# Deep Learning

Ian Goodfellow
Yoshua Bengio
Aaron Courville

# Contents

# Website

www.deeplearningbook.org

This book is accompanied by the above website. The website provides a variety of supplementary material, including exercises, lecture slides, corrections of mistakes, and other resources that should be useful to both readers and instructors.