Navarurh Kumar   1 ▼

🏠    Course Homepage    Individual Final Report   **SUBMIT TURNITIN ASSIGNMENT**

# Submit Turnitin Assignment

Congratulations - your submission is complete! This is your digital recei
copy of this receipt from within the Document Viewer.

**Author:**
Navarurh Kumar

**Assignment title:**
Individual Final Report

**Submission title:**
Final Term Project Report

**File name:**
NXK180010_BUAN6320_FinalTermProject.pdf

**File size:**
517.12K

**Page count:**
17

**Word count:**
5063

**Character count:**
24613

**Submission date:**
12-Dec-2018 08:25PM (UTC-0600)

**Submission ID:**
1056138766

«    Page 1    »

We take your privacy very seriously. We do not share your details for marketing purposes with any externa
be shared with our third party partners ONLY so that we may offer our service.

Return to assignment list