UTC	departed
2	SEA	2015-01-01 08:21:00 UTC	IAD	null	null	2015-01-01 08:38:00 UTC	wheelsoff
3	SEA	2015-01-01 08:21:00 UTC	IAD	2015-01-01 12:48:00 UTC	22.0	2015-01-01 12:48:00 UTC	arrived
4	KOA	2015-01-01 10:11:00 UTC	SEA	2015-01-01 15:45:00 UTC	40.0	2015-01-01 15:45:00 UTC	arrived
5	SEA	2015-01-01 16:43:00 UTC	PSP	null	null	2015-01-01 16:43:00 UTC	departed

As expected, we see three events for the SEA-IAD flight, one at departure, the next at wheelsoff, and the third at arrival. The arrival delay is known only at arrival.

BigQuery is a columnar database, so a query that selects all fields:

```
SELECT
  *
FROM
  dsongcp.flights_simevents
ORDER BY EVENT_TIME ASC
```

will be inefficient. However, we do need all of the event data in order to send out event notifications. Therefore, we traded off storage for speed by adding an extra column called EVENT_DATA to our BigQuery table and populated it in our Dataflow pipeline as follows:

```
def create_event_row(fields):
    feadict = dict(fields) # copy
    feadict['EVENT_DATA'] = json.dumps(fields)
    return feadict
```

Then, our query to pull the events could simply be as follows:

```
SELECT
  EVENT_TYPE,
  EVENT_TIME,
  EVENT_DATA
FROM
  dsongcp.flights_simevents
WHERE
  EVENT_TIME >= TIMESTAMP('2015-05-01 00:00:00 UTC')
  AND EVENT_TIME < TIMESTAMP('2015-05-03 00:00:00 UTC')
ORDER BY
  EVENT_TIME ASC
```

```
LIMIT  
2
```

The result looks like this:

Row	EVENT_TYPE	EVENT_TIME	EVENT_DATA
1	wheelsoff	2015-05-01 00:00:00 UTC	{"FL_DATE": "2015-04-30", "UNIQUE_CARRIER": "DL", "ORIGIN_AIRPORT_SEQ_ID": "1295302", "ORIGIN": "LGA", "DEST_AIRPORT_SEQ_ID": "1330303", "DEST": "MIA", "CRS_DEP_TIME": "2015-04-30T23:29:00", "DEP_TIME": "2015-04-30T23:35:00", "DEP_DELAY": 6.0, "TAXI_OUT": 25.0, "WHEELS_OFF": "2015-05-01T00:00:00", "CRS_ARR_TIME": "2015-05-01T02:53:00", "CANCELLED": false, "DIVERTED": false, "DEP_AIRPORT_TZOFFSET": -14400.0, "ARR_AIRPORT_TZOFFSET": -14400.0, "EVENT_TYPE": "wheelsoff", "EVENT_TIME": "2015-05-01T00:00:00"}
2	departed	2015-05-01 00:00:00 UTC	{"FL_DATE": "2015-04-30", "UNIQUE_CARRIER": "DL", "ORIGIN_AIRPORT_SEQ_ID": "1295302", "ORIGIN": "LGA", "DEST_AIRPORT_SEQ_ID": "1320402", "DEST": "MCO", "CRS_DEP_TIME": "2015-04-30T23:55:00", "DEP_TIME": "2015-05-01T00:00:00", "DEP_DELAY": 5.0, "CRS_ARR_TIME": "2015-05-01T02:45:00", "CANCELLED": false, "DIVERTED": false, "DEP_AIRPORT_TZOFFSET": -14400.0, "ARR_AIRPORT_TZOFFSET": -14400.0, "EVENT_TYPE": "departed", "EVENT_TIME": "2015-05-01T00:00:00"}

This table will serve as the source of our events; it is from such a query that we will simulate streaming flight data.

Publishing an Event Stream to Cloud Pub/Sub

Now that we have the source events from the raw flight data, we are ready to simulate the stream. Streaming data in Google Cloud Platform is typically published to Cloud Pub/Sub, a serverless real-time messaging service. Cloud Pub/Sub provides reliable delivery and can scale to more than a million messages per second. Unless you are using [Cloud Pub/Sub Lite](#) (which is a single-zone service that is built for low-cost operation), Pub/Sub stores copies of messages in multiple zones to provide “at least once” guaranteed delivery to subscribers, and there can be many simultaneous subscribers.

Our simulator will read from the events table in BigQuery (populated in the previous section) and publish messages to Cloud Pub/Sub. Essentially, we will walk through the flight event records, getting the notification time from each, and simulate publishing those events as they happen.

Speed-Up Factor

However, we'll also use a mapping between the event notification time (arrival or departure time based on event) and the current system time. Why? Because it is inefficient to always simulate the flight events at real-time speeds. Instead, we might want to run through a day of flight data in an hour (as long as the code that processes these events can handle the increased data rate). At other times, we may be running our event-processing code in a debugging environment that is slower and so we might want to slow down the simulation. I will refer to this ratio between the actual time and simulation time as the *speed-up factor*—the speed-up factor will be greater than 1 if we want the simulation to be faster than real time and less than 1 if we want it to be slower than real time.

Based on the speed-up factor, we'll have to do a linear transformation of the event time to system time. If the speed-up factor is 1, a 60-minute difference between the start of the simulation in event time and the current record's timestamp should be encountered 60 minutes after the start of the simulation. If the speed-up factor is 60, a 60-minute difference in event time translates to a 1-minute difference in system time, and so the record should be published a minute later. If the event time clock is ahead of the system clock, we sleep for the necessary amount of time so as to allow the simulation to catch up.

The simulation consists of four steps (see also [Figure 4-5](#)):

- Run the query to get the set of flight event records to publish.
- Iterate through the query results.
- Accumulate events to publish as a batch.
- Publish accumulated events and sleep as necessary.

Even though this is an ETL pipeline, the need to process records in strict sequential order and sleep in between makes this ETL pipeline a poor fit for Cloud Dataflow. Instead, we'll implement this as a pure Python program. The problem with this choice is that the simulation code is not fault tolerant—if the simulation fails, it will not automatically restart and definitely will not start from the last successfully notified event.

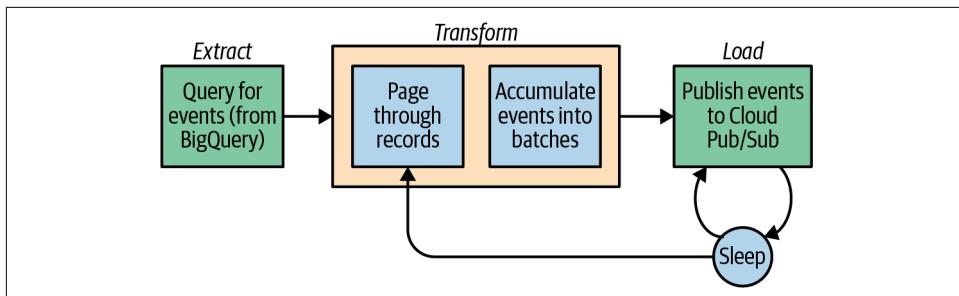


Figure 4-5. The four steps of simulation.

The simulation code that we are writing is only for quick experimentation with streaming data. Hence, I will not take the extra effort needed to make it fault-tolerant. If we had to do so, we could make the simulation fault-tolerant by starting from a BigQuery query that is bounded in terms of a time range with the start of that time range automatically inferred from the last-notified record in Cloud Pub/Sub. Then, we could launch the simulation script from a Docker container and use Cloud Run or Google Kubernetes Engine to automatically restart the simulation if the simulation code fails.

Get Records to Publish

The BigQuery query is parameterized by the start and end time of the simulation and can be invoked through the Google Cloud API for Python (see `04_streaming/simulate/simulate.py` in the GitHub repository):

```

bqclient = bq.Client(args.project)
querystr = """
SELECT
    EVENT_TYPE,
    EVENT_TIME AS NOTIFY_TIME,
    EVENT_DATA
FROM
    dsongcp.flights_simevents
WHERE
    EVENT_TIME >= TIMESTAMP('{}')
    AND EVENT_TIME < TIMESTAMP('{}')
ORDER BY
    EVENT_TIME ASC
"""
rows = bqclient.query(querystr.format(args.startTime,
                                         args.endTime))

```

This, however, is a bad idea. Do you see why?

It's because we are getting the start time and end time from the command line of the simulation script and directly passing it into BigQuery. This is called *SQL injection*,

which can lead to security problems.¹⁷ A better approach is to use *parameterized queries*—the BigQuery query contains the parameters marked as `@startTime`, etc., and the Python query function takes the definitions via the job configuration parameter:

```
bqclient = bq.Client(args.project)
querystr = """
SELECT
    EVENT_TYPE,
    EVENT_TIME AS NOTIFY_TIME,
    EVENT_DATA
FROM
    dsongcp.flights_simevents
WHERE
    EVENT_TIME >= @startTime
    AND EVENT_TIME < @endTime
ORDER BY
    EVENT_TIME ASC
"""

job_config = bq.QueryJobConfig(
    query_parameters=[
        bq.ScalarQueryParameter("startTime", "TIMESTAMP", args.startTime),
        bq.ScalarQueryParameter("endTime", "TIMESTAMP", args.endTime),
    ]
)
rows = bqclient.query(querystr, job_config=job_config)
```

The query function returns an object (called `rows` in the preceding snippet) that we can iterate through:

```
for row in rows:
    # do something
```

What do we need to do for each of the rows? We'll need to iterate through the records, build a batch of events, and publish each batch. Let's see how each of these steps is done.

How Many Topics?

As we walk through the query results, we need to publish events to Cloud Pub/Sub. We have three choices in terms of the architecture:

- We could publish all the events to a single topic. However, this can be wasteful of network bandwidth if we have a subscriber that is interested only in the `wheels off` event. Such a subscriber will have to parse the incoming event, decode the `EVENT_TYPE` file in the JSON, and discard events in which they are not interested.

¹⁷ It opens a door to someone passing in queries that could, for example, delete a table. This [XKCD cartoon](#) is famous for highlighting the issue.

- We could publish all the events to a single topic, but add attributes to each message. For example, to publish an event with two attributes `event_type` and `carrier`, we'd do:

```
publisher.publish(topic, event_data,
    event_type='departed', carrier='AA')
```

Then, the subscriber could ask for **server-side filtering** based on an attribute or combination of attributes when creating the subscription:

```
subscriber.create_subscription(request={..., "filter":
    "attributes.carrier='AS' AND attributes.event_type='arrived'"})
```

- Create a separate topic per event type (i.e., an `arrived` topic, a `departed` topic, and a `wheelsoff` topic).

Option 1 is the simplest, and should be your default choice unless you will have many subscribers that are interested only in subsets of the event stream. If you will have subscribers that are interested in subsets of the event stream, choose between Options 2 and 3.

Option 2 adds software complexity. Option 3 adds infrastructure complexity. I suggest choosing Option 3 when you have only one attribute, and that attribute has only a handful of options. This limits infrastructure complexity while keeping publisher and subscriber code simple. Choose Option 2 when you have many attributes each with many possible values because Option 3 in such a scenario will lead to an explosion in the number of topics.

Iterating Through Records

We will choose to have a separate topic per event type (i.e., an `arrived` topic, a `departed` topic, and a `wheelsoff` topic), so we create three topics:¹⁸

```
for event_type in ['wheelsoff', 'arrived', 'departed']:
    topics[event_type] = publisher.topic_path(args.project, event_type)
    try:
        publisher.get_topic(topic=topics[event_type])
        logging.info("Already exists: {}".format(topics[event_type]))
    except:
        logging.info("Creating {}".format(topics[event_type]))
        publisher.create_topic(name=topics[event_type])
```

After creating the topics, we call the `notify()` method passing along the rows read from BigQuery:

```
# notify about each row in the dataset
programStartTime = datetime.datetime.utcnow()
```

¹⁸ See `04_streaming/simulate/simulate.py` in the GitHub repository.

```

simStartTime = datetime.datetime.strptime(args.startTime,
                                         TIME_FORMAT).replace(tzinfo=pytz.UTC)
notify(publisher, topics, rows, simStartTime,
       programStartTime, args.speedFactor)

```

Building a Batch of Events

The `notify()` method consists of accumulating the rows into batches, publishing a batch, and sleeping until it is time to publish the next batch:

```

def notify(publisher, topics, rows, simStartTime, programStart, speedFactor):
    # sleep computation
    def compute_sleep_secs(notify_time):
        time_elapsed = (datetime.datetime.utcnow() -
                        programStart).seconds
        sim_time_elapsed = (notify_time - simStartTime).seconds / speedFactor
        to_sleep_secs = sim_time_elapsed - time_elapsed
        return to_sleep_secs

    tonotify = {}
    for key in topics:
        tonotify[key] = list()

    for row in rows:
        event, notify_time, event_data = row

        # how much time should we sleep?
        if compute_sleep_secs(notify_time) > 1:
            # notify the accumulated tonotify
            publish(publisher, topics, tonotify, notify_time)
            for key in topics:
                tonotify[key] = list()

        # recompute sleep, since notification takes a while
        to_sleep_secs = compute_sleep_secs(notify_time)
        if to_sleep_secs > 0:
            logging.info('Sleeping {} seconds'.format(to_sleep_secs))
            time.sleep(to_sleep_secs)

        tonotify[event].append(event_data)
    # left-over records; notify again
    publish(publisher, topics, tonotify, notify_time)

```

There are a few points to be made here. First is that we work completely in UTC so that the time difference computations make sense. Second, we always compute whether to sleep by looking at the time difference since the start of the simulation. If we simply keep moving a pointer forward, errors in time will accumulate. Finally, note that we check whether the sleep time is more than a second the first time, so as to give records time to accumulate. If, when you run the program, you do not see any sleep, your speed-up factor is too high for the capability of the machine running the simulation code and the network between that machine and Google Cloud Platform.

Slow down the simulation, get a larger machine, or run it behind the Google firewall (such as in Cloud Shell or on a Compute Engine instance).

Publishing a Batch of Events

The `notify()` method that we saw in the previous code example has accumulated the events in between sleep calls. Even though it appears that we are publishing one event at a time, the publisher actually maintains a separate batch for each topic:

```
def publish(publisher, topics, allevents):
    for key in topics: # 'departed', 'arrived', etc.
        topic = topics[key]
        events = allevents[key]
        logging.info('Publishing {} {} events'.format(len(events), key))
        for event_data in events:
            publisher.publish(topic, event_data.encode())
```

Note that Cloud Pub/Sub does not guarantee the order in which messages will be delivered, especially if the subscriber lets a huge backlog build up. Out-of-order messages will happen, and downstream subscribers will need to deal with them. Cloud Pub/Sub guarantees “at least once” delivery and will resend the message if the subscriber does not acknowledge a message in time. I will use Cloud Dataflow to ingest from Cloud Pub/Sub, and Cloud Dataflow deals with both these issues (out-of-order and duplication) transparently.

We can try out the simulation by typing the following:

```
python3 simulate.py --startTime '2015-05-01 00:00:00 UTC' \
--endTime '2015-05-04 00:00:00 UTC' --speedFactor=60
```

This will simulate three days of flight data (the end time is exclusive) at 60 times real-time speed and stream the events into three topics on Cloud Pub/Sub.¹⁹ Because the simulation starts off from a BigQuery query, it is quite straightforward to limit the simulated events to just a single airport or to airports within a latitude/longitude bounding box.

In this section, we looked at how to produce an event stream and publish those events in real time. Throughout this book, we can use this simulator and these topics for experimenting with how to consume streaming data and carry out real-time analytics.

¹⁹ At 60 times real-time speed, the 3 days of flight data will take over an hour to complete. Hopefully, that's enough time to complete the rest of the chapter. If not, just restart the simulator. If you get done early, hit Ctrl-C to stop the simulator.

Real-Time Stream Processing

Now that we have a source of streaming data that includes location information, let's look at how to build a real-time dashboard. [Figure 4-6](#) presents the reference architecture for many solutions on Google Cloud Platform.²⁰

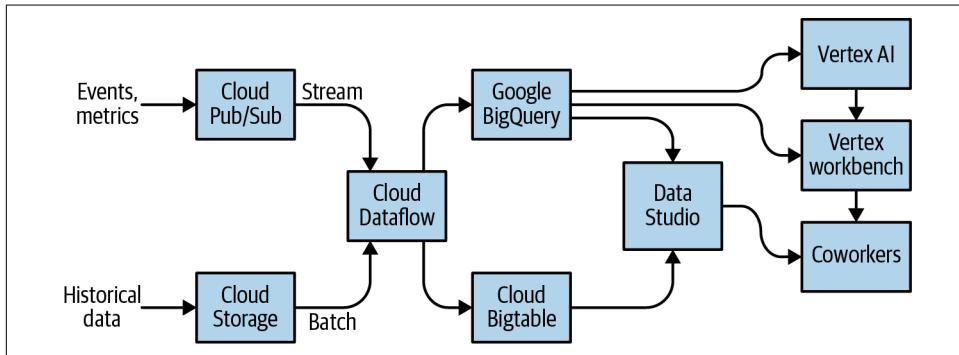


Figure 4-6. Reference architecture for data processing on Google Cloud Platform.

In the previous section, we set up a real-time stream of events into Cloud Pub/Sub that we can aggregate in Cloud Dataflow and write to BigQuery. Data Studio can connect to BigQuery and provide a real-time, interactive dashboard. Let's get started.

Streaming in Dataflow

When we carried out the time correction of the raw flight data, we were working off a complete BigQuery flights table, processing them in Cloud Dataflow, and writing the events table into BigQuery. Processing a finite, bounded input dataset is called *batch processing*.

Here, though, we need to process events in Cloud Pub/Sub that are streaming in. The dataset is *unbounded*. Processing an unbounded set of data is called *stream processing*. Fortunately, the code to do stream processing in Apache Beam is identical to the code to do batch processing.

We could simply receive the events from Cloud Pub/Sub similarly to how we read data from a CSV file:²¹

```
topic_name = "projects/{}/topics/arrived".format(project)
events = (pipeline
    | 'read' >> beam.io.ReadFromPubSub(topic=topic_name)
```

²⁰ For an example, see the [reference architecture](#) to analyze games on mobile devices.

²¹ See [04_streaming/realtime/avg01.py](#) in the GitHub repository.

```
| 'parse' >> beam.Map(lambda s: json.loads(s))  
)
```

The only change we have to do is to turn on the streaming flag in the Dataflow options:

```
argv = [  
    ...  
    '--streaming',  
]
```

We can stream the read-in events to BigQuery using code similar to what we used in batch processing:

```
schema = 'FL_DATE:date,...,EVENT_TYPE:string,EVENT_TIME:timestamp'  
(events  
    | 'bqout' >> beam.io.WriteToBigQuery(  
        'dsongcp.streaming_events', schema=schema,  
        create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED  
    )  
)
```

In the preceding code, we subscribe to a topic in Cloud Pub/Sub and begin reading from it. As each message streams in, we parse the message, convert it to a TableRow in BigQuery, and then write it out. Indeed, if this is all we need, we can simply use the Google-provided Dataflow template that goes from [Pub/Sub to BigQuery](#).

But let's say that we want to read both the arrived events and the departed events and write them to the same BigQuery table. We can do that quite simply in Beam:

```
events = {}  
for event_name in ['arrived', 'departed']:  
    topic_name = "projects/{}/topics/{}".format(project, event_name)  
    events[event_name] = (pipeline  
        | 'read:{}'.format(event_name) >>  
            beam.io.ReadFromPubSub(topic=topic_name)  
        | 'parse:{}'.format(event_name) >> beam.Map(  
            lambda s: json.loads(s))  
)  
  
all_events = (events['arrived'], events['departed']) | beam.Flatten()
```

Flattening the two sets of events concatenates them into a single collection. We then write out `all_events` to BigQuery.

To try this code out, we need to run the simulator we wrote in the previous section so that the simulator can publish events to the Pub/Sub topics. To start the simulation, start the Python simulator that we developed in the previous section:

```
python simulate.py --startTime '2015-05-01 00:00:00 UTC'  
--endTime '2015-05-04 00:00:00 UTC' --speedFactor 30
```

The simulator will send events from May 1, 2015, to May 3, 2015, at 30 times real-time speed, so that an hour of data is sent to Cloud Pub/Sub in two minutes. You can do this from Cloud Shell or from your local laptop. (If necessary, run `install_packages.sh` to install the necessary Python packages and `gcloud auth application-default login` to give the application the necessary credentials to execute queries.)

In another terminal, start `avg01.py` to read the stream of events and write them out to BigQuery. We can then query the dataset in BigQuery even as the events are streaming in. The BigQuery UI may not even show this streaming table yet, but it can be queried:

```
SELECT * FROM dsongcp.streaming_events  
ORDER BY EVENT_TIME DESC  
LIMIT 5
```

Windowing a Pipeline

Although we could do just a straight data transfer, I'd like to do more. When I populate a real-time dashboard of flight delays, I'd like the information to be aggregated over a reasonable interval—for example, I want a moving average of flight delays and the total number of flights over the past 60 minutes at every airport. So, rather than simply take the input received from Cloud Pub/Sub and stream it out to BigQuery, I'd like to carry out time-windowed analytics on the data as I'm receiving it and write those analytics to BigQuery.²² Cloud Dataflow can help us do this.

While we may be averaging over 60 minutes, how often should we compute this 60-minute average? It might be advantageous, for example, to use a sliding window and compute this 60-minute average every five minutes.

Streaming Aggregation

The key difference between batch aggregation and streaming aggregation is the unbounded nature of the data in stream processing. What does an operation like “max” mean when the data is unbounded? After all, whatever our maximum at this point in time, a large value could come along in the stream at a later point.

A key concept when aggregating streaming data is that of a window that becomes the scope for all aggregations. Here, we apply a time-based sliding window on the pipeline. From now on, all grouping, aggregation, and so on is within that time window and there is a separate maximum, average, etc. in each time window:

²² If you wanted to write the raw data that is received to BigQuery, you could do that, too, of course—that is what is shown in the previous code snippet. In this section, I assume that we need only the aggregate statistics over the past hour.

```

stats = (all_events
    | 'byairport' >> beam.Map(by_airport)
    | 'window' >> beam.WindowInto(
        beam.window.SlidingWindows(60 * 60, 5 * 60))
    | 'group' >> beam.GroupByKey()
    | 'stats' >> beam.Map(lambda x: compute_stats(x[0], x[1]))
)

```

Let's walk through the preceding code snippet carefully.

The first thing we do is to take all the events and apply the `by_airport` transformation to the events:

```
| 'byairport' >> beam.Map(by_airport)
```

What this does is to pull out the origin airport for departed events and destination airport for arrival events:

```

def by_airport(event):
    if event['EVENT_TYPE'] == 'departed':
        return event['ORIGIN'], event
    else:
        return event['DEST'], event

```

Next, we apply a sliding window to the event stream. The window is of 60 minutes duration, applied every 5 minutes:

```
| 'window' >> beam.WindowInto(
    beam.window.SlidingWindows(60 * 60, 5 * 60))
```

Then, we apply a `GroupByKey`:

```
| 'group' >> beam.GroupByKey()
```

What's the key?

In the `by_airport` function mentioned previously, we made the airport the key and the entire event object the value. So, the `GroupByKey` groups events by airport.

But the `GroupByKey` is not just by airport. Because we have already applied a sliding window, there is a separate group created for each time window. So, each group now consists of 60 minutes of flight events that arrived or departed at a specific airport.

It is on these groups that we call the `compute_stats` function in the last `Map` of the snippet:

```
| 'stats' >> beam.Map(lambda x: compute_stats(x[0], x[1]))
```

The `compute_stats` function takes the airport and list of events at that airport, and then computes some statistics:

```

def compute_stats(airport, events):
    arrived = [event['ARR_DELAY'] for event in events
               if event['EVENT_TYPE'] == 'arrived']

```

```

avg_arr_delay = float(np.mean(arrived))
    if len(arrived) > 0 else None

departed = [event['DEP_DELAY'] for event in events
            if event['EVENT_TYPE'] == 'departed']
avg_dep_delay = float(np.mean(departed))
    if len(departed) > 0 else None

num_flights = len(events)
start_time = min([event['EVENT_TIME'] for event in events])
latest_time = max([event['EVENT_TIME'] for event in events])

return {
    'AIRPORT': airport,
    'AVG_ARR_DELAY': avg_arr_delay,
    'AVG_DEP_DELAY': avg_dep_delay,
    'NUM_FLIGHTS': num_flights,
    'START_TIME': start_time,
    'END_TIME': latest_time
}

```

In the preceding code, we pull out the arrived events and compute the average arrival delay. Similarly, we compute the average departure delay on the departed events. We also compute the number of flights in the time window at this airport and return all these statistics.

The statistics are then written out to BigQuery using code that should look familiar by now:

```

stats_schema = ','.join(
    ['AIRPORT:string,AVG_ARR_DELAY:float,AVG_DEP_DELAY:float',
     'NUM_FLIGHTS:int64,START_TIME:timestamp,END_TIME:timestamp'])
(stats
| 'bqout' >> beam.io.WriteToBigQuery(
    'dsongcp.streaming_delays', schema=stats_schema,
    create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED
)
)

```

As with the previous section, we can start the simulator, and then start *avg02.py*. When I did this, the resulting aggregations were getting produced every 5 minutes, but within each 5 minute period, the events being notified about covered a 150 minute range (because the 30x simulation processes 150 minutes of data in 5 minutes).

The stream processing engine was applying the sliding windows based on the time on a wall clock. We, however, want it to apply the window based on the timestamp within the images.

How do we do that?

Using Event Timestamps

We have to add an attribute at the time we publish the events (in *simulate.py*):

```
publisher.publish(topic, event_data.encode(), EventTimeStamp=timestamp)
```

Then, in our Beam pipeline, when read from Pub/Sub, we should tell the pipeline to disregard the publish time in Pub/Sub and use this attribute of the message as the timestamp instead:

```
| 'read:{}'.format(event_name) >> beam.io.ReadFromPubSub(  
    topic=topic_name, timestamp_attribute='EventTimeStamp')
```

With this change, when I run the query:

```
SELECT * FROM dsongcp.streaming_delays  
WHERE AIRPORT = 'ATL'  
ORDER BY END_TIME DESC
```

I get rows approximately 5 minutes apart as expected:

Row	AIRPORT	AVG_ARR_DELAY	AVG_DEP_DELAY	NUM_FLIGHTS	START_TIME	END_TIME
1	ATL	35.72222222222222	13.666666666666666	48	2015-05-01 02:24:00 UTC	2015-05-01 03:21:00 UTC
2	ATL	35.25	8.717948717948717	59	2015-05-01 02:15:00 UTC	2015-05-01 03:12:00 UTC
3	ATL	38.666666666666664	9.882352941176471	52	2015-05-01 02:19:00 UTC	2015-05-01 03:12:00 UTC
4	ATL	38.473684210526315	5.916666666666667	55	2015-05-01 02:15:00 UTC	2015-05-01 03:08:00 UTC
5	ATL	35.111111111111114	5.53125	50	2015-05-01 02:15:00 UTC	2015-05-01 03:03:00 UTC

The reported times are not exactly 5 minutes apart because the reported times correspond to the earliest/latest flight in Atlanta within the time window. Note also that the length of the time window is approximately an hour.

It is likely, however, that Cloud Shell or your local laptop will struggle to keep up with the event stream. We need to be executing this pipeline in Dataflow in a serverless way.

Executing the Stream Processing

To run the Beam pipeline in Cloud Dataflow, all I have to do is to change the runner (see *avg03.py* in the GitHub repository):

```
argv = [
    '--project={0}'.format(project),
    '--job_name=ch04avgdelay',
    '--streaming',
    ...
    '--runner=DataflowRunner'
]
```

Before we start this pipeline, though, it is a good idea to delete the rows already written to the BigQuery table by *avg02.py* in the previous section. The easiest way to do this is to run the following SQL DML command to *truncate* the table:

```
TRUNCATE TABLE dsongcp.streaming_delays
```

Running *avg03.py* will launch off a Dataflow job. If you now browse to the Cloud Platform console, to the Cloud Dataflow section, you will see that a new streaming job has started and that the pipeline looks like that shown in [Figure 4-7](#).

The pipeline processes flight events as they stream into Pub/Sub, aggregates them into time windows, and streams the resulting statistics into BigQuery.

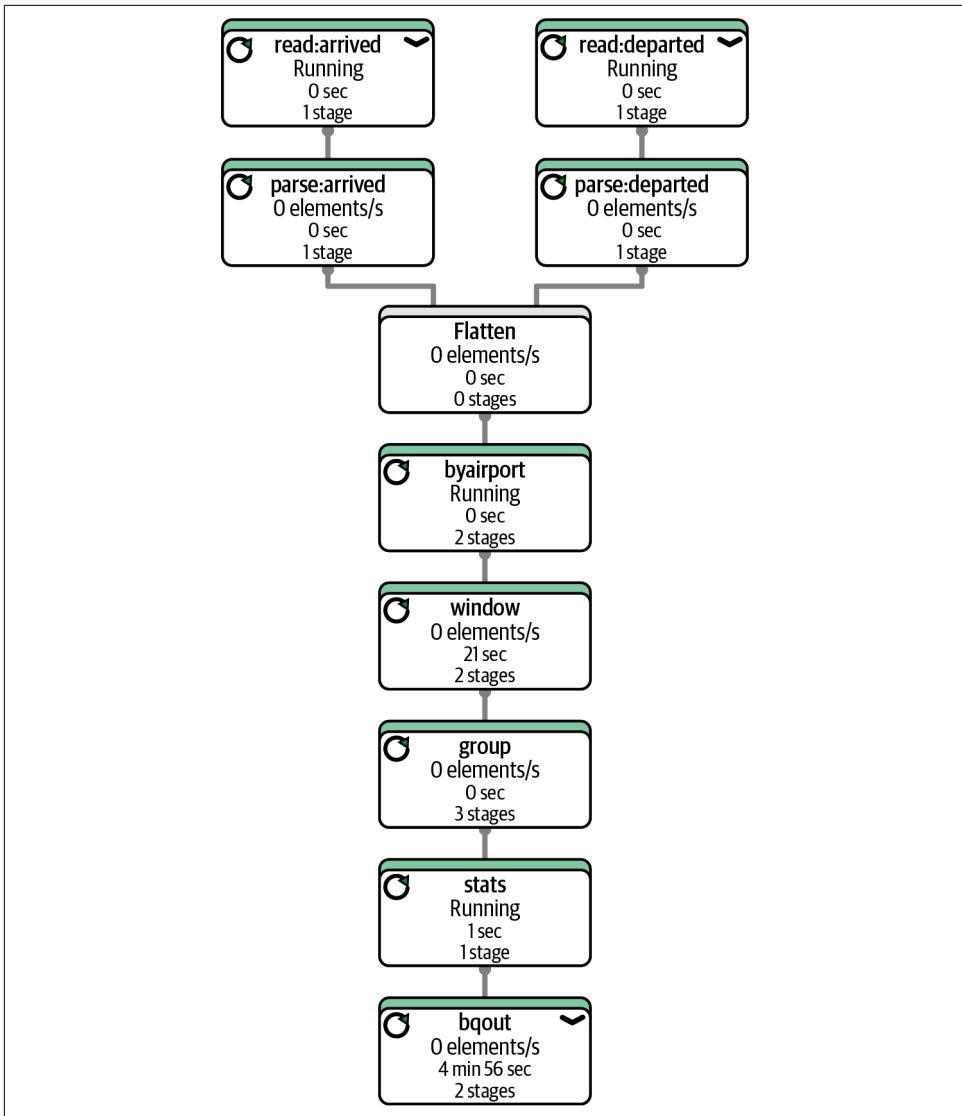


Figure 4-7. The streaming pipeline to compute delay statistics in real time at each airport.

Analyzing Streaming Data in BigQuery

Two minutes after the launch of your program,²³ the first set of data will make it into BigQuery. You can query for the statistics for a specific airport from the BigQuery console using the same query as before:

```
SELECT * FROM dsongcp.streaming_delays
WHERE AIRPORT = 'ATL'
ORDER BY END_TIME DESC
```

The cool thing is that we can do this querying even as the data is streaming! How would we get the latest data for all airports? We could get all the data for each airport, order it by time, and take the latest:

```
SELECT
    AIRPORT,
    ARRAY_AGG(
        STRUCT(AVG_ARR_DELAY, AVG_DEP_DELAY, NUM_FLIGHTS, END_TIME)
        ORDER BY END_TIME DESC LIMIT 1) AS a
FROM dsongcp.streaming_delays d
GROUP BY AIRPORT
```

The results look something like this:

Row	AIRPORT	a.AVG_ARR_DELAY	a.AVG_DEP_DELAY	a.NUM_FLIGHTS	a.END_TIME
1	BUR	-6.8	-5.66666666666667	8	2015-05-01 03:26:00 UTC
2	HNL	17.11111111111111	-3.7777777777777777	18	2015-05-01 03:46:00 UTC
3	CVG	-7.75	null	4	2015-05-01 03:48:00 UTC
4	PHL	5.636363636363637	16.5	13	2015-05-01 03:48:00 UTC
5	IND	40.6	null	5	2015-05-01 03:45:00 UTC

Queries like these on streaming data will be useful when we begin to build our dashboard. For example, the first query will allow us to build a time series chart of delays at a specific airport. The second query will allow us to build a map of average delays across the country.

²³ Recall that we are computing aggregates over 60 minutes every 5 minutes. Cloud Dataflow treats the first “full” window as happening 65 minutes into the simulation. Because we are simulating at 30 times speed, this is two minutes on your clock.

Real-Time Dashboard

Now that we have streaming data in BigQuery and a way to analyze it as it is streaming in, we can create a dashboard that shows departure and arrival delays in context. Two maps can help explain our contingency table-based model to end users: current arrival delays across the country and current departure delays across the country.

To pull the data to populate these charts, we need to add a BigQuery data source in Data Studio. Although Data Studio supports specifying the query directly in the user interface, it is much better to create a view in BigQuery and use that view as a data source in Data Studio. BigQuery views have a few advantages over queries that you type into Data Studio: they tend to be reusable across reports and visualization tools, there is only one place to change if an error is detected, and BigQuery views map better to access privileges (Cloud Identity and Access Management roles) based on the columns they need to access.

Here is the query that I used to create the view:

```
CREATE OR REPLACE VIEW dsongcp.airport_delays AS
WITH delays AS (
    SELECT d.*,
           a.LATITUDE, a.LONGITUDE
      FROM dsongcp.streaming_delays d
     JOIN dsongcp.airports a USING(AIRPORT)
     WHERE a.AIRPORT_IS_LATEST = 1
)
SELECT
    AIRPORT,
    CONCAT(LATITUDE, ',', LONGITUDE) AS LOCATION,
    ARRAY_AGG(
        STRUCT(AVG_ARR_DELAY, AVG_DEP_DELAY, NUM_FLIGHTS, END_TIME)
        ORDER BY END_TIME DESC LIMIT 1) AS a
   FROM delays
  GROUP BY AIRPORT, LONGITUDE, LATITUDE
```

This is slightly different from the second query in the previous section in that it also adds the location of the airport by joining against the airports table.

Having saved the view in BigQuery, we can create a data source for the view in Data Studio, just as we did in the previous chapter:

- Visit <https://datastudio.google.com>.
- Create a BigQuery data source, point it to the `airport_delays` view, and connect to it.
- Change the location field from Text to a Geo | Latitude, Longitude, then click Create Report.
- Add a Geo Chart to the report.

- Specify the location field as the geo dimension (see [Figure 4-8](#)).
- Specify average departure delay as the dimension and United States as the zoom level.
- Change the style so that the color bar includes all areas.
- Repeat for the arrival delay.

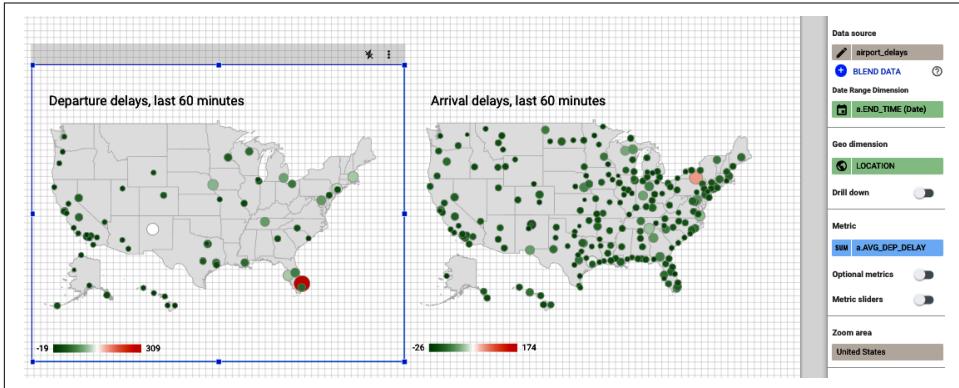


Figure 4-8. Dashboard of latest flight data from across the United States.

It is worth reflecting on what we did in this section. We processed streaming data in Cloud Dataflow, creating 60-minute moving averages that we streamed into BigQuery. We then created a view in BigQuery that would show the latest data for each airport, even as it was streaming in. We connected that to a dashboard in Data Studio. Every time the dashboard is refreshed, it pulls new data from the view, which in turn dynamically reflects the latest data in BigQuery.

Summary

In this chapter, we discussed how to build a real-time analysis pipeline to carry out streaming analytics and populate real-time dashboards. In this book, we are using a dataset that is not available in real time. Therefore, we simulated the creation of a real-time feed so that I could demonstrate how to build a streaming ingest pipeline. Building the simulation also gives us a handy test tool—no longer do we need to wait for an interesting event to happen. We can simply play back a recorded event!

In the process of building out the simulation, we realized that time handling in the original dataset was problematic. Therefore, we improved the handling of time in the original data and created a new dataset with UTC timestamps and local offsets. This is the dataset that we will use going forward.

We also looked at the reference architecture for handling streaming data in Google Cloud Platform. First, receive your data in Cloud Pub/Sub so that the messages can

be received asynchronously. Process the Cloud Pub/Sub messages in Cloud Dataflow, computing aggregations on the data as needed, and stream either the raw data or aggregate data (or both) to BigQuery. We worked with all three Google Cloud products (Cloud Pub/Sub, Cloud Dataflow, and BigQuery) using the Google Cloud client libraries in Python. However, in none of these cases did we ever need to create a virtual machine ourselves—these are all serverless and autoscaled offerings. We thus were able to concentrate on writing code, letting the platform manage the rest.

Suggested Resources

The Apache Beam website has [interactive coding exercises](#), called Katas, that provide an excellent hands-on way to learn streaming concepts and how to implement them using Beam.

Dataflow templates are prewritten Apache Beam pipelines that are handy for data migration. In the chapter, we mentioned the Dataflow template for ingesting data from Pub/Sub into BigQuery. Dataflow templates also exist for non-Google sources. For example, there is a Dataflow connector from SAP HANA to BigQuery, as described in the 2017 Google blog post [“Using Apache Beam and Cloud Dataflow to Integrate SAP HANA and BigQuery”](#) by Babu Prasad Elumalai and Mark Shalda. That particular connector is written in Java.

This tutorial walks you through the process of [creating your own Dataflow template](#). Any Dataflow pipeline can be made into a template for easy sharing and convenient launch.

In this 2021 article, “Processing Billions of Events in Real Time at Twitter,” Twitter engineers Lu Zhang and Chukwudiuto Malife describe how Twitter [processes 400 billion events](#) in real time using Dataflow.

Interactive Data Exploration with Vertex AI Workbench

In every major field of study, there is usually a towering figure who did much of the seminal work and blazed the trail for what that particular discipline would evolve into. Classical physics has Newton, relativity has Einstein, game theory has John Nash, and so on. When it comes to computational statistics (the field of study that develops computationally efficient methods for carrying out statistical operations), the towering figure is John W. Tukey. At Bell Labs, he collaborated with John von Neumann on early computer designs soon after World War II—famously, Tukey was responsible for coining the word *bit*. Later, at Princeton (where he founded its statistics department), Tukey collaborated with James Cooley to develop the fast Fourier transform, one of the first examples of using divide-and-conquer to address a formidable computational challenge.

While Tukey was responsible for many “hard science” mathematical and engineering innovations, some of his most enduring work is about the distinctly softer side of science. Unsatisfied that most of statistics overemphasized confirmatory data analysis (i.e., statistical hypothesis testing such as paired *t*-tests),¹ Tukey developed a variety of approaches to do what he termed *exploratory data analysis* (EDA)² and many practical statistical approximations. It was Tukey who developed the box plot, jack-knifing, range test, median-median regression, and so on and gave these eminently practical methods a solid mathematical grounding by motivating them in the context of simple conceptual models that are applicable to a wide variety of datasets. In this chapter, we follow Tukey’s approach of carrying out exploratory data analysis to identify

¹ See [this Wikipedia article](#) for an excellent overview of statistical hypotheses testing.

² John Tukey. *Exploratory Data Analysis*. Addison-Wesley, 1977.

important variables, discover underlying structure, develop parsimonious models, and use those models to identify unusual patterns and values.



All of the code snippets in this chapter are available in the folder [05_bqnotebook](#) of the book's GitHub repository. See the *README.md* file in that directory for instructions on how to do the steps described in this chapter.

Exploratory Data Analysis

Ever since Tukey introduced the world to the value of EDA in 1977, the traditional first step for any data scientist has been to analyze raw data by using a variety of graphical techniques. This is not a fixed set of methods or plots; rather, it's an approach meant to develop insight into a dataset and enable the development of robust statistical models. Specifically, as a data scientist, you should do the following:

- Test any underlying assumptions, such as that a particular value will always be present or will always lie within a certain range. For example, as discussed in [Chapter 2](#), the distribution of the distance between any pair of airports might help verify whether the distance is the same across the dataset or whether it reflects the actual flight distance.
- Use intuition and logic to identify important variables. Verify that these variables are, indeed, important as expected. For example, plotting the relationship between departure delay and arrival delay might validate assumptions about the recording of these variables.
- Discover underlying structures in the data (i.e., relationships between important variables and situations such as the data falling into specific statistical regimes). It might be useful to examine whether the season (summer versus winter) has an impact on how often flight delays can be made up.
- Develop a parsimonious model—a simple model with explanatory power, that you can use to hypothesize about what reasonable values in the data look like. If there is a simple relationship between departure delay and arrival delay, values of either delay that are far off the trendline might warrant further examination.
- Detect outliers, anomalies, and other inexplicable data values. This depends on having that parsimonious model. Thus, further examination of outliers from the simple trend between departure and arrival delays might lead to the discovery that such values off the trendline correspond to rerouted flights.
- Discover any potential overarching data quality problems such as the issues we found in [Chapter 3](#) with time being recorded without being UTC.

To carry out exploratory data analysis, it is necessary to load the data in a form that makes interactive analysis possible. In this chapter, we load data into Google BigQuery, explore the data in Vertex AI Workbench, carry out quality control based on what we discover about the dataset, build a new model, and evaluate the model to ensure that it is better than the model we built in [Chapter 4](#). As we go about loading the data and exploring it and move on to building models and evaluating them, we'll discuss a variety of considerations that come up, from security to pricing.

Both exploratory data analysis and the dashboard creation discussed in [Chapter 3](#) involve the creation of graphics. However, the steps differ in two ways—in terms of the purpose and in terms of the audience. The aim of dashboard creation is to crowd-source insight into the working of models from end users and is, therefore, primarily about presenting an explanation of the models to end users. In [Chapter 3](#), I recommended doing it very early in your development cycle, but that advice was more about Agile development and getting feedback early than about statistical rigor.³ The aim of EDA is for you, the data engineer, to develop insights about the data before you delve into developing sophisticated models. The audience for EDA is typically other members of your team and yourself, not end users. In some cases, especially if you uncover strange artifacts in the data, the audience could be the data engineering team that produces the dataset you are working with. For example, when we discovered the problem that the times were being reported in local time, with no UTC offsets, we could have relayed that information back to the US Bureau of Transportation Statistics (BTS).⁴ In any case, the assumption is that the audience for an EDA graphic is statistically sophisticated. Although you probably would not include a violin plot in a dashboard meant for end users,⁵ you would have no compunctions about using it in an EDA chart that is meant for data scientists.

Doing exploratory data analysis on large datasets poses a few challenges. To test that a particular value will always be present, for example, you would need to check every row of a tabular dataset, and if that dataset is many millions of rows, these tests can take hours. An interactive ability to quickly explore large datasets is indispensable. On Google Cloud Platform, BigQuery provides the ability to run Cloud SQL queries on unindexed datasets (i.e., your raw data) in a matter of seconds even if the datasets

³ James Shore. *The Art of Agile Development*, 2nd ed. With Diana Larsen, Gitte Klitgaard, and Shane Warden. O'Reilly Media, 2021.

⁴ I did confirm with the BTS via email that the times in the dataset were, indeed, in the local time zone of the corresponding airport. The BTS being a government agency that was making the data freely available, I didn't broach the matter of them producing a separate dataset with UTC timestamps. In a contractual relationship with a vendor, however, this is the type of change you might request as a result of EDA.

⁵ A violin plot is a way of visualizing probability density functions. See the [seaborn documentation](#) for examples of violin plots.

are in the terabyte scale. Therefore, in this chapter, we load the flight data into BigQuery.

Anscombe's Quartet

The statistician Francis Anscombe illustrated that graphs are essential to good statistical analysis using a very powerful example. All four of the datasets shown in [Figure 5-1](#) have the same mean, variance, linear fit, and correlation (to two decimal places) but are obviously quite different from one another.

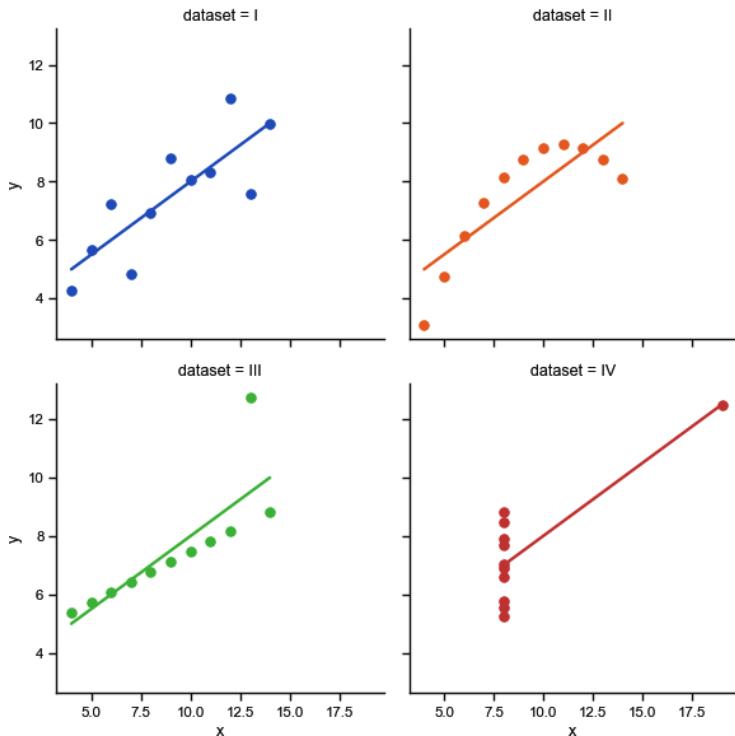


Figure 5-1. Anscombe's Quartet. Figure from “Graphs in Statistical Analysis” by Francis Anscombe in American Statistician 27, no. 1 (1973): 17–21 as recreated in the seaborn documentation.

Anscombe used the quartet to emphasize that summary statistics are not a substitute for graphic data—it's particularly important to graph outliers to develop a holistic understanding of the data.

Identifying outliers and underlying structure typically involves using univariate and bivariate plots.⁶ The graphs themselves can be created using Python plotting libraries—I use a mixture of *Matplotlib* and *seaborn* to do this. The networking overhead of moving data between storage and the graphics engine can become prohibitive if I carry out data exploration on my laptop—for exploration of large datasets to be interactive, we need to bring the analytics closer to the data. Therefore, we will want to use a cloud computer (not my laptop) to carry out graph generation. Generating graphics on a cloud computer poses a display challenge because the graphs are being generated on a Compute Engine instance that is *headless*—that is, it has no input or output devices. A Compute Engine instance has no keyboard, mouse, or monitor, and frequently has no graphics card.⁷ Such a machine is accessed purely through a network connection. Fortunately, desktop programs and interactive read-eval-print loops (REPLs) are no longer necessary to create visualizations. Instead, notebook servers such as *Jupyter* have become the standard way that data scientists create graphs and disseminate executable reports. On Google Cloud Platform, Vertex AI Workbench provides a fully managed way to run Jupyter notebooks that connect to Google Cloud Platform services.

Exploration with SQL

Let's start in the BigQuery console by exploring the time-corrected dataset that we created in BigQuery in [Chapter 4](#):

```
SELECT
    ORIGIN,
    AVG(DEP_DELAY) AS dep_delay,
    AVG(ARR_DELAY) AS arr_delay,
    COUNT(ARR_DELAY) AS num_flights
FROM
    dsongcp.flights_tzcorr
GROUP BY
    ORIGIN
```

⁶ A *univariate graph* is a graph of just one variable. For example, we might plot the histogram of arrival delays. A *bivariate graph* is a graph of two variables. For example, we might plot the median taxi-out time by airport.

⁷ Special [general purpose graphics processing unit \(GPU\) instances exist](#) that are used for high-performance computing applications, but for generating the graphs in this chapter, CPU instances are sufficient.

The result consists of 322 airports (the order you get might be different):

Row	ORIGIN	dep_delay	arr_delay	num_flights
1	OTZ	5.209103840682787	6.562952243125903	691
2	HPN	11.782807151007983	9.087898089171965	7850
3	SJU	9.8362379921783	2.506036485508635	26257
4	ANC	3.2497373643048966	-0.4801384732734849	17043
5	CVG	8.826792206581548	5.244408048666357	21370

Let's look at just the major airports, which we can define as airports that have on average more than 10 flights a day. To do this we can filter by airports that have a sufficient number of flights:

```
WITH all_airports AS (
    SELECT
        ORIGIN,
        AVG(DEP_DELAY) AS dep_delay,
        AVG(ARR_DELAY) AS arr_delay,
        COUNT(ARR_DELAY) AS num_flights
    FROM
        dsongcp.flights_tzcorr
    GROUP BY
        ORIGIN
)
SELECT * FROM all_airports WHERE num_flights > 3650
ORDER BY dep_delay DESC
```

We are thresholding the number of flights at 3,650 because there are 365 days in the dataset. The result, when I did it, was:

Row	ORIGIN	dep_delay	arr_delay	num_flights
1	ORD	13.305085522847	7.596119952650316	304120
2	EWR	13.182294215975096	3.9227994696288535	107849
3	BWI	12.893989460498512	6.768316724436742	92320
4	LGA	12.764120915158792	5.043357442317552	103281
5	IAD	12.23048266485387	4.505307971508886	36643

It makes sense that airports that serve major American cities experience the worst departure delays (ORD serves Chicago, EWR and LGA serve New York City, and BWI and IAD serve Washington, DC).

What if we restrict this analysis to January, reducing the number of flights threshold to 310 since there are 31 days in January?

```
WITH all_airports AS (
    SELECT
```

```

        ORIGIN,
        AVG(DEP_DELAY) AS dep_delay,
        AVG(ARR_DELAY) AS arr_delay,
        COUNT(ARR_DELAY) AS num_flights
    FROM
        dsongcp.flights_tzcorr
    WHERE EXTRACT(MONTH FROM FL_DATE) = 1
    GROUP BY
        ORIGIN
)
SELECT * FROM all_airports WHERE num_flights > 310
ORDER BY dep_delay DESC

```

Now, we get a somewhat stranger set of airports:

Row	ORIGIN	dep_delay	arr_delay	num_flights
1	ASE	20.86779661016949	16.988095238095244	588
2	ORD	19.96205128205124	17.016131923283723	22316
3	JAC	18.787172011661802	16.096209912536445	343
4	SBN	18.491891891891886	16.326975476839234	367
5	FAT	18.12554744525547	17.63823529411766	680

I don't recognize four of the five airports on this list, and I'm a rather frequent traveler. A short Google Search later, I learned that ASE is a ski resort (Aspen, Colorado) as is JAC (Jackson Hole, Wyoming). This makes sense—ski resorts are open only in winter, have to load up bulky baggage, and probably suffer more weather-related delays.

Using the average delay to characterize airports is not ideal, though. What if most flights to Aspen were actually on time but a few highly delayed flights (perhaps flights delayed by several hours) are skewing the average? I'd like to see a distribution function of the values of arrival and departure delays. BigQuery itself cannot help us with graphs—instead, we need to tie the BigQuery backend to a graphical, interactive exploration tool. Data scientists tend to use Jupyter Notebooks for EDA, so I'll use Vertex AI Workbench, which offers fully managed Jupyter Notebooks.

Reading a Query Explanation

Before I move on to Notebooks, though, we want to see if there are any red flags regarding query performance on our table in BigQuery. In the BigQuery console, there is a tab (next to Results) labeled “Execution details.” [Figure 5-2](#) shows the explanation of the January query.

The screenshot shows the BigQuery web interface. At the top, there are buttons for RUN, SAVE, SHARE, SCHEDULE, and MORE. Below that, a code editor displays a SQL query:

```

4   ORIGIN,
5   AVG(DEP_DELAY) AS dep_delay,
6   AVG(ARR_DELAY) AS arr_delay,
7   COUNT(ARR_DELAY) AS num_flights
8
9   FROM
10  dsongcp.flights_tzcorr
11 WHERE

```

A processing location indicator shows "Processing location: US". Below the code editor is a "Query results" section. It includes a "SAVE RESULT" button. The results are presented in a table with four columns: Row, ORIGIN, dep_delay, arr_delay, and num_flights. The data shows two rows:

Row	ORIGIN	dep_delay	arr_delay	num_flights
1	MEI	23.914691943127963	22.564285714285717	420
2	ASE	23.631280158955438	20.561003420752563	3508

Figure 5-2. The explanation of a query in BigQuery.

Our query has been executed in three stages. Expand each of the stages to see their details:

- The first stage (see Figure 5-3) pulls the origin, departure delay, arrival delay, and date for each flight and filters the result by looking at the month. Then, it groups them by origin, computes averages on each shard of data, and writes them to `_stage00_output`. `_stage00_output` is organized by the hash of the ORIGIN.
- The second stage (see Figure 5-4) reads the fields organized by ORIGIN, computes the average delays and count (but starting from the SHARD averages), and filters the result to ensure that the count is greater than 310. Note that the query has been optimized a bit—my WHERE clause was actually outside the WITH statement, but it has been moved here so as to minimize the amount of data written out to `_stage01_output`.
- The third stage (see Figure 5-5) simply sorts the rows by departure delay and writes to the output.

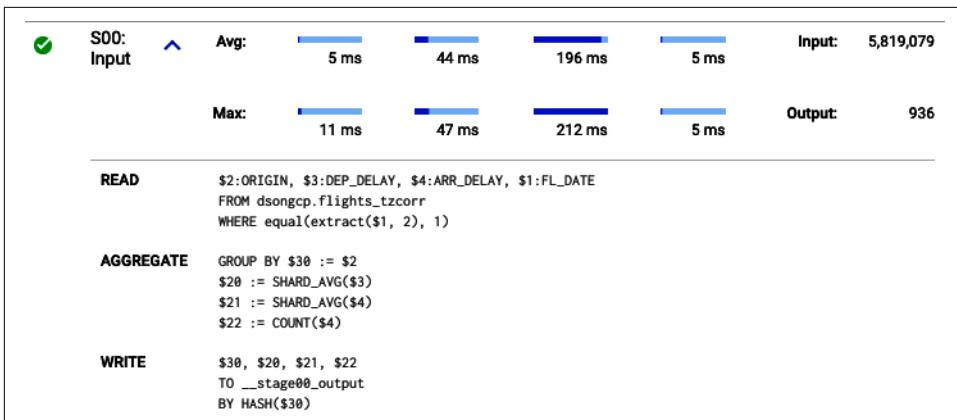


Figure 5-3. The first stage.

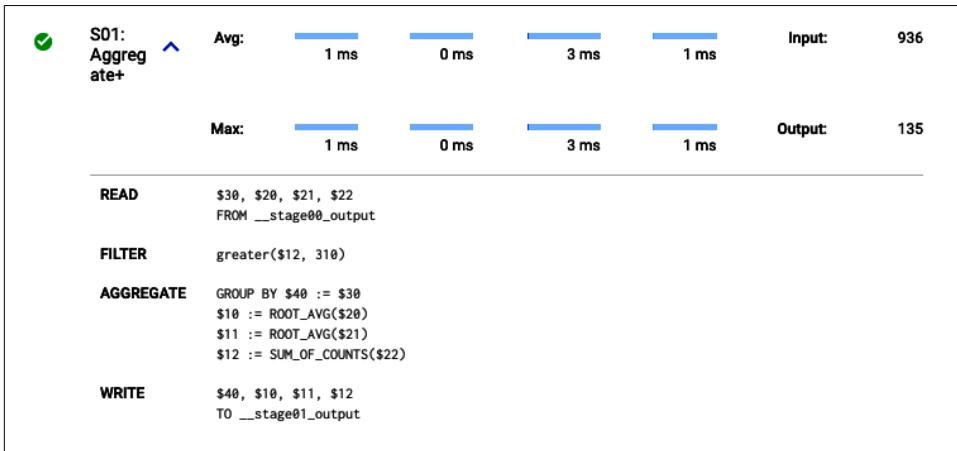


Figure 5-4. The second stage.

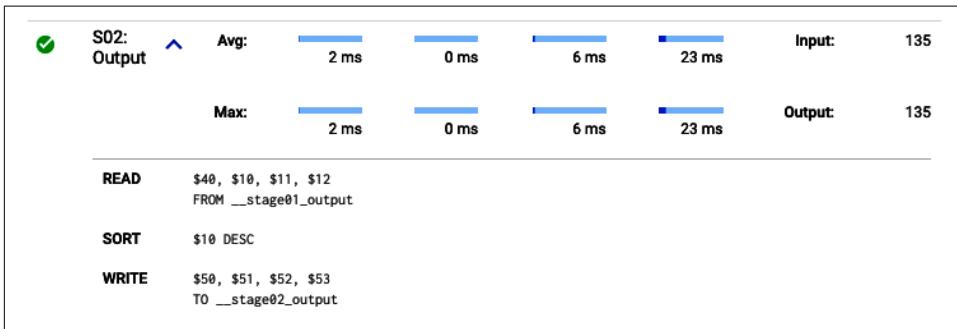


Figure 5-5. The third stage.

Based on the preceding second-stage optimization, is it possible to write the query itself in a better way? Yes, by using the `HAVING` keyword:

```
SELECT
    ORIGIN,
    AVG(DEP_DELAY) AS dep_delay,
    AVG(ARR_DELAY) AS arr_delay,
    COUNT(ARR_DELAY) AS num_flights
FROM
    dsongcp.flights_tzcorr
WHERE EXTRACT(MONTH FROM FL_DATE) = 1
GROUP BY
    ORIGIN
HAVING num_flights > 310
ORDER BY dep_delay DESC
```

In the rest of this chapter, I will use this form of the query that avoids the `WITH` statement. By using the `HAVING` keyword, we are not relying on the query optimizer to minimize the amount of data written to `__stage01_output`.

What do the times in the graphics mean? Each stage (see [Figure 5-6](#)) is broken into four steps: wait, read, compute, and write. The average and maximum time spent in each of these steps by the BigQuery workers is reported. So, in another example shown in [Figure 5-6](#), BigQuery spends an average of 353 milliseconds ($37 + 115 + 195 + 6$) in this stage. A worker could spend as much as 608 milliseconds ($49 + 308 + 242 + 9$) in it, though.

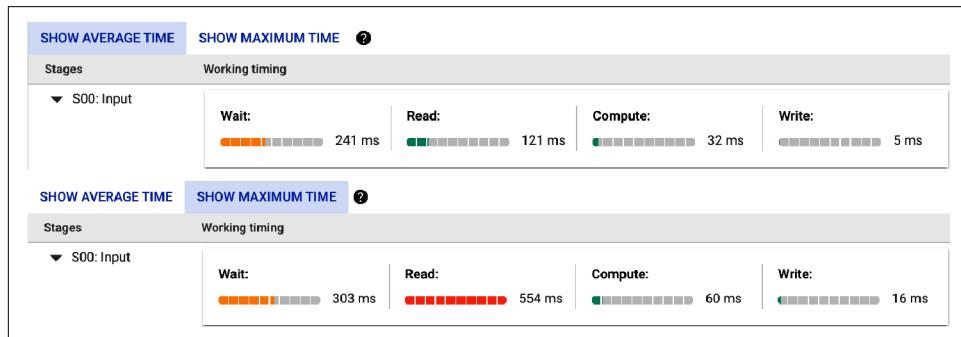


Figure 5-6. Average and maximum time spent by BigQuery workers by steps. Each part of each stage in the query explanation is depicted by a color bar that represents the fraction of time spent in that part.

The length of the bar is the time taken by the most time-consuming step in the stage—so the length of the bar corresponds to 308 ms and the color in each bar is the fraction of that time spent in this step. In other words, the bars are all normalized to the time taken by the longest step (wait, read, compute, or write). A large difference between the average and the maximum (as in the read step of [Figure 5-6](#)) indicates a

skew—there are some workers who are doing a lot more work than others. Sometimes, it is inherent to the query, but at other times, it might be possible to rework the storage or partitions to reduce such skew.⁸

The wait time is the time spent waiting for the necessary compute resources to become available—a very high value here indicates a job that could not be immediately scheduled on the cluster. A high wait time could occur if you are using the BigQuery flat-rate plan and someone else in your organization is already consuming all of the paid-for capacity. The solution would be to run your job at a different time, make your job smaller, or negotiate with the group that is using the available resources. The read step reads the data required at this stage from the table or from the output of the previous stage. A high value of read time indicates that you might consider reworking the query so that most of the data is read in the initial stages. The compute step carries out the computation required—if you run into high values here, consider whether you can carry out some of the operations in postprocessing or if you could omit the use of user-defined functions (UDFs).⁹ The write step writes to temporary storage or to the response and is mainly a function of the amount of data being written out in each stage—optimization of this step typically involves moving filtering options to occur in the innermost query (or earliest stage), although as we saw earlier, the BigQuery optimizer can do some of this automatically.

For all three stages in our query, the read step is what takes the most amount of time, indicating that our query is I/O bound and that the basic cost of reading the data is what dominates the query. It is clear from the numbers in the input column (6 million to 936 to 135) that we are already doing quite well at funneling the data through and processing most of the data in earlier stages. We also noticed from the explanation that BigQuery has already optimized things by moving the filtering step to the earliest possible stage—there is no way to move it any earlier because it is not possible to filter on the number of flights until that value is computed. On the other hand, if this is a frequent sort of filter, it might be helpful to add a table indicating the traffic at each airport and join with this table instead of computing the aggregate each time. It might also be possible to achieve an approximation to this by adding a column indicating some characteristic (such as the population) of the metropolitan area that each airport serves. For now, without any idea of the kinds of airports that the typical user of this dataset will be interested in, there is little to be done. We are determined to process all the data, and processing all the data requires time spent reading that data.

⁸ Reducing the skew is not just about reducing the time taken by the workers who have a lot to do. You should also see if you can rework the query to combine the work carried out by the underworked workers, so that you have fewer workers overall.

⁹ BigQuery supports UDFs in JavaScript, but the excessive use of UDFs can slow down your query, and certain UDFs can be [high-compute queries](#).

If we don't need statistics from all the data, we could consider sampling the data and computing our statistics on that sample instead.

Exploratory Data Analysis in Vertex AI Workbench

Data scientists have moved en masse to using notebooks because notebooks greatly streamline the workflow of developing, visualizing, collaborating, and publishing in science.¹⁰ The contrast between the user experience of a Jupyter Notebook and the way exploratory data analysis was carried out a few years ago is stark. Take, for example, [Figure 5-7](#), which appears in one of my papers about [different ways to track storms in weather radar images](#). Just this single graphic required wrangling multiple languages (C++, R, Java), concepts (distributed programming, statistics), data formats (CSV, PNG, LaTeX, PDF), and collaboration mechanisms (FTP, email)!¹¹ Today, I'd do it all in a Jupyter Notebook with the big data analysis carried out in BigQuery or Dataflow.

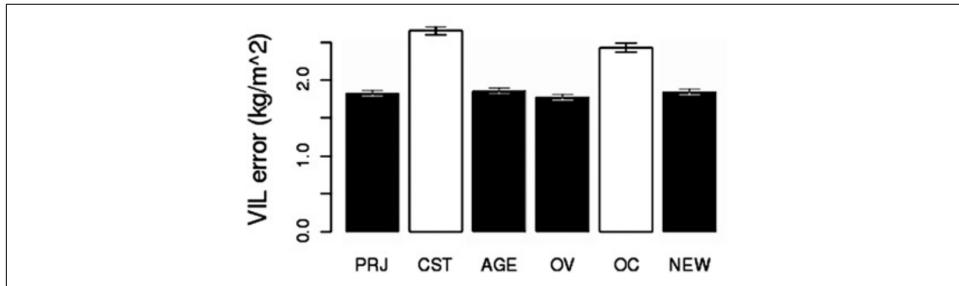


Figure 5-7. Graph created using a complex workflow.

¹⁰ Jupyter Notebooks are a realization of the *literate programming* concept envisioned by Donald Knuth in 1984. By allowing data scientists to interweave statistical and business logic in code and the output of that code within a literate statistical programming paradigm, notebooks foster replicability and reuse.

¹¹ If you are interested in the gory details: [Figure 5-7](#) was created by running the methods in question (PRJ, CST, etc.) on a large dataset of weather radar imagery and computing various evaluation metrics (VIL error in the graphic). For performance reasons, this was done in C++. The metric for each pair of images was written out to a text file (a different text file for each method), and it is the aggregates of those metrics that are reported in [Figure 5-7](#). The text files had to be wrangled from the cluster of machines on which they were written out, combined by key (the method used to track storms), and then aggregated. This code, essentially a MapReduce operation, was written in Java. The resulting aggregate files were read by an R program that ranked the methods, determined the appropriate shading, and wrote out an image file in PNG format. These PNG images were incorporated into a LaTeX report, and a compiler run to create a shareable document in PDF from the LaTeX and PNG sources. It was this PDF of the paper that we could disseminate to interested colleagues. If the colleague then suggested a change, we'd go through the process all over again. The ordering of the programs—C++, followed by Java, followed by R, followed by LaTeX, followed by attaching the PDF to an email—was nontrivial, and there were times when we skipped something in between, resulting in incorrect graphics or text that didn't match the graphs.

Jupyter Notebooks

JupyterLab is open source software that provides an interactive scientific computing experience in a variety of languages, including Python, R, Julia, and Scala. The key unit of work is a Jupyter Notebook, which is a document that contains code, visualizations, and explanatory text. The code in the document is executed on and served from a web server that runs JupyterLab.

A key issue with notebooks is how to manage the web servers that serve them. Running the notebook server on our laptop will work, but is not ideal. Instead, we want the notebook server to run on a cloud machine for the following reasons:

- Because the code is executed on the notebook server, we want the notebook server to be able to handle bigger datasets. It is easier to get a more powerful machine on the public cloud than it is to upgrade one's laptop. You can simplify managing the provisioning of notebook servers by running them on demand in the public cloud. This way, we can stop the machines when we leave for the day, instead of paying for machines even when we are not at work.
- Data science workloads such as machine learning require heavy, repetitive computation—typically, we scale up such workloads using GPUs. However, GPUs are expensive and become superseded by better hardware rather quickly. Ideally, you want to be able to add/remove GPUs on demand from these machines, rather than pay for GPUs all the time that we are using notebooks.
- A common pattern is to develop on small datasets and basic hardware and then, once we have the code working, to execute the code on large datasets on a more powerful machine. This ability to change the infrastructure is possible on the public cloud.
- We might even schedule the execution of these jobs periodically, or in response to an event such as the arrival of new data.

Once we say that we are going to run notebooks ephemerally on hardware that depends on the computations that we are doing, lifecycle management becomes quite important. Vertex AI Workbench on Google Cloud gives us a fully managed notebook experience.

To start a fully managed notebook in Google Cloud, visit the GCP web console, navigate to Vertex AI Workbench, and choose the tab for a Google-managed notebook. Then, create a notebook with the name dsongcp-ch5. Look at the Advanced settings and note that it's possible to add a GPU if we want (see [Figure 5-8](#)). Note also there is a default time period after which the notebook server will automatically shut down. We can always restart it by opening the notebook from the GCP console.

Vertex AI Workbench provides a hosted version of JupyterLab on Google Cloud; it knows how to authenticate against Google Cloud so as to provide easy access to

Cloud Storage, BigQuery, Cloud Dataflow, Vertex AI Training, and so on. A few minutes after you launch a managed notebook, the console shows you a JupyterLab link. When you are done using the Notebooks instance, you can manually delete the instance from the web console. You can also stop the instance when you are not using it—you won’t be charged for the CPU resources, which are the bulk of the cost, although resources like disks will continue to be charged for. As shown in [Figure 5-8](#), this shutdown can be made to happen automatically after an idle time period that you specify.

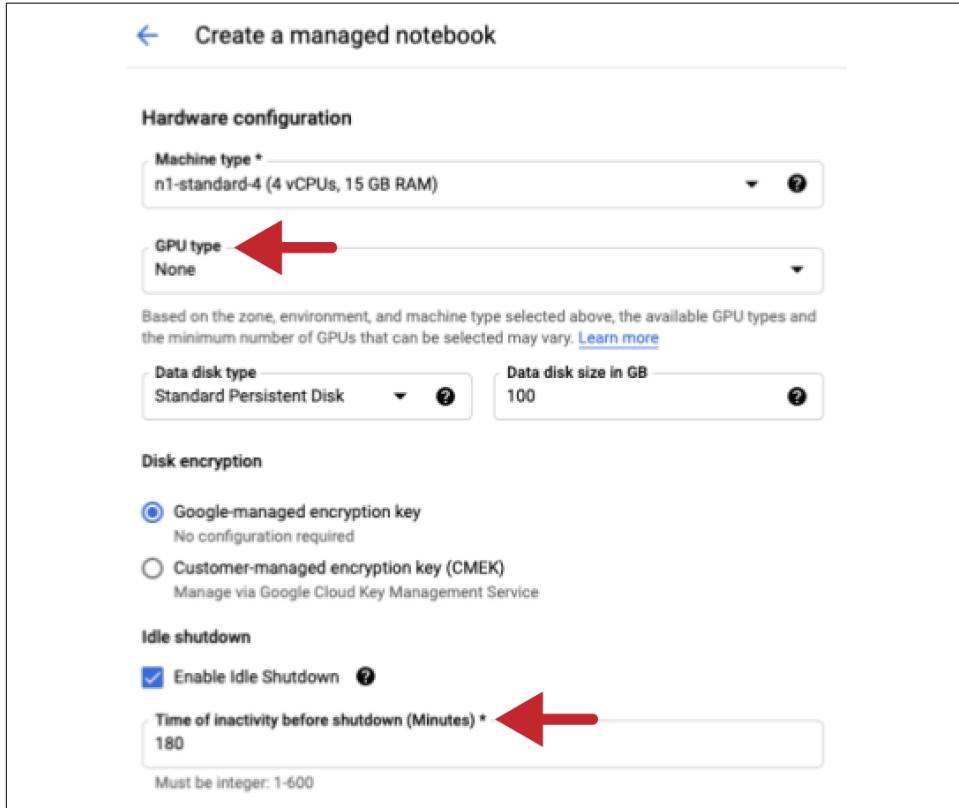


Figure 5-8. Options when creating a managed notebook in Vertex AI Workbench.

Creating a Notebook

After you have launched the Notebooks instance and navigated to JupyterLab, you can create a new Python notebook from the launcher menu that it starts up with. Alternately, navigate to the folder in which you want this notebook to appear and select File > New Notebook.

A notebook contains two key types of cells: a *markdown* cell is a cell with text content,¹² whereas a *code* cell has source code. When you run a markdown cell, the cell contents are formatted nicely, whereas when you run a code cell, the notebook displays the output of *print* statements from that cell.

For example, suppose that you type into a notebook the contents of [Figure 5-9](#) (with the first cell a markdown cell, and the second cell a Python cell).

You can run the cell your cursor is in by clicking Run (or use the keyboard shortcut Ctrl/Cmd + Shift + Enter). You could also run all cells by clicking “Run all cells.” When you click this, the cells are evaluated, and the results rendered.

The screenshot shows a Jupyter Notebook interface. At the top, there is a header bar with the title "Ch 5. Interactive Data Analysis". Below the header, there are two cells. The first cell is a markdown cell containing text about carrying out interactive data analysis using BigQuery and Vertex AI Workbench. The second cell is a code cell containing Python code to print the values of variables `a` and `b`. A red arrow points from the code cell down to the rendered output, which is the text "Ch 5. Interactive Data Analysis" followed by the printed output "a=3 b=8".

```
# Ch 5. Interactive Data Analysis

This notebook introduces carrying out interactive data analysis of data in BigQuery using a Jupyter Notebook managed by Vertex AI Workbench.

This cell, for example, is a mark-down cell. Which is why you are seeing text. The cell that follows is a Python code cell. The output of that cell is whatever is printed out from it.

[1]: a = 3
      b = a + 5
      print("a={} b={}".format(a,b))
```

Run

Ch 5. Interactive Data Analysis

```
This notebook introduces carrying out interactive data analysis of data in BigQuery using a Jupyter Notebook managed by Vertex AI Workbench.

This cell, for example, is a mark-down cell. Which is why you are seeing text. The cell that follows is a Python code cell. The output of that cell is whatever is printed out from it.

[1]: a = 3
      b = a + 5
      print("a={} b={}".format(a,b))

a=3 b=8
```

Figure 5-9. What you type into the notebook (top) and what is rendered (bottom).

¹² The *README.md* files in the repository are another example of markdown files. Markdown is ubiquitous enough that it is worth learning the key syntax from a [cheat sheet](#).

Note that the markdown has been converted into a visual document, the Python code has been evaluated, and the resulting output printed out.

Jupyter Commands

You can `git clone` the repository for this book within the notebook environment by typing:

```
!git clone https://github.com/GoogleCloudPlatform/data-science-on-gcp
```

You could have also used the git icon in Vertex AI Workbench or run the preceding command from a terminal. Regardless of how you interact with git, get into the habit of practicing source-code control on changed notebooks.

The use of the exclamation point (when you type `!git` into a code cell) is an indication to Jupyter that the line is not Python code, but is instead a shell command. If you have multiple lines of a shell command, you can start a cell with `%bash`, for example:

```
%%bash
wget tensorflow ...
pip install ...
```

Installing Packages

Which Python packages are already installed in Notebooks, and which ones will we have to install? One way to check which packages are installed is to type the following:

```
%pip freeze
```

This lists the Python packages installed. Another option is to add in imports for packages and see if they work. Let's do that with packages that I know that we'll need:

```
import matplotlib.pyplot as plt
import seaborn as sb
import pandas as pd
import numpy as np
```

NumPy is the canonical numerical Python library that provides efficient ways of working with arrays of numbers. *Pandas* is an extremely popular data analysis library that provides a way to do operations such as group by and filter on in-memory dataframes. *Matplotlib* is a Matlab-inspired module to create graphs in Python. *seaborn* provides extra graphics capabilities built on top of the basic functionality provided by Matplotlib. All these are open source packages that are installed in Vertex AI Workbench by default.

Had I needed a package that was not already installed, I could have installed it using `pip`. For example, to install the `pytz` package that we used in [Chapter 4](#), execute this code within a cell:

```
%pip install pytz
```

Often, you will need to restart the Python kernel for the new package to be picked up (you can do this using the Restart Kernel button on the notebook ribbon user interface).

Jupyter Magic for Google Cloud

When we used `%%bash` in the previous section, we were using a Jupyter *magic*, a syntactic element that marks what follows as special.¹³ This is how Jupyter can support multiple interpreters or engines. Jupyter knows what language a cell is in by looking at the magic at its beginning. For example, try typing the following into a code cell:

```
%%html  
This cell will print out a <b> HTML </b> string.
```

You should see the HTML rendering of that string being printed out on evaluation, as depicted in [Figure 5-10](#).

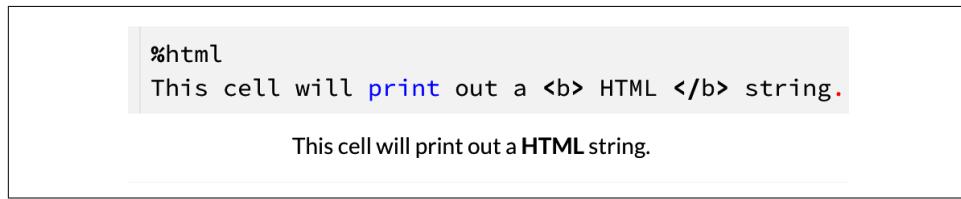


Figure 5-10. Jupyter magic for HTML rendering.

Jupyter mágics provide a mechanism to run a wide variety of languages and ways to add some more. The BigQuery Python package has added a few mágics to make the interaction with Google Cloud Platform convenient.

For example, you can run a query on your BigQuery table using the `%%bigquery` magic environment that comes with Vertex AI Workbench:

```
%%bigquery  
SELECT  
    COUNTIF(arr_delay >= 15)/COUNT(arr_delay) AS frac_delayed  
FROM dsongcp.flights_tzcorr
```

If you get the fraction of flights that are delayed, as shown in [Figure 5-11](#), all is well.

If not, look at the error message and carry out appropriate remedial actions. You might need to authenticate yourself, set the project you are working in, or change permissions on the BigQuery table.

¹³ Mágics use a syntax element that is not valid in the underlying language. In the case of notebooks in Python, this is the % symbol.

```
[7]: %%bigrquery
SELECT
    COUNTIF(arr_delay >= 15)/COUNT(arr_delay) AS frac_delayed
FROM dsongcp.flights_tzcorr

Query complete after 0.00s: 100%|██████████| 1/1 [00:00<00:00, 1030.0
4query/s]
Downloading: 100%|██████████| 1/1 [00:01<00:00, 1.99s/rows]

[7]: frac_delayed
0      0.186111
```

Figure 5-11. The `%%bigrquery` magic environment that comes with Vertex AI Workbench.

The fact that we refer to `%%bigrquery` as a Jupyter magic should indicate that this is not pure Python—you can execute this only within a notebook environment. The magic, however, is simply a wrapper function for Python code.¹⁴ If there is a piece of code that you’d ultimately want to run outside a notebook (perhaps as part of a scheduled script), it’s better to use the underlying Python and not the magic pragma.¹⁵

```
sql = """
SELECT
    COUNTIF(arr_delay >= 15)/COUNT(arr_delay) AS frac_delayed
FROM dsongcp.flights_tzcorr
"""

from google.cloud import bigrquery
bq = bigrquery.Client()
df = bq.query(sql).to_dataframe()
print(df)
```

One way to use the underlying Python is to use the `google.cloud.bigrquery` package—this allows us to use code independent of the notebook environment. This is, of course, the same `bigrquery` package in the Cloud Client Library that we used in Chapters 2 and 4. The client library includes interconnections between BigQuery results and NumPy/Pandas to simplify the creation of graphics.

Exploring Arrival Delays

Now that we have a notebook up and running, let’s use it to do exploratory analysis of arrival delays because this is the variable we want to be able to predict.

¹⁴ See [GitHub](#) for the Python code being wrapped.

¹⁵ See `05_bqnotebook/exploration.ipynb` in the GitHub repository.

Basic Statistics

To pull the arrival delays corresponding to the model created in [Chapter 3](#) (i.e., of the arrival delay for flights that depart more than 10 minutes late), we can do the following:

```
%%bigquery df
SELECT ARR_DELAY, DEP_DELAY
FROM dsongcp.flights_tzcorr
WHERE DEP_DELAY >= 10
```

This code uses the `%%bigquery` magic to run the SQL statement and stores the result set into a Pandas dataframe named `df`. Recall that in [Chapter 4](#), we did this using the Google Cloud Platform API.

After we have the dataframe, getting fundamental statistics about the two columns returned by the query is as simple as this:

```
df.describe()
```

This gives us the mean, standard deviation, minimum, maximum, and quartiles of the arrival and departure delays given that departure delay is more than 10 minutes, as illustrated in [Figure 5-12](#) (see the WHERE clause of the query):

The screenshot shows a Jupyter Notebook cell with the command `[12]: df.describe()`. The output displays statistical summary data for the `ARR_DELAY` and `DEP_DELAY` columns. The data is presented in a table format with two columns: `ARR_DELAY` and `DEP_DELAY`. The rows include count, mean, std, min, 25%, 50% (median), 75%, and max. All values are in scientific notation.

	ARR_DELAY	DEP_DELAY
count	1.286778e+06	1.294778e+06
mean	4.611797e+01	5.094516e+01
std	6.360700e+01	6.151423e+01
min	-7.800000e+01	1.000000e+01
25%	1.100000e+01	1.700000e+01
50%	2.700000e+01	3.000000e+01
75%	5.900000e+01	6.100000e+01
max	1.971000e+03	1.988000e+03

Figure 5-12. Getting the fundamental statistics of a Pandas dataframe.

Plotting Distributions

Beyond just the statistical capabilities of Pandas, we can also pass the Pandas dataframes and underlying NumPy arrays to plotting libraries like `seaborn`. For example, to plot a violin plot of our decision surface from [Chapter 3](#) (i.e., of the arrival delay for flights that depart more than 10 minutes late), we can do the following:

```
sns.set_style("whitegrid")
ax = sns.violinplot(data=df, x='ARR_DELAY', inner='box', orient='h')
ax.axes.set_xlim(-50, 300);
```

This produces the graph shown in [Figure 5-13](#).

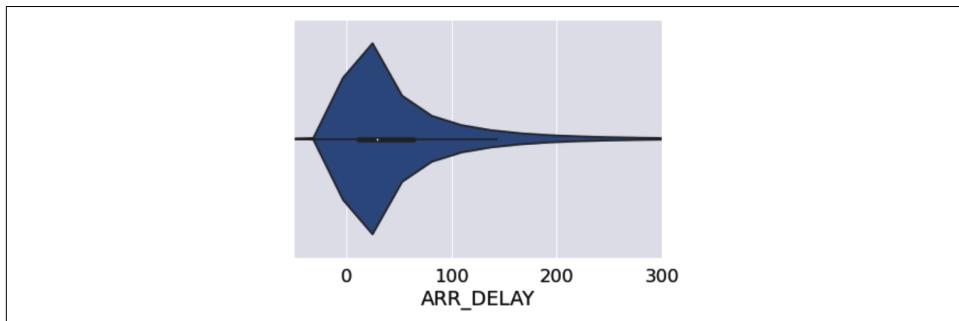


Figure 5-13. Violin plot of arrival delay.

A violin plot is a *kernel density plot*,¹⁶ that is, it is an estimate of the probability distribution function (PDF).¹⁷ We see that, even though the distribution peaks around 10 minutes (which is the mode), deviations around this peak are skewed toward larger delays than smaller ones. Importantly, we also notice that there is only one peak—the distribution is not, for example, bimodal.

Let's compare the violin plot for flights that depart more than 10 minutes late with the violin plot for flights that depart less than 10 minutes late and zoom in on the x-axis close to our 15-minute threshold. First, we pull all of the delays using the following:

```
%%bq df
SELECT ARR_DELAY, DEP_DELAY
FROM dsongcp.flights_tzcorr
```

In this query, I have dropped the `WHERE` clause. Instead, we will rely on Pandas to do the thresholding. I can now create a new column in the Pandas dataframe that is either `True` or `False` depending on whether the flight departed less than 10 minutes late:

```
df['ontime'] = df['DEP_DELAY'] < 10
```

We can graph this new Pandas dataframe using *seaborn*:

¹⁶ A kernel `density plot` is just a smoothed histogram—the challenge lies in figuring out how to smooth the histogram while balancing interpretability against the loss of information. Here, I'm just letting *seaborn* use its default settings for the smoothing bandwidth.

¹⁷ See the discussion of the PDF in [Chapter 1](#).

```
ax = sns.violinplot(data=df, x='ARR_DELAY', y='ontime',
                     inner='box', orient='h')
ax.set_xlim(-50, 200)
```

The difference between the previous violin plot and this one is the inclusion of the `ontime` column. This results in a violin plot (see [Figure 5-14](#)) that illustrates how different flights that depart 10 minutes late are from flights that depart early.

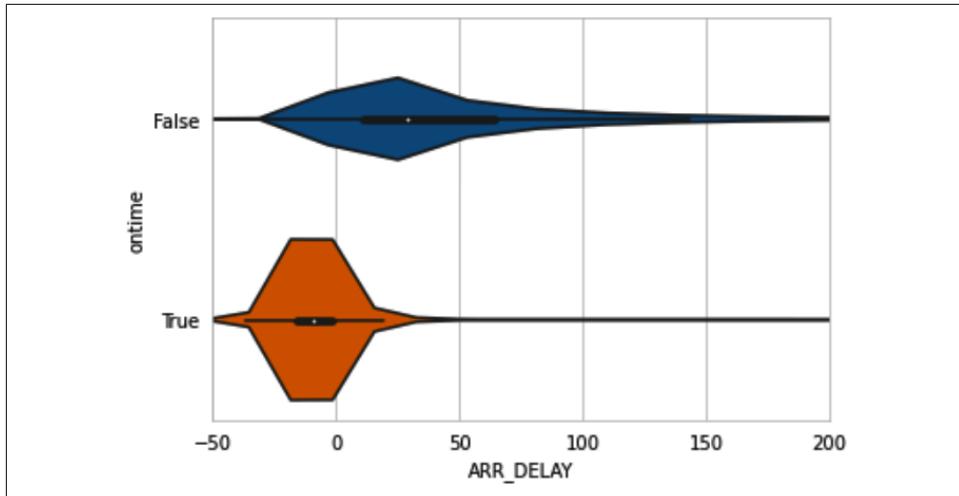


Figure 5-14. Difference between violin plots of all late flights (top) versus on-time flights (bottom).



The angular peak of the top violin plot indicates that the `seaborn` default smoothing was too coarse. You can fix this by passing in a `gridsize` parameter:

```
ax = sns.violinplot(data=df, x='ARR_DELAY', y='ontime',
                     inner='box', orient='h', gridsize=1000)
```

But doing so will make the computation take much longer. The [notebook](#) in the GitHub repository shows what the result looks like with greater smoothing.

As we discussed in [Chapter 3](#), it is clear that the 10-minute threshold separates the dataset into two separate statistical regimes, so that the typical arrival delay for flights that depart more than 10 minutes late is skewed toward much higher values than for flights that depart more on time. We can see this in [Figure 5-14](#), both from the shape of the violin plot and from the box plot that forms its center. Note how centered the

on-time flights are versus the box plot (the dark line in the center) for delayed flights.¹⁸

However, the extremely long, skinny tail of the violin plot is a red flag—it is an indication that the dataset might pose modeling challenges. Let's investigate what is going on.

Quality Control

We can continue writing queries in the notebook, but doing so on the BigQuery console gives me immediate feedback on syntax and logic errors. So, I switch over to [the BigQuery console](#) and type in my first query:

```
SELECT
    AVG(ARR_DELAY) AS arrival_delay
FROM
    dsongcp.flights_tzcorr
GROUP BY
    DEP_DELAY
ORDER BY
    DEP_DELAY
```

This should give me the average arrival delay associated with every value of departure delay (which, in this dataset, is stored as an integer number of minutes). I got back more than one thousand rows. Are there really more than one thousand unique values of DEP_DELAY? What's going on?

Oddball values

To look at this further, let's add more elements to my initial query:

```
SELECT
    DEP_DELAY,
    AVG(ARR_DELAY) AS arrival_delay,
    COUNT(ARR_DELAY) AS numflights
FROM
    dsongcp.flights_tzcorr
GROUP BY
    DEP_DELAY
ORDER BY
    DEP_DELAY
```

The resulting table explains what's going on. The first few rows have only a few flights each:

¹⁸ I created this second, zoomed-in violin plot by adding `ax.set_xlim(-50, 50)`.

Row	DEP_DELAY	arrival_delay	numflights
1	null	null	0
2	-82.0	-80.0	1
3	-68.0	-87.0	1
4	-61.0	-77.0	1
5	-56.0	-26.0	1

However, departure delay values of a few minutes have hundreds of thousands of flights:

56	0.0	-5.100334305600409	328442
57	1.0	-4.188285855693881	159619
58	2.0	-3.2246696399075128	121080
59	3.0	-2.1957821784079146	104177
60	4.0	-1.2101860730716607	92813

Oddball values that are such a small proportion of the data can probably be ignored. Moreover, if the flight does really leave 82 minutes early, I'm quite sure that you won't be on the flight, and if you are, you know that you will make the meeting. There is no reason to complicate our statistical modeling with such odd values.

Outlier removal: Big data is different

How can you remove such outliers? There are two ways to filter the data: one would be based on the departure delay variable itself, keeping only values that met a condition such as this:

```
WHERE dep_delay > -15
```

A second method would be to filter the data based on the number of flights:

```
WHERE numflights > 300
```

The second method—using a quality-control filter that is based on removing data for which we have insufficient examples—is preferable.

This is an important point that gets at the key difference between statistics on “normal” datasets and statistics on big data. Although I agree that the term *big data* has become completely hyped, people who claim that big data is just data are missing a key point—the fundamental approach to problems becomes different when datasets grow sufficiently large. The way we detect outliers is just one such example.

For a dataset that numbers in the hundreds to thousands of examples, you would filter the dataset and remove values outside, say, $\mu \pm 3\sigma$ (where μ is the mean and σ the

standard deviation).¹⁹ We can find out what the range would be by running a BigQuery query on the table:

```
SELECT
    AVG(DEP_DELAY) - 3*STDDEV(DEP_DELAY) AS filtermin,
    AVG(DEP_DELAY) + 3*STDDEV(DEP_DELAY) AS filtermax
FROM
    dsongcp.flights_tzcorr
```

This yields the range $[-102, 121]$ minutes so that the WHERE clause would become as follows:

```
WHERE dep_delay BETWEEN -102 AND 121
```

Of course, a filter that retains values in the range $\mu \pm 3\sigma$ is based on an implicit assumption that the distribution of departure delays is Gaussian. We can avoid such an assumption by using percentiles, perhaps by omitting the top and bottom 5% of values:

```
SELECT
    APPROX_QUANTILES(DEP_DELAY, 20)
FROM
    dsongcp.flights_tzcorr
```

This would lead us to retain values in the range $[-9, 66]$. Regardless of how we find the range, though, the range is based on an assumption that unusually high and low values are outliers.

On datasets that number in the hundreds of thousands to millions of examples, thresholding your input data based on value is dangerous because you can very well be throwing out valuable nuance—if there are sufficient examples of a delay of 150 minutes, it is worth modeling such a value regardless of how far off the mean it is. Customer satisfaction and “long-tail” business strategies might hinge on our systems coping well with usually small or large values. There is, therefore, a world of difference between filtering our data using:

```
WHERE dep_delay > -15
```

versus filtering it using:

```
WHERE numflights > 370
```

The first method imposes a threshold on the input data and is viable only if we are sure that a departure delay of less than -15 minutes is absurd. The second method, on the other hand, is based on how often certain values are observed—the larger our dataset grows, the less unusual any particular value becomes.

¹⁹ In a Gaussian distribution, 99.7% of values lie within three standard deviations of the mean. It's a handy way to identify outliers.

The term *outlier* is, therefore, somewhat of a misnomer when it comes to big data. An outlier implies a range within which values are kept, with outliers being values that lie outside that range. Here, we are keeping data that meets a criterion involving frequency of occurrence—any value is acceptable as long as it occurs often enough in our data.²⁰

How Far Can You Trust Quality Flags in Datasets?

Many datasets include metadata about their data quality. These might even be on a row-by-row basis. Should you discard any rows whose quality is marked as being bad?

The quality flags that you find in many datasets are themselves highly suspect. Many of them are set with no basis on a holistic understanding of the environment (“the instrument was left unshielded”) but simply based on statistical analysis of the data values. The statistical techniques used are often carried over from the days of small datasets. So, if you can do your own analysis based on frequency of occurrence, you should.

Of course, if you don’t have time to examine the data, the quality flag in the dataset is better than nothing. Take it into consideration! However, the way to take it into consideration is to treat it as one more input to your model and not to discard the supposedly bad values.

In our flights dataset, we will trust flags such as whether the flight was canceled or diverted (those reflect an understanding of the environment) but carry out our own statistical analysis of values such as departure delay based on occurrence frequency.

Filtering data on occurrence frequency

To filter the dataset based on frequency of occurrence, we first need to compute the frequency of occurrence and then threshold the data based on it. We can accomplish this by using a HAVING clause:

```
SELECT
    DEP_DELAY,
    AVG(ARR_DELAY) AS arrival_delay,
    STDDEV(ARR_DELAY) AS stddev_arrival_delay,
    COUNT(ARR_DELAY) AS numflights
FROM
    dsongcp.flights_tzcorr
GROUP BY
```

²⁰ In this dataset, floating-point numbers have already been discretized. For example, arrival delays have been rounded to the nearest minute. If this is not the case, you will have to discretize continuous data before computing the frequency of occurrence.

```

    DEP_DELAY
HAVING
  numflights > 370
ORDER BY
  DEP_DELAY

```

Why threshold the number of flights at 370? This number derives from a guideline called the *three-sigma rule*,²¹ which is traditionally the range within which we consider “nearly all values”²² to lie. If we assume (for now; we’ll verify it soon) that at any departure delay, arrival delays are normally distributed, we can talk about things that are true for “almost every flight” if our population size is large enough. Because 99.73% of values in a Gaussian distribution lie within the three-sigma bounds, filtering our dataset so that we have at least $1 / (1 - 0.9973) = 370$ examples of each input value is a rule of thumb that achieves this.²³

How different would the results be if we were to choose a different threshold? We can look at the number of flights that are removed by different quality-control thresholds by looking at the slope of a linear model between arrival delay and departure delay using this query:

```

CREATE TEMPORARY FUNCTION linear_fit(NUM_TOTAL INT64, THRESH INT64)
RETURNS STRUCT<thresh INT64, num_removed INT64, ln FLOAT64>
AS (((
  SELECT AS STRUCT
    THRESH,
    (NUM_TOTAL - SUM(numflights)) AS num_removed,
    AVG(arrival_delay * numflights) / AVG(dep_delay * numflights) AS ln
  FROM
  (
    SELECT
      DEP_DELAY,
      AVG(ARR_DELAY) AS arrival_delay,
      STDDEV(ARR_DELAY) AS stddev_arrival_delay,
      COUNT(ARR_DELAY) AS numflights
    FROM

```

²¹ For a normal distribution (at each departure delay, the number of flights is in the hundreds to thousands, so usual statistical thinking applies), 68.27% of values lie in the $\mu \pm \sigma$ range, 95.45% of values lie in the $\mu \pm 2\sigma$ range, and 99.73% of values lie in the $\mu \pm 3\sigma$ range. That last range is termed the three-sigma rule. For more information, see the [Encyclopedia of Mathematics entry](#) for the three-sigma rule.

- ²² Traditions, of course, are different in different fields and often depend on how much data you can reasonably collect in that field. In business statistics, this three-sigma rule is quite common. In the social sciences and in medicine, two-sigma is the typical significance threshold. Meanwhile, when the Higgs boson discovery announcement was made, the significance threshold to classify it as a true discovery and not just a statistical artifact was five-sigma or 1 in 3.5 million (see the [blog at Scientific American](#)).
- ²³ It might appear fishy that this number is independent of the size of the dataset, but if you think about it, this rule of thumb has to be such that we are less likely to discard outliers the more data we have. The larger the dataset, the more likely it is that there will be 370 instances of any particular condition. Corner cases on small datasets will have enough company on very large datasets.

```

        dsongcp.flights_tzcorr
    GROUP BY
        DEP_DELAY
    )
    WHERE numflights > THRESH
)
;

```

Running this function for various different thresholds on `numflights` (see `exploration.ipynb` in the GitHub repository), we get the following results:

Row	stats.thresh	stats.num_removed	stats.lm
1	1000	175873	0.25
2	500	143801	0.34
3	370 (three-sigma rule)	135518	0.36
4	300	129835	0.38
5	200	123640	0.40
6	100	115471	0.43
7	22 (two-sigma rule)	108247	0.45
8	10	106958	0.46
9	5	106319	0.46

As you can see, the slope varies extremely slowly as we remove fewer and fewer flights by decreasing the threshold. Thus, the differences in the model created for thresholds of 300, 370, or 500 are quite minor. However, that model is quite different from that created if the threshold were 5 or 10. The order of magnitude of the threshold matters, but perhaps not the exact value.

Arrival Delay Conditioned on Departure Delay

Now that we have a query that cleans up oddball values of departure delay from the dataset, we can take the query over to the Jupyter Notebook to continue our exploratory analysis and to develop a model to help us make a decision on whether to cancel our meeting.

In [Chapter 3](#), we built a simple model based on simply thresholding the departure delay. Here, however, we see that there are many flights for each value of departure delay. Given a certain departure delay, what arrival delays are likely?

Distribution of arrival delays

I simply copy and paste from the BigQuery console to the notebook and give the Pandas dataframe a name, as shown here:

```

%%bigquery depdelay
SELECT

```

```

    DEP_DELAY,
    AVG(ARR_DELAY) AS arrival_delay,
    STDDEV(ARR_DELAY) AS stddev_arrival_delay,
    COUNT(ARR_DELAY) AS numflights
  FROM
    dsongcp.flights_tzcorr
 GROUP BY
    DEP_DELAY
 HAVING numflights > 370
 ORDER BY DEP_DELAY

```

We can display the first five rows of the dataframe using [:5]:

```
depdelay[:5]
```

The result is:

Row	DEP_DELAY	arrival_delay	stddev_arrival_delay	numflights
1	-23.0	-23.888646288209607	11.432163250582196	458
2	-22.0	-23.22748815165877	12.590133374822704	633
3	-21.0	-22.29978118161926	11.558312559289162	914
4	-20.0	-21.40782122905028	12.066489232808147	1432
5	-19.0	-20.430769230769243	11.910133697086701	1950

Let's plot this data to see what insight we can get. Even though we have been using *seaborn* so far, Pandas itself has plotting functions built in:

```
ax = depdelay.plot(kind='line', x='DEP_DELAY',
                    y='arrival_delay', yerr='stddev_arrival_delay')
```

This yields the plot shown in [Figure 5-15](#).

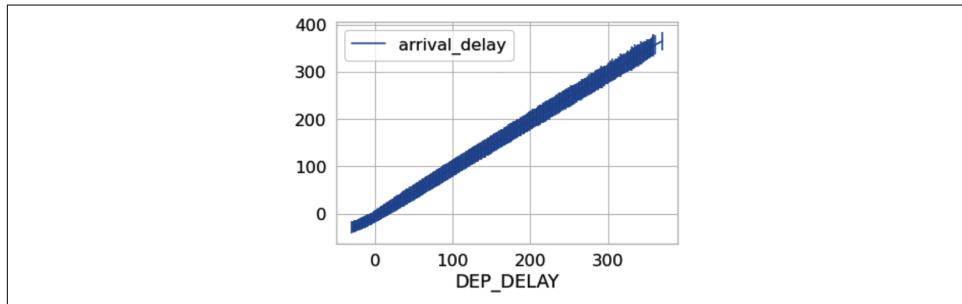


Figure 5-15. Relationship between departure delay and arrival delay.

It certainly does appear as if the relationship between departure delay and arrival delay is quite linear. The width of the standard deviation of the arrival delay is also pretty constant, on the order of 10 minutes.

Applying a probabilistic decision threshold

Recall from [Chapter 1](#) that our decision criteria are 15 minutes and 30%. If the plane is more than 30% likely to be delayed (on arrival) by more than 15 minutes, we want to send a text message asking to postpone the meeting. At what departure delay does this happen?

By computing the standard deviation of the arrival delays corresponding to each departure delay, we implicitly assumed that arrival delays are normally distributed. For now, let's continue with that assumption. I can examine a [complementary cumulative distribution table](#) and find where 0.3 happens. From the table, this happens at $Z = 0.52$.

Let's now go back to Jupyter to plug this number into our dataset:

```
Z_30 = 0.52
depdelay['arr_delay_30'] = (Z_30 * depdelay['stddev_arrival_delay']) \
    + depdelay['arrival_delay']
ax = plt.axes()
depdelay.plot(kind='line', x='DEP_DELAY', y='arr_delay_30',
               ax=ax, ylim=(0,30), xlim=(0,30), legend=False)
ax.set_xlabel('Departure Delay (minutes)')
ax.set_ylabel('> 30% prob of this Arrival Delay (minutes)');

x = np.arange(0, 30)
y = np.ones_like(x) * 15
ax.plot(x, y, 'r');
```

The plotting code yields the plot depicted in [Figure 5-16](#).

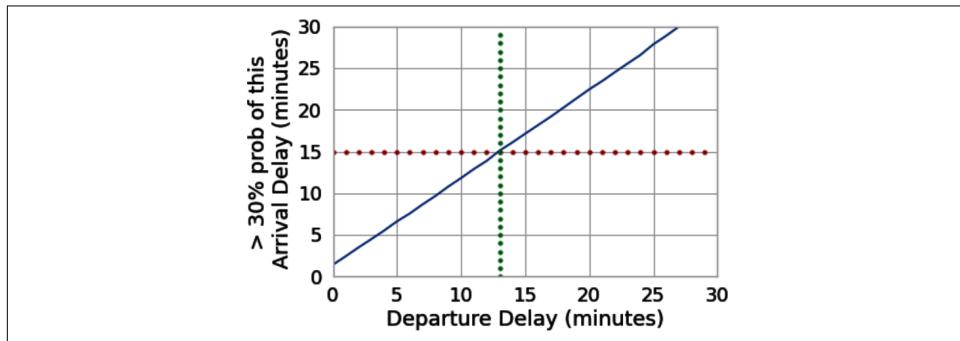


Figure 5-16. Choosing the departure delay threshold that results in a 30% probability of an arrival delay of < 15 minutes.

Looking up the x-axis value corresponding to the decision threshold of 15 minutes (see dotted lines in [Figure 5-16](#)). It appears that our decision criteria translate to a departure delay of 13 minutes. If the departure delay is 13 minutes or more, the aircraft is more than 30% likely to be delayed by 15 minutes or more.

Empirical probability distribution function

The analysis in the previous section used the number 0.52, which assumes that the distribution of flights at each departure delay is normally distributed. What if we drop that assumption? We then will need to empirically determine the 30% likelihood at each departure delay. Happily, we do have at least 370 flights at each departure delay (the joys of working with large datasets!), so we can simply compute the 30th percentile for each departure delay.

We can compute the 30th percentile in BigQuery by discretizing the arrival delays corresponding to each departure delay into 100 bins and picking the arrival delay that corresponds to the 70th bin:

```
SELECT
    DEP_DELAY,
    APPROX_QUANTILES(ARR_DELAY, 101)[OFFSET(70)] AS arrival_delay_30th,
    COUNT(ARR_DELAY) AS numflights
FROM
    dsongcp.flights_tzcorr
GROUP BY
    DEP_DELAY
HAVING numflights > 370
ORDER BY DEP_DELAY
```

The function `APPROX_QUANTILES()` takes the `ARR_DELAY` and divides it into $N + 1$ bins (here we specified $N = 101$).²⁴ The first bin is the approximate minimum, the last bin the approximate maximum, and the rest of the bins are what we'd traditionally consider the bins. Hence, the 70th percentile is the 71st element of the result. The `[]` syntax finds the n th element of that array—`OFFSET(70)` will provide the 71st element because `OFFSET` is zero-based.²⁵ Why 70 and not 30? Because we want the arrival delay that could happen with 30% likelihood and this implies the larger value.

The results of this query provide the empirical 30th percentile threshold for every departure delay:

Row	DEP_DELAY	arrival_delay_30th	numflights
1	-23.0	-20.0	458
2	-22.0	-19.0	633
...			
39	15.0	14.0	38835
40	16.0	15.0	35771

²⁴ This function computes [the approximate quantiles](#) because computing the exact quantiles on large datasets, especially of floating-point values, can be very expensive in terms of space. Instead, most big data databases use some variant of [Greenwald and Khanna's algorithm](#) to compute approximate quantiles.

²⁵ Had I used `ORDINAL` instead of `OFFSET`, it would have been 1-based.

Row	DEP_DELAY	arrival_delay_30th	numflights
41	17.0	16.0	33964
...			

Plugging the query back into the Jupyter Notebook, we can avoid the Z-lookup and Z-score calculation associated with Gaussian distributions. We get the chart shown in [Figure 5-17](#).

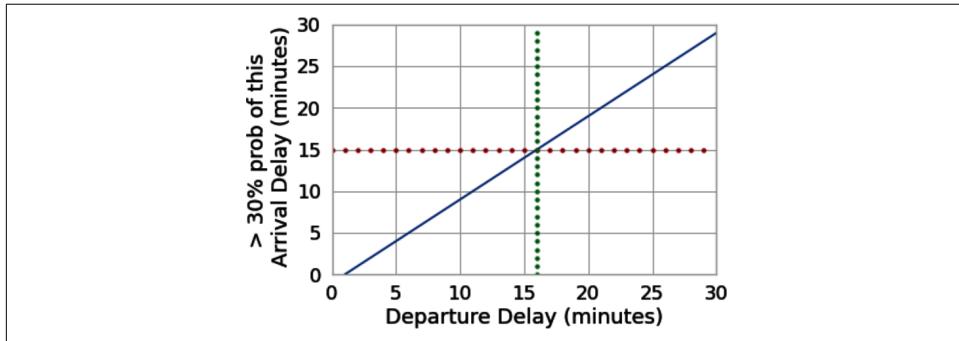


Figure 5-17. Departure delay threshold that results in a 30% likelihood of an arrival delay of < 15 min.

The answer is...

From the chart in [Figure 5-16](#), our decision threshold, without the assumption of normal distribution, is 16 minutes. If a flight is delayed by more than 16 minutes, there is a greater than 30% likelihood that the flight will arrive more than 15 minutes late.

Recall that the aforementioned threshold is conditioned on rather conservative assumptions—you are going to cancel the meeting if there is more than 30% likelihood of being late by 15 minutes. What if you are a bit more audacious in your business dealings, or if this particular customer will not be annoyed by a few minutes' wait? What if you won't cancel the meeting unless there is a greater than 70% chance of being late by 15 minutes? The good thing is that it is easy enough to come up with a different decision threshold for different people and different scenarios out of the same basic framework.

Another thing to notice is that the addition of the actual departure delay in minutes has allowed us to make a better decision than going with just the contingency table. Using just the contingency table, we would cancel meetings whenever flights were just 10 minutes late. Using the actual departure delay and a probabilistic decision framework, we were able to avoid canceling our meeting unless flights were delayed by 16 minutes or more.

Evaluating the Model

But how good is this advice? How many times will my advice to cancel or not cancel the meeting be the correct one? Had you asked me that question, I would have hemmed and hawed—I don't know how accurate the threshold is because we have no independent sample. Let's address that now—as our models become more sophisticated, an independent sample will be increasingly required.

There are two broad approaches to finding an independent sample:

- Collect new data. For example, we could go back to BTS and download 2016 data and evaluate the recommendation on that dataset.
- Split the 2015 data into two parts. Create the model on the first part (called the *training set*), and evaluate it on the second part (called the *test set*).

The second approach is more common because datasets tend to be finite. In the interest of being practical here, let's do the same thing even though, in this instance, we could go back and get more data.²⁶

When splitting the data, we must be careful. We want to ensure that both parts are representative of the full dataset (and have similar relationships to what you are predicting), but at the same time we want to make sure that the testing data is independent of the training data. To understand what this means, let's take a few reasonable splits and talk about why they won't work.

Random Shuffling

We might split the data by randomly shuffling all the rows in the dataset and then choosing the first 70% as the training set, and the remaining 30% as the test set. In BigQuery, you could do that using the `RAND()` function:

```
SELECT
    ORIGIN, DEST,
    DEP_DELAY,
    ARR_DELAY
FROM
    dsongcp.flights_tzcorr
WHERE
    RAND() < 0.7
```

The `RAND()` function returns a value between 0 and 1, so approximately 70% of the rows in the dataset will be selected by this query. However, there are several problems with using this sampling method for machine learning:

²⁶ By the time we get to [Chapter 10](#), I will have based so many decisions on 2015–2018 data that I will get 2019 data to act as a truly independent test set.

- It is not nearly as easy to get the 30% of the rows that were not selected to be in the training set to use as the test dataset.
- The `RAND()` function returns different things each time it is run, so if you run the query again, you will get a different 70% of rows. In this book, we are experimenting with different machine learning models, and this will play havoc with comparisons between models if each model is evaluated on a different test set.
- The order of rows in a BigQuery result set is not guaranteed—it is essentially the order in which different workers return their results. So, even if you could set a random seed to make `RAND()` repeatable, you'll still not get repeatable results. You'd have to add an `ORDER BY` clause to explicitly sort the data (on an `ID` field, a field that is unique for each row) before doing the `RAND()`. This is not always going to be possible.

Further, on this particular dataset, random shuffling is problematic for another reason. Flights on the same day are probably subject to the same weather and traffic factors. Thus, the rows in the training set and test sets will not be independent if we simply shuffle the data. This consideration is relevant only for this particular dataset—shuffling the data and taking the first 70% will work for other datasets that don't have this interrow dependence, as long as you have an `id` field.

We could split the data such that Jan–Sep 2015 is training data and Oct–Dec is testing data. But what if delays can be made up in summer but not in winter? This split fails the representativeness test. Neither the training dataset nor the test dataset will be representative of the entire year if we split the dataset by months.

Splitting by Date

The approach that we will take is to find all the unique days in the dataset, shuffle them, and use 70% of these days as the training set and the remainder as the test set. For repeatability, I will store this division as a table in BigQuery.

The first step is to get all the unique days in the dataset:

```
SELECT
  DISTINCT(FL_DATE) AS FL_DATE
FROM
  dsongcp.flights_tzcorr
ORDER BY
  FL_DATE
```

The next step is to select a random 70% of these to be our training days:

```
SELECT
  FL_DATE,
  IF(ABS(MOD(FARM_FINGERPRINT(CAST(FL_DATE AS STRING)), 100)) < 70,
    'True', 'False') AS is_train_day
```

```

FROM (
  SELECT
    DISTINCT(FL_DATE) AS FL_DATE
  FROM
    dsongcp.flights_tzcorr)
ORDER BY
  FL_DATE

```

In the preceding query, the hash value of each of the unique days from the inner query is computed using the FarmHash library and the `is_train_day` field is set to True if the last two digits of this hash value are less than 70:²⁷

Row	FL_DATE	is_train_day
1	2015-01-01	True
2	2015-01-02	False
3	2015-01-03	False
4	2015-01-04	True
5	2015-01-05	True

The final step is to save this result as a table in BigQuery:

```

CREATE OR REPLACE TABLE dsongcp.trainday AS
...

```

We can join with this table whenever we want to pull out training rows. For your convenience, the preceding query is in the GitHub repository in the file `trainday.txt`, so you can simply do:

```
cat trainday.txt | bq query --nouse_legacy_sql
```

In some chapters, we won't be using BigQuery. Just in case we aren't using BigQuery, I will also export the table as a CSV file—we can do this on the web console, but we can also script it:

```
bq extract dsongcp.trainday gs://${BUCKET}/flights/trainday.csv
```

Training and Testing

Now, I can go back and edit my original query to carry out the percentile using only data from my training days. To do that, I will change this string in my original query:

```

FROM
  dsongcp.flights_tzcorr

```

to:

²⁷ See [the Google Open Source Blog](#) for a description and [GitHub](#) for the code.

```

FROM
  dsongcp.flights_tzcorr
JOIN dsongcp.trainday USING(FL_DATE)
WHERE is_train_day = 'True'

```

Now, the percentile is computed out only on days for which `is_train_day` is `True`.

The code to create the plot remains the same. On running it, the threshold (the x-axis value of the intersection point) remains consistent, as depicted in [Figure 5-18](#).

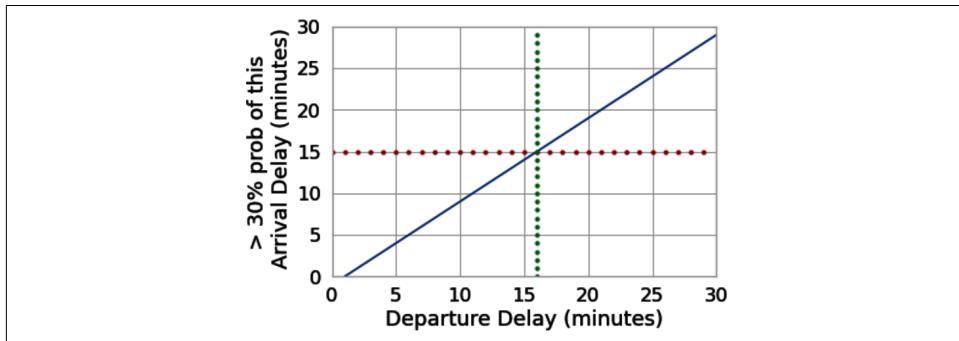


Figure 5-18. The departure delay threshold remains consistent with earlier methods.

This is gratifying because we get the same answer—16 minutes—after creating the empirical probabilistic model on just the training data.

Let's formally evaluate how well our recommendation of 16 minutes does in terms of predicting an arrival delay of 15 minutes or more. To do that, we have to find the number of times that we would have wrongly canceled a meeting or missed a meeting. We can compute these numbers using this query on days that are not training days:

```

SELECT
  SUM(IF(DEP_DELAY < 16
        AND arr_delay < 15, 1, 0)) AS correct_nocancel,
  SUM(IF(DEP_DELAY < 16
        AND arr_delay >= 15, 1, 0)) AS wrong_nocancel,
  SUM(IF(DEP_DELAY >= 16
        AND arr_delay < 15, 1, 0)) AS wrong_cancel,
  SUM(IF(DEP_DELAY >= 16
        AND arr_delay >= 15, 1, 0)) AS correct_cancel
FROM (
  SELECT
    DEP_DELAY,
    ARR_DELAY
  FROM
    dsongcp.flights_tzcorr
  JOIN dsongcp.trainday USING(FL_DATE)
  WHERE is_train_day = 'False'
)

```

Note that unlike when I was computing the decision threshold, I am not removing outliers (i.e., thresholding on 370 flights at a specific departure delay) when evaluating the model—outlier removal is part of my training process, and the evaluation needs to be independent of that. The second point to note is that this query is run on days that are not in the training dataset. Running this query in BigQuery, I get:

Row	correct_nocancel	wrong_nocancel	wrong_cancel	correct_cancel
1	1259740	66081	52827	217669

We will cancel meetings corresponding to a total of $52,827 + 217,669$ or around 270k flights. What fraction of the time are these recommendations correct? We can do the computation in the notebook (assuming that the dataframe is named eval):

```
print(eval['correct_nocancel'] /  
      (eval['correct_nocancel'] + eval['wrong_nocancel']))  
print(eval['correct_cancel'] /  
      (eval['correct_cancel'] + eval['wrong_cancel']))
```

Figure 5-19 presents the results.

```
[36]: print(df_eval['correct_nocancel'] /  
           (df_eval['correct_nocancel'] + df_eval['wrong_nocancel']))  
print(df_eval['correct_cancel'] /  
      (df_eval['correct_cancel'] + df_eval['wrong_cancel']))  
  
0    0.947403  
dtype: float64  
0    0.8187  
dtype: float64
```

Figure 5-19. Computing accuracy on independent test dataset.

It turns out when I recommend that you not cancel your meeting, I will be correct 95% of the time, and when I recommend that you cancel your meeting, I will be correct 82% of the time.

Why is this not 70%? Because the populations are different. In creating the model, we found the 70th percentile of arrival delay given a specific departure delay. In evaluating the model, we looked at the dataset of all flights. One's a marginal distribution, and the other's the full one. Another way to think about this is that the 95% figure is padded by all the departure delays of more than 20 minutes when canceling the meeting is an easy call.

We could evaluate right at the decision boundary by changing our scoring function:

```
SELECT  
SUM(IF(DEP_DELAY = 15  
      AND arr_delay < 15, 1, 0)) AS correct_nocancel,
```

```

SUM(IF(DEP_DELAY = 15
       AND arr_delay >= 15, 1, 0)) AS wrong_nocancel,
SUM(IF(DEP_DELAY = 16
       AND arr_delay < 15, 1, 0)) AS wrong_cancel,
SUM(IF(DEP_DELAY = 16
       AND arr_delay >= 15, 1, 0)) AS correct_cancel
...

```

If we do that, evaluating only at departure delays of 15 and 16 minutes, the contingency table and ratios look like those in [Figure 5-20](#).

```

%%bigquery eval
SELECT
  SUM(IF(DEP_DELAY = 15
         AND arr_delay < 15, 1, 0)) AS correct_nocancel,
  SUM(IF(DEP_DELAY = 15
         AND arr_delay >= 15, 1, 0)) AS wrong_nocancel,
  SUM(IF(DEP_DELAY = 16
         AND arr_delay < 15, 1, 0)) AS wrong_cancel,
  SUM(IF(DEP_DELAY = 16
         AND arr_delay >= 15, 1, 0)) AS correct_cancel
FROM (
  SELECT
    DEP_DELAY,
    ARR_DELAY
  FROM
    dsongcp.flights_tzcorr
  JOIN dsongcp.trainday USING(FL_DATE)
  WHERE is_train_day = 'False'
)

```

Query complete after 0.00s: 100%|██████████| 3/3 [00:00<00:00, 1313.46query/s]
 Downloading: 100%|██████████| 1/1 [00:00<00:00, 1.06rows/s]

```

eval.head()

  correct_nocancel  wrong_nocancel  wrong_cancel  correct_cancel
0              7684            2935           6787          2942

print(eval['correct_nocancel'] / (eval['correct_nocancel'] + eval['wrong_nocancel']))
print(eval['correct_cancel'] / (eval['correct_cancel'] + eval['wrong_cancel']))

0    0.723609
dtype: float64
0    0.302395
dtype: float64

```

Figure 5-20. Evaluating only at marginal decisions.

As expected, we are correct to not cancel the meeting 72% of the time, close to our target of 70%. We chose the departure delay threshold of 16 minutes on the training dataset because we expected to be 70% correct in not canceling if we do so, and now we've proved on an independent dataset that this is the case. This model achieves the 70% correctness measure that was our target but does so by canceling fewer flights than the contingency table-based model of [Chapter 3](#).

Summary

In this chapter, we began to carry out exploratory data analysis. To be able to interactively analyze our large dataset, we loaded the data into BigQuery, which gave us the ability to carry out queries on millions of rows in a matter of seconds. We required sophisticated statistical plotting capabilities, and we obtained that by using a Jupyter Notebook in the form of Vertex AI Workbench.

In terms of the model itself, we were able to use nonparametric estimation of the 30th percentile of arrival delays, at each departure delay, to pick the departure delay threshold. We discovered that doing this allows us to cancel fewer meetings while attaining the same target correctness. We evaluated our decision threshold on an independent set of flights by dividing our dataset into two parts—a training set and a testing set—based on randomly partitioning the distinct days that comprise our dataset.

Suggested Resources

To learn how to carry out EDA using Python libraries like Matplotlib, NumPy, and Pandas, read the O'Reilly Media book *Hands-On Exploratory Data Analysis with Python* by Suresh Kumar Mukhiya and Usman Ahmed. For a more theoretically grounded introduction to the topic, consider taking the [online course on EDA](#) from Johns Hopkins.

Peruse and work through the [gallery of seaborn plots](#) so that you are familiar with the different ways of visualizing data that are available.

While Jupyter is great for EDA, it is not a great environment for developing stand-alone Python programs. Ultimately, you will want to refactor your notebook code into functions and move them into a Python package. At that point, use a proper integrated development environment like [PyCharm](#) or [Visual Studio Code](#). This is exactly what we will do in this book. To learn this workflow, read this 2018 article by Florian Wilhelm, “[Working Efficiently with JupyterLab Notebooks](#)”. You don't need to install Jupyter because Vertex AI Workbench manages that for you, but the rest of Florian's advice applies.

Bayesian Classifier with Apache Spark on Cloud Dataproc

Having become accustomed to running queries in BigQuery where there were no clusters to manage, I'm dreading going back to configuring and managing Hadoop clusters. But I did promise you a tour of data science on the cloud, and in many companies, Hadoop plays an important role in that.

In this chapter, we tackle the next stage of our data science problem, by creating a Bayes model to predict the likely arrival delay of a flight. We will do this through an integrated workflow that involves BigQuery and Spark SQL.



All of the code snippets in this chapter are available in the folder [06_dataproc](#) of the book's [GitHub repository](#). See the *README.md* file in that directory for instructions on how to do the steps described in this chapter.

MapReduce and the Hadoop Ecosystem

MapReduce was described in [a paper by Jeff Dean and Sanjay Ghemawat](#) as a way to process large datasets on a cluster of machines. They showed that many real-world tasks can be decomposed into a sequence of two types of functions: `map` functions that process key-value pairs to generate intermediate key-value pairs, and `reduce` functions that merge all the intermediate values associated with the same key. A flexible and general-purpose framework can run programs that are written following this MapReduce model on a cluster of commodity machines. Such a MapReduce framework will take care of many of the details that make writing distributed system applications so difficult—the framework, for example, will partition the input data

appropriately, schedule running the program across a set of machines, and handle job or machine failures.

How MapReduce Works

Imagine that you have a large set of documents and you want to compute word frequencies on that dataset. Before MapReduce, this was an extremely difficult problem. One approach you might take would be to scale up—that is, to get an extremely large, powerful machine.¹ The machine will hold the current word frequency table in memory, and every time a word is encountered in the document, this word frequency table will be updated. Here it is in pseudocode:

```
wordcount(Document[] docs):
    wordfrequency = {}
    for each document d in docs:
        for each word w in d:
            wordfrequency[w] += 1
    return wordfrequency
```

We can make this a multithreaded solution by having each thread work on a separate document, sharing the word frequency table between the threads, and updating this in a thread-safe manner. You will at some point, though, run into a dataset that is beyond the capabilities of a single machine. At that point, you will want to scale out, by dividing the documents among a cluster of machines. Each machine on the cluster then processes a fraction of the complete document collection. The programmer implements two methods, `map` and `reduce`:

```
map(String docname, String content):
    for each word w in content:
        emitIntermediate(w, 1)

reduce(String word, Iterator<int> intermediate_values):
    int result = 0;
    for each v in intermediate_values:
        result += v;
    emit(result);
```

The framework manages the orchestration of the maps and reduces and interposes a group-by-key in between (i.e., it's the framework that makes these calls—not the programmer):

```
wordcount(Document[] docs):
    for each doc in docs:
        map(doc.name, doc.content)
    group-by-key(key-value-pairs)
```

¹ In [Chapter 2](#), we discussed scaling up, scaling out, and data in situ from the perspective of data center technologies. This background is useful to have here.

```

for each key in key-values:
    reduce(key, intermediate_values)

```

To improve speed in an environment in which network bisection bandwidth (see [Chapter 2](#)) is low,² the documents are stored on local drives attached to the compute instance. The map operations are then scheduled by the MapReduce infrastructure in such a way that each map operation runs on a compute instance that already has the data it needs (this assumes that the data has been presharded on the cluster), as shown in [Figure 6-1](#).

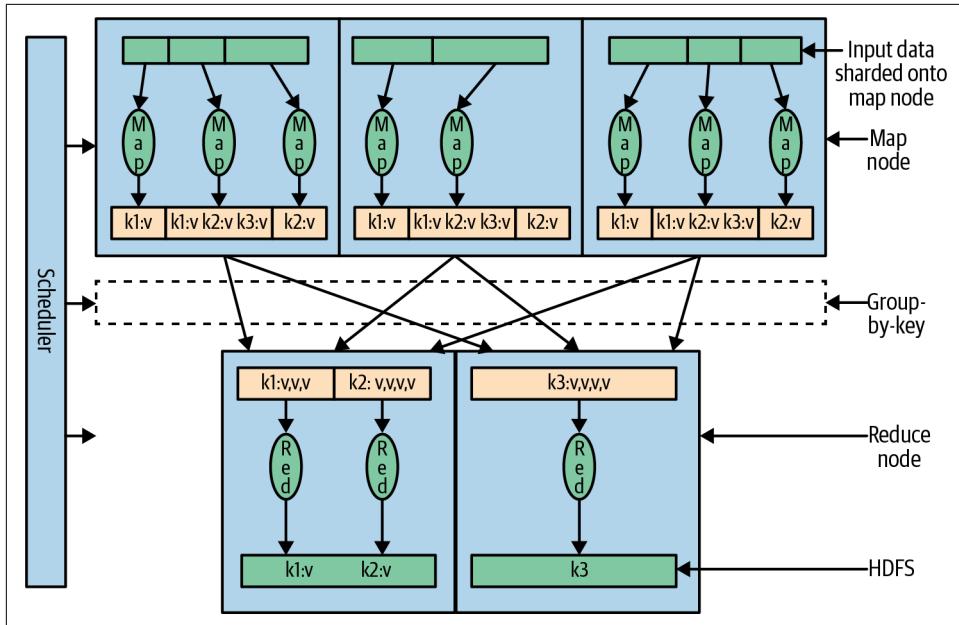


Figure 6-1. MapReduce is an algorithm for distributed processing of datasets in which the data are presharded onto compute instances such that each map operation can access the data it needs using a local filesystem call.

As the diagram indicates, there can be multiple map and reduce jobs assigned to a single machine. The key capability that the MapReduce framework provides is the orchestration and massive group-by-key after the map tasks complete and before the reduce jobs can begin.

² See [Slide 6 of Dean and Ghemawat's original presentation](#)—the MapReduce architecture they proposed assumes that the cluster has limited bisection bandwidth and local, rather slow drives.

Apache Hadoop

When Dean and Ghemawat published the MapReduce paper, they did not make Google's MapReduce implementation open source.³ **Hadoop** is open source software that was created from parts of **Apache Nutch**, an open source web crawler created by Doug Cutting based on a couple of Google papers. Cutting modeled the distributed file system in his crawler on Google's descriptions of the **Google File System** (a predecessor of the **Colossus filesystem** that is in use within Google Cloud Platform today) and the data processing framework on Dean and Ghemawat's MapReduce paper. These two parts were then factored out into Hadoop in 2006 as the Hadoop Distributed File System (HDFS) and the MapReduce engine.

Hadoop today is managed by the Apache Software Foundation. It is a framework that runs applications using the MapReduce algorithm, enabling these applications to process data in parallel on a cluster of commodity machines. Apache Hadoop provides Java libraries necessary to write MapReduce applications (i.e., the `map` and `reduce` methods) that will be run by the framework. In addition, it provides a scheduler, called **YARN**, and a distributed file system (HDFS). To run a job on Hadoop, the programmer submits a job by specifying the location of the input and output files (typically, these will be in HDFS) and uploading a set of Java classes that provide the implementation of the `map` and `reduce` methods.

Google Cloud Dataproc

Normally, the first step in writing Hadoop jobs is to get a Hadoop installation going. This involves setting up a cluster, installing Hadoop on it, and configuring the cluster so that the machines all know about one another and can communicate with one another in a secure manner. Then, you'd start the YARN and MapReduce processes and finally be ready to write some Hadoop programs.

On Google Cloud, **Google Cloud Dataproc** makes it convenient to spin up a Hadoop cluster that is capable of running MapReduce, Pig, Hive, Presto, and Spark.

If you are using Spark, Dataproc offers a fully managed, serverless Spark environment—you can simply submit a Spark program and Dataproc will execute it. In this way, Dataproc is to Apache Spark what Dataflow is to Apache Beam. In fact, Dataproc and Dataflow share backend services. At the time I'm writing this chapter (December 2021), this serverless execution environment in Dataproc supports only Spark, although there are plans to expand it to other frameworks commonly used in Hadoop clusters.

³ Now, research papers from Google are often accompanied by open source implementations—Kubernetes, Apache Beam, TensorFlow, and Inception are examples.

Even if you are not using Spark, Dataproc will still reduce the toil associated with running Hadoop workloads in several ways:

- Dataproc ties into Cloud identity and access management (IAM), Cloud Logging, etc., so that you don't have to manage security or logging on a cluster-by-cluster basis.
- It is autoscaling and will shrink or grow to accommodate your workloads, so you don't have to manage and provision machines yourself.
- It reads directly off Cloud Storage, so you don't have to manage the storage yourself.
- It offers a metadata service so that, even if you run clusters only for the duration of the job, Hive jobs can have persistent metadata.

We can create a fully configured Dataproc cluster by using the following single `gcloud` command:⁴

```
gcloud dataproc clusters create ch6cluster \
    --enable-component-gateway \
    --region us-central1 --zone us-central1-a \
    --master-machine-type n1-standard-4 \
    --master-boot-disk-size 500 --num-workers 2 \
    --worker-machine-type n1-standard-4 \
    --worker-boot-disk-size 500 --image-version 2.0 \
    --properties dataproc:dataproc.personal-auth.user=$EMAIL \
    --optional-components JUPYTER --project $PROJECT \
    --scopes https://www.googleapis.com/auth/cloud-platform
```

A minute or so later, the Cloud Dataproc cluster is created, all ready to go. The `--num-workers`, `--worker-machine-type`, and `--master-machine-type` parameters specify the hardware configuration of the cluster. The `scopes` parameter indicates what Cloud IAM roles this cluster's service account should have. For example, to create a cluster that will run programs that will need to administer Cloud Bigtable and invoke BigQuery queries, you could specify the scope as follows:

```
--scopes=https://www.googleapis.com/auth/bigtable.admin,bigquery
```

Here, I'm allowing the cluster to work with all Google Cloud Platform products. Cloud Dataproc allows you to specify an image version, so that any work you carry out is repeatable. Leave out `--image-version` to use the latest stable version. The `--enable-component-gateway` parameter creates readily accessible, but secure, https proxy endpoints for various services running on the cluster. Besides the standard

⁴ As discussed in [Chapter 3](#), the `gcloud` command makes a REST API call, so this can be done programmatically. You could also use the Google Cloud Platform web console. This command is available in the GitHub repository as `05_dataproc/create_cluster.sh`.

Hadoop services, we also want Jupyter, and so we specify it as an optional component. If your data (that will be processed by the cluster) is in a single-region bucket on Google Cloud Storage, you should create your cluster in that same zone to take advantage of the high bisection bandwidth within a Google data center; that's what the `--zone` specification does.

Although the cluster creation command supports a `--bucket` option to specify the location of a staging bucket to store such things as configuration and control files, best practice is to allow Cloud Dataproc to determine its own staging bucket. This allows you to keep your data separate from the staging information needed for the cluster to carry out its tasks. Cloud Dataproc will create a separate bucket in each geographic region, choose an appropriate bucket based on the zone in which your cluster resides, and reuse such Cloud Dataproc-created staging buckets between cluster create requests if possible.⁵

Because we specified `--enable-component-gateway`, we can verify that Hadoop is running by visiting the Cloud Dataproc section of the Google Cloud Platform web console and accessing the HDFS NameNode web interface (from the Web Interfaces section of the cluster details). You should be able to see the list of data nodes.



If you want to use Secure Shell (SSH) to connect to the cluster, you can, but you'd have to give the master node an external IP in order to do so. This is generally not a good idea. Instead, interact with the cluster through the available web interfaces. Later in this chapter, I'll show you how to install software on startup so that you don't need to SSH into the cluster to install software.

Need for Higher-Level Tools

The word count example is embarrassingly parallel, and therefore trivial to implement in terms of a single map and a single reduce operation. However, it is nontrivial to cast more complex data processing algorithms into sequences of map and reduce operations. Higher-level solutions are called for, and as different organizations implemented add-ons to the basic Hadoop framework and made these additions available as open source, the Hadoop ecosystem was born.

[Apache Pig](#) provided one of the first ways to simplify the writing of MapReduce programs to run on Hadoop. Apache Pig requires you to write code in a language called *Pig Latin*; these programs are then converted to sequences of MapReduce programs, and these MapReduce programs are executed on Apache Hadoop. Because Pig Latin

⁵ To find the name of the staging bucket created by Cloud Dataproc, run `gcloud dataproc clusters describe`.

(sometimes just referred to as Pig) comes with a command-line interpreter, it is very conducive to interactive creation of programs meant for large datasets. At the same time, it is possible to save the interactive commands and execute the script on demand. This provides a way to achieve both embarrassingly parallel data analysis and data flow sequences consisting of multiple interrelated data transformations. Pig can optimize the execution of MapReduce sequences, thus allowing the programmer to express tasks naturally without worrying about efficiency.

Apache Hive provides a mechanism to project structure onto data that is already in distributed storage. With the structure (essentially a table schema) projected onto the data, it is possible to query, update, and manage the dataset using SQL. Typical interactions with Hive use a command-line tool or a Java Database Connectivity (JDBC) driver.

Pig and Hive both rely on the distributed storage system to store intermediate results. **Apache Spark**, on the other hand, takes advantage of in-memory processing and a variety of other optimizations. Because many data pipelines start with large, out-of-memory data, but quickly aggregate it to something that can be fit into memory, Spark can provide dramatic speedups when compared to Pig and as well as speedups for Spark SQL when compared to Hive.⁶ In addition, because Spark (like Pig and Big-Query) optimizes the directed acyclic graph (DAG) of successive processing stages, it can provide gains over handwritten Hadoop operations. With the growing popularity of Spark, a variety of machine learning, data mining, and streaming packages have been written for it. Hence, in this chapter, we focus on Spark solutions. Cloud Dataproc, though, provides an execution environment for Hadoop jobs regardless of the abstraction level (i.e., whether you submit jobs in Hadoop, Pig, Hive, or Spark). All these software packages are installed by default on Cloud Dataproc.

Jobs, Not Clusters

We will look at how to submit jobs to the Cloud Dataproc clusters shortly, but after you are done with the cluster, delete it by using the following:

```
gcloud dataproc clusters delete ch6cluster
```

You can even set the cluster up so that it is automatically deleted if it's idle for a specific time duration.

This is not the typical Hadoop workflow—if you are used to an on-premises Hadoop installation, you might have set up the cluster a few months ago and it has remained up since then. The better practice on Google Cloud Platform, however, is to delete the cluster after you are done. The reasons are twofold. First, it typically takes less than

⁶ Hive has been sped up in recent years through the use of a new application framework ([Tez](#)) and a variety of optimizations for [long-lived queries](#).

two minutes to start a cluster. Because cluster creation is fast and can be automated, it is wasteful to keep unused clusters around—you are paying for the cluster regardless of whether you are running anything useful on them. Second, one reason that on-premises Hadoop clusters are kept always on is because the data is stored on HDFS. Although you can use HDFS in Cloud Dataproc (recall that we looked at HDFS NameNode to get the status of the Hadoop cluster), it is not recommended. Instead, it is better to keep your data on Google Cloud Storage and directly read from Cloud Storage in your MapReduce jobs—the original MapReduce practice of assigning map processes to nodes that already have the necessary data came about in an environment in which network bisection speeds were low. On the Google Cloud Platform, for which network bisection speeds are on the order of a petabit per second, the best practice has changed. Instead of sharding your data onto HDFS, keep your data on Cloud Storage and read the data into an ephemeral cluster, as demonstrated in [Figure 6-2](#).

Because of the high network speed that prevails within the Google data center, reading from Cloud Storage is competitive with HDFS in terms of speed for sustained reads of large files (the typical Hadoop use case). If your use case involves frequently reading small files, reading from Cloud Storage could be slower than reading from HDFS. However, even in this scenario, you can counteract this lower speed by simply creating more compute nodes—because storage and compute are separate, you are not limited to the number of nodes that happen to have the data. Because Hadoop clusters tend to be underutilized, you will often save money by creating an ephemeral cluster many times the size of an always-on cluster with an HDFS filesystem. Getting the job done quickly with a lot more machines and deleting the cluster when you are done is often the more frugal option (you should measure this on your particular workflow, of course, and estimate the cost of different scenarios).⁷ This method of operating with short-lived clusters is also quite conducive to the use of [preemptible instances](#)—you can create a cluster with a given number of standard instances and many more preemptible instances, thus getting a lower cost for the full workload.

⁷ For a tool to estimate costs quickly, go to the [Google Cloud Pricing Calculator](#).

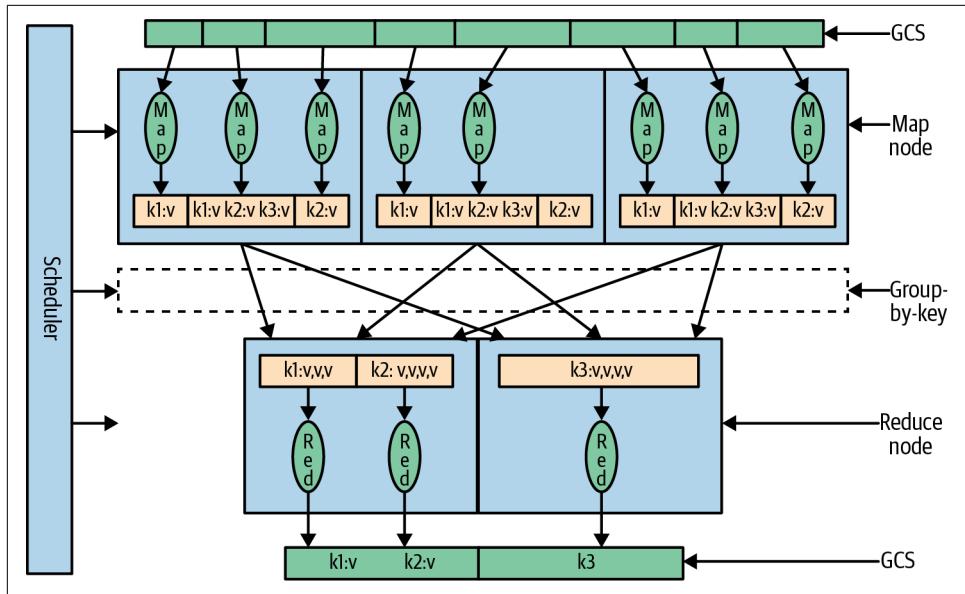


Figure 6-2. Because network bisection speeds on Google Cloud are on the order of a petabit per second, best practice is to keep your data on Cloud Storage and simply spin up short-lived compute nodes to perform the map operations. These nodes will read the data across the network. In other words, there is no need to preshard the data.

Preinstalling Software

Creating and deleting clusters on demand is fine if you want a plain, vanilla Hadoop cluster, but what if you need to install specific software on the individual nodes?

There are two approaches. One is to create your own custom Docker images and ask Dataproc to use those:

```
gcloud dataproc clusters create --image=...
```

You can create these images starting from an existing Dataproc base image and adding any other packages you require in your Dockerfile.

The second option is to use *initialization actions*. These are simply startup executables, stored on Cloud Storage, that will be run on the nodes of the cluster. For example, suppose that we want a specific Python package:

- Create a script to carry out whatever software we want preinstalled:⁸

```
#!/bin/bash

# Things to do on both Master and Worker
apt-get -y update
apt-get install python-dev
apt-get install python-pip
pip install --upgrade google-api-python-client

ROLE=$(/usr/share/google/get_metadata_value attributes/dataproc-role)
if [[ "${ROLE}" == 'Master' ]]; then
    cd home/dataproc
    git clone https://github.com/GoogleCloudPlatform/data-science-on-gcp
fi
```

Now, when the cluster is created, the specified packages will exist on all the nodes and the GitHub repository will exist on the Master node.

- Save the script on Cloud Storage:

```
#!/bin/bash
BUCKET=cloud-training-demos-ml
ZONE=us-central1-a
INSTALL=gs://$BUCKET/flights/dataproc/install_on_cluster.sh

# upload install file
gsutil cp install_on_cluster.sh $INSTALL
```

- Supply the script to the cluster creation command:⁹

```
gcloud dataproc clusters create \
    --num-workers=2 \
    ...
    --initialization-actions=$INSTALL \
    ch6cluster
```

Some components, like Jupyter, are already available for installation in Cloud Dataproc. For Jupyter, we could get away with just specifying it as one of the `--optional-components` to be installed.

⁸ One efficient use of initialization actions is to preinstall all the third-party libraries you might need, so that they don't have to be submitted with the job. This script is `06_dataproc/install_on_cluster.sh`.

⁹ This script is in the GitHub repository of this book as `06_dataproc/create_cluster.sh`.

Quantization Using Spark SQL

So far, we have used only one variable in our dataset—the departure delay—to make our predictions of the arrival delay of a flight. However, we know that the distance the aircraft needs to fly must have some effect on the ability of the pilot to make up for delays en route. The longer the flight, the more likely it is that small delays in departure can be made up in flight. So, let's build a statistical model that uses two variables—the departure delay and the distance to be traveled.

One way to do this is to put each flight into one of several bins, as shown in [Table 6-1](#).

Table 6-1. Quantizing distance and departure delay to carry out Bayesian classification over two variables

	< 10 min	10–12 min	12–15 min	> 15 min
< 100 miles	For example:			
	<ul style="list-style-type: none">Arrival Delay \geq 15 min: 150 flightsArrival Delay < 15 min: 850 flights85% of flights have arrival delay < 15 minutes			
100–500 miles				
> 500 miles				

For each bin, I can look at the number of flights within the bin that have an arrival delay of more than 15 minutes and the number of flights with an arrival delay of less than 15 minutes, and then determine which category is higher. The majority vote then becomes our prediction for that entire bin. Because our threshold for decisions is 70% (recall that we want to cancel the meeting if there is a 30% likelihood that the flight will be late), we'll recommend canceling the meeting for flights that fall into a bin if the fraction of arrival delays of less than 15 minutes is less than 0.7. This method is called *Bayesian classification*, and the statistical model is simple enough that we can build it from scratch with a few lines of code.

The probability that the flight will be late given that the distance x_0 is 120 miles and the departure delay x_1 is 8 minutes is called the *conditional probability*,¹⁰ written as $P(C_{\text{late}} \mid x_0, x_1)$. Within each bin, we are calculating the conditional probability $P(C_{\text{ontime}} \mid x_0, x_1)$ and $P(C_{\text{late}} \mid x_0, x_1)$ where (x_0, x_1) is the pair of predictor variables (mileage and departure delay) and C_k is one of two classes depending on the value of the arrival delay of the flight. Because the probability of a specific value of a continuous variable is zero, we need to estimate the probability over an interval, and, in this case, the intervals are given by the bins. Thus, to estimate $P(C_{\text{ontime}} \mid x_0, x_1)$,

¹⁰ For a good, intuitive introduction to conditional probability, see [Statistics How To](#).

we find the bin that (x_0, x_1) falls into and use that as the estimate of $P(C_{ontime})$. If this is less than 70%, our decision will be to cancel the meeting.

Of all the ways of estimating a conditional probability, the way we are doing it—by divvying up the dataset based on the values of the variables—is the easiest, but it will work only if we have large enough populations in each of the bins. This method of directly computing the probability tables works with two variables, but will it work with 20 variables? How likely is it that there will be enough flights for which the departure airport is TUL, the distance is about 350 miles, the departure delay is about 10 minutes, the taxi-out time is about 4 minutes, and the hour of day that the flight departs is around 7 a.m?

As the number of variables increases, we will need more sophisticated methods in order to estimate the conditional probability. A scalable approach that we can employ if the predictor variables are independent is a method called *Naive Bayes*. In the Naive Bayes approach, we compute the probability tables by taking each variable in isolation (i.e., computing $P(C_{ontime} | x_0)$ and $P(C_{ontime} | x_1)$ separately) and then multiplying them to come up with $P(C_k | x_i)$.¹¹ However, for just two variables, for a dataset this big, we can get away with binning the data and directly estimating the conditional probability.

JupyterLab on Cloud Dataproc

Developing the Bayesian classification from scratch requires being able to interactively carry out development. Although we could spin up a Cloud Dataproc cluster, connect to it via SSH, and do development on the Spark read–eval–print loop (REPL), it would be better to use JupyterLab and get a notebook experience similar to how we worked with BigQuery in [Chapter 5](#).

Among the web interfaces that we enabled with `--enable-component-gateway` was that for JupyterLab. Hence, we can connect to it similar to the way we connected to the HDFS NameNode, from the Google Cloud web console, in the Web Interfaces part of the cluster details section.

In Jupyter, use the File Browser on the left to navigate to `/home/dataproc` on the local disk and open the notebook `06_dataproc/quantization.ipynb` in the clone of the course repository that you find there.

¹¹ The exact calculation involves dividing by a scaling factor so that the outcome is the probability. See the [Wikipedia entry on the Naive Bayes classifier](#) for more details on the mathematics.

Independence Check Using BigQuery

Before we can get to computing the proportion of delayed flights in each bin, we need to decide how to quantize the delay and distance. What we do not want are bins with very few flights—in such cases, statistical estimates will be problematic. In fact, if we could somehow spread the data somewhat evenly between bins (using *quantization*), it would be ideal.

For simplicity, we would like to choose the quantization thresholds for distance and for departure delay separately, but we can do this only if they are relatively independent. Let's verify that this is the case. Cloud Dataproc is integrated with the managed services on Google Cloud Platform, so even though we have our own Hadoop cluster, we can still call out to BigQuery from the notebook that is running on Cloud Dataproc. Using BigQuery, Pandas, and seaborn as we did in [Chapter 5](#), here's what the query looks like:

```
sql = """
SELECT DISTANCE, DEP_DELAY
FROM dsongcp.flights_tzcorr
WHERE RAND() < 0.001 AND dep_delay > -20
    AND dep_delay < 30 AND distance < 2000
"""
df = bq.query(sql).to_dataframe()
sns.set_style("whitegrid")
g = sns.jointplot(x=df['DISTANCE'], y=df['DEP_DELAY'], kind="hex",
                   height=10, joint_kws={'gridsize':20})
```

The query samples the full dataset, pulling in 1/1,000 of the flights' distance and departure delay fields (that lie within reasonable ranges) into a Pandas dataframe. This sampled dataset is sent to the seaborn plotting package and a hexbin plot is created. The resulting graph is shown in [Figure 6-3](#).

Each hexagon of a hexagonal bin plot is colored based on the number of flights in that bin, with darker hexagons indicating more flights. It is clear that at any distance, a wide variety of departure delays is possible and for any departure delay, a wide variety of distances is possible. The distribution of distances and departure delays in turn is similar across the board. There is no obvious trend between the two variables—in [Figure 6-3](#), note that the Pearson correlation coefficient is 0.07. This indicates that we can treat the two variables as independent.

The distribution plots at the top and right of the center panel of the graph show how the distance and departure delay values are distributed. This will affect the technique that we can use to carry out quantization. Note that the distance is distributed relatively uniformly until about 1,000 miles, beyond which the number of flights begins to taper off. The departure delay, on the other hand, has a long tail and is clustered around -5 minutes. We might be able to use equispaced bins for the distance variable (at least in the 0- to 1,000-mile range), but for the departure delay variable, our bin

size must be adaptive to the distribution of flights. In particular, our bin size must be wide in the tail areas and relatively narrow where there are lots of points.

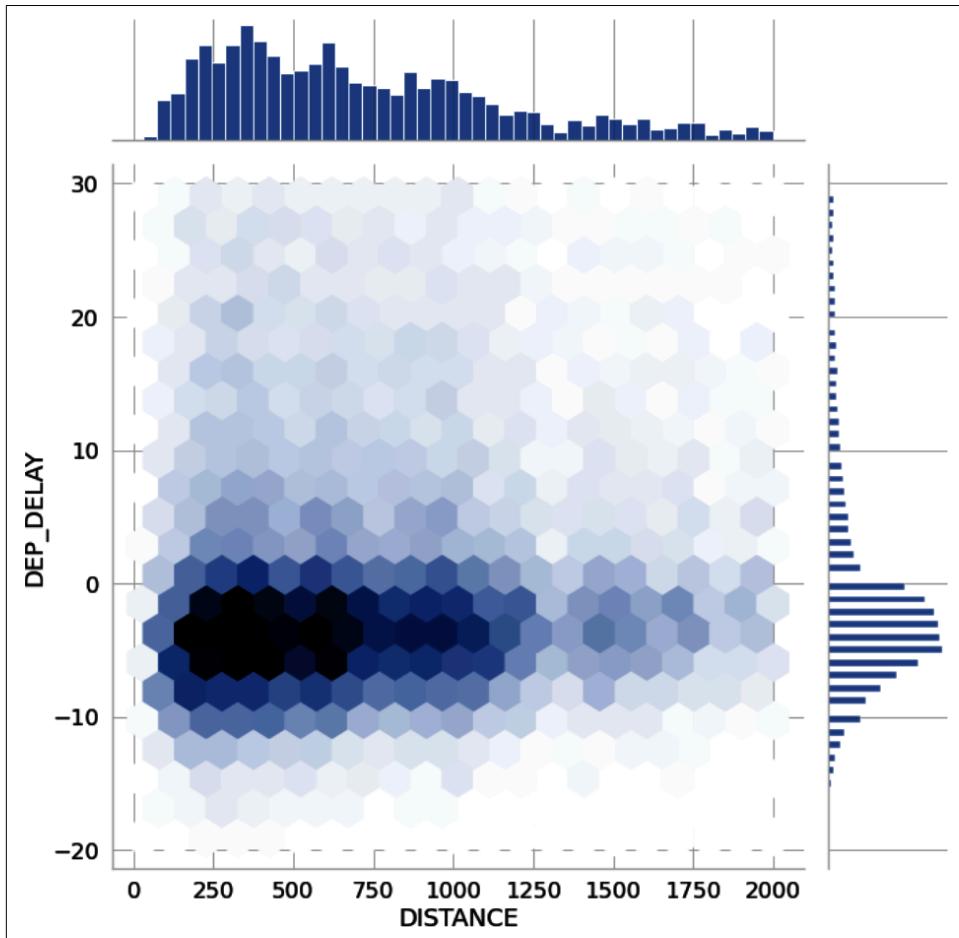


Figure 6-3. The hexbin plot shows the joint distribution of departure delay and the distance flown. You can use such a plot to verify whether the fields in question are independent.

There is one issue with the hexbin plot in Figure 6-3: we have used data that we are not allowed to use. Recall that our model must be developed using only the training data. While we used it only for some light exploration, it is better to be systematic about excluding days that will be part of our evaluation dataset from all model development. To do that, we need to join with the `traindays` table and retain only days for which `is_train_day` is `True`. We could do that in BigQuery, but even though Cloud Dataproc is integrated with other Google Cloud Platform services, invoking BigQuery from a Hadoop cluster feels like a cop-out. So, let's try to recreate the same

plot as before, but this time using Spark SQL, and this time using only the training data.

Spark SQL in JupyterLab

A Spark session can be created by typing the following into a code cell:

```
from pyspark.sql import SparkSession
spark = SparkSession \
    .builder \
    .appName("Bayes classification using Spark") \
    .getOrCreate()
```

With the `spark` variable in hand, we can read in the time-corrected JavaScript Object Notation (JSON) files that we wrote to Google Cloud Storage in [Chapter 4](#):

```
inputs = 'gs://{}//flights/tzcorr/all_flights-*'.format(BUCKET)
flights = spark.read.json(inputs)
```

Even though we do want to ultimately read all the flights and create our model from all of the data, we will find that development goes faster if we read a fraction of the dataset. So, let's change the input from `all_flights-*` to `all_flights-00000-*`:

```
inputs = 'gs://{}//flights/tzcorr/all_flights-00000-*'.format(BUCKET))
```

Because I had 26 JSON files, doing this change means that I will be processing just the first file, and we will notice an increase in speed of 26 times during development. Of course, we should not draw any conclusions from processing such a small sample other than that the code works as intended.¹² After the code has been developed on 4% of the data, we'll change the string so as to process all the data and increase the cluster size so that this is also done in a timely manner. Doing development on a small sample on a small cluster ensures that we are not underutilizing a huge cluster of machines while we are developing the code.

With the `flights` dataframe created as shown previously, we can employ SQL on the dataframe by creating a temporary view (it is available only within this Spark session):

```
flights.createOrReplaceTempView('flights')
```

¹² Just as an example, the Google Cloud Dataflow job that wrote out this code could have ordered the JSON file by date, and in that case, this file will contain only the first 14 days of the year.

Now, we can employ SQL to query the `flights` view, for example by doing this:

```
results = spark.sql('SELECT COUNT(*) FROM flights WHERE dep_delay > -20 AND distance < 2000')
results.show()
```

On my development subset, this yields the following result:

```
+-----+
|count(1)|
+-----+
| 59665 |
+-----+
```

This is just about right to comfortably fit in memory, but even if it were somewhat larger I dare not go any smaller than 2%–3% of the data, even in development.

To create the `traindays` dataframe, we can follow the same steps, but for a CSV file this time:

```
traindays = spark.read \
    .option("header", "true") \
    .option("inferSchema", "true") \
    .csv('gs://{}//flights/trainday.csv'.format(BUCKET))
traindays.createOrReplaceTempView('traindays')
```

A quick check illustrates that `traindays` has been read, and the column names and types are correct:

```
results = spark.sql('SELECT * FROM traindays')
results.head(5)
```

This yields the following:

```
[Row(FL_DATE='2015-01-01', is_train_day=True),
 Row(FL_DATE='2015-01-02', is_train_day=False),
 Row(FL_DATE='2015-01-03', is_train_day=False),
 Row(FL_DATE='2015-01-04', is_train_day=True),
 Row(FL_DATE='2015-01-05', is_train_day=True)]
```

To restrict the `flights` dataframe to contain only training days, we can do a SQL join:

```
statement = """
SELECT
    f.FL_DATE AS date,
    CAST(distance AS FLOAT) AS distance,
    dep_delay,
    IF(arr_delay < 15, 1, 0) AS ontime
FROM flights f
JOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day AND
    f.dep_delay IS NOT NULL
ORDER BY
```

```

    f.dep_delay DESC
"""
flights = spark.sql(statement)

```

Now, we can use the `flights` dataframe for the hexbin plots after clipping the x-axis and y-axis to reasonable limits:

```

df = flights[(flights['distance'] < 2000) & \
             (flights['dep_delay'] > -20) & \
             (flights['dep_delay'] < 30)]

```

When we drew the hexbin plot in the previous section, we sampled the data to 1/1,000, but that was because we were passing in a Pandas dataframe to seaborn. This sampling was done so that the Pandas dataframe would fit into memory. However, whereas a Pandas dataframe must fit into memory, a Spark dataframe does not. As of this writing (i.e., December 2021), though, there is no way to directly plot a Spark dataframe either—you must convert it to a Pandas dataframe; therefore, we will still need to sample it, at least when we are processing the full dataset.

Because there are about 50,000 rows on 1/25 of the data, we expect the full dataset to have about 6 million rows. Let's sample this down to about 100,000 records, which would be about 0.02 of the dataset:

```

pdf = df.sample(False, 0.02, 20).toPandas()
g = sns.jointplot(x=pdf['distance'], y=pdf['dep_delay'], kind="hex",
                   height=10, joint_kws={'gridsize':20})

```

This yields a hexbin plot that is not very different from the one we ended up with in the previous section. The conclusion—that we need to create adaptive-width bins for quantization—still applies. Just to be sure, though, this is the point at which I'd repeat the analysis on the entire dataset to ensure our deductions are correct had I done only the Spark analysis. However, we did do it on the entire dataset in BigQuery, so let's move on to creating adaptive bins.

Histogram Equalization

To choose the quantization thresholds for the departure delay and the distance in an adaptive way (wider thresholds in the tails and narrower thresholds where there are a lot of flights), we will adopt a technique from image processing called *histogram equalization*.¹³

Low-contrast digital images have histograms of their pixel values distributed such that most of the pixels lie in a narrow range. Take, for example, the photograph in Figure 6-4.¹⁴

¹³ For examples of histogram equalization applied to improve the contrast of images, go to [OpenCV.org](#).

¹⁴ Photograph by the author.



Figure 6-4. Original photograph of the pyramids of Giza used to demonstrate histogram equalization.

As depicted in [Figure 6-5](#), the histogram of pixel values in the Pyramids image is clustered around two points: the dark pixels in the shade, and the bright pixels in the sun.



Figure 6-5. Histogram of pixel values in photograph of the pyramids.

Let's remap the pixel values such that the full spectrum of values is present in the image, so that the new histogram looks like that shown in [Figure 6-6](#).

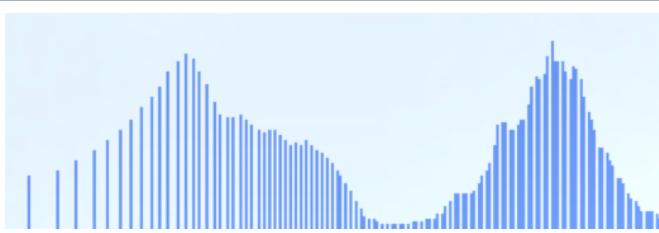


Figure 6-6. Histogram of pixels after remapping the pixels to occupy the full range.

The remapping is of pixel values and has no spatial component. For example, all pixel values of 125 in the old image might be changed to a pixel value of 5 in the new image

regardless of where they are in terms of horizontal and vertical position. [Figure 6-7](#) presents the remapped image.



Figure 6-7. Note that after the histogram equalization, the contrast in the image is enhanced.

What we implicitly did was to remap the pixel values, such that each section of the spectrum from black to white now has approximately the same number of pixels (whereas previously they were all in the gray middle). Histogram equalization has helped to enhance the contrast in the image and bring out finer details. Look, for example, at the difference in the rendering of the sand in front of the pyramid or of the detail of the midsection of Khafre's pyramid (the tall one in the middle).¹⁵

How is this relevant to what we want to do? We are also looking to remap values when we seek to find quantization intervals. In the image case, the output had the same range as the input. But in our flight example, we'd like to remap a distance value of 422 miles to a quantized value of perhaps 3. As in histogram equalization, we want the bin values to be uniformly distributed. We can, therefore, apply the same technique as is employed in the image processing filter to achieve this.

What we want to do is to divide the spectrum of distance values into, say, five bins. The first bin will contain all values in $[0, d_0]$, the second will contain values in $[d_0, d_1]$, and so on, until the last bin contains values in $[d_4, \infty)$. Histogram equalization requires that d_0, d_1 , and so on be such that the number of flights in each bin is approximately equal—that is, for the data to be uniformly distributed after quantization. As in the example photograph of the pyramids, it won't be perfectly uniform because the

¹⁵ Not the tallest, though. Khufu's pyramid (the tall pyramid in the forefront) is taller and larger, but has been completely stripped of its alabaster topping and is situated on slightly lower ground.

input values are also discrete. However, the goal is to get as close to an equalized histogram as possible.

With histogram equalization, at any specific departure delay, the number of flights at each distance and delay bin should remain large enough that our conclusions are statistically valid. Assuming independence and 6 million total flights, if we divvy up the data into 100 bins (10 bins per variable), we will have about 60,000 flights in each bin. That's probably still okay, but let's be safe and divvy up the data into just five bins each. Divvying up the data into five bins implies a probability range of 0, 0.2, ..., 0.8 or five probabilistic thresholds:

```
np.arange(0, 1.0, 0.2)
```

Finding thresholds that make the two quantized variables uniformly distributed is quite straightforward using the approximate quantiles method discussed in [Chapter 5](#). There is an `approxQuantile()` method available on the Spark dataframe also:

```
distthresh = flights.approxQuantile('distance',
                                     list(np.arange(0, 1.0, 0.2)), 0.02)
delaythresh = flights.approxQuantile('dep_delay',
                                      list(np.arange(0, 1.0, 0.2)), 0.02)
```

On the development dataset, here's what the distance thresholds turn out to be:

```
[130.0, 370.0, 621.0, 1009.0]
```

The zeroth percentile is essentially the minimum. The next ones are the 25th percentile, median, and 75th percentile. In order to have all bin boundaries, we can tack on infinity at the end:

```
distthresh[-1] = float('inf')
```

Other than setting the policy (histogram equalization), we don't need to be in the business of choosing distance thresholds. This automation is important because it allows us to dynamically update thresholds if necessary on the most recent data,¹⁶ taking care to use the same set of thresholds in prediction as was used in training.

We can similarly quantize the departure delay thresholds into equal boundaries, and we get:

```
[-22.0, -5.0, -3.0, 0.0, inf]
```

Unfortunately, this variable is not as well-behaved as the distance—more than 75% of flights depart on-time or early, so the really interesting delayed departures are all hidden in the last bin. This is going to be a problem that we will fix shortly.

¹⁶ If, for example, there is a newfangled technological development that enables pilots to make up time in the air better or, more realistically, if new regulations prompt airline schedulers to start padding their flight-time estimates.

Bayesian Classification

Now that we have the quantization thresholds, what we need to do is find out the recommendation (whether to cancel the meeting) for each bin based on whether 70% of flights in that bin are on time or not.

Bayes in Each Bin

We can find the flights that belong to the m th distance bin and n th delay bin by slicing the full set of flights:

```
bdf = flights[(flights['distance'] >= distthresh[m])  
    & (flights['distance'] < distthresh[m+1])  
    & (flights['dep_delay'] >= delaythresh[n])  
    & (flights['dep_delay'] < delaythresh[n+1])]
```

Once we do that, we can compute the fraction of flights that arrive on time for this bin:

```
ontime_frac = (bdf.agg(F.sum('ontime')).collect()[0][0] /  
                bdf.agg(F.count('ontime')).collect()[0][0])
```

Looping through the first on-time fractions for the first few bins, we immediately notice a problem:

m	n	ontime_frac
0	0	0.9853403141361257
0	1	0.9847756410256411
0	2	0.9753028890959925
0	3	0.6346045989904655
1	0	0.9721913236929922
1	1	0.9650856389986825
1	2	0.9711299153807864
1	3	0.5715380684721513

The on-time fraction is nearly 100% for all the delay bins except the largest value for n . This makes perfect sense because only the last departure delay bin has any delayed flights.

We'll have to fix this—one way to do so is to hand-select the departure delay bins. Because we already looked at thresholding the departure delay in [Chapter 3](#), we know that the interesting range is between 10 and 20 minutes and that departure delays are reported in integer minutes. So, we simply need to try delay variables of 10, 11, 12, ..., 20 minutes.

So, let's change the delay thresholds and look at them in increments of one minute:

```
delaythresh = range(10, 20)
```

To find the delay threshold for each distance threshold where the value is closest to the 0.70 decision boundary (see *quantization.ipynb* in the GitHub repository):

```

df['score'] = abs(df['frac_ontime'] - 0.7)
bayes = (df.sort_values(['score']).groupby('dist_thresh')
         .head(1).sort_values('dist_thresh'))

```

The resulting model is a *lookup table* that consists of a delay threshold for each distance bin:

Distance bin	delay_thresh
[130, 370]	17
[370, 621]	13
[621, 1009]	17
[1009, inf]	18

If the departure delay is greater than the threshold corresponding to how far the flight is, then we will cancel the meeting because we expect the flight to be late. We are finding the delay beyond which we need to cancel flights for each distance and saving just that threshold. This makes it quite easy to productionize the model—just write out the preceding table as a CSV file, perhaps, and ask the developer of the application to apply the appropriate threshold based on the lookup table.

For example, what is the appropriate decision for a flight with a distance of 800 miles that departs 16 minutes late? The flight falls into the [621, 1009] bin. For such flights, we need to cancel the meeting only if the flight departs 17 or more minutes late—a shorter departure delay is something that can be made up en route.

Evaluating the Model

How well does this model do? To evaluate the model, we have to look at flights that were not used in creating the model. The held-out days are obtained by looking for:

```
t.is_train_day == 'False'
```

We can compute the contingency table values for any given bin using:

```

SELECT
    ROUND(SUM(IF(dep_delay < {2:f} AND arr_delay < 15, 1, 0))/COUNT(*), 2)
        AS correct_nocancel,
    ROUND(SUM(IF(dep_delay >= {2:f} AND arr_delay < 15, 1, 0))/COUNT(*), 2)
        AS false_positive,
    ROUND(SUM(IF(dep_delay < {2:f} AND arr_delay >= 15, 1, 0))/COUNT(*), 2)
        AS false_negative,
    ROUND(SUM(IF(dep_delay >={2:f} AND arr_delay >= 15, 1, 0))/COUNT(*), 2)
        AS correct_cancel,
    COUNT(*) AS total_flights
FROM flights f
JOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day == 'False' AND

```

```

f.distance >= {0:f} AND f.distance < {1:f}
""".format( distthresh[m], distthresh[m+1],
    bayes[
        bayes['dist_thresh'] == distthresh[m]
        ]['delay_thresh'].values[0] )

```

When I did this, I got the results shown in [Figure 6-8](#). We can not put too much stock in this model, though, because it was trained on just 1/25 of the data. Let's fix that next.

bin	correct_nocancel	false_positive	false_negative	correct_cancel	total_flights
130-370 miles	0.82	0.02	0.02	0.13	3131
370-621 miles	0.77	0.03	0.03	0.15	3626
621-1009 miles	0.8	0.03	0.03	0.14	3782
1009-100000 miles	0.8	0.04	0.05	0.11	7624

Figure 6-8. Results of evaluating the Bayes model.

Dynamically Resizing Clusters

The thresholds in the previous section have been computed on about 1/25 of the data (recall that our input was only one shard: `all-flights-00000-of-*`). So, we should find the actual thresholds that we will want to use by repeating the processing on all of the training data at hand. To do this in a timely manner, we will also want to increase our cluster size. Fortunately, we don't need to bring down our Cloud Data-proc cluster in order to add more nodes.

Let's add machines to the cluster so that it has 20 workers, 15 of which are *secondary* and so are heavily discounted in price:¹⁷

```
gcloud dataproc clusters update ch6cluster\  
--num-secondary-workers=15 --num-workers=5 --region=us-central1
```

The secondary machines are preemptible. These machines are provided by Google Cloud Platform at a large (fixed) discount to standard Google Compute Engine instances in return for users' flexibility in allowing the machines to be taken away at very short notice.¹⁸ They are particularly helpful on Hadoop workloads because Hadoop is fault-tolerant and can deal with machine failure—it will simply reschedule those jobs on the machines that remain available. Using preemptible machines on your jobs is a frugal choice—here, the five standard workers are sufficient to finish the task in a reasonable time. However, the availability of 15 more machines means that our task could be completed four times faster and much more inexpensively than if we have only standard machines in our cluster.¹⁹

We can navigate to the Google Cloud Platform console in a web browser and check that our cluster now has 20 workers, as illustrated in Figure 6-9.

A screenshot of the Google Cloud Platform Clusters page. At the top, there is a header with the word "Clusters", a "CREATE CLUSTER" button, and a refresh icon. Below the header is a table with three columns: "Name", "Zone", and "Total worker nodes". There is one row in the table, representing the cluster "ch6cluster" located in the zone "us-central1-a" with a total of 20 worker nodes.

Name	Zone	Total worker nodes
ch6cluster	us-central1-a	20

Figure 6-9. The cluster now has 20 workers.

Now, go to the JupyterLab Notebook and change the input variable to process the full dataset. Next, in the JupyterLab Notebook, click Kernel > “Restart Kernel and Clear

¹⁷ Trying to increase the number of workers might have you hitting against (soft) quotas on the maximum number of CPUs, drives, or addresses. If you hit any of these soft quotas, request an increase from the Google Cloud Platform console's [section on Compute Engine quotas](#). Besides the necessary CPU quota, you may need to ask for an increase in [Persistent Disk](#) and [In-use IP addresses](#). Because a Cloud Dataproc cluster is in a single region, these are *regional* quotas. See the [documentation on resource quotas in Compute Engine](#) for details. In [Chapter 7](#), I had to ask for additional CPUs, and the process of getting a quota increased is explained there as well. If you are in an organization where increasing the quota is a bureaucratic process, ask for the larger quota you will need for [Chapter 7](#) now.

¹⁸ Less than a minute's notice as of this writing in December 2021.

¹⁹ If the preemptible instances cost 20% of a standard machine (as they do as of this writing in December 2021), the 15 extra machines cost us only as much as three standard machines.

All Outputs” to avoid mistakenly using a value corresponding to the development dataset. Then, click on Run > “Run all Cells.”

All the graphs and charts are updated. After we have the results, we can resize the cluster back to something smaller so that we are not wasting cluster resources:

```
gcloud dataproc clusters update ch6cluster\  
--num-secondary-workers=0 --num-workers=2
```

On the full dataset, the lookup table is:

Distance bin	delay_thresh
[31, 328]	14
[328, 541]	15
[541, 802]	15
[802, inf]	17

Note that the thresholds changed (the quantiles are different once we add the remaining 95% of information). The delay threshold also changes quite smoothly as the distance increases. The behavior matches our intuition that we can be tolerant of longer delays on longer flights.

Comparing to Single Threshold Model

Yes, but is this better than the single, universal threshold that we used in [Chapter 5](#)? How well does this new two-variable model perform? We can modify the evaluation BigQuery query from [Chapter 5](#) to add in a distance criterion and supply the appropriate threshold for that distance:

```
SELECT  
    SUM(IF(DEP_DELAY = 14  
        AND arr_delay < 15,  
        1,  
        0)) AS wrong_cancel,  
    SUM(IF(DEP_DELAY = 14  
        AND arr_delay >= 15,  
        1,  
        0)) AS correct_cancel  
FROM (  
    SELECT  
        DEP_DELAY,  
        ARR_DELAY  
    FROM  
        dsongcp.flights_tzcorr f  
    JOIN  
        dsongcp.trainday t  
    ON  
        f.FL_DATE = t.FL_DATE  
    WHERE
```

```
t.is_train_day = 'False'
AND f.DISTANCE < 328)
```

In this query, 14 minutes is the newly determined threshold for distances under 328 miles and the WHERE clause is now limited to flights over distances of less than 328 miles. The result is:

Row	wrong_cancel	correct_cancel
1	1244	582

This indicates that we cancel meetings when it is correct to do so $582 / (582 + 1,244)$ or 32% of the time—remember that our goal was 30%. Similarly, we can do the other four distance categories. Both with this model and with a model that took into account only the departure delay (as in [Chapter 5](#)), we are able to get reliable predictions—canceling meetings when the flight has more than a 30% chance of being delayed.

When we have two models that perform equally well on the primary metric, it is possible that we can see if they differ on a secondary metric that also matters to us. Even if two models have the same reliability (of 30%), a model that allows us to achieve that reliability while canceling fewer meetings would be preferable. Or perhaps there is a certain category of meetings that are more important than others. The more complex model is worthwhile if we end up canceling fewer important meetings or if we can be more fine-grained in our decisions (i.e., change which meetings we cancel). If our secondary metric is the total number of meetings canceled, we can compute the sum of `correct_cancel` and `wrong_cancel` over all flights. In the case of using only the departure delay variable, we used a threshold of 16 minutes, and we would have canceled 270k meetings. How about now? Let's look at the total number of flights in the test set that would cause us to cancel our meetings:

```
SELECT
  SUM(IF(DEP_DELAY >= 14 AND DISTANCE < 328, 1, 0)) +
  SUM(IF(DEP_DELAY >= 15 AND DISTANCE >= 328 AND DISTANCE < 541, 1, 0)) +
  SUM(IF(DEP_DELAY >= 15 AND DISTANCE >= 541 AND DISTANCE < 802, 1, 0)) +
  SUM(IF(DEP_DELAY >= 17 AND DISTANCE >= 802, 1, 0))
AS cancel
FROM (
  SELECT
    DEP_DELAY,
    ARR_DELAY,
    DISTANCE
  FROM
    dsongcp.flights_tzcorr f
  JOIN
    dsongcp.trainday t
  ON
    f.FL_DATE = t.FL_DATE
```

```
WHERE
t.is_train_day = 'False')
```

This turns out to be 275k. It appears, then, that our simpler univariate model got us pretty much the same results as this more complex model using two variables.²⁰ However, the decision surfaces are different—in the single-variable threshold, we cancel meetings whenever the flight is delayed by 16 minutes or more. However, when we take into account distance, we cancel more meetings corresponding to shorter flights (threshold is now 14–15 minutes) and cancel fewer meetings corresponding to longer flights (threshold is now 17 minutes). One reason to use this two-variable Bayes model over the one-variable threshold determined empirically is to make such finer-grained decisions. This might or might not be important—it comes down to whether longer flights are typically those corresponding to more important meetings.

Why did we not get an improvement in the number of canceled meetings? Perhaps the round-off in the delay variables (they are rounded off to the nearest minute) has hurt our ability to locate more precise thresholds. Also, maybe the extra variable would have helped if I'd used a more sophisticated model—direct evaluation of conditional probability on relatively coarse-grained quantization bins is a very simple method. In [Chapter 7](#), we explore a more complex approach.

Orchestration

So far, in this chapter, we have developed and run the Spark jobs interactively. Once you have done so, you will want to operationalize the job and run it routinely. Although you can productionize a Jupyter Notebook using tools such as [Papermill](#), I recommend that you convert the code into a Python program that you can execute in a standalone way.

²⁰ Occam's razor suggests that we should pick the simpler model whenever its performance is comparable to a more complex model—it costs money and time to collect the data corresponding to additional features and keep them quality-controlled. Famously, although Netflix paid out \$1 million to a team that developed a better recommendation algorithm, it also [announced that it had no plans](#) to put that algorithm into production because of the combination of additional engineering effort and changes to Netflix's business model. We have to balance this desire for simplicity against the additional expressive power offered by the more complex model. For example, greater granularity might raise interesting questions to get you greater performance later—perhaps we can investigate the reason behind the dip in the second quartile, identify the city pairs driving this degradation, and impose rules that address the scenario.



Although you can just copy-paste the code cells out of a Jupyter Notebook into a Python file, a more systematic way is to convert the Jupyter Notebook to a Python program using a tool called **nbconvert**. After that, we can do minor editing to get rid of the display cells (such as plotting the hexbin plots). The Python program corresponding to the Jupyter Notebook that we've developed so far is in the GitHub repository as *bayes_in_spark.py*.

Submitting a Spark Job

We already have a Dataproc cluster that we have been using during development—our Jupyter Notebook is running on a Dataproc cluster. We can submit our Python program to this cluster from Cloud Shell:

```
gcloud dataproc jobs submit pyspark \  
  --cluster ch6cluster --region $REGION \  
  bayes_in_spark.py \  
  -- \  
  --bucket $BUCKET --debug
```

This is the approach you'd take if you have an already running cluster and wish to submit Spark jobs to it. However, this will require ensuring that the cluster has enough resources to handle your job. If you happen to submit your job at the same time as some other team that is already utilizing the cluster, your job might run very slowly. The solution for this problem is to use job-specific clusters.

Workflow Template

I recommend that you create ephemeral clusters, run jobs on them, and then delete them when you are done. Instead of doing them manually, you can automate it using a *workflow template*:

```
TEMPLATE=ch6eph  
MACHINE_TYPE=n1-standard-4  
CLUSTER=ch6eph  
  
gcloud dataproc --quiet workflow-templates create $TEMPLATE
```

The first step of the template will be to create a cluster of the appropriate size and set it to be a managed cluster so that it gets deleted once all the steps in the template are complete:

```
gcloud dataproc workflow-templates set-managed-cluster $TEMPLATE \  
  --master-machine-type $MACHINE_TYPE \  
  --worker-machine-type $MACHINE_TYPE \  
  --initialization-actions $STARTUP_SCRIPT \  
  --num-preemptible-workers=3 --num-workers 2 \  
  --
```

```
--image-version 2.0 \
--cluster-name $CLUSTER
```

Then, you can add jobs to the template. For example, to run a Pig program, we'd do:

```
gcloud dataproc workflow-templates add-job \
  pig gs://$BUCKET/bayes_final.pig \
  --step-id create-report \
  --workflow-template $TEMPLATE \
  -- --bucket=$BUCKET
```

Finally, instantiate the template to run the jobs and delete the cluster once done:

```
gcloud dataproc workflow-templates instantiate $TEMPLATE
```

One key change that we have to make to our program is to ensure that the output of the Spark program goes to a Cloud Storage location (rather than a local file on disk) because the cluster will be deleted once the job is complete:

```
bayes.to_csv('gs://${BUCKET}/flights/bayes.csv'.format(BUCKET),
             index=False)
```

Cloud Composer

If your Dataproc job is part of a larger data pipeline, you will typically write the data pipeline in Apache Airflow. Cloud Composer provides a fully managed experience for Airflow on Google Cloud.

Within your Airflow graph, you can [launch the workflow template](#) using an Airflow operator:

```
start_template_job = DataprocInstantiateWorkflowTemplateOperator(
    ...
)
```

Your considerations might change, however, if your company owns a Hadoop cluster on premises. In that case, you will submit Spark jobs to that long-lived cluster and will typically be concerned with ensuring that the cluster is not overloaded.

For the scenario in which you own an on-premises cluster, you might want to consider using a public cloud as a spillover in those situations for which there are more jobs than your cluster can handle. You can achieve this by monitoring YARN jobs and sending such spillover jobs to Cloud Dataproc. Cloud Composer provides the necessary plug-ins to be able to do this, but discussing how to set up such a hybrid system is beyond the scope of this book.

Autoscaling

When we created the workflow template, we specified the number of workers in the cluster. When we were developing the Spark program, we resized the cluster to add workers when we were ready to create the model on the full dataset. In both scenar-

ios, we have to know what size of cluster we need. This can be difficult for new workloads and for spiky jobs.

On a production system, it is possible to tell Dataproc to *autoscale*—the autoscaler monitors the cluster and, when it sees that the machines are getting maxed out, it adds more workers. When it sees machines on the cluster being idle, it shuts down a few workers. To do this, specify an autoscaling policy when creating the Dataproc cluster:

```
gcloud dataproc clusters create ch6cluster \
    --autoscaling-policy=ch6policy \
    ...
```

The autoscaling policy is specified in a YAML file with the syntax:

```
workerConfig:
  minInstances: 3
  maxInstances: 10
  weight: 1
secondaryWorkerConfig:
  minInstances: 0
  maxInstances: 20
  weight: 1
basicAlgorithm:
  cooldownPeriod: 2m
  yarnConfig:
    scaleUpFactor: 0.05
    scaleDownFactor: 1.0
    scaleUpMinWorkerFraction: 0.0
    scaleDownMinWorkerFraction: 0.0
    gracefulDecommissionTimeout: 1h
```

Note how the range of the number of primary and secondary workers is specified, as is the rate by which workers are scaled up. Once the autoscaling policy is specified, it is registered using `gcloud`:

```
gcloud dataproc autoscaling-policies import ch6policy \
    --source=filepath/filename.yaml \
    --region=region
```

Then, submit jobs to the cluster as and when you need them to be run:

```
gcloud dataproc jobs submit pyspark \
    --cluster=ch6cluster bayes_final.py \
    --bucket=$BUCKET
```

The cluster will be autoscaled based on the resources needed by the job subject to the limits specified in the policy.

Serverless Spark

However, even autoscaling requires a cluster to be running.

An even better approach would be if we can simply submit a Spark program to a Dataproc service, and the service starts the cluster, runs the job, autoscales it if necessary, and deletes the cluster. We'd like our Spark job to be *serverless*.

To do so, we put the script itself on Cloud Storage:

```
gsutil cp bayes_on_spark.py gs://$BUCKET/
```

and submit the job using gcloud:

```
gcloud beta dataproc batches submit pyspark \
--project=$(gcloud config get-value project) \
--region=$REGION \
gs://${BUCKET}/bayes_on_spark.py \
-- \
--bucket ${BUCKET} --debug
```

See *submit_serverless.sh* in the GitHub repository for details.

Once we do this, Dataproc takes care of all the infrastructure details. It runs our job and puts the lookup table in Cloud Storage. We can see the status of the job and examine its logs using the GCP console (see [Figure 6-10](#)).

The screenshot shows the GCP Dataproc section of the web console. On the left, there's a sidebar with categories: Jobs on Clusters (Clusters, Jobs, Workflows, Autoscaling policies), Serverless (Batches, Component exchange), and Utilities. The 'Batches' section under 'Serverless' is currently selected and highlighted in blue. In the main area, a specific job entry is displayed with the following details:

Batch ID	d129d3ad635444fb8d29757f9b070c3b
Batch UUID	01d2d000-40e1-4f04-8aa7-be509ad92ac7
Resource type	Batch
Status	Succeeded
DETAILS	
Start time	Oct 12, 2021, 11:22:20 AM
Elapsed time	9 min 34 sec
Region	us-central1
Batch type	PySpark
Main python file	gs://ai-analytics-solutions-dsongcp/bayes_on_spark.py
Arguments	--bucket ai-analytics-solutions-dsongcp

Figure 6-10. To view the status of the serverless Spark job, visit the Dataproc section of the GCP web console.

When I did this, I got:

```
dist_thresh,delay_thresh,frac_ontime,score
31.0,14.0,0.6874006359300477,0.012599364069952212
328.0,15.0,0.6958465263550009,0.004153473644999073
544.0,15.0,0.7054307116104869,0.005430711610486916
802.0,17.0,0.6874393150322182,0.012560684967781732
```

This matches the lookup table that we got when we ran the notebook on the full dataset:

Distance bin	delay_thresh
[31, 328]	14
[328, 541]	15
[541, 802]	15
[802, inf]	17

Because running the job involves only a single `gcloud` command, it is possible to schedule the Spark program to run every month, or whenever a new month of data is received by updating the Cloud Run and Cloud Scheduler solution that we created in [Chapter 2](#).

Summary

In this chapter, we explored how to create a two-variable Bayes model to provide insight as to whether to cancel a meeting based on the likely arrival delay of a flight. We quantized the two variables (distance and departure delay), created a conditional probability lookup table, and examined the on-time arrival percentage in each bin. We carried out the quantization using histogram equalization and on-time arrival percentage computations in Spark.

Upon discovering that equalizing the full distribution of departure delays resulted in a very coarse sampling of the decision surface, we chose to go with the highest possible resolution in the crucial range of departure delay.²¹ However, to ensure that we would have statistically valid groupings, we also made our quantization thresholds coarser in distance. On doing this, we discovered that the probability of the arrival delay being less than 15 minutes varied rather smoothly. Because of this, our conditional probability lookup reduced to a table of thresholds that could be applied cleanly using IF-THEN rules.

On evaluating the two-variable model, we found that we would be canceling about the same number of meetings as with the single-variable model while retaining the

²¹ One-minute increments in the range (10, 20).

same overall accuracy. We hypothesize that the improvement isn't higher because the departure delay variable has already been rounded off to the nearest minute, limiting the scope of any improvement we can make.

In terms of tooling, we created a three-node Cloud Dataproc cluster for development and resized it on the fly to 20 workers when our code was ready to be run on the full dataset. Cloud Dataproc goes a long way toward making this a low-touch endeavor—we saw that it is possible to create and schedule ephemeral jobs because Dataproc provides a serverless experience for Spark. The reason that Dataproc can create, resize, and delete the clusters for our job is that our data is held not in HDFS, but on Google Cloud Storage. We carried out development in JupyterLab, which gives us an interactive notebook experience. We also found that we were able to integrate BigQuery and Spark SQL into our workflow on the Hadoop cluster.

Finally, we converted the notebook into a Python Spark program that can be run routinely in production. We explored different ways of doing this: submitting the Spark to an already running cluster, creating a Workflow template, using Cloud Composer, and running Spark in a serverless way. Of these, serverless Spark involves the least amount of fiddling around with infrastructure. If you can do serverless, do serverless.

Suggested Resources

The most common reason that organizations use Dataproc today is that they used to have Hadoop clusters on premises. On-premises Hadoop workloads are typically lifted-and-shifted to Dataproc, optimized to take advantage of cloud computing (for example, using ephemeral clusters), and then modernized over time to BigQuery and Dataflow. This [technical guide](#) steps you through the considerations in moving the data, migrating the jobs, and connecting various types of clients and security tools.

One common question is whether to use Dataflow or Dataproc—both these products support data ingest and data processing. The flowchart in [Figure 6-11](#), from the Google Cloud documentation, suggests that this depends on whether you want to use Hadoop tools like Spark.

As the flowchart suggests, the second common reason that organizations use Dataproc is that they wish to have some degree of control over their infrastructure. Perhaps they have tasks that have to be carried out on premises for regulatory reasons. When doing so it is important to adopt [best practices](#) for storage, compute, and operations. I would add a third consideration here—Dataflow is much better at streaming than any alternative on Hadoop. My colleague Grace Mollison has [collected these flow charts](#) on her website.

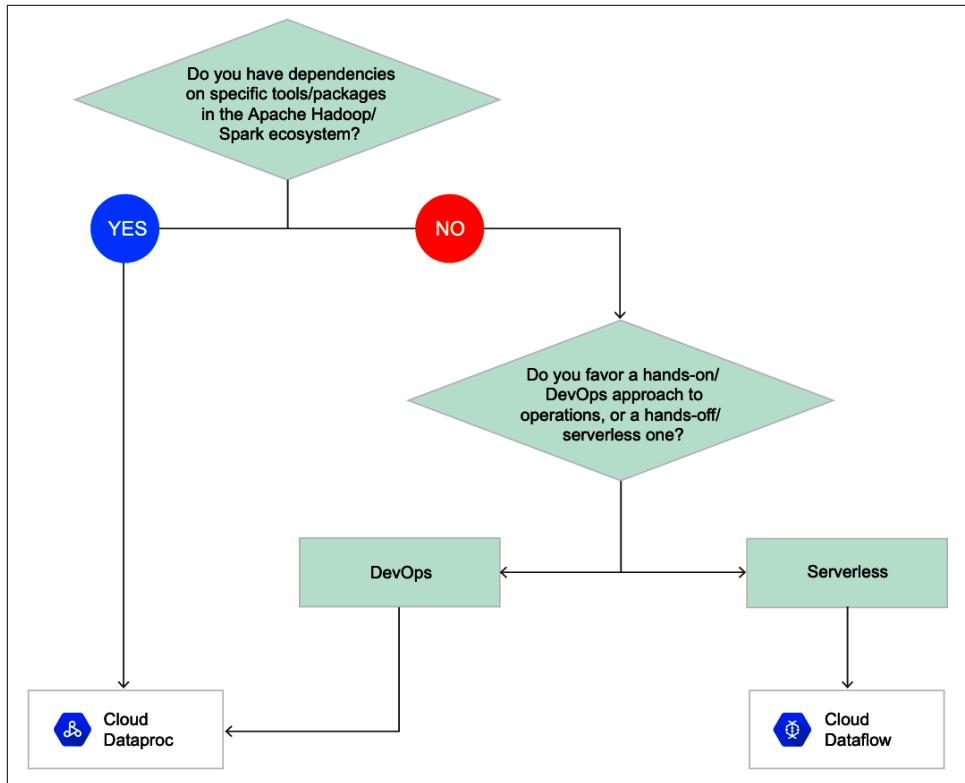


Figure 6-11. Choosing between Dataproc and Dataflow.

Logistic Regression Using Spark ML

In [Chapter 6](#), we created a model based on two variables—distance and departure delay—to predict the probability that a flight will be more than 15 minutes late. We found that we could get a finer-grained decision if we used a second variable (distance) instead of using just one variable (departure delay).

Why not use all the variables in the dataset? Or at least many more of them? In particular, I'd like to use the `TAXI_OUT` variable—if it is too high, the flight will be stuck on the runway waiting for the airport tower to allow the plane to take off, and so the flight is likely to be delayed. The Naive Bayes approach in [Chapter 6](#) was quite limiting in terms of being able to incorporate additional variables. As we add variables, we would need to continue slicing the dataset into smaller and smaller bins. We would then find that many of our bins would contain very few samples, resulting in decision surfaces that would not be well behaved. Remember that, after we binned the data by distance, we found that the departure delay decision boundary was quite well behaved—departure delays above a certain threshold were associated with the flight not arriving on time. Our simplification of the Bayesian classification surface to a simple threshold that varied by bin would not have been possible if the decision boundary had been noisier.¹ The more variables we use, the more bins we will have, and this good behavior will begin to break down. This sort of breakdown in good behavior as the number of variables (or *dimensions*) increases is called the *curse of dimensionality*; it affects many statistical and machine learning techniques, not just the quantization-based Bayesian approach of [Chapter 6](#).

¹ We also put our thumb on the scale a little. Recall that we quantized the distance variable quite coarsely. Had we quantized distance into many more bins, there would have been fewer flights in each bin.



All of the code snippets in this chapter are available in the folder [07_sparkml](#) of the [GitHub repository](#). See the *README.md* file in that directory for instructions on how to do the steps described in this chapter.

Logistic Regression

One way to address the breakdown in behavior as the number of variables increases is to change the approach from that of directly evaluating the probability based on the input dataset. Instead, we could attempt to fit a smooth function on the variables in the dataset (a multidimensional space) to the probability of a flight arriving late (a single-dimensional space) and use the value of that function as the estimate of the probability. In other words, we could try to find a function f , such that:

$$P(Y) \approx f(x_0, x_1, \dots, x_{n-1})$$

In our case, x_0 could be the departure delay, x_1 the taxi-out time, x_2 the distance, and so on. Each row will have different values for the x 's and represent different flights. The idea is that we have a function that will take these x 's and somehow transform them to provide a good estimate of the probability that the flight corresponding to that row's input variables is on time.

How Logistic Regression Works

One of the simplest transformations of a multidimensional space to a single-dimensional one is to compute a weighted sum of the input variables, as demonstrated here:

$$L = w_0x_0 + w_1x_1 + \dots + w_{n-1}x_{n-1} + b$$

The w 's (called the weights) and the constant b (called the intercept) are constants, but we don't initially know what they are. We need to find "good" values of the w 's and b such that the weighted sum for any row closely approximates either 1 (when the flight is on time) or 0 (when the flight is late). Because this process of finding good values is averaged over a large dataset, the value L is the prediction that flights with a specific departure delay, taxi-out time, and so on will arrive on time. If that number is 0.8, we would like it to be that 80% of such flights would arrive on time and 20% would be late. In other words, rather than L being simply 1 or 0, we'd like it to be the probability that the flight will arrive on time.

There is a problem, though. The preceding weighted sum cannot function as a probability. This is because the linear combination (L) can take any value, whereas a probability will need to lie between 0 and 1. One common solution for this problem is to transform the linear combination using the *logistic* function:

$$P(Y) = \frac{1}{1 + e^{-L}}$$

Fitting a logistic function of a linear combination of variables to binary outcomes (i.e., finding “good” values for the w ’s and b such that the estimated $P(Y)$ are close to the actual recorded outcome of the flight being on time) is called *logistic regression*.

In machine learning, the original linear combination, L , which lies between $-\infty$ and ∞ , is called the *logit*. You can see that if the logit adds up to ∞ , e^{-L} will be 0 and so, $P(Y)$ will be 1. If the original linear combination adds up to $-\infty$, then e^{-L} will be ∞ and so, $P(Y)$ will be 0. Therefore, Y could be an event such as the flight arriving on time and $P(Y)$, the probability of that event happening. Because of the transformation, $P(Y)$ will lie between 0 and 1 as required for anything to be a probability.

If $P(Y)$ is the probability, the logit, L , is given by the following:

$$\log_e \frac{P(Y)}{1 - P(Y)}$$

The *odds* is the ratio of the probability of the event happening, $P(Y)$, to the probability of the event not happening, $1 - P(Y)$. Therefore, the logit can also be interpreted as the log-odds where the base of the logarithm is e .

Figure 7-1 depicts the relationship between the logit, the probability, and the odds.

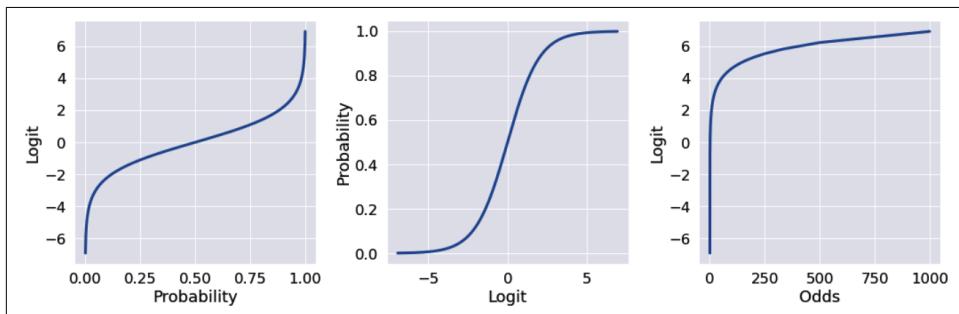


Figure 7-1. The relationships among the probability, the logit, and the odds.

Spend some time looking at the graphs in [Figure 7-1](#) and gaining an intuitive understanding for what the relationships mean. For example, see if you can answer this set of questions (feel free to sketch out the curves as you do your reasoning):²

- At equal odds (i.e., the flight is as likely to be late as not), what is the logit? How about if the odds are predominantly in favor of a flight being on time?
- At what probability does the logit function change most rapidly?
- Where is the gradient (rate of change) of the logit function slowest?
- Which logit change, from 2 to 3 or from 2 to 1, will result in a greater change in probability?
- How does the value of the intercept, b , affect the answer to Question 4?
- Suppose the intercept is zero. If all the input variables double, what happens to the logit?
- If the logit value doubles, what happens to the probability? How does this depend on the original value of the logit?
- What logit value does it take to provide a probability of 0.95? How about 0.995?
- How extreme do the input variables have to be to attain probabilities that are close to zero or one?

Many practical considerations in classification problems in machine learning derive from this set of relationships. So, as simple as these curves are, it is important that you understand their implications.

The name *logistic regression* is a little confusing—regression is normally a way of fitting a real-valued number, whereas classification is a way of fitting to a categorical outcome. Here, we fit the observed variables (the x 's) to a logit (which is real-valued)—this is the regression that is being referred to in the name. However, we then use the logistic function that has no free parameters (no weights to tune) to transform the

² Answers (but don't take my word for it): (1) At equal odds, the probability is $\frac{1}{2}$ and the logit is zero. If the odds are predominantly in favor of the flight being on time, the probability of on-time arrival is nearly 1.0 and the logit has a very large positive value. (2) The logit function changes fastest at probabilities near zero and one. (3) The gradient of the logit function is slowest near a probability of $\frac{1}{2}$. (4) The change from 2 to 1 will result in a greater change in probability. Near probabilities of zero and one, larger logit changes are required to have the same impact on probability. As you get nearer to a probability of $\frac{1}{2}$, smaller logit changes suffice. (5) The intercept directly impacts the value of the logit, so it moves the first curve upwards or downwards. (6) The logit doubles. (7) If the original probability is near zero or one, doubling of the logit has negligible impact (look at what happens if the logit changes from 4 to 8, for example). If the original probability is between about 0.3 and 0.7, the relationship is quite linear: doubling of the logit ends up causing a proportional increase in probability. (8) About 3, and about 5. (9) Doubling of the logit value at 0.95 is required to reach 0.995. Lots of “energy” in the input variables is required to move the probability needle at the extremes.

real-valued number to a probability. Overall, therefore, logistic regression functions as a classification method.

Spark ML Library

Given a multivariable dataset, Spark has the ability to carry out logistic regression and give us the optimal weight for each variable. Spark's logistic regression module will give us the w 's and b if we show it a bunch of x 's and the corresponding Y 's. The logistic regression module is part of Apache Spark's machine learning library, MLLib, which you can program against in Java, Scala, Python, or R. Spark MLLib (colloquially known as Spark ML) includes implementations of many canonical machine learning algorithms: decision trees, random forests, alternating least squares, k-means clustering, association rules, support vector machines, and so on. Spark can execute on a Hadoop cluster; thus, it can scale to large datasets.

The problem we are solving—to find a set of weights that optimizes model predictions based on known outcomes—is an instance of a *supervised* learning problem. In supervised learning problems, the actual answers, called *labels*, need to be known for some dataset. As illustrated in [Figure 7-2](#), first you ask the machine (here, Spark) to learn (the w 's) from data (the x 's) that has labels (the Y 's). This is called *training*.

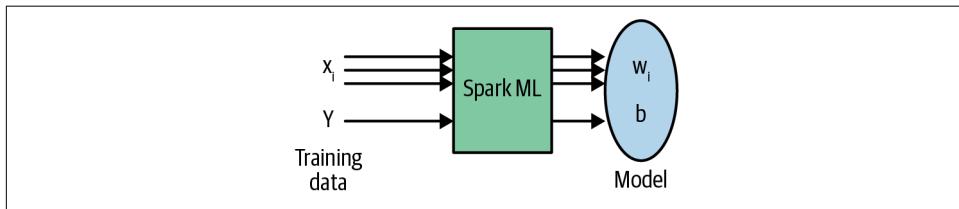


Figure 7-2. In supervised learning, the machine (here, Spark) learns a set of parameters (the w 's and b 's) from training data that consists of inputs (x 's) and their corresponding labels (Y 's).

The learned set of weights, along with the original equation (the logistic function of a linear combination of x 's), is called a *model*. After you have learned a model from the training data, you can save it to a file. Then, whenever you want to make a prediction for a new flight, you can re-create the model from that file, pass in the x 's in the same order, and compute the logistic function and obtain the estimate $P(Y)$. This process, called *prediction*, might be carried out in real time in response to a request that includes the input variables, whereas the training of the model might happen less frequently, as shown in [Figure 7-3](#).

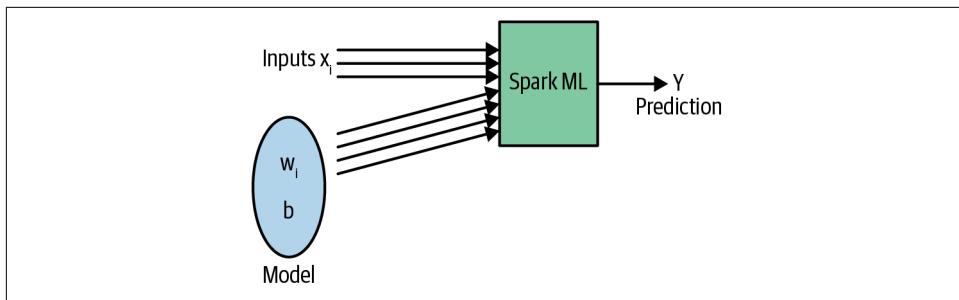


Figure 7-3. You can use the learned set of weights to predict the value of Y for new data (x 's).

Of course, the prediction can be carried outside of Spark—all that we'd need to do is compute the weighted sum of the inputs and the intercept, and then compute the logistic function of the weighted sum. In general, though, you should seek to use the same libraries for prediction as you use for training. This helps to mitigate *training-serving skew*, the situation that we talked about in [Chapter 2](#) in which the input variables in prediction are subtly different from those used in training and which leads to poorly performing models.

Getting Started with Spark Machine Learning

To run Spark conveniently, I will continue to use the Cloud Dataproc cluster that I launched in the previous chapter. In case you deleted it, start a new one using:

```
cd 06_dataproc
./create_cluster.sh bucketname region
```

Even though we are going to develop the logistic regression code in a Jupyter Notebook, we should keep in mind that our end goal is to run the machine jobs routinely. To achieve that goal, it is important to keep a notebook with the final machine learning workflow and export this notebook into a standalone program. You can submit the standalone program to the cluster whenever the machine learning needs to be repeated on new datasets.

As in [Chapter 6](#), we open the notebook from the Web Interfaces section of the Google Cloud console. Within the notebook, we start out creating a `SparkContext` variable `sc` and the `SparkSession` variable `spark`:³

```
from pyspark.sql import SparkSession
from pyspark import SparkContext
sc = SparkContext('local', 'logistic')
spark = SparkSession \
```

³ See `logistic_regression.ipynb` in the GitHub repository for this book.

```
.builder \
.appName("Logistic regression w/ Spark ML") \
.getOrCreate()
```

After we do this, any line that works in the interactive shell will also work when launched from the notebook *logistic_regression.ipynb* or the script *logistic.py*. The application name (*logistic*) will show up in logs when the script is run.

Spark Logistic Regression

The logistic regression implementation, L-BFGS, is in *pyspark.mllib* and is named for the initials of the independent inventors (Broyden, Fletcher, Goldfarb, and Shanno)⁴ of a popular, iterative, fast-converging optimization algorithm. The L-BFGS algorithm is used by Spark to find the weights that minimize the *logistic loss* function:

$$\sum \log(1 + e^{-yL})$$

over the training dataset, where y , the training label, is either -1 or 1 and L is the logit computed from the input variables, the weights, and the intercept.

So, let's begin by adding `import` lines for the Python classes that we'll need:

```
from pyspark.mllib.classification import LogisticRegressionWithLBFGS
from pyspark.mllib.regression import LabeledPoint
```

Knowing the details of the logistic formulation and loss function used by Spark is important. This is because different machine learning libraries use equivalent (but not identical) versions of these formulae. For example, another common approach taken by machine learning frameworks is to minimize the *cross-entropy*:

$$\sum -y \log P(Y) - (1 - y) \log (1 - P(Y))$$

Here, y is the training label, and $P(Y)$ is the probabilistic output of the model. In that case, the training label will need to be 0 or 1 . I won't go through the math, but even though the two loss functions look very different, minimizing the logistic loss and minimizing the cross-entropy loss turn out to be equivalent.

Rather confusingly, the Spark documentation notes that “a binary label y is denoted as either $+1$ (positive) or -1 (negative), which is convenient for the formulation. However, the negative label is represented by 0 in *spark.mllib* instead of -1 , to be consistent with multiclass labeling.”⁵ In other words, Spark ML uses the logistic loss

⁴ The L stands for low-memory—this is a limited-memory variant of the BFGS algorithm.

⁵ See the section on loss functions in [the Spark documentation](#).

function, but requires that the labels we provide be 0 or 1. There really is no substitute for reading the documentation!

To summarize, the preprocessing that you might need to do to your input data depends on the formulation of the loss function employed by the machine learning framework you are using. Suppose that we decide to do this:

$$\begin{aligned}y &= 0 \text{ if arrival delay } \geq 15 \text{ minutes} \\y &= 1 \text{ if arrival delay } < 15 \text{ minutes}\end{aligned}$$

Because we have mapped on-time flights to 1, the machine learning algorithm (after training) will predict the probability that the flight is on time.

Creating a Training Dataset

First, let's read in the list of training days. To do that, we need to read *trainday.csv* from Cloud Storage, remembering that the comma-separated value (CSV) file has a header that will help with inferring its schema:

```
traindays = spark.read \
    .option("header", "true") \
    .csv('gs://{}/flights/trainday.csv'.format(BUCKET))
```

For convenience, I'll make this a Spark SQL view, as well:

```
traindays.createOrReplaceTempView('traindays')
```

We can print the first few lines of this file:

```
spark.sql("SELECT * from traindays LIMIT 5").show()
```

This obtains the following eminently reasonable result:

```
+-----+-----+
| FL_DATE|is_train_day|
+-----+-----+
|2015-01-01|      True|
|2015-01-02|     False|
|2015-01-03|     False|
|2015-01-04|      True|
|2015-01-05|      True|
+-----+-----+
```

While we're developing the code (on my minimal Hadoop cluster that is running Jupyter), it would be easier to read only a small part of the dataset. Hence, we define the `inputs` variable to be just one of the shards:

```
inputs = 'gs://{}//flights/tzcorr/all_flights-00000-*'.format(BUCKET)
```

After we have developed all of the code, we can change the inputs to the full dataset:

```
#inputs = 'gs://{}//flights/tzcorr/all_flights-*'.format(BUCKET) # FULL
```

For now, though, let's leave the latter line commented out. We can read in the flights dataset as we did in [Chapter 6](#):

```
flights = spark.read.json(inputs)
flights.createOrReplaceTempView('flights')
```

Training will need to be carried out on the flights that were on days for which `is_train_day` is True:

```
trainquery = """
SELECT
    f.*
FROM flights fJOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day == 'True'
"""
traindata = spark.sql(trainquery)
```

Dealing with corner cases

Let's verify that `traindata` does contain the data we need. We can look at the first few (here, the first two) rows of the dataframe using the following:

```
traindata.head(2)
```

The result seems quite reasonable:

```
[Row(ARR_AIRPORT_LAT=33.43416667, ARR_AIRPORT_LON=-112.01166667,
ARR_AIRPORT_TZOFFSET=-25200.0, ARR_DELAY=-16.0, ARR_TIME='2015-07-28T18:20:00',
CANCELLED=False, CRS_ARR_TIME='2015-07-28T18:36:00',
CRS_DEP_TIME='2015-07-28T17:05:00', DEP_AIRPORT_LAT=33.9425,
DEP_AIRPORT_LON=-118.40805556, DEP_AIRPORT_TZOFFSET=-25200.0, DEP_DELAY=-3.0,
DEP_TIME='2015-07-28T17:02:00', DEST='PHX', DEST_AIRPORT_SEQ_ID='1410702',
DISTANCE='370.00', DIVERTED=False, FL_DATE='2015-07-28', ORIGIN='LAX',
ORIGIN_AIRPORT_SEQ_ID='1289203', TAXI_IN=6.0, TAXI_OUT=14.0,
UNIQUE_CARRIER='AA', WHEELS_OFF='2015-07-28T17:16:00',
WHEELS_ON='2015-07-28T18:14:00'),
```

Date fields are dates, and airport codes are reasonable, as are the latitudes and longitudes. But eyeballing is no substitute for truly verifying that all of the values exist.

So, let's restrict the query to fields we want:

```
SELECT
    DEP_DELAY, TAXI_OUT, ARR_DELAY, DISTANCE
FROM flights f
...
```

Knowing that the four variables we are interested in are all floats, we can ask Spark to compute simple statistics over the full dataset:

```
traindata.describe().show()
```

The `describe()` method computes column-by-column statistics, and the `show()` method causes those statistics to be printed. We now get the following:⁶

summary	DEP_DELAY	TAXI_OUT	ARR_DELAY	DISTANCE
count	259692	259434	258706	275062
mean	13.178	16.9658	9.7319	802.3747
stddev	41.8886	10.9363	45.0384	592.254
min	-61.0	1.0	-77.0	31.0
max	1587.0	225.0	1627.0	4983.0

Notice anything odd?

Notice the count statistic. There are 275,062 `DISTANCE` values, but only 259,692 `DEP_DELAY` values, and even fewer `TAXI_OUT` values. What is going on? This is the sort of thing that you will need to chase down to find the root cause. In this case, the reason has to do with flights that are scheduled but never leave the gate and flights that depart the gate but never take off. Similarly, there are flights that take off (and have a `TAXI_OUT` value) but are diverted and do not have an `ARR_DELAY`. In the data, these are denoted by `NULL`, and Spark's `describe()` method doesn't count `NULLs`.

We don't want to use canceled and diverted flights for training either. One way to tighten up the selection of our training dataset would be to simply remove `NULLs`, as shown here:

```
trainquery = """
SELECT
    DEP_DELAY, TAXI_OUT, ARR_DELAY, DISTANCE
FROM flights f
JOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day == 'True' AND
    f.dep_delay IS NOT NULL AND
    f.arr_delay IS NOT NULL
"""

traindata = spark.sql(trainquery)
traindata.describe().show()
```

Running this gets us a consistent value of the count across all the columns:

summary	DEP_DELAY	TAXI_OUT	ARR_DELAY	DISTANCE
count	258706	258706	258706	258706

⁶ Your results might be different because the actual flight records held in your first shard (recall that the input is `all_flights-00000-*`) are possibly different.

However, I strongly encourage you not to do this. Removing NULLs is merely fixing the symptom of the problem. What we really want to do is to address the root cause. In this case, you'd do that by removing flights that have been canceled or diverted, and fortunately, we do have this information in the data. So, we can change the query to be the following:

```
trainquery = """
SELECT
    DEP_DELAY, TAXI_OUT, ARR_DELAY, DISTANCE
FROM flights f
JOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day == 'True' AND
    f.CANCELLED == 'False' AND
    f.DIVERTED == 'False'
"""

traindata = spark.sql(trainquery)
traindata.describe().show()
```

This, too, yields the same counts as when we threw away the NULLs, thereby demonstrating that our diagnosis of the problem was correct.

Discovering corner cases and problems with an input dataset at the time we begin training a machine learning model is quite common. In this case, I knew this problem was coming and was careful to select the CANCELLED and DIVERTED columns to be part of my input dataset (in [Chapter 2](#)). In real life, you will need to spend quite a bit of time troubleshooting this, potentially adding new logging operations to your ingest code to uncover the reason that underlies a simple problem. What you should not do is to simply throw away bad values.



Bad values (like NULL) are usually a symptom of a problem. Investigate the issue. Don't simply discard bad values.

Creating training examples

Now that we have the training data, we can look at the documentation for `LogisticRegressionModel` to determine the format of its input. The documentation indicates that each row of the training data needs to be transformed to a `Labeled Point` whose documentation in turn indicates that its constructor requires a label and an array of features, all of which need to be floating-point numbers.

Let's create a method that will convert each data point of our dataframe into a *training example* (an example is a combination of the input features and the true answer):

```

def to_example(raw_data_point):
    return LabeledPoint(
        float(raw_data_point['ARR_DELAY'] < 15), # on-time?
        [
            raw_data_point['DEP_DELAY'],
            raw_data_point['TAXI_OUT'],
            raw_data_point['DISTANCE'],
        ]
    )

```

Note that we have created a label and an array of features. Here, the features consist of three numeric fields that we pass in as-is. It is good practice to create a separate method that takes the raw data and constructs a training example because this allows us to fold in other operations as well. For example, we can begin to do preprocessing of the feature values, and having a method to construct training examples allows us to reuse the code between training and evaluation.

After we have a way to convert each raw data point into a training example, we need to apply this method to the entire training dataset. We can do this by mapping the dataset row by row:

```
examples = traindata.map(to_example)
```

Training the Model

Now that we have a dataframe in the requisite format, we can ask Spark to fit the training dataset to the labels:

```
lrmodel = LogisticRegressionWithLBFGS.train(examples, intercept=True)
```

We'd have specified `intercept = False` if we believed that when all $x = 0$, the prediction needed to be 0. We have no reason to expect this, so we ask the model to find a value for the intercept.

When the `train()` method completes, the `lrmodel` will have the weights and intercept, and we can print them out:

```
print lrmodel.weights, lrmodel.intercept
```

This yields the following:⁷

```
[ -0.164, -0.132, 0.000294] 5.1579
```

The `weights` is an array, one for each variable. These numbers, plus the formula for logistic regression, are enough to set up code for the model in any language we choose. Remember that in our labeled points, 0 indicated late arrivals and 1 indicated

⁷ Because of random seeds used in the optimization process, and different data in shards, your results will be different.

on-time arrivals. So, applying these weights to the departure delay, taxi-out time, and flight distance of a flight will yield the probability that the flight will be on time.

In this case, it appears that the departure delay has a weight of -0.164 . The negative sign indicates that the higher the departure delay, the lower the probability that the flight will be on time (which sounds about right). On the other hand, the sign on the distance is positive, indicating that higher distances are associated with more on-time behavior. Even though we are able to look at the weights and reason with them on this dataset, such reasoning will begin to break down if the variables are not independent. If you have highly correlated input variables, the magnitudes and signs of the weights are very hard to interpret.

Let's try out a prediction:

```
lrmodel.predict([6.0,12.0,594.0])
```

The result is 1—that is, the flight will be on time when the departure delay is 6 minutes, the taxi-out time is 12 minutes, and the flight is for 594 miles. Let's change the departure delay from 6 minutes to 36 minutes:

```
lrmodel.predict([36.0,12.0,594.0])
```

The result now is 0—the flight won't arrive on time.

But wait a minute. We want the output to be a probability, not 0 or 1 (the final label). To do that, we can remove the implicit threshold of 0.5:

```
lrmodel.clearThreshold()
```

With the thresholding removed, we get probabilities. The probability of arriving late increases as the departure delay increases.

By keeping two of the variables constant, it is possible to study how the probability varies as a function of one of the variables. For example, at a departure delay of 20 minutes and a taxi-out time of 10 minutes, this is how the distance affects the probability that the flight is on time:

```
dist = np.arange(10, 2000, 10)
prob = [lrmodel.predict([20, 10, d]) for d in dist]
plt.plot(dist, prob)
```

Figure 7-4 shows the plot.

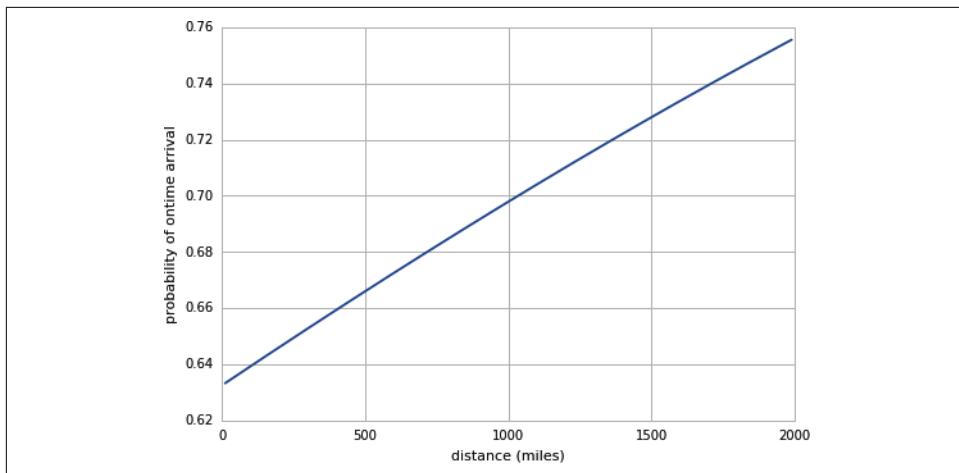


Figure 7-4. How the distance of the flight affects the probability of on-time arrival. According to our model, longer flights tend to have higher likelihoods of arriving on time, but the effect (0.63 to 0.76) is rather minor.

As you can see, the effect is relatively minor. The probability increases from about 0.63 to about 0.76 as the distance changes from a very short hop to a cross-continental flight. On the other hand, if we hold the taxi-out time and distance constant and examine the dependence on departure delay, we see a more dramatic impact (see Figure 7-5):

```
delay = np.arange(-20, 60, 1)
prob = [lrmrmodel.predict([d, 10, 500]) for d in delay]
ax = plt.plot(delay, prob)
```

Although the probabilities are useful to be able to plot the behavior of the model in different scenarios, we do want a specific decision threshold. Recall that we want to cancel the meeting if the probability of the flight arriving on time is less than 70%. So, we can change the decision threshold:

```
lrmrmodel.setThreshold(0.7)
```

Now, the predictions are 0 or 1, with the probability threshold set at 0.7.

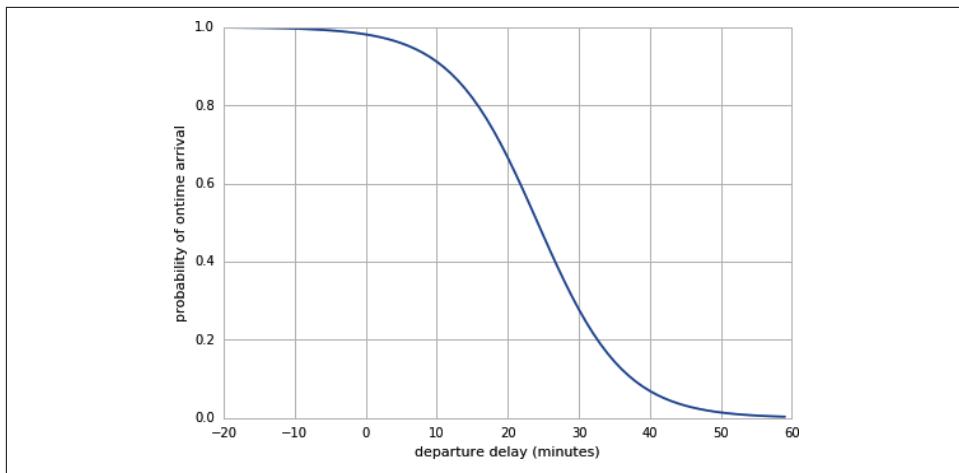


Figure 7-5. How the departure delay of the flight affects the probability of on-time arrival. The effect of the departure delay is rather dramatic.

Predicting Using the Model

Now that we have a trained model, we can save it to Cloud Storage and retrieve it whenever we need to make a prediction. To save the model, we provide a location on Cloud Storage:

```
MODEL_FILE='gs://' + BUCKET + '/flights/sparkmloutput/model'
lrmodel.save(sc, MODEL_FILE)
```

To retrieve the model, we load it from the same location:

```
from pyspark.mllib.classification import LogisticRegressionModel
lrmodel = LogisticRegressionModel.load(sc, MODEL_FILE)
lrmodel.setThreshold(0.7)
```

Note that we must take care to set the decision threshold; it is not part of the model.

Now, we can use the `lrmodel` variable to carry out predictions:

```
print lrmodel.predict([36.0,12.0,594.0])
```

Obviously, this code could be embedded into a Python web application to create a prediction web service or API.

A key point to realize is that whereas model training in Spark is distributed and requires a cluster, model prediction is a pretty straightforward mathematical computation. Model training is a batch operation that requires the ability to scale out to multiple processors, but online prediction requires fast computation on a single processor. When the machine learning model is relatively small (as in our logistic

regression workflow),⁸ hardware optimizations (like graphics processing units [GPUs]) are not needed in the training stage. Thus, when choosing how to resize the cluster to run our machine learning training job over the full dataset, it is more cost-effective to simply add more CPUs.

When would you need GPUs in machine learning? GPUs are potentially needed in the prediction stage for small models. GPUs become useful in prediction even for small models if the system needs to provide for low latency and a high number of queries per second (QPS). Of course, had we been training a deep learning model for image classification with hundreds of layers, GPUs would have been called for both in training and in prediction.

Evaluating a Model

Now that we have a trained model, we can evaluate its performance on the test days, a set of days that were not used in training (we created this test dataset in [Chapter 5](#)). To do that, we change the query to pull out the test days:

```
testquery = trainquery.replace(\n    "t.is_train_day == 'True'", "t.is_train_day == 'False'\")\nprint testquery
```

Here is the resulting query:

```
SELECT\n    DEP_DELAY, TAXI_OUT, ARR_DELAY, DISTANCE\nFROM flights f\nJOIN traindays t\nON f.FL_DATE == t.FL_DATEWHERE\n    t.is_train_day == 'False' AND\n    f.CANCELLED == 'False' AND\n    f.DIVERTED == 'False'
```

We then carry out the same ML pipeline as we did during training:

```
testdata = spark.sql(testquery)\nexamples = testdata.map(to_example)
```

Note that we are able to reuse the function `to_example` to go from the raw data to the training examples.

As soon as we have the dataframe `examples`, we can have the model predict the label given the set of features for each row:

```
labelpred = examples.map(lambda p: \\n        (p.label, lrmodel.predict(p.features)))
```

⁸ Very large deep neural networks, such as those used for image classification, are another story. Such models can have hundreds of layers, each with hundreds of weights. Here, we have three weights—four if you count the intercept.

The `map` function applies the `predict` method to each row of features and creates a dataframe that contains the true label and the model prediction for each row.

To evaluate the performance of the model, we first find out how many flights we will need to cancel and how accurate we are in terms of flights we cancel and flights we don't cancel:

```
def eval(labelpred):
    cancel = labelpred.filter(lambda (label, pred): pred == 1)
    nocancel = labelpred.filter(lambda (label, pred): pred == 0)
    corr_cancel = cancel.filter(lambda (label, pred): \
        label == pred).count()
    corr_nocancel = nocancel.filter(lambda (label, pred): \
        label == pred).count()
    return {'total_cancel': cancel.count(), \
            'correct_cancel': float(corr_cancel)/cancel.count(), \
            'total_nocancel': nocancel.count(), \
            'correct_nocancel': float(corr_nocancel)/nocancel.count()\
    }
```

Here's what the resulting statistics turn out to be:

```
{'correct_cancel': 0.7917474551623849, 'total_nocancel': 115949,
 'correct_nocancel': 0.9571363271783284, 'total_cancel': 33008}
```

As discussed in [Chapter 5](#), the reason the correctness percentages are not 70% is because the 70% threshold is on the marginal distribution—the accuracy percentages here are computed on the total distribution and so are padded by the easy decisions. However, let's go back and modify the evaluation function to explicitly print out statistics around the decision threshold—this is important to ensure that we are, indeed, making a probabilistic decision.

We should clear the threshold so that the model returns probabilities and then carry out the evaluation twice: once on the full dataset, and next on only those flights that fall near the decision threshold of 0.7:

```
lrmodel.clearThreshold() # so it returns probabilities
labelpred = examples.map(lambda p: \
    (p.label, lrmodel.predict(p.features)))
print eval(labelpred)
# keep only those examples near the decision threshold
labelpred = labelpred.filter(lambda (label, pred):\
    pred > 0.65 and pred < 0.75)
print eval(labelpred)
```

Of course, we must change the evaluation code to work with probabilities instead of with categorical predictions. The four variables now become as follows:

```
cancel = labelpred.filter(lambda (label, pred): pred < 0.7)
nocancel = labelpred.filter(lambda (label, pred): pred >= 0.7)
corr_cancel = cancel.filter(lambda (label, pred): \
    label == int(pred >= 0.7)).count()
```

```
corr_nocancel = nocancel.filter(lambda (label, pred): \
                                label == int(pred >= 0.7)).count()
```

When run, the first set of results remains the same, and the second set of results now yields this:

```
{'correct_cancel': 0.30886504799548276, 'total_noncancel': 2224,
 'correct_noncancel': 0.7383093525179856, 'total_cancel': 1771}
```

Note that we are correct about 74% of the time in our decision to go ahead with a meeting (our target was 70%). Although useful to verify that the code is working as intended, the actual results are meaningless because we are not running on the entire dataset, just one shard of it. Therefore, the final step is to export the code from the notebook, remove the various `show()` and `plot()` functions, and create a submittable script.⁹

We can submit the script to the Cloud Dataproc cluster from a laptop with the Cloud SDK installed, from Cloud Shell, or from the Cloud Dataproc section of the [Google Cloud Platform web console](#). Before we submit the script, we need to resize the cluster, so that we can process the entire dataset on a larger cluster than the one on which we did development. So, what we need to do is to increase the size of the cluster, submit the script, and decrease the size of the cluster after the script is done. Alternatively, we can use an autoscaling Dataproc cluster or serverless Spark as discussed in [Chapter 6](#).¹⁰

Running `logistic.py` on a more powerful cluster, we create a model and then evaluate it on the test days. We obtain these results from the logs:

```
All flights: {'total_cancel': 291895, 'correct_cancel': 0.8122920913342127,
 'total_noncancel': 1304422, 'correct_noncancel': 0.9642401002129679}
```

Looking at these overall results, notice that we are canceling about 292k meetings. We are accurate 96.4% of the time we don't cancel a meeting and 81.2% of the time when we decide to cancel. Remember that we canceled 275k meetings when using the Bayesian classifier in [Chapter 6](#) and were correct 83% of the time that we decided to cancel. Based on our secondary criterion, the Naive Bayes approach in [Chapter 6](#) with only two variables is better than the logistic regression approach in this chapter with three variables! There is, in machine learning, no substitute for experimentation to see what works.

How about the results on the marginal distribution (i.e., on flights for which our probability is near 0.7)? These, too, are in the logs. They indicate that we are, indeed, making an appropriately probabilistic decision—our decision to go ahead with a meeting is correct 73% of the time:

⁹ See `logistic.py` in `07_sparkml`.

¹⁰ The script `submit_spark.sh` in `07_sparkml` uses the serverless Spark approach.

```
Flights near decision threshold: {'total_cancel': 15084,  
'correct_cancel': 0.3325377883850438, 'total_nocancel': 18441,  
'correct_nocancel': 0.7279431701100808}
```

A close examination of the preceding numbers indicates what is going on—we are correct to cancel 33% of the time while in [Chapter 6](#), we were correct to cancel only 30% of the time. This extra 3% allows us to cancel more meetings. It is clear, then, that simply looking at the number of meetings we cancel is not a good secondary criterion of model performance. Instead, we will find a different criterion by comparing the performance of the model against ideal performance across all thresholds. We will do that in the next section.

Feature Engineering

Still, it is unclear whether we really needed all three variables in the logistic regression model. Any variable you include in a machine learning model brings with it an increased danger of overfitting. Known as the [principle of parsimony](#) (often referred to as Occam's razor),¹¹ the idea is that it is preferable to use a simpler model to a more complex model that has similar accuracy—the fewer variables we have in our model, the better it is.

One manifestation of the principle of parsimony is that there are practical considerations involved with every new variable in a machine learning model. A hand-coded system of rules can deal with the presence or absence of values (for example, if the variable in question was not collected) for specific variables relatively easily—simply write a new rule to handle the case. On the other hand, machine learning models require the presence of enough data when that variable value is absent. Thus, a machine learning model will often be unusable if all of its variable values are not present in the data. Even if a variable value is present in some new data, it might be defined differently or computed differently, and we might have to go through an expensive retraining effort in order to use it. Thus, extra variables pose issues around the applicability of a machine learning model to new situations. We should attempt to use as few variables as possible.

Experimental Framework

It is also unclear whether the three variables we chose are the ones that matter most. Perhaps we could use more variables from the dataset besides the three we are using. In machine learning terminology, the inputs to the model are called *features*. The features could be different from the raw input variables because the input variables could be transformed in some way before being provided to the model. The process

¹¹ William of Ockham (later spelled Occam), who was a friar in medieval England, actually wrote “*Pluralitas non est ponenda sine necessitate*,” which translates to “Entities should not be multiplied unnecessarily.”

of designing the transforms that are carried out on the raw input variables is called *feature engineering*.

To test whether a feature provides value to the model, we need to build an experimental framework. We could begin with one feature (departure delay, for example) and test whether the incorporation of a new feature (perhaps the distance) improves model performance. If it does, we keep it and try out one more feature. If not, we discard it and try the next feature on our list. This way, we get to select a subset of features and ensure that they do matter. Another approach is to train a model with all possible features, remove a feature, and retrain. If performance doesn't go down, leave the feature out. At the end, as before, we will be left with a subset of features that matter. The second approach is preferable because it allows us to capture interactions—perhaps a feature by itself doesn't matter, but its presence alongside another is powerful. Choosing the set of features through a systematic process is called *feature selection*.

For both feature engineering and feature selection, it is important to devise an experimental framework to test out our hypotheses as to whether a feature is needed. On which dataset should we evaluate whether a feature is important? We cannot evaluate how much it improves accuracy on the training dataset itself because the model might be fitting noise in the training data. Instead, we need an independent dataset in order to carry out feature selection. However, we cannot use the test dataset because if we use the test dataset in our model creation process, it is no longer independent and cannot be used as a good indicator of model performance. Therefore, we will split the training dataset itself into two parts—one part will be used for training, whereas the other will be held out and used to evaluate different models.

Figure 7-6 shows what our experimentation framework is going to look like.

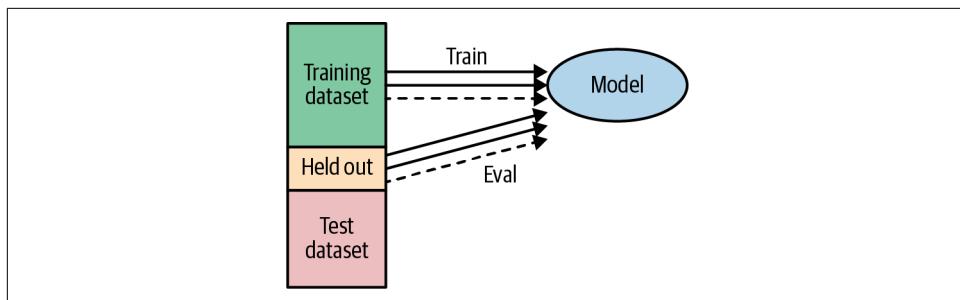


Figure 7-6. We often split a dataset into three parts. The training dataset is used to tune the weights of a model, and the held-out dataset is used to evaluate the impact of model changes, such as feature selection. The independent test dataset is used only to gauge the performance of the final, selected model.

First, we break the full dataset into two parts and keep one part of it for the final evaluation of models (we did this in [Chapter 6](#) when we created the `traindays` dataset). This part is called the test dataset and is what we have been using for end-of-chapter evaluations. However, when we are creating several models and need to choose among them, we cannot use the test dataset. Therefore, we will split our original training dataset itself into two parts. We'll retain the larger part of it for actual training and use the held-out portion to evaluate the model. In [Figure 7-6](#), for example, we are deciding whether to use the third input variable. We do this based on whether that feature improves model performance enough.

Choosing a metric

What metric should we evaluate in order to choose between two models? *Not* the number of canceled flights! It is easy to game metrics computed from the contingency table because it is possible to change the probability threshold to get a wide range of accuracy, precision, or recall metrics.¹² A contingency table-based metric is a good way to understand the performance of a model, but not to choose between two models unless care is taken to ensure that the measure is not tunable by changing the threshold. One way to do this would be to, for example, compare precision at a fixed recall rate, but you should do this only if that fixed recall is meaningful. In our problem, however, it is the probability that is meaningful, not the precision or recall.

Hence, it is not possible to fix either of them, and we are left comparing two pairs of numbers.

Another way to avoid the problem that the metric can be gamed is to use a measure that uses the full distribution of probabilities that are output by the model. When carrying out feature selection or any other form of hyperparameter tuning, we could use a measure such as the logistic loss or cross-entropy that conveys this full distribution. As a simpler, more intuitive measure that nevertheless uses the full distribution of probabilities, let's use the root mean squared error (RMSE) between the true labels and the probabilistic predictions:

```
totsqe = labelpred.map(  
    lambda data: (data[0] - data[1]) * (data[0] - data[1])  
).sum()  
rmse = np.sqrt(totsqe / float(cancel.count() + nocancel.count()))
```

¹² The Machine Learning Crash Course from Google has a [good explanation of these metrics](#). Remember that we need to threshold the probabilistic output of the model to get the entries in the contingency table. A correct cancel, for example, is the situation that the flight arrived more than 15 minutes late and the predicted probability of on-time arrival was less than 0.7. The metrics evaluated on the contingency table are extremely sensitive to this choice of threshold. Different models will be different at thresholds of 0.65, 0.68, or 0.70, especially for models whose performance is quite comparable. If, for example, we want the overall correct cancel percentage to be 80%, we can change the threshold to get this. We can also change the threshold to get an overall correct cancel percentage of 20% if that is what we desire.

What is “enough” of an improvement when it comes to RMSE? There are no hard-and-fast rules. We need the model performance to improve enough to outweigh the drawbacks involved with additional inputs and the loss of agility and model runtime speed that additional input features entail. Here, I will choose to use 0.5% as my threshold. If the model performance, based on some metric we decide upon, isn’t reduced by at least 0.5% by removing a variable, I won’t use the extra variable.

Creating the held-out dataset

Because the held-out dataset is going to be used only for model evaluation and only within Spark, we do not need to create the held-out dataset the same way we created the test dataset in [Chapter 5](#). For example, we do not need to store the held-out days as a separate dataset that can be read from multiple frameworks.¹³ However, the principle of repeatability still applies—every time we run our Spark program, we should get the same set of held-out days. Otherwise, it will not be possible to compare the performance of different models (because evaluation metrics depend on the dataset on which they are evaluated).

After I read the `traindays` dataset, I will add in a new temporary column called `holdout` that will be initialized from a random array:¹⁴

```
from pyspark.sql.functions import rand
SEED = 13
traindays = traindays.withColumn("holdout", rand(SEED) > 0.8) # 20%
traindays.createOrReplaceTempView('traindays')
```

I am passing in a seed so that I get the exact same array (and hence the same set of held-out days) every time I run this Spark code.

The first few rows of the `traindays` table are now as follows:

```
Row(FL_DATE=u'2015-01-01', is_train_day=u'True', holdout=False),
Row(FL_DATE=u'2015-01-02', is_train_day=u'False', holdout=True),
Row(FL_DATE=u'2015-01-03', is_train_day=u'False', holdout=False),
Row(FL_DATE=u'2015-01-04', is_train_day=u'True', holdout=False),
Row(FL_DATE=u'2015-01-05', is_train_day=u'True', holdout=True),
```

Note that we have both `is_train_day` and `holdout`—obviously, we are not going to be holding out any *test* data, so the query to pull training samples is as follows:

```
SELECT
*
```

¹³ Recall that we stored the training versus test days both as a BigQuery table and as a CSV file on Cloud Storage. We had to save the `traindays` dataset to persistent storage because otherwise we would have run into the problem of different hash function implementations in Spark, Pig, and Tensorflow. There would have been no way to evaluate model performance on the same dataset.

¹⁴ See `experimentation.ipynb` and `experiment.py` in `07_sparkml`.

```

FROM flights f
JOIN traindays t
ON f.FL_DATE == t.FL_DATE
WHERE
    t.is_train_day == 'True' AND
    t.holdout == False AND
    f.CANCELLED == 'False' AND
    f.DIVERTED == 'False'

```

After we have trained the model, we evaluate it not on the test data as before, but on the held-out data:

```

evalquery = trainquery.replace("t.holdout == False",
                             "t.holdout == True")

```

After we have developed this code, we can export it out of the notebook into a stand-alone script so that it can be submitted conveniently.

Feature Selection

Let's use this experimental framework and the held-out dataset to decide whether all three input variables are important. As explained earlier, we will remove one variable at a time and check the RMSE on the evaluation dataset. Conveniently, this involves changing only the `to_example()` method from:

```

def to_example(raw_data_point):
    features = [
        fields['DEP_DELAY'],
        fields['DISTANCE'],
        fields['TAXI_OUT'],
    ]
    return LabeledPoint(
        float(fields['ARR_DELAY'] < 15), #ontime
        features)

```

to:

```

def to_example(raw_data_point):
    features = [
        # fields['DEP_DELAY'],
        fields['DISTANCE'],
        fields['TAXI_OUT'],
    ]
    return LabeledPoint(
        float(fields['ARR_DELAY'] < 15), #ontime
        features)

```

Creating a large cluster

In the last chapter, running the logistic regression program script with serverless Spark took 15 minutes. Of this, nearly two minutes were for the cluster to start and an additional three minutes for it to autoscale to a sufficient size.

We are going to run many variations of this job, and because we don't quite know what possibilities we want to try, we want jobs that start immediately and finish quickly. Saving five minutes by having an already sized cluster ready to go will speed up our experimentation by 33% for our use case.

Therefore, let's create a cluster consisting of 50 machines that have 8 vCPUs each:¹⁵

```
#!/bin/bash
gcloud dataproc clusters create ch7cluster \
--enable-component-gateway \
--region ${REGION} --zone ${REGION}-a \
--master-machine-type n1-standard-4 \
--master-boot-disk-size 500 \
--num-workers 30 --num-secondary-workers 20 \
--worker-machine-type n1-standard-8 \
--worker-boot-disk-size 500 \
--project $PROJECT \
--scopes https://www.googleapis.com/auth/cloud-platform
```

Unfortunately, when I tried it, the script immediately exited with an error:

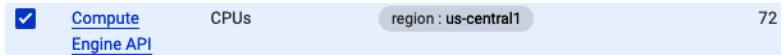
```
ERROR: (gcloud.dataproc.clusters.create) INVALID_ARGUMENT:
Insufficient 'CPUS' quota. Requested 404.0, available 67.0.
```

I didn't have the necessary quota—I needed 404 CPUs, but had only 67 available.

Increasing quota

The number of CPUs that I'm allowed to use in a region is an example of a *soft quota*. Soft quotas are meant to guard against human error and billing surprises. They are quite easy to change.

To increase my quota, I visited [the Quotas page on the GCP web console](#) and found the Compute Engine CPUs quota for the region I am interested in:



My quota is 72, but I must already be using 5, which is why the error indicated that only 67 are available.

I then clicked on “Edit quotas,” requested 500, and explained why:

¹⁵ See `create_large_cluster.sh`.

Compute Engine API

Quota: CPUs - us-central1

Current limit: 72

Enter a new quota limit. Your request will be sent to your service provider for approval. [?](#)

New limit *

500

Request description *

Needed to demonstrate experimentation in Spark.

Your description will be sent to your service provider and is used to evaluate your request. It's useful to include the intent of the quota usage, future growth plans, region or zone spread, and any additional requirements or dependencies.

[DONE](#)

Amazingly enough, the approval for my quota increase arrived in my email inbox within a minute of my requesting it. However, the email asked me to wait 15 minutes for the quota increase to get propagated throughout.

Once the quotas page showed that my quota increase had gone through, I was able to rerun the cluster creation command successfully.

Autoscale up and down

A 400-CPU machine cluster feels a bit dangerous, though—what if we forget and leave it running for months on end? Let's add an autoscaling policy to this cluster (we could have done it when creating the cluster too, but creating the cluster first and then adding the autoscaling policy allows us to verify that we have the quota to go as high as needed).

First, I create a policy file:

```
workerConfig:  
  minInstances: 2  
  maxInstances: 30  
secondaryWorkerConfig:  
  maxInstances: 20  
basicAlgorithm:  
  cooldownPeriod: 15m  
yarnConfig:  
  scaleUpFactor: 0.05  
  scaleDownFactor: 1.0  
  gracefulDecommissionTimeout: 1h
```

Then, I update the policy of the cluster that I just created:

```

gcloud dataproc autoscaling-policies import experiment-policy \
--source=autoscale.yaml --region=$REGION

gcloud dataproc clusters update ch7cluster \
--autoscaling-policy=experiment-policy --region=$REGION

```

The difference between this and serverless Spark is that once this cluster is started, we can simply submit jobs to it. There is no need to wait for the cluster to start at the beginning of each experiment. At the same time, the cluster will remain at the auto-scaled size, and as long as we submit the next job to it within 15 minutes, we can continue using the machines already spun up. After an hour of no activity, the cluster is decommissioned.

Yes, this is all a bit hacky. BigQuery and its near-instantaneous spin-up, autoscaling, and spin-down are so much nicer! Still, we are in Spark world here, and this is pretty flexible for being in Spark world.

When you're using a large cluster and submitting jobs one after the other, it is a good idea to monitor the Hadoop nodes to ensure that all the machines are available and that jobs are being spread over the entire cluster. You can access the monitoring details from the Dataproc section of the GCP web console. You should see an uptick in CPU usage during the regression part of the code and use of all 50 node managers as demonstrated in [Figure 7-7](#).

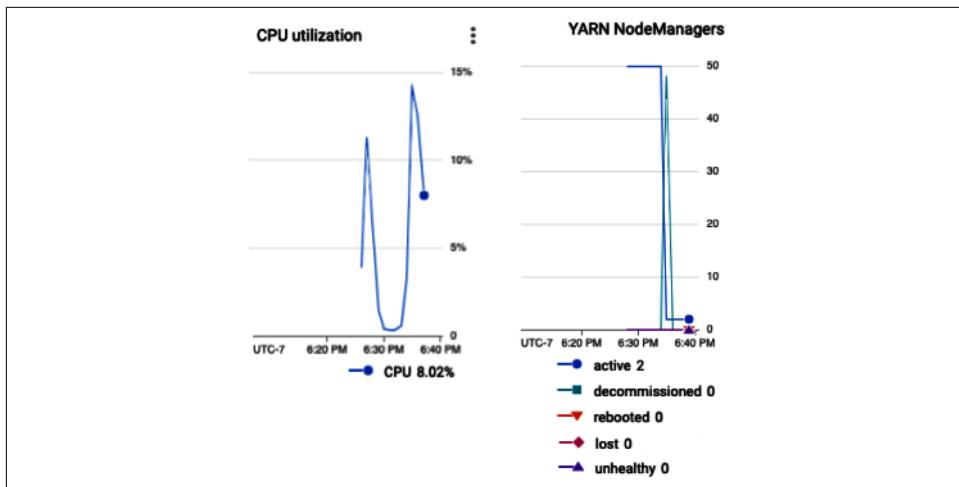


Figure 7-7. Monitor the CPU usage and use of all the nodes during the regression phase to ensure that all the machines are being used effectively.

In my case, monitoring revealed that simply increasing the number of workers didn't actually spread out the processing to all the nodes. This is because Spark estimates the number of partitions based on the raw input data size (here just a few gigabytes), so

that estimate is probably too low for our 50-worker cluster. Hence, I modified the reading code to add in an explicit repartitioning step:

```
traindata = spark.sql(trainquery).repartition(1000)
```

And similarly, when reading the evaluation data:

```
evaldata = spark.sql(evalquery).repartition(1000)
```

After this bit of optimization, I am able to run an experiment in about 10 minutes, as I expected, by backing out the startup and upscaling time from the serverless Spark runtime. The longer the experiment, the less of a benefit all this fine tuning provides. We are better off leaving it to serverless Spark.

Removing features

Now that we have a faster way of running the job, we can carry out our experiment of removing one variable at a time (i.e., in Experiment 3, we use `DEP_DELAY` and `TAXI_OUT` as input features):

Experiment #	Variables	RMSE	Percent increase in RMSE
1	DEP_DELAY DISTANCE TAXI_OUT	0.205	N/A
2	Remove DEP_DELAY	0.361	76%
3	Remove DISTANCE	0.207	1%
4	Remove TAXI_OUT	0.227	11%

It is clear that some variables are dramatically more important than others, but also that all three variables do carry information and you should not discard them. Contrary to my assumption in [Chapter 6](#), the distance has the least amount of impact—had I known this earlier, [Chapter 6](#) might have involved Bayesian selection of departure delay and taxi-out time!

Feature Transformations

In the previous section, we carried out feature selection to determine that the distance variable was the least important of the variables included in the model. However, there is another possibility why including the distance variable did not affect the RMSE greatly—it could be that the distance variable ranges over a wider range (thousands of miles), whereas the time intervals range to only a few minutes. Therefore, the distance variable might be overpowering the other variables—assuming that both variables are given equal weights, the impact of the distance variable will be magnified 1,000-fold over that of the time interval. In an effort to reduce the impact of distance on both the sum and the gradient, the optimization process might move the distance weight to 0. And thus, the distance might end up not playing much of a role.

This is less likely in the case of logistic regression because it is a linear model, and gradients in linear models can be scaled more effectively. However, having all the variables have similar magnitudes will become important as our models become more complex.

Scaling

Another reason we'd like to scale all the input values so that they all have similar (and small) magnitudes is that the initial, random weights in optimization routines are often set to be in the range -1 to 1 , and this is where the optimizer starts to search. Thus, starting with all the variables having less than unit magnitude can help the optimizer converge more efficiently and effectively. Following are common choices for scaling functions:

- To map the minimum value of the variable to -1 and the maximum to 1 . This involves scanning the dataset to find the minimum and maximum within each column—Spark has a class called `AbsScaler` that will do this for us (but it requires an extra pass through the data). However, the choices of the minimum and maximum need not be exact, so we can use the data exploration that we carried out in previous chapters to do an approximate scaling. As long as we scale the variables the same way during training and during prediction (for example, if we scale linearly such that a distance of 30 maps to -1.0 and a distance of 6,000 maps to 1.0 both in training and during prediction), the precise minimum and maximum don't matter.
- To map the mean of the variable within the column to 0 and values one standard deviation away to -1 or 1 . The tails of the distribution will map to quite large values, but such values will also be rarer. This serves to emphasize unusual values while linearly scaling the common ones.

Let's use Option 1, which is to do linear scaling between an approximate minimum and maximum. Let's assume that departure delays are in the range $(-30, 30)$, distance in the range $(0, 2,000)$, and taxi-out times in the range $(0, 20)$. These are approximate, but quite reasonable, and using these reasonable values helps us to avoid being overly affected by outliers in the data.¹⁶ To map these ranges to $(-1, 1)$, this is the transformation we would carry out on the input variables inside the `to_example()` function that converts a row into a training example:

```
def to_example(raw_data_point):
    return LabeledPoint(\
```

¹⁶ This presupposes that you have done exploratory data analysis, and understand your data. If you don't know what a reasonable range is, then you could use the 5th and 95th percentiles of the column values in the training data.

```

        float(raw_data_point['ARR_DELAY'] < 15), #ontime \
[ \
    raw_data_point['DEP_DELAY'] / 30, \
    (raw_data_point['DISTANCE'] / 1000) - 1, \
    (raw_data_point['TAXI_OUT'] / 10) - 1, \
]
)

```

After making these changes to scale the three variables, and running the experiment again, I see that the RMSE is unaffected—scaling doesn't make a difference:

Experiment #	Variables	RMSE	Percent improvement
1 (values from experiment #1 repeated for convenience)	Raw values of DEP_DELAY DISTANCE TAXI_OUT	0.20537	N/A
5	Scaled values of the three variables	0.20537	0

Clipping

Another possible preprocessing that we could carry out is called *clipping*. Values that are beyond what we would consider reasonable are clamped to the bounds. For example, we could treat distances of more than 2,000 miles as 2,000 miles, and departure delays of more than 30 minutes as 30 minutes. This allows the optimization algorithm to focus on the part of the data where the majority of the data lies, and not be pulled off the global minimum by outliers. Also, some error metrics can be susceptible to outliers, so it is worth doing one more experiment after clipping the input variables.

Adding clipping to scaled variables is straightforward:

```

def to_example(raw_data_point):
    def clip(x):
        if x < -1:
            return -1
        if x > 1:
            return 1
        return x
    return LabeledPoint(
        float(raw_data_point['ARR_DELAY'] < 15), #ontime \
[ \
    clip(raw_data_point['DEP_DELAY'] / 30), \
    clip((raw_data_point['DISTANCE'] / 1000) - 1), \
    clip((raw_data_point['TAXI_OUT'] / 10) - 1), \
]
)

```

Carrying out the experiment on clipped variables, and adding in the RMSE to the table, we see these results:

Experiment #	Transform	RMSE	Percent improvement	Keep transform?
1 (repeated for convenience)	Raw values of DEP_DELAY DISTANCE TAXI_OUT	0.20537	N/A	N/A
6	Scaled	0.20537	None	No
7	Clipped	0.20538	Negative	No

It turns out that neither scaling nor clipping matters for this algorithm (logistic regression) in this framework (Spark ML). In general, though, experimenting with different preprocessing transforms should be part of your workflow. It could have a dramatic impact.

Feature Creation

So far, we have tried out the three numeric predictors in our dataset. Why did I pick only the numeric fields in the dataset? Because the logistic regression model is, at heart, just a weighted sum. We can quite easily multiply and add numeric values (actually not all numeric values but those that are continuous),¹⁷ but what does it mean to use a timestamp such as 2015-03-13-11:00:00 as an input variable to the logistic regression model?

We cannot simply convert such a timestamp to a numeric quantity, such as the day number within the year, and then add it to the model. One rule of thumb is that to use a value as input to a model, you should have at least 5 to 10 instances of that value appearing in the data. Columns that are too specific to a particular row or handful of rows can cause the model to become *overfit*—an overfit model is one that will perform extremely well on the training dataset (essentially, it will memorize exactly what happened at each historical timestamp—for example, that flights in the Midwest were delayed on May 11, 2015) but then not work well on new data. It won’t work well on new data because the timestamps in that data (such as May 11, 2018) will not have been observed.¹⁸

Therefore, we need to do smarter things with attributes such as timestamps so that they are both relevant to the problem and not too specific. We could, for example, use the hour of day as an attribute in the model. The hour of day might matter—most airports become busy in the early morning and early evening because many flights

¹⁷ For example, you cannot add and multiply employee ID numbers even if they are numeric. Employee ID numbers are not continuous.

¹⁸ The reason that we stratify our dataset into training and validation is to catch nonobvious instances of overfitting. If the training loss is very low, but the validation error is high, it is likely that the model has overfit. It is then up to you to diagnose the reason!

are scheduled around the daily work schedules of business travelers. In addition, delays accumulate over the day because an aircraft that is delayed in arriving is also delayed on takeoff.

Suppose that we extract the hour of day from the timestamp. Given a timestamp such as 2015-03-13-11:00:00, what is the hour? It's 11, of course, but the 11 is in the time zone corresponding to the UTC time zone. On the other hand, we care about the time zone at the American airport because many airports are busier early in the morning and late in the evening. This is one instance for which it is local time that matters. Thus, to extract the hour of day, we will need to correct by the time zone offset and then extract the hour of day. The *feature* hour of day is computed from two input variables—the departure time and the time zone offset.

It is worth pausing here and clarifying that I am making a distinction between the words *input* and *feature*—the timestamp is the input, the hour of day is the feature. What the client application provides when it wants a prediction is the input, but what the ML model is trained on is the feature. The feature could be (as in the case of scaling of input variables) a transformation of an input variable. In other cases, as with the hour of day, it could be a combination of multiple input variables. The `to_example()` method is the method that converts inputs (`raw_data_point`) to *examples* (where each example is a tuple of features and a label). Different machine learning APIs will ask for inputs, features, or examples, and it is good to be clear on what exactly the three terms mean.

Given a departure timestamp and a time zone offset,¹⁹ we can compute the hour in local time using the time-handling code we discussed in [Chapter 4](#):

```
def to_example(raw_data_point):
    def get_local_hour(timestamp, correction):
        import datetime
        TIME_FORMAT = '%Y-%m-%dT%H:%M:%S'
        t = datetime.datetime.strptime(timestamp, TIME_FORMAT)
        d = datetime.timedelta(seconds=correction)
        t = t + d
        return t.hour
    return LabeledPoint(
        float(raw_data_point['ARR_DELAY'] < 15), #ontime \
        [ \
            raw_data_point['DEP_DELAY'], \
            raw_data_point['TAXI_OUT'], \
            get_local_hour(raw_data_point['DEP_TIME'], \
                           raw_data_point['DEP_AIRPORT_TZOFFSET']) \
        ])
```

¹⁹ The time zone offset is a float, and must be added to the schema as such.

There is one potential problem with treating the hour of day as a straightforward number. Hour 22 and hour 2 are only 4 hours apart, and it would be good to capture this somehow. An elegant way to work with periodic variables in machine learning is to convert them to two features— $\sin(\theta)$ and $\cos(\theta)$, where θ in this case would be the angle of the hour hand in a 24-hour clock:²⁰

```
def to_example(raw_data_point):
    def get_local_hour(timestamp, correction):
        import datetime
        TIME_FORMAT = '%Y-%m-%dT%H:%M:%S'
        t = datetime.datetime.strptime(timestamp, TIME_FORMAT)
        d = datetime.timedelta(seconds=correction)
        t = t + d
        theta = np.radians(360 * t.hour / 24.0)
        return [np.sin(theta), np.cos(theta)]

    features = [ \
        raw_data_point['DEP_DELAY'], \
        raw_data_point['TAXI_OUT'], \
    ]
    features.extend(get_local_hour(raw_data_point['DEP_TIME'],
                                   raw_data_point['DEP_AIRPORT_TZOFFSET']))
    return LabeledPoint(\n        float(raw_data_point['ARR_DELAY'] < 15), #ontime \
    features)
```

This encoding of a periodic variable using the sin and cos makes it two features. These two features capture the information present in the periodic variable, but do not distort the distance between two values.

Another approach would be to *bucketize* the hour. For example, we could group hours 20 to 23 and hours 0 to 5 as “night,” hours 6 to 9 as “morning,” and so on. Obviously, bucketizing takes advantage of what human experts know about the problem. We suspect that the behavior of flight delays changes depending on the time of day—long taxi-out times during busy hours are probably built into the scheduled arrival time, but a plane that experiences a long taxi-out because the towing vehicle broke down and had to be replaced will almost definitely arrive late. Hence, our bucketizing of the hour of day relies on our intuition of what the busy hours at an airport are:

```
def get_category(hour):
    if hour < 6 or hour > 20:
        return [1, 0, 0] # night
    if hour < 10:
        return [0, 1, 0] # morning
    if hour < 17:
```

²⁰ The distribution of hours in the dataset, we assume, follows the Fisher–von Mises distribution, which describes points distributed on an n -dimensional sphere. If $n = 2$, this reduces to points on a unit circle. An hour hand is just that.

```

        return [0, 0, 1] # mid-day
    else:
        return [0, 0, 0] # evening
def get_local_hour(timestamp, correction):
    ...
    return get_category(t.hour)

```

You might find two things odd about the preceding snippet:

- We are returning binary numbers, for example [1, 0, 0]. The first number in the triplet captures whether the hour is at night or not. The second captures whether it is in the morning or not. In essence, we convert the hour variable which is numeric and continuous into four independent variables, each of which is 1 or 0.
- But we don't have four binary numbers, one for each class. We have only three! Note that the vector corresponding to the last category is [0, 0, 0] and not [0, 0, 0, 1], as you might have expected. This is because we don't want the four features to always add up to one—that would make them linearly dependent. This trick of dropping the last column keeps the values independent. Assuming that we have N categories, bucketizing will make the hour variable into N – 1 features.

How do we know which of these methods—the raw value, the sin/cos trick, or bucketizing—works best for the hour of day? We don't, so we need to run an experiment and choose (note that we are using the departure delay, distance, and taxi-out, and now adding a new variable to see if it helps):

Experiment #	Transform	RMSE
1 (repeated for convenience)	Without hour	0.20537
8	raw hour	0.20536
9	sin(theta) cos(theta)	0.20535
10	bucketize	0.20538

The fact that we aren't able to reduce the RMSE by adding the hour information suggests that the hour of day is already adequately captured by the scheduling rules used by the airlines. It can be tempting to simply throw away the hour information, but we should follow our systematic process of keeping all our variables and then discarding one variable at a time—it is possible that the hour of day doesn't matter now, but it might after we include some other variables. So, for now, let's pick one of the possibilities arbitrarily—I will use the bucketed hour as the way to create a feature from the timestamp. Of course, I could have created additional features from the timestamp input variable—day of week, season, and so on.

Spark ML supports a rich set of feature transforms—it is a good idea to go through that list and learn the types of variables for which they are meant. Knowing the tools

available to you is a prerequisite to be able to call on them when appropriate. If this is the first time you are encountering machine learning, you might be surprised by how much we are relying on our experimental framework. Running experiments like this, though, is the unglamorous work that lies behind most machine learning applications. Although my approach in this chapter has required careful record keeping, a better way would be if our machine learning framework would provide structure around experiments, not just single training operations. Spark ML **provides this functionality** via the `CrossValidator` tool, but even that still requires quite a bit of scaffolding.

There is another problem, though. The runtime increased from 10 minutes to 15 minutes with the addition of the hour variable. This doesn't bode well—we have a lot more features we want to try to include.

Categorical Variables

How about using the airport code as a predictor? What we are doing when using the airport code as a predictor is that we are asking the ML algorithm to learn the idiosyncrasies of each airport. I remember, for example, sitting on the taxiway at New York's LaGuardia airport for nearly 45 minutes and then being surprised when the flight arrived in Dallas a few minutes ahead of schedule! Apparently, a 45-minute taxi-out time in New York is quite common and nothing to be worried about.

To use timestamp information, we extracted a numeric part of the timestamp—the hour—and used it in our model. We tried using it in raw form, as a periodic variable, and as a bucketized set of categories. We cannot use that approach here because there is no numeric component to the letters DFW or LGA. So how can we use the airport code as an input to our model?

The trick here is to realize that bucketizing the hour was a special case of making the variable categorical. A more bludgeon-like, but often effective, approach is to do *one-hot encoding*. Essentially, the hour variable is made into 24 features. The 11th feature is 1.0 and the rest of the features 0.0 when the hour is 11, for example. One-hot encoding is the standard way to deal with *categorical* features (i.e., features for which there is no concept of magnitude or ranking between different values of the variable).²¹ This is the way we'll need to encode the departure airport if we were minded to include it in our model—we'd essentially have one feature per airport, so that the DFW feature would be 1, and the rest of the features 0 for flights that departed from DFW.

²¹ It is not the case that all strings are categorical and all numeric columns are continuous. To use my previous example, an employee ID might be numeric but is categorical. On the other hand, student grades (A+, A, A-, B+, B, etc.) are strings but can easily be translated to a continuous variable.

Unlike bucketing hours, though, we need to find all the possible airport codes (called the *vocabulary*) and assign a specific binary column to them. For example, we might need to assign DFW to the 143rd column. Fortunately, we don't need to write the code. One-hot encoding is available as a prebuilt feature transformation in Spark; we can add a new column of vectors to the `traindata` dataframe using this code:

```
def add_categorical(df):
    from pyspark.ml.feature import OneHotEncoder, StringIndexer
    indexer = StringIndexer(inputCol='ORIGIN',
                            outputCol='origin_index')
    index_model = indexer.fit(df) # ❶
    indexed = index_model.transform(df) # ❷
    encoder = OneHotEncoder(inputCol='origin_index',
                            outputCol='origin_onehot')
    return encoder.transform(indexed) # ❸
traindata = add_categorical(traindata)
```

- ❶ Create an index from origin airport codes (e.g., DFW) to an origin index (e.g., 143).
- ❷ Transform the dataset so that all flights with `ORIGIN=DFW` have `origin_index=143`.
- ❸ One-hot encode the index into a binary vector that is used as input to training.

During evaluation, the same change needs to be made to the dataset, except that the index model will need to be reused from training (so that DFW continues to map to 143). In other words, we need to save the `index_model` and carry out the last three lines before prediction. So we modify the `add_categorical()` method to:²²

```
index_model = None
def add_categorical(df, train=False):
    from pyspark.ml.feature import OneHotEncoder, StringIndexer
    if train:
        indexer = StringIndexer(inputCol='ORIGIN',
                                outputCol='origin_index')
        index_model = indexer.fit(df)
        indexed = index_model.transform(df)
        encoder = OneHotEncoder(inputCol='origin_index',
                                outputCol='origin_onehot')
        return encoder.transform(indexed)
    traindata = add_categorical(traindata, train=True)
    ...
    evaldata = add_categorical(evaldata)
```

²² See `experiment.py`.