

ML Concepts

Home

Crash Course

≡

Filter

Linear regression (70 min)

Logistic regression (35 min)

Classification (70 min)

Data

Working with numerical data (85 min)

Working with categorical data (50 min)

Datasets, generalization, and overfitting (105 min)

Introduction (5 min)

Data characteristics (10 min)

Labels (10 min)

Imbalanced datasets (10 min)

Building the criminal dataset

Home

>

Products

>

Machine Learning

>

ML Concepts

>

Crash Course

Was this helpful?

👍

👎

Datasets, generalization, and overfitting

🔖

📄

[Send feedback](#)

On this page

Introduction

🕒

Estimated module length: 105 minutes

🎓

Learning objectives

- Identify four different characteristics of data and datasets.
- Identify at least four different causes of data unreliability.
- Determine when to discard missing data and when to impute it.
- Differentiate between direct and derived labels.
- Identify two different ways to improve the quality of human-rated labels.
- Explain why to subdivide a dataset into a training set, validation set, and test set; identify a potential problem in data splits.
- Explain overfitting and identify three possible causes for it.
- Explain the concept of regularization. In particular, explain the following:
 - Bias versus variance (adaptation to outliers...)
 - L_2 regularization, including Lambda (regularization rate)
 - Early stopping
- Interpret different kinds of loss curves; detect convergence and overfitting in loss curves.

✓

Prerequisites:

This module assumes you are familiar with the concepts covered in the following modules:

- [Introduction to Machine Learning](#)
- [Linear regression](#)
- [Working with numerical data](#)
- [Working with categorical data](#)

Introduction

This module begins with a leading question. Choose one of the following answers:

If you had to prioritize improving one of the following areas in your machine learning project, which would have the most impact?

Improving the quality of your dataset

✓

Data trumps all. The quality and size of the dataset matters much more than which shiny algorithm you use to build your model.

Correct answer.

Applying a more clever loss function to training your model

❑

And here's an even more leading question:

Take a guess: In your machine learning project, how much time do you typically spend on data preparation and transformation?

Less than half of the project time

❑

More than half of the project time

✓

Yes, ML practitioners spend the majority of their time constructing datasets and doing feature engineering.

Correct answer.

In this module, you'll learn more about the characteristics of machine learning datasets, and how to prepare your data to ensure high-quality results when training and evaluating your model.

Help Center

←

Previous

Test your knowledge (10 min)

Next

Data characteristics (10 min)

→

Was this helpful?

👍

👎

[Send feedback](#)

Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 4.0 License](#), and code samples are licensed under the [Apache 2.0 License](#). For details, see the [Google Developers Site Policies](#). Java is a registered trademark of Oracle and/or its affiliates.

Last updated 2024-10-09 UTC.

Connect

Blog

Instagram

LinkedIn

X (Twitter)

YouTube

Programs

Google Developer Groups

Google Developer Experts

Accelerators

Women Techmakers

Google Cloud & NVIDIA

Developer consoles

Google API Console

Google Cloud Platform Console

Google Play Console

Firebase Console

Actions on Google Console

Cast SDK Developer Console

Chrome Web Store Dashboard

Google Home Developer Console

Google for Developers

Android

Chrome

Firebase

Google Cloud Platform

Google AI

All products

Terms | Privacy

Sign up for the Google for Developers newsletter

[Subscribe](#)

🌐

English ▾