ORACLE
University

# Oracle Cloud Infrastructure Generative AI Professional

Student Guide

S1107455GC10

# Table of Contents