Charu C. Aggarwal

# Neural Networks and Deep Learning

A Textbook

Springer

Charu C. Aggarwal
IBM T. J. Watson Research Center
International Business Machines
Yorktown Heights, NY, USA

# Contents