



BÀI GIẢNG NHẬP MÔN KHOA HỌC DỮ LIỆU

Trần Quang Quý*

Hà Thị Thanh†

Thái Nguyên, tháng 06 năm 2023

Mục lục

1	Mở đầu	5
2	Giới thiệu chung	6
2.1	Tổng quan về Khoa học dữ liệu	6
2.1.1	Lịch sử	6
2.1.2	Dữ liệu hóa	6
2.1.3	Xu hướng tương lai	6
2.2	Ngành Khoa học dữ liệu	7
2.2.1	Khái niệm	7
2.2.2	Các cấp độ	7
2.3	Ngôn ngữ lập trình	9
2.3.1	Ngôn ngữ Python	9
2.3.2	Ngôn ngữ R	11
2.3.3	Power BI	12
3	Cài đặt môi trường Anaconda	15
3.1	Anaconda và Jupyter Notebook	15
3.1.1	Tìm hiểu Anaconda	15
3.1.2	Tìm hiểu Jupyter Notebook	17
3.1.3	Cài đặt Anaconda và Jupyter Notebook	17
3.2	Thao tác dữ liệu với Python	17
3.2.1	Dữ liệu Data Frame	19
3.2.2	Đọc dữ liệu	20

*Trường Đại học Công nghệ Thông tin & Truyền thông - ICTU, tqquy@ictu.edu.vn

†Trường Đại học Công nghệ Thông tin & Truyền thông - ICTU, httthanh@ictu.edu.vn

3.2.3	Thao tác dữ liệu	21
3.3	Kết chương	33
4	Thống kê mô tả	33
4.1	Dữ liệu và tham số mô tả	33
4.1.1	Số liệu thống kê	33
4.1.2	Các tham số mô tả	34
4.1.3	Tổng thể và mẫu	34
4.2	Chuẩn bị dữ liệu	35
4.2.1	Thu thập và chuẩn bị dữ liệu	35
4.2.2	Tập dữ liệu Adult	36
4.3	Tính toán tham số	41
4.3.1	Giá trị trung bình	41
4.3.2	Giá trị phương sai	42
4.3.3	Giá trị trung vị	43
4.4	Phân phối dữ liệu	44
4.4.1	Hàm khối xác suất	46
4.4.2	Hàm phân phối tích lũy	47
4.5	Dữ liệu ngoại lai	48
4.5.1	Khái niệm	48
4.5.2	Tìm giá trị ngoại lai	49
4.6	Bất đối xứng	51
4.6.1	Độ lệch	51
4.6.2	Hệ số lệch trung vị Pearson	53
5	Phân phối xác suất	53
5.1	Phân phối nhị thức	53
5.1.1	Khái niệm	53
5.1.2	Tính toán phân phối nhị thức	54
5.2	Phân phối Poisson	55
5.2.1	Khái niệm	55
5.2.2	Tính toán phân phối Poisson	55
5.3	Phân phối đều liên tục	56
5.3.1	Khái niệm	56
5.3.2	Tính toán phân phối đều	57
5.4	Hàm phân phối mũ	57
5.4.1	Khái niệm	57
5.4.2	Tính toán phân phối mũ	57
5.5	Phân phối chuẩn	57
5.5.1	Khái niệm	57
5.6	Phân phối Chi-squared	59
5.6.1	Khái niệm	59
5.7	Phân phối Student	59
5.8	Phân phối F	60
6	Ước lượng khoảng	61
6.1	Khái niệm	61

6.2	Ví dụ	62
7	Ước lượng quần thể	62
7.0.1	Ước tính trung bình cho quần thể	63
7.0.2	Ước tính khi biết phương sai	63
8	Trực quan dữ liệu	65
8.1	Tổng quan về trực quan	65
8.1.1	Khái niệm	65
8.1.2	Các bước chính trực quan dữ liệu	66
8.2	Thư viện ggplot2	67
8.2.1	Cấu trúc lớp	67
8.2.2	Lớp ggplot	67
8.2.3	Lớp geom	68
8.2.4	Trực quan nhóm	72
8.2.5	Cơ chiều dữ liệu	73
8.2.6	Trực quan bằng facet	74
8.3	Ánh xạ dữ liệu	75
8.3.1	Ánh xạ trong hàm	75
8.3.2	Ánh xạ ngoài hàm	76
8.4	Biểu đồ thống kê	77
8.4.1	Xây dựng biểu đồ	77
8.4.2	Biểu đồ Histogram	78
9	Khám phá dữ liệu EDA	84
9.1	Thư viện dplyr	84
9.1.1	Giới thiệu về toán tử pipe	85
9.1.2	Toán tử Forward Pipe	85
9.1.3	Toán tử T pipe	87
9.1.4	Toán tử Assigning pipe	88
9.1.5	Toán tử Backward pipe	90
9.2	Các hàm chính trong dplyr	90
9.2.1	Hàm Select	91
9.2.2	Hàm filter	91
9.2.3	Hàm mutate	92
9.2.4	Hàm summarize và group_by	92
9.3	Phân tích dữ liệu gapminder	93
9.3.1	Lựa chọn cột	93
9.3.2	Bổ sung cột	95
9.3.3	Thống kê nhóm	96
9.3.4	Tần số xuất hiện	98
9.4	Trực quan dữ liệu	99
9.4.1	Biểu đồ dây	99
9.4.2	Biểu đồ mật độ xác suất	101
9.4.3	Nâng cao	102
9.5	Biến đổi chiều	103

10	Hệ số tương quan	105
10.1	Khái niệm	105
10.2	Hệ số tương quan Pearson	105
10.2.1	Karl Pearson	106
10.2.2	Ý nghĩa của hệ số Pearson	106
10.2.3	Công thức tính hệ số tương quan Pearson	107
10.3	Hệ số tương quan Spearman	107
10.3.1	Charles Edward Spearman	108
10.3.2	Tính hệ số tương quan Spearman	108
10.4	Hệ số tương quan Kendall	108
10.5	Áp dụng:	109
11	Hồi quy tuyến tính	113
11.1	Khái niệm	113
11.2	Bài toán	113
11.3	Phân tích toán học	114
11.3.1	Dạng của Hồi quy tuyến tính	114
11.3.2	Sai số dự đoán	114
11.3.3	Hàm mất mát	114
11.3.4	Nghiệm cho bài toán	115
11.4	Thực hành trên Python	115
11.4.1	Hiển thị dữ liệu	116
11.4.2	Nghiệm theo công thức	117
11.4.3	Nghiệm theo thư viện scikit-learn	119
11.5	Hồi quy đa thức	119
11.6	Hạn chế của hồi quy tuyến tính	120
12	Kiểm định giả thuyết thống kê	120
12.1	Khái niệm	120
12.2	Mô hình chung	124
12.3	Kiểm định Chi bình phương	125
12.3.1	Bài toán ví dụ	127
12.4	Kiểm định t-test	131
12.4.1	Phân phối t-student	131
12.4.2	Bài toán	132
12.4.3	Kiểm định	133
12.4.4	One-Sample t-Test	135
12.5	Kiểm định t-test trong R	136
12.5.1	Kiểm định t-test 1 mẫu	136
13	Khoa học dữ liệu trong kinh doanh	141
13.1	Khai phá luật kết hợp	141
13.2	Thuật toán Apriori	144
13.2.1	Mô tả thuật toán Apriori	145
13.2.2	Bài tập với Apriori	146
13.3	Tập dữ liệu Retail	147
13.3.1	Load thư viện	147

13.3.2	Bổ sung thuộc tính	149
13.3.3	Ma trận thưa	150
13.4	Tạo luật kết hợp	152
13.4.1	Các giới hạn	154
13.4.2	Loại bỏ quy tắc	154
13.4.3	Tương quan	154
13.5	Trực quan luật kết hợp	156
13.5.1	Scatter-Plot	156
13.5.2	Graph-Based Visualizations	158
13.5.3	Individual Rule	159
14	Tài liệu tham khảo	160
15	Phụ lục: Phân tổng hợp các code Python và R được sử dụng	161

1 Mở đầu

Data science còn được biết đến với tên gọi là Khoa học dữ liệu. Đúng với tên gọi của nó, về mặt bản chất, đây chính là công việc thu thập và phân tích dữ liệu. Data science là một lĩnh vực liên ngành mà trong đó, những bộ dữ liệu được xử lý, sắp xếp và giải mã bằng các mô hình thống kê hay phương pháp toán học.

Các kỹ sư khoa học dữ liệu sẽ sử dụng tất cả các điểm dữ liệu khác nhau thu thập được, từ đó tạo ra một mô hình dữ liệu hoặc thuật toán để áp dụng cho từng mục đích cụ thể mang tính chiến lược. Có 2 mục đích chính:

- Sử dụng dữ liệu để phân tích chuyên sâu về một vấn đề nào đó và đưa ra các giải pháp, hoặc dự đoán cho tương lai.
- Xây dựng các mô hình dữ liệu để tạo ra các sản phẩm, hoặc tính năng công nghệ nào đó.

Có thể nói, ứng dụng của Data Science là vô cùng quan trọng khi số lượng dữ liệu được tạo ra mỗi ngày là rất lớn, tăng lên theo cấp số nhân, và ngày càng trở nên phức tạp. Điều đó khiến việc xử lý dữ liệu thông qua các công cụ như Excel không còn khả thi. Thay vào đó, các kỹ sư khoa học dữ liệu phải sử dụng đến ngôn ngữ lập trình cao cấp hơn, phổ biến nhất là Python và R, để hiểu dữ liệu.

Trong nội dung học phần “Nhập môn khoa học dữ liệu”, tiếng Anh là Introduction to Data Science sẽ trình bày những hiểu biết cũng như những kiến thức về thống kê, phân tích dữ liệu và một số nội dung khác liên quan tới Khoa học dữ liệu. Hai ngôn ngữ được sử dụng chính trong quyển bài giảng là ngôn ngữ Python và ngôn ngữ R.

2 Giới thiệu chung

2.1 Tổng quan về Khoa học dữ liệu

2.1.1 Lịch sử

Chắc chắn bạn đã có kinh nghiệm về khoa học dữ liệu dưới nhiều hình thức. Điều này có thể khẳng định được! Ví dụ đơn giản là khi bạn đang tìm kiếm thông tin trên web bằng cách sử dụng công cụ tìm kiếm hoặc hỏi đường trên điện thoại di động của mình, bạn đang tương tác với các sản phẩm khoa học dữ liệu. Khoa học dữ liệu đã hỗ trợ giải quyết một số nhiệm vụ hàng ngày phổ biến nhất của chúng ta trong vài năm. Thuật ngữ Khoa học dữ liệu không phải là mới mẻ, mà nó là sản phẩm của Toán học và Thống kê từ xưa trong một thời gian dài. Thống kê là một ngành khoa học lâu đời đứng trên vai những người khổng lồ của thế kỷ 18 như Pierre Simon Laplace (1749–1827) và Thomas Bayes (1701–1761)[1]

2.1.2 Dữ liệu hóa

Tính mới của khoa học dữ liệu không bắt nguồn từ kiến thức khoa học mới nhất, mà bắt nguồn từ một sự thay đổi mang tính đột phá trong xã hội của chúng ta do sự phát triển của công nghệ gây ra: dữ liệu hóa - **Datification**. Dữ liệu hóa là quá trình kết xuất thành các khía cạnh dữ liệu của thế giới chưa từng được định lượng trước đây. Ở cấp độ một thực thể nào đó, danh sách các khái niệm được dữ liệu hóa rất dài và vẫn đang tiếp tục phát triển: mạng lưới kinh doanh, quyển sách chúng ta đang đọc, những bộ phim chúng ta thích, thức ăn chúng ta ăn, hoạt động thể chất, mua sắm, hành vi lái xe của chúng ta, v.v. Ở cấp độ kinh doanh, các công ty đang xác định dữ liệu bán cấu trúc (semi-structured) mà trước đây đã từng bị loại bỏ không được sử dụng như: nhật ký hoạt động web, hoạt động mạng máy tính, tín hiệu máy móc, v.v. Dữ liệu phi cấu trúc (unstructured), chẳng hạn như báo cáo bằng văn bản, e-mail hoặc bản ghi âm, hiện đang được lưu trữ không chỉ cho mục đích lưu trữ mà còn được phân tích[2].

2.1.3 Xu hướng tương lai

Khoa học dữ liệu (Data science) là ngành khoa học về việc khai phá, quản trị và phân tích dữ liệu để dự đoán các xu hướng trong tương lai và đưa ra các quyết định, chiến lược hành động. Khoa học dữ liệu (Data science) gồm ba phần chính: tạo và quản trị dữ liệu, phân tích dữ liệu, và áp dụng kết quả phân tích thành những hành động có giá trị. Việc phân tích và sử dụng dữ liệu dựa vào ba nguồn tri thức: toán học (thống kê toán học - Mathematical Statistics), công nghệ thông tin (máy học - Machine Learning) và tri thức của lĩnh vực ứng dụng cụ thể.

Số liệu lớn (Big Data) đã cách mạng hóa các công ty và đem lại cho họ một lợi thế cạnh tranh. Các công ty này cần những người chuyên môn, thành thạo trong việc xử lý, quản lý, phân tích và hiểu xu hướng trong dữ liệu. Chính vì thế mà ngành Khoa học dữ liệu (Data science) càng ngày càng trở thành xu hướng và được săn đón.

2.2 Ngành Khoa học dữ liệu

2.2.1 Khái niệm

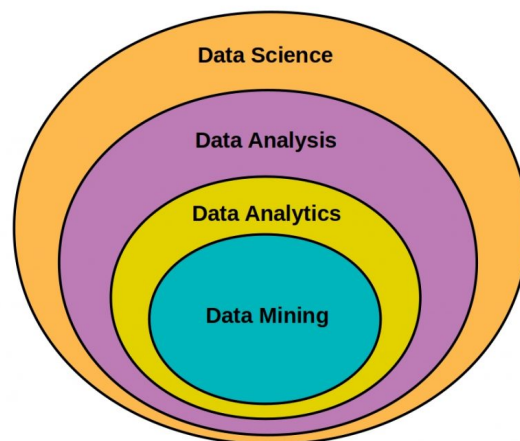
Có nhiều quan điểm khác nhau về ngành khoa học dữ liệu (data science) nhưng về cơ bản thì đây là lĩnh vực sử dụng các phương pháp khoa học, thuật toán và hệ thống công nghệ thông tin để giải mã ‘thông tin’ đằng sau các dữ liệu ngẫu nhiên, từ đó biến các báo cáo thô thành các thông tin có giá trị.

Những thông tin này có thể ứng dụng đa dạng vào nhiều lĩnh vực khác nhau như giúp ngành y tế phân tích hình ảnh y khoa, phân nhóm khách hàng trong ngành tài chính hay dự báo rủi ro cho nhóm ngành sản xuất. Dưới sự hỗ trợ của khoa học dữ liệu, nhiều vấn đề trước giờ vẫn ‘tốn nhân công, tốn nguồn lực’ để giải quyết, nay hoàn toàn có thể xử lý nhanh chóng, hiệu quả hơn, đặc biệt khi công nghệ ngày càng hoàn thiện và phổ biến..

2.2.2 Các cấp độ

Khoa học dữ liệu là một ngành lớn và có những phân cấp nhỏ hơn để giải quyết những vấn đề cụ thể, chi tiết trong các lĩnh vực khác nhau:

- Cấp thứ 1- data mining: khoanh vùng và tìm kiếm dữ liệu đầu vào dưới sự hỗ trợ của công nghệ và thuật toán thống kê.
- Cấp thứ 2- data analytics: Cách con người ứng dụng công nghệ phân tích để biến những dữ liệu đã khai thác thành các nhóm thông tin cụ thể.
- Ở cấp thứ 3- data analysis: hệ thống hóa và xây dựng cấu trúc cho những thông tin đầu ra của cấp thứ 2. Những thông tin này sẽ được đào sâu phân tích, kết nối và diễn giải thành các thông tin thông dụng, mang tính ứng dụng cao.

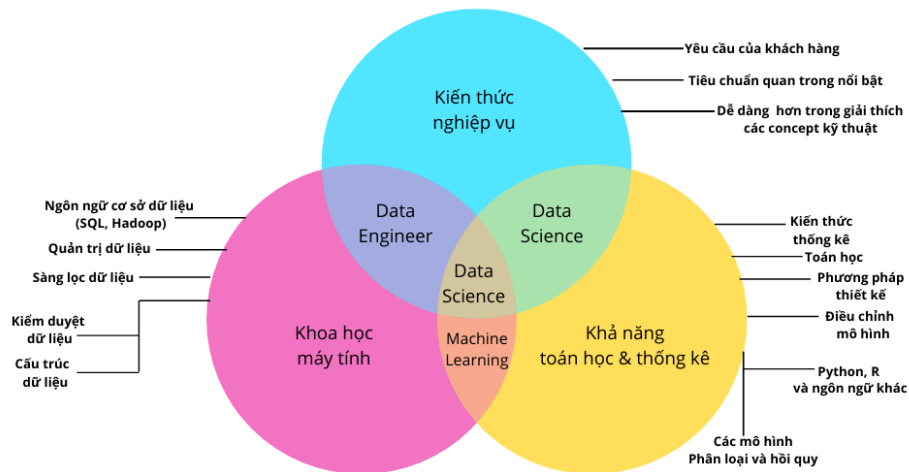


Hình 1: Các cấp độ trong Khoa học dữ liệu

Data science ngày càng được ứng dụng trong nhiều lĩnh vực khi công nghệ ngày càng cải tiến khiến cho tính chính xác của dữ liệu gia tăng. Trong đó tính ứng dụng của khoa học dữ liệu đặc biệt rõ nét trong kinh doanh vì các quyết định kinh doanh sẽ có được cơ sở chắc chắn và giúp sớm loại bỏ các quyết định nhiều rủi ro.

Các nhóm kỹ năng cần thiết của một nhà khoa học dữ liệu bao gồm Phân tích (Analytics), Lập trình (Programming), và Kiến thức chuyên ngành (Domain Knowledge). Chính vì thế, nếu theo học ngành Khoa học dữ liệu, chúng ta sẽ được học một số các môn chuyên ngành như:

- Thống kê ứng dụng (Applied Statistics)
- Nhập môn Khoa học máy tính (Introduction to Computer Science)
- Lập trình cùng Python, R hay SQL (Programming with Python/R/SQL)
- Trực quan hóa dữ liệu (Data Visualization)
- Xác suất (Probability)
- Khai phá dữ liệu (Data Mining)
- Công cụ trực quan hóa dữ liệu - Tableau, Power BI

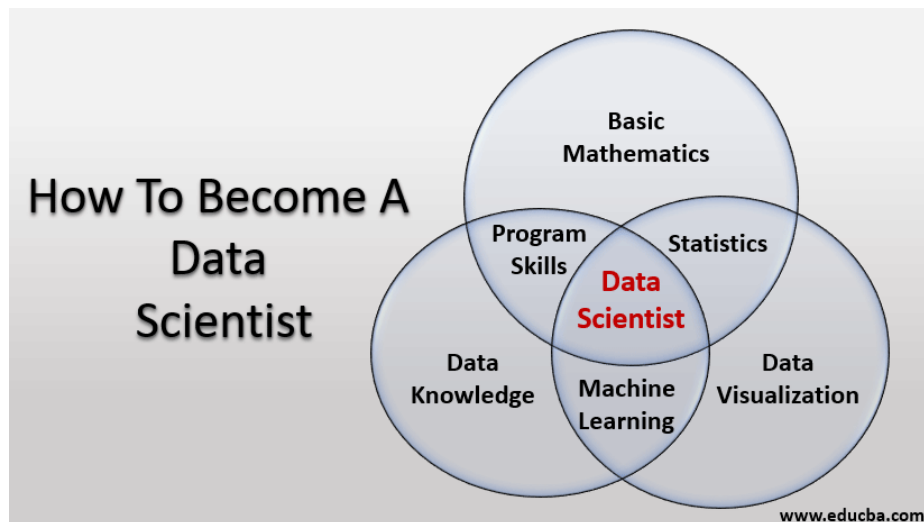


Hình 2: Tổng quan chung về Data Science

Để trở thành một nhà Khoa học dữ liệu - Data Scientist, chúng ta cần có kiến thức tổng quan như miêu tả ở Hình 2 và Hình 3.

Nếu theo đuổi ngành khoa học dữ liệu trong trường học, sau đây là những kỹ năng về kỹ thuật mà chúng ta cần học. Nói cách khác, để có thể tìm kiếm cơ hội việc làm và phát triển trong ngành data science thì chúng ta không thể bỏ qua những kỹ năng quan trọng này:

- Phân tích thống kê và tính toán
- Xử lý tập dữ liệu lớn
- Sắp xếp dữ liệu
- Machine learning
- Deep learning



Hình 3: Kỹ năng để trở thành Data Scientist

- Trực quan hoá dữ liệu
- Toán học
- Lập trình (các ngôn ngữ lập trình Python, Java, C/C++, v.v.)
- Big data
- Thống kê

Ngoài ra, ngành khoa học dữ liệu nói chung và một data scientist nói riêng không thể thiếu những công cụ hỗ trợ phân tích dữ liệu như SQL, R, Spark, SAS, Hive, v.v.

2.3 Ngôn ngữ lập trình

Có rất nhiều ngôn ngữ (công cụ) giúp chúng ta thao tác với dữ liệu và khám phá dữ liệu. Gần đây, hai ngôn ngữ như Python và R đã chứng minh được sức mạnh như những ngôn ngữ mang tính kịch bản (script) để xử lý dữ liệu.

2.3.1 Ngôn ngữ Python

Python là ngôn ngữ lập trình máy tính bậc cao thường được sử dụng để xây dựng trang web và phần mềm, tự động hóa các tác vụ và tiến hành phân tích dữ liệu. Python là ngôn ngữ có mục đích chung, nghĩa là nó có thể được sử dụng để tạo nhiều chương trình khác nhau và không chuyên biệt cho bất kỳ vấn đề cụ thể nào.

Một vài sự thật thú vị về Python:

- Python được phát triển vào cuối những năm 1980 bởi Guido van Rossum tại Viện Nghiên cứu Quốc gia về Toán học và Khoa học Máy tính ở Hà Lan với tư cách là người kế thừa ngôn ngữ ABC có khả năng xử lý và giao tiếp ngoại lệ.
- Python có nguồn gốc từ các ngôn ngữ lập trình như ABC, Modula 3, small talk, Algol-68.

- Van Rossum đã chọn tên Python cho ngôn ngữ mới từ một chương trình truyền hình, Monty Python's Flying Circus.
- Trang Python là một tệp có phần mở rộng .py chứa có thể là sự kết hợp của Thẻ HTML và tập lệnh Python.
- Vào tháng 12 năm 1989, người sáng tạo đã phát triển trình thông dịch python đầu tiên như một sở thích, và sau đó vào ngày 16 tháng 10 năm 2000, Python 2.0 được phát hành với nhiều tính năng mới.
- Vào ngày 3 tháng 12 năm 2008, Python 3.0 được phát hành với nhiều thử nghiệm hơn và bao gồm các tính năng mới.
- Python là một ngôn ngữ kịch bản mã nguồn mở.
- Python là mã nguồn mở, có nghĩa là bất kỳ ai cũng có thể tải xuống miễn phí từ trang chủ và sử dụng nó để phát triển các chương trình. Mã nguồn của nó có thể được truy cập và sửa đổi theo yêu cầu trong dự án.
- Python là một trong những ngôn ngữ chính thức tại Google.

Ứng dụng của Python:

Python thường được sử dụng để phát triển trang web và phần mềm, tự động hóa tác vụ, phân tích dữ liệu và trực quan hóa dữ liệu. Vì tương đối dễ học, Python đã được nhiều người không phải là lập trình viên như kế toán và nhà khoa học áp dụng cho nhiều công việc hàng ngày, chẳng hạn như tổ chức tài chính.

Phân tích dữ liệu và học máy:

Python đã trở thành một yếu tố chính trong khoa học dữ liệu, cho phép các nhà phân tích dữ liệu và các chuyên gia khác sử dụng ngôn ngữ này để thực hiện các phép tính thống kê phức tạp, tạo trực quan hóa dữ liệu, xây dựng thuật toán học máy, thao tác và phân tích dữ liệu cũng như hoàn thành các nhiệm vụ khác liên quan đến dữ liệu.

Python có thể xây dựng nhiều dạng trực quan hóa dữ liệu khác nhau, chẳng hạn như biểu đồ đường và thanh, biểu đồ hình tròn, biểu đồ 3D. Python cũng có một số thư viện cho phép các lập trình viên viết chương trình để phân tích dữ liệu và học máy nhanh hơn và hiệu quả hơn, như TensorFlow và Keras.

Phát triển web:

Python thường được sử dụng để phát triển back-end của trang web hoặc ứng dụng—những phần mà người dùng không nhìn thấy. Vai trò của Python trong phát triển web có thể bao gồm gửi dữ liệu đến và đi từ máy chủ, xử lý dữ liệu và giao tiếp với cơ sở dữ liệu, định tuyến URL và đảm bảo tính bảo mật. Python cung cấp một số khuôn khổ để phát triển web. Những cái thường được sử dụng bao gồm Django và Flask.

Một số công việc phát triển web sử dụng Python bao gồm kỹ sư phụ trợ, nhà phát triển Python, kỹ sư phần mềm và kỹ sư DevOps.