

Phase 2 BERT Model Training Guideline

ST1506: DSDA FINAL YEAR PROJECT

COMPANY PROFILING TEXT MINING

August 2021

Synopsis

The purpose of this documentation is to demonstrate a detailed breakdown on the content of the ipython file for model training, so that users will understand the rationale behind the content.

This documentation is for data scientists who want to see the usage of BERT in this specific context and possibly improve from the current work.

Table of Contents

Breakdown Of Notebook Contents.....	4
Hardware Rendering.....	4
Loading necessary libraries	4
Exploratory Data Analysis	6
Data Pre-processing	10
Model Training.....	15
Model Evaluation	18

Breakdown Of Notebook Contents

Hardware Rendering

1. When the user chooses to run the code using GPU in Google Colab, the user can run the notebook as per usual so long as the runtime has been set to GPU in Google Colab.
2. However, if the user chooses to run the code using TPU, the user firstly must initialize itself first. An output will be rendered, which will state the number of logical TPU devices the TPU cluster comprises of, which are capable of parallel processing: This process also defines a distribution strategy for distributed training over the number of logical TPU devices.

```
2.1. Check if TPU or GPU is utilized.
We need to check if we are currently using TPU or GPU.

import tensorflow as tf
# if test_gpu_device() is True:
#     print('Using GPU')
# else:
#     print('Using TPU')
#     resolver = tf.distribute.cluster_resolver.TFClusterResolver()
#     config = tf.distribute.cluster_resolver.TFConfigProto()
#     config.device_policy = 'TPU'
#     config.num_replicas_in_sync = 8
#     tpu_strategy = tf.distribute.experimental.TPUStrategy(resolver)
#     strategy = tpu_strategy
# else:
#     strategy = tf.distribute.Strategy(resolver)

# print('Using GPU')
# print('No GPU Available', tpu_strategy.experimental_local_physical_device())
tpu_strategy = tpu_strategy

# Using TPU
tf.distribute.experimental.initialize(tf.config.experimental.get_device_values('/TPU:0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000,1001,1002,1003,1004,1005,1006,1007,1008,1009,1010,1011,1012,1013,1014,1015,1016,1017,1018,1019,1020,1021,1022,1023,1024,1025,1026,1027,1028,1029,1030,1031,1032,1033,1034,1035,1036,1037,1038,1039,1040,1041,1042,1043,1044,1045,1046,1047,1048,1049,1050,1051,1052,1053,1054,1055,1056,1057,1058,1059,1060,1061,1062,1063,1064,1065,1066,1067,1068,1069,1070,1071,1072,1073,1074,1075,1076,1077,1078,1079,1080,1081,1082,1083,1084,1085,1086,1087,1088,1089,1090,1091,1092,1093,1094,1095,1096,1097,1098,1099,1100,1101,1102,1103,1104,1105,1106,1107,1108,1109,1110,1111,1112,1113,1114,1115,1116,1117,1118,1119,1120,1121,1122,1123,1124,1125,1126,1127,1128,1129,1130,1131,1132,1133,1134,1135,1136,1137,1138,1139,1140,1141,1142,1143,1144,1145,1146,1147,1148,1149,1150,1151,1152,1153,1154,1155,1156,1157,1158,1159,1160,1161,1162,1163,1164,1165,1166,1167,1168,1169,1170,1171,1172,1173,1174,1175,1176,1177,1178,1179,1180,1181,1182,1183,1184,1185,1186,1187,1188,1189,1190,1191,1192,1193,1194,1195,1196,1197,1198,1199,1200,1201,1202,1203,1204,1205,1206,1207,1208,1209,1210,1211,1212,1213,1214,1215,1216,1217,1218,1219,1220,1221,1222,1223,1224,1225,1226,1227,1228,1229,1230,1231,1232,1233,1234,1235,1236,1237,1238,1239,1240,1241,1242,1243,1244,1245,1246,1247,1248,1249,1250,1251,1252,1253,1254,1255,1256,1257,1258,1259,1260,1261,1262,1263,1264,1265,1266,1267,1268,1269,1270,1271,1272,1273,1274,1275,1276,1277,1278,1279,1280,1281,1282,1283,1284,1285,1286,1287,1288,1289,1290,1291,1292,1293,1294,1295,1296,1297,1298,1299,1300,1301,1302,1303,1304,1305,1306,1307,1308,1309,1310,1311,1312,1313,1314,1315,1316,1317,1318,1319,1320,1321,1322,1323,1324,1325,1326,1327,1328,1329,1330,1331,1332,1333,1334,1335,1336,1337,1338,1339,1340,1341,1342,1343,1344,1345,1346,1347,1348,1349,1350,1351,1352,1353,1354,1355,1356,1357,1358,1359,1360,1361,1362,1363,1364,1365,1366,1367,1368,1369,1370,1371,1372,1373,1374,1375,1376,1377,1378,1379,1380,1381,1382,1383,1384,1385,1386,1387,1388,1389,1390,1391,1392,1393,1394,1395,1396,1397,1398,1399,1400,1401,1402,1403,1404,1405,1406,1407,1408,1409,1410,1411,1412,1413,1414,1415,1416,1417,1418,1419,1420,1421,1422,1423,1424,1425,1426,1427,1428,1429,1430,1431,1432,1433,1434,1435,1436,1437,1438,1439,1440,1441,1442,1443,1444,1445,1446,1447,1448,1449,1450,1451,1452,1453,1454,1455,1456,1457,1458,1459,1460,1461,1462,1463,1464,1465,1466,1467,1468,1469,1470,1471,1472,1473,1474,1475,1476,1477,1478,1479,1480,1481,1482,1483,1484,1485,1486,1487,1488,1489,1490,1491,1492,1493,1494,1495,1496,1497,1498,1499,1500,1501,1502,1503,1504,1505,1506,1507,1508,1509,1510,1511,1512,1513,1514,1515,1516,1517,1518,1519,1520,1521,1522,1523,1524,1525,1526,1527,1528,1529,1530,1531,1532,1533,1534,1535,1536,1537,1538,1539,1540,1541,1542,1543,1544,1545,1546,1547,1548,1549,1550,1551,1552,1553,1554,1555,1556,1557,1558,1559,1560,1561,1562,1563,1564,1565,1566,1567,1568,1569,1570,1571,1572,1573,1574,1575,1576,1577,1578,1579,1580,1581,1582,1583,1584,1585,1586,1587,1588,1589,1590,1591,1592,1593,1594,1595,1596,1597,1598,1599,1600,1601,1602,1603,1604,1605,1606,1607,1608,1609,1610,1611,1612,1613,1614,1615,1616,1617,1618,1619,1620,1621,1622,1623,1624,1625,1626,1627,1628,1629,1630,1631,1632,1633,1634,1635,1636,1637,1638,1639,1640,1641,1642,1643,1644,1645,1646,1647,1648,1649,1650,1651,1652,1653,1654,1655,1656,1657,1658,1659,1660,1661,1662,1663,1664,1665,1666,1667,1668,1669,1670,1671,1672,1673,1674,1675,1676,1677,1678,1679,1680,1681,1682,1683,1684,1685,1686,1687,1688,1689,1690,1691,1692,1693,1694,1695,1696,1697,1698,1699,1700,1701,1702,1703,1704,1705,1706,1707,1708,1709,1710,1711,1712,1713,1714,1715,1716,1717,1718,1719,1720,1721,1722,1723,1724,1725,1726,1727,1728,1729,1730,1731,1732,1733,1734,1735,1736,1737,1738,1739,1740,1741,1742,1743,1744,1745,1746,1747,1748,1749,1750,1751,1752,1753,1754,1755,1756,1757,1758,1759,1760,1761,1762,1763,1764,1765,1766,1767,1768,1769,1770,1771,1772,1773,1774,1775,1776,1777,1778,1779,1780,1781,1782,1783,1784,1785,1786,1787,1788,1789,1790,1791,1792,1793,1794,1795,1796,1797,1798,1799,1800,1801,1802,1803,1804,1805,1806,1807,1808,1809,1810,1811,1812,1813,1814,1815,1816,1817,1818,1819,1820,1821,1822,1823,1824,1825,1826,1827,1828,1829,1830,1831,1832,1833,1834,1835,1836,1837,1838,1839,1840,1841,1842,1843,1844,1845,1846,1847,1848,1849,1850,1851,1852,1853,1854,1855,1856,1857,1858,1859,1860,1861,1862,1863,1864,1865,1866,1867,1868,1869,1870,1871,1872,1873,1874,1875,1876,1877,1878,1879,1880,1881,1882,1883,1884,1885,1886,1887,1888,1889,1890,1891,1892,1893,1894,1895,1896,1897,1898,1899,1900,1901,1902,1903,1904,1905,1906,1907,1908,1909,1910,1911,1912,1913,1914,1915,1916,1917,1918,1919,1920,1921,1922,1923,1924,1925,1926,1927,1928,1929,1930,1931,1932,1933,1934,1935,1936,1937,1938,1939,1940,1941,1942,1943,1944,1945,1946,1947,1948,1949,1950,1951,1952,1953,1954,1955,1956,1957,1958,1959,1960,1961,1962,1963,1964,1965,1966,1967,1968,1969,1970,1971,1972,1973,1974,1975,1976,1977,1978,1979,1980,1981,1982,1983,1984,1985,1986,1987,1988,1989,1990,1991,1992,1993,1994,1995,1996,1997,1998,1999,2000,2001,2002,2003,2004,2005,2006,2007,2008,2009,2010,2011,2012,2013,2014,2015,2016,2017,2018,2019,2020,2021,2022,2023,2024,2025,2026,2027,2028,2029,2030,2031,2032,2033,2034,2035,2036,2037,2038,2039,2040,2041,2042,2043,2044,2045,2046,2047,2048,2049,2050,2051,2052,2053,2054,2055,2056,2057,2058,2059,2060,2061,2062,2063,2064,2065,2066,2067,2068,2069,2070,2071,2072,2073,2074,2075,2076,2077,2078,2079,2080,2081,2082,2083,2084,2085,2086,2087,2088,2089,2090,2091,2092,2093,2094,2095,2096,2097,2098,2099,2100,2101,2102,2103,2104,2105,2106,2107,2108,2109,2110,2111,2112,2113,2114,2115,2116,2117,2118,2119,2120,2121,2122,2123,2124,2125,2126,2127,2128,2129,2130,2131,2132,2133,2134,2135,2136,2137,2138,2139,2140,2141,2142,2143,2144,2145,2146,2147,2148,2149,2150,2151,2152,2153,2154,2155,2156,2157,2158,2159,2160,2161,2162,2163,2164,2165,2166,2167,2168,2169,2170,2171,2172,2173,2174,2175,2176,2177,2178,2179,2180,2181,2182,2183,2184,2185,2186,2187,2188,2189,2190,2191,2192,2193,2194,2195,2196,2197,2198,2199,2200,2201,2202,2203,2204,2205,2206,2207,2208,2209,2210,2211,2212,2213,2214,2215,2216,2217,2218,2219,2220,2221,2222,2223,2224,2225,2226,2227,2228,2229,2230,2231,2232,2233,2234,2235,2236,2237,2238,2239,2240,2241,2242,2243,2244,2245,2246,2247,2248,2249,2250,2251,2252,2253,2254,2255,2256,2257,2258,2259,2260,2261,2262,2263,2264,2265,2266,2267,2268,2269,2270,2271,2272,2273,2274,2275,2276,2277,2278,2279,2280,2281,2282,2283,2284,2285,2286,2287,2288,2289,2290,2291,2292,2293,2294,2295,2296,2297,2298,2299,2300,2301,2302,2303,2304,2305,2306,2307,2308,2309,2310,2311,2312,2313,2314,2315,2316,2317,2318,2319,2320,2321,2322,2323,2324,2325,2326,2327,2328,2329,2330,2331,2332,2333,2334,2335,2336,2337,2338,2339,2340,2341,2342,2343,2344,2345,2346,2347,2348,2349,2350,2351,2352,2353,2354,2355,2356,2357,2358,2359,2360,2361,2362,2363,2364,2365,2366,2367,2368,2369,2370,2371,2372,2373,2374,2375,2376,2377,2378,2379,2380,2381,2382,2383,2384,2385,2386,2387,2388,2389,2390,2391,2392,2393,2394,2395,2396,2397,2398,2399,2400,2401,2402,2403,2404,2405,2406,2407,2408,2409,2410,2411,2412,2413,2414,2415,2416,2417,2418,2419,2420,2421,2422,2423,2424,2425,2426,2427,2428,2429,2430,2431,2432,2433,2434,2435,2436,2437,2438,2439,2440,2441,2442,2443,2444,2445,2446,2447,2448,2449,2450,2451,2452,2453,2454,2455,2456,2457,2458,2459,2460,2461,2462,2463,2464,2465,2466,2467,2468,2469,2470,2471,2472,2473,2474,2475,2476,2477,2478,2479,2480,2481,2482,2483,2484,2485,2486,2487,2488,2489,2490,2491,2492,2493,2494,2495,2496,2497,2498,2499,2500,2501,2502,2503,2504,2505,2506,2507,2508,2509,2510,2511,2512,2513,2514,2515,2516,2517,2518,2519,2520,2521,2522,2523,2524,2525,2526,2527,2528,2529,2530,2531,2532,2533,2534,2535,2536,2537,2538,2539,2540,2541,2542,2543,2544,2545,2546,2547,2548,2549,2550,2551,2552,2553,2554,2555,2556,2557,2558,2559,2560,2561,2562,2563,2564,2565,2566,2567,2568,2569,2570,2571,2572,2573,2574,2575,2576,2577,2578,2579,2580,2581,2582,2583,2584,2585,2586,2587,2588,2589,2590,2591,2592,2593,2594,2595,2596,2597,2598,2599,2600,2601,2602,
```

2. In Section 2.4., we will load the training, and the validation dataset as Pandas Dataframes. The training dataset will be used as the training dataset for the model, while the validation dataset will serve as a test dataset.

[illegible]

Figure 3: Loading the excel files and preview of the training dataset

Exploratory Data Analysis

1. First, we will identify the various columns in the training dataset to get a better understanding of the dataset. As shown below are the columns identified in the Dataframe.

```
[ ] 1 # see the row headers of the entire pandas dataframe first
    2 list(companies.columns)

['Company_ID',
 'Company',
 'Country',
 'PIC',
 'Sector',
 'Subsector',
 'Archetype',
 'Valuechain',
 'Websites',
 'Company Profile Information',
 'Remarks']
```

Figure 4: List of columns in the training dataset

2. Here, we identify the number of records, as well as the number of unique labels.

```
# get the total number of records in the dataframe
df_count = df_train['Company_ID'].count()

# get count of unique contries where companies are based in
df_countCountry = df_train['Country'].nunique()

# get count of total unique sectors where companies are from
df_countSector = df_train['Sector'].nunique()

# get count of total unique subsector where companies are from
df_countsubSector = df_train['Subsector'].nunique()

# get count of total unique valuechain where companies are from
df_countValuechain = df_train['Valuechain'].nunique()

# get count of total unique archetypes
df_countArchetype = df_train['Archetype'].nunique()

print('Total number of records:', df_count)
print('Total number of countries:', df_countCountry)
print('Total number of sectors:', df_countSector)
print('Total number of subsectors:', df_countsubSector)
print('Total number of valuechain:', df_countValuechain)
print('Total number of archetypes:', df_countArchetype)

Total number of records: 9600
Total number of countries: 14
Total number of sectors: 16
Total number of subsectors: 37
Total number of valuechain: 18
Total number of archetypes: 94
```

Figure 5: Number of records, and labels in interest

- We proceed to identify the unique Archetype labels present in the dataset.

```
print('List of unique archetype:\n{}'.format(df_archetype))

List of unique archetype:
others                2161
building_material_manufacturer    573
buildings & industrial_contractor    496
consumer discretionary distributor    358
cni_service providers    279
...
MIDSTREAM                2
tisp - tower                2
tisp - fiber cable        2
building_material_manufacturer    1
metals and mining        1
Name: Archetype, Length: 94, dtype: int64
```

Figure 6: Unique archetype breakdown

- Here, we drop the unnecessary columns from the dataset, so that they do not interfere with the EDA process.

```
# declare the list of the row names that are redundant
rows_to_drop = ['Company_ID', 'PIC', 'Websites', 'Remarks']

# use a conditional expression to filter out those rows
df_filteredCompanies = df_train.drop(labels=rows_to_drop, axis=1)

df_filteredCompanies
```

	Company	Country	Sector	Subsector	Archetype	Valuechain	Company Profile Information
0	AKSORN CHAROEN TAT ACT. CO.,LTD.	THAILAND	TMT	media	media_aggregator/distributor	Midstream	For over 80 years of experience in creating an...
1	DONGGUAN SHENG YIA CLEANING APPLIAME CO.LTD	CHINA	TMT	consumer electronics	consumer electronics_distributor	Downstream	Yatai's main products cover various cleaning m...
2	EXIS TECH SDN. BHD.	MALAYSIA	oos	others	others	NaN	In the beginning, it started off by providing ...
3	BEI JING ESTRABA IMPORT AND EXPORT	CHINA	NaN	NaN	NaN	NaN	NaN
4	Aztech Electronics Pte Ltd	SINGAPORE	TMT	consumer electronics	consumer electronics_distributor	Downstream	Being a turnkey, one-stop integrated solutions...
...
9595	SOONBEE INVESTMENT HOLDINGS PTE. LT D	SINGAPORE	OOS	others	others	NaN	SOONBEE INVESTMENT HOLDINGS is an ACRA-registe...
9596	SICHUAN CHINA RAILWAY WENRUI REAL E	CHINA	NaN	NaN	NaN	NaN	NaN
9597	GROUP INDUSTRIES PTE LTD	SINGAPORE	CNI	building_material	building_material_manufacturer	MIDSTREAM	Group Industries Pte Ltd was founded in 1986. ...
9598	SIAM NISSAN EASTERN CO.,LTD.	THAILAND	IND	auto & mec	auto dealer	Trader	Auto Dealership - Retail
9599	FAT INVESTMENT HONG KONG LIMITED	HONG KONG	NaN	NaN	NaN	NaN	NaN

Figure 7: New training dataset after dropping columns

- Records with all NaN values are extracted, shown, and discarded, as they cannot be used in the training process.

```
# find all the rows with nan data in sector, subsector, archetype and valuechain
cols_to_check = ['Sector', 'Subsector', 'Archetype', 'Valuechain', 'Company Profile Information']
empty = df_filteredCompanies[df_filteredCompanies[cols_to_check].isnull().all(1)]
empty
```

	Company	Country	Sector	Subsector	Archetype	Valuechain	Company Profile Information
3	BEI JING ESTRABA IMPORT AND EXPORT	CHINA	NaN	NaN	NaN	NaN	NaN
12	CHANGTU COUNTRY LONGXING FERTILIZER CO.,LTD	CHINA	NaN	NaN	NaN	NaN	NaN
28	ZIBOBOSHANHONGLIWEI MOTOR CO.,LTD	CHINA	NaN	NaN	NaN	NaN	NaN
48	BEIJING DUO MEIDUO SHIYOU PRODUCTS SALES CO.,...	CHINA	NaN	NaN	NaN	NaN	NaN
63	TRUSVEST SDN. BHD.	MALAYSIA	NaN	NaN	NaN	NaN	NaN
...
9576	REAL CHARM INVESTMENT LIMITED	HONG KONG	NaN	NaN	NaN	NaN	NaN
9585	JUJIN INTERNATIONAL TRADE (SHANGHAI)	CHINA	NaN	NaN	NaN	NaN	NaN
9592	QITAI TIANSHAN CEMENT CO., LTD	CHINA	NaN	NaN	NaN	NaN	NaN
9596	SICHUAN CHINA RAILWAY WENRUI REAL E	CHINA	NaN	NaN	NaN	NaN	NaN
9599	FAT INVESTMENT HONG KONG LIMITED	HONG KONG	NaN	NaN	NaN	NaN	NaN

1149 rows x 7 columns

Figure 8: Records with NaN values

6. The records with valid data are then kept.

```
# now we get the dataset that are valid
df_valid = pd.concat([df_filteredCompanies, empty, empty]).drop_duplicates(keep=False)
df_valid
```

	Company	Country	Sector	Subsector	Archetype	Valuechain	Company Profile Information
0	AKSORN CHAROEN TAT ACT. CO.,LTD.	THAILAND	TMT	media	media_aggregator/distributor	Midstream	For over 80 years of experience in creating an...
1	DONGGUAN SHENGYA CLEANING APPLIAME CO LTD	CHINA	TMT	consumer electronics	consumer electronics_distributor	Downstream	Yatai's main products cover various cleaning m...
2	EXIS TECH SDN. BHD.	MALAYSIA	oos	others	others	NaN	In the beginning, it started off by providing ...
4	Aztech Electronics Pte Ltd	SINGAPORE	TMT	consumer electronics	consumer electronics_distributor	Downstream	Being a turnkey, one-stop integrated solutions...
5	TONGDUN INTERNATIONAL PTE LTD	SINGAPORE	tmt	it_services	it_services	midstream	Tongdun Technology is a professional third-par...
...
9593	BANGKOK PATANA SCHOOL	THAILAND	CG	retail n distribution	consumer discretionary distributor	Downstream	Bangkok Patana School is Thailand's original B...
9594	INDRATI LINES PTE LTD	SINGAPORE	OOS	others	others	NaN	Indra is one of the leading global technology ...
9595	SOONBEE INVESTMENT HOLDINGS PTE. LT D	SINGAPORE	OOS	others	others	NaN	SOONBEE INVESTMENT HOLDINGS is an ACRA-registe...
9597	GROUP INDUSTRIES PTE LTD	SINGAPORE	CNI	building_material	building_material_manufacturer	MIDSTREAM	Group Industries Pte Ltd was founded in 1986. ...
9598	SIAM NISSAN EASTERN CO.,LTD.	THAILAND	IND	auto & mec	auto dealer	Trader	Auto Dealership - Retail

8447 rows x 7 columns

Figure 9: Dataframe with the valid results

7. Here, we will turn the numbers seen in the EDA into graphs, so that it is easier to spot the trend. Shown below is a bar chart of the breakdown of labels in Sector. It seems that OOS, CNI and REH make up the majority of the dataset.

We also plotted a pie chart of the countries found in the dataset. We can see that most of the portions comes from Thai companies, followed by China and Malaysia.

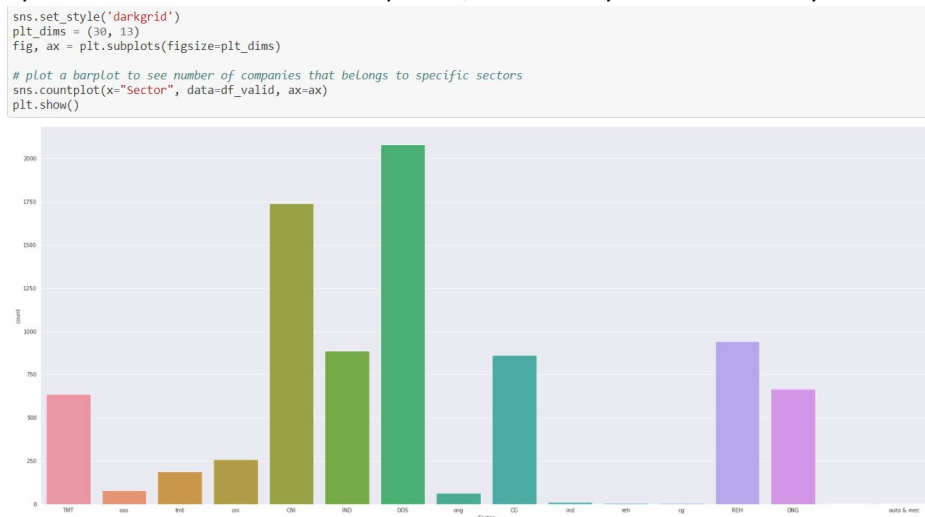


Figure 10: Bar Chart of the various Sector counts

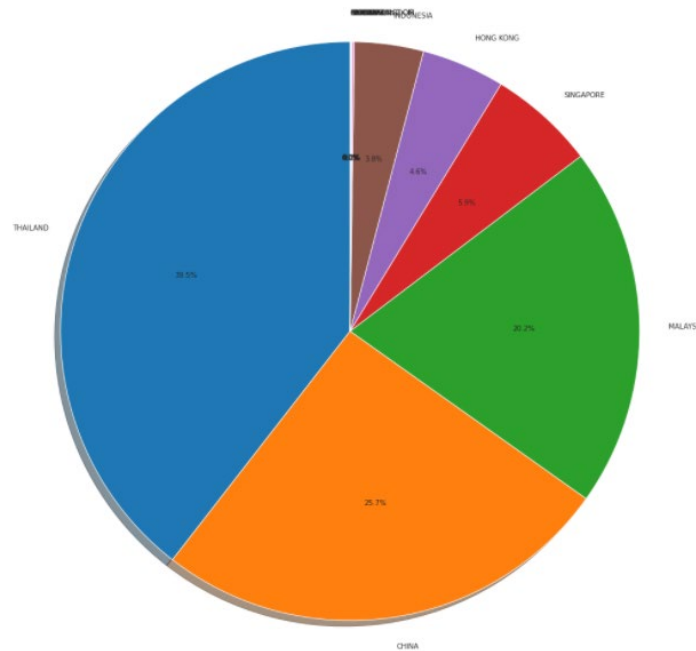


Figure 11: Pie chart breakdown of the countries the companies originate from

8. We would like to see examples of how the worded company profile description looks like. Therefore, we printed the 1st 10 companies text description for analysis.

```
# configure pandas dataframe to let us see the entire company description IN FULL
pd.set_option('display.max_colwidth', None)

# get the 1st 10 results and observe
df_valid.loc[0:10, 'Company Profile Information']
```

0
For over 80 years of experience in creating and developing high-quality learning materials has enabled us to provide world-class educational innovation to meet the needs of all teachers, students, institutions and educational authorities.

1
Yatai's main products cover various cleaning machinery, cleaning agents, cleaning tools, stone maintenance and other cleaning products; cleaning solutions and services include product technical consultation, product customization, employee training, maintenance and so on.

2
In the beginning, it started off by providing technical support for test handlers, then moving into module design and production. Its first, in-house designed, full-fledged handler was introduced in 2008. Since then, the company has designed and produced a wide range of turnkey and pick-and-place solutions for its customers all over the world.

4
Being a turnkey, one-stop integrated solutions provider based in Singapore, Aztech is equipped with state-of-the-art equipment, R&D, design, manufacturing and packaging capabilities to deliver a seamless, unified experience. Each and every time. Always striving towards the edge of technology for more than 34 years, we have been building capabilities to serve clients' manufacturing needs, including the consumer electronics, telecommunications, healthcare, LED lighting, automotive and technology start-up market segments.

5
Tongdun Technology is a professional third-party intelligent risk management and decision-making service provider headquartered in Hangzhou, Zhejiang. By integrating artificial intelligence into business scenarios, Tongdun Technology offers solutions in intelligent user analysis, intelligent risk management, intelligent antifraud and intelligent operation to clients from various industries including financial industry, internet business, logistics, healthcare, retail, smart cities and government bodies. Over 10,000 corporate clients have chosen Tongdun Technology's products and services.

6
For over 30 years, EPMSTIGA has been servicing the construction and the oil and gas industries with distinction. We try our hardest to create win-win situations and value-add to our clients with every transaction.

7
Sales of hydraulic equipment and parts

8
For Excelkos, it all started with a far-sighted vision to deliver quality industrial products and solutions to meet the increasing demands of various industries. The vision became a reality in 1988 when the Excelkos Group was established to distribute and provide quality rubber industrial chemicals.

9
Kum Eng Huat is the authorised dealer for Osram and Philips, as well as a trusted supplier of lighting solutions for major projects with the local government.

10
Wholesale of other building materials, not elsewhere classified
Name: Company Profile Information, dtype: object

Figure 12: First 10 company text descriptions

Data Pre-processing

1. To ensure that there are no escape characters such as (`\n`) and (`\t`), we will first replace those with a space (' ') to ensure that it does not interfere with the model learning process later on.

We also used a distribution plot to see the most frequent word length among the entire dataset. From the frequency distribution plot, we can assume that most company descriptions have around 700 to 800 word descriptions.

4.1. Removing `\n` and `\t`

Now, we will like to standardize all the paragraphs such that they are homogenous, before we tokenize the paragraph

```
# get rid of the \n found in the respective descriptions
df_valid = df_valid.replace(["\n", "\t"], ' ', regex=True)

# now plot a distribution plot to see the word length distribution
sns.distplot(df_valid["length"], kde=True)
plt.title('Word Length Distribution')
plt.show()
```

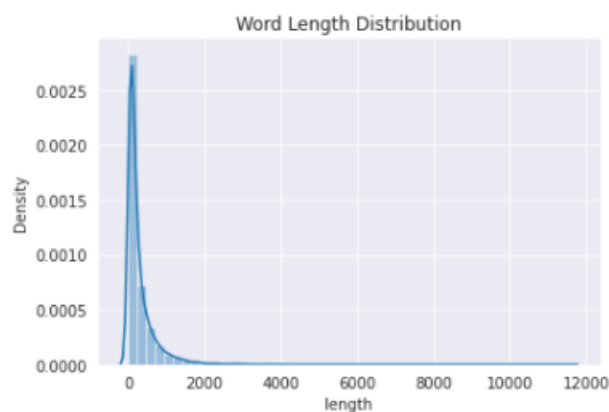


Figure 13: Removing special characters and getting the word length distribution

2. We then proceed on to populate NaN cells with spaces to proceed on with text tokenization later on.

```
# fill na with space instead of others
df_valid.fillna(" ", inplace=True)
df_valid
```

	Company	Country	Sector	Subsector	Archetype	Valuechain	Company Profile Information	length
0	AKSORN CHAROEN TAT ACT. CO.,LTD.	THAILAND	TMT	media	media_aggregator/distributor	Midstream	For over 80 years of experience in creating and developing high-quality learning materials has enabled us to provide world-class educational innovation to meet the needs of all teachers, students, institutions and educational authorities.	238
1	DONGGUAN SHENGYA CLEANING APPLIAME CO.LTD	CHINA	TMT	consumer electronics	consumer electronics_distributor	Downstream	Yatai's main products cover various cleaning machinery, cleaning agents, cleaning tools, stone maintenance and other cleaning products; cleaning solutions and services include product technical consultation, product customization, employee training, maintenance and so on.	273
2	EXIS TECH SDN. BHD.	MALAYSIA	oos	others	others		In the beginning, it started off by providing technical support for test handlers, then moving into module design and production. Its first, in-house designed, full-fledged handler was introduced in 2008. Since then, the company has designed and produced a wide range of turret and pick-and-place solutions for its customers all over the world.	344

Figure 14: A snippet off populating NaN cells to space

3. Section 4.5. of the notebook retrieves the 4 labels, namely sector, subsector, archetype and valuechain, and assign unique labels to the respective company description. These unique labels will be used later as the training labels when training the model later.

4.5. Assigning tags

In this section, we will be assigning tags to every row, so that we can make use of the given keywords as the training dataset labels.

```
# Programmatically assign tags to each definition
sector_keywords = pd.read_excel('./sector_master_definition.xlsx')
df_keywords = sector_keywords[['Sector', 'Subsector', 'Archetype', 'Value Chain', 'Sector Keywords']]

# capitalise all tags
df_keywords['Value Chain'] = df_keywords['Value Chain'].str.upper()
df_keywords.fillna(' ', inplace=True)
df_keywords['Sector Keywords'] = df_keywords['Sector Keywords'].str.upper()
df_keywords['Sector Keywords'].replace(' ', '[]', inplace=True)

# save unique tags, sorted for consistency across runs
sector = np.sort(df_keywords['Sector']).unique()
subsector = np.sort(df_keywords['Subsector']).unique()
archetype = np.sort(df_keywords['Archetype']).unique()
valuechain = np.sort(df_keywords['Value Chain']).unique()
print(len(sector), len(subsector), len(archetype), len(valuechain))
tag_counts = [len(sector), len(subsector), len(archetype), len(valuechain)]

# assign number tag list to each row
taglist = []
for index, row in df_keywords.iterrows():
    temp = []

    temp.append(np.where(sector == row['Sector'])[0][0])
    temp.append(np.where(subsector == row['Subsector'])[0][0])
    temp.append(np.where(archetype == row['Archetype'])[0][0])
    temp.append(np.where(valuechain == row['Value Chain'])[0][0])

    taglist.append(temp)

# assign completed taglist to column in dataframe
df_keywords['list_tag'] = taglist

7 32 92 9
```

Figure 15: Assigning tags to each company description

4. We redid data pre-processing for the training dataset for homogeneity.

```
# process data for homogeneity
df_valid['Valuechain'] = df_valid['Valuechain'].str.split().str.join(' ')
df_valid['Valuechain'] = df_valid['Valuechain'].str.upper()
df_valid['Sector'] = df_valid['Sector'].str.upper()
df_valid['Valuechain'].replace(' ', ' ', inplace=True)

taglist_df = []
# process tags for records
for index, row in df_valid.iterrows():
    temp = []

    try: # for error handling
        temp.append(np.where(sector == row['Sector'])[0][0])
        temp.append(np.where(subsector == row['Subsector'])[0][0])
        temp.append(np.where(archetype == row['Archetype'])[0][0])
        temp.append(np.where(valuechain == row['Valuechain'])[0][0])

        taglist_df.append(temp)
    except:
        # drop data if not valid
        df_valid.drop(index, inplace=True)

df_valid['list_tag'] = taglist_df

df_valid.shape

(8423, 9)
```

Figure 16: Pre-Processing df_valid

- Similar to point 4, we also did a similar data pre-processing to the test dataset under Section 4.7. in the notebook.

4.7. Preprocessing data for validation testset.

We will perform data preprocessing for the validation testset.

```
df_test.replace('NaN', np.NaN, inplace=True)

# drop unnecessary columns
df_test.drop(rows_to_drop, axis=1, inplace=True)

# replace newline characters in validation data
df_test = df_test.replace('\n', ' ', regex=True)

# fill in NaN values in validation data
df_test.fillna(' ', inplace=True)

# change dtype of validation data columns
for i in columns_to_convert:
    df_valid[i] = df_valid[i].astype(str)

# add tags to validation data
taglist_df = []
# process tags for records
for index, row in df_test.iterrows():
    temp = []

    try: # for error handling
        temp.append(np.where(sector == row['Sector'])[0][0])
        temp.append(np.where(subsector == row['Subsector'])[0][0])
        temp.append(np.where(archetype == row['Archetype'])[0][0])
        temp.append(np.where(valuechain == row['Valuechain'])[0][0])

    taglist_df.append(temp)
except Exception as e:
    # drop data if not valid
    print(row.name, e, '\n')
    df_test.drop(index, inplace=True)

df_test['list_tag'] = taglist_df
```

Figure 17: Data pre-processing for the test dataset

- We went on to ensure that the dtypes for the respective columns are in string and not other dtypes.

We made use of HuggingFace tokenizer function. We first load BertConfig, which is the configuration of the Bert model, to bert-base-uncased to instantiate a BERT model according to the specified arguments. We then modify one of the parameter, changing the output hidden states to be False. We then proceed on to load the BERT tokenizer Fast, which essentially tokenizer faster than its BERT tokenizer due to the code being implemented as a “fast” BERT tokenizer.

```
# we will have to ensure all the dtype of the respective columns are in string and not float for spacy to handle properly, so now we will
attempt to convert all into strings
columns_to_convert = ['Sector', 'Subsector', 'Archetype', 'Valuechain', 'Company Profile Information']

for i in columns_to_convert:
    df_valid[i] = df_valid[i].astype(str)

# instead of using spacy tokenizer, we will make use of the preprocessing and tokenizer in the python package from huggingface
from transformers import TFBertModel, BertConfig, BertTokenizerFast

# use the uncased bert model
config = BertConfig.from_pretrained('bert-base-uncased')

# set the output_hidden_states to false since we have 4 different classes with their respective labels AND we want to define our own hidden
layers
config.output_hidden_states = False

# load the BERT tokenizer first. BertTokenizerFast tokenize faster than BertTokenizer as it uses Rust
tokenizer = BertTokenizerFast.from_pretrained(pretrained_model_name_or_path='bert-base-uncased', config=config)
```

Figure 18: Loading BERT tokenizer from HuggingFace

7. We will proceed to randomise our dataset for more randomness.

```
# randomise dataset first here
df_train_rand = df_valid.sample(frac=1).reset_index(drop=True)
df_test_rand = df_test.sample(frac=1).reset_index(drop=True)

# drop the length column section from the train_rand df since its not necessary
df_train_rand.drop(['length'], axis=1, inplace=True)

# define one hot encode function
def one_hot(arr, n_cat):
    output = []
    for n in arr:
        result = np.zeros(n_cat)
        result[n] = 1

        output.append(result)

    return np.array(output, dtype=int)
```

Figure 19: Randomise dataset and further dropping unnecessary columns

8. We will now split the training dataset (df_train_rand) and test dataset (df_test_rand) into X_train and y_train, X_test and y_test respectively. The y labels are the label for the respective company descriptions.

```
# split into data and labels
X_train = tokenizer(text=df_train_rand['Company Profile Information'].to_list(),
                    add_special_tokens=True, # add special tokens like [SEP] and others
                    max_length=100, # this is the max length of the sentence-to-be-token
                    truncation=True,
                    padding=True,
                    return_tensors='tf', # to return it as tf tensors to feed into keras API
                    return_token_type_ids=False,
                    return_attention_mask=True, # generate attention mask
                    verbose=True)

X_test = tokenizer(text=df_test_rand['Company Profile Information'].to_list(),
                  add_special_tokens=True,
                  max_length=100, # we use the max length that BERT can handle
                  truncation=True,
                  padding=True,
                  return_tensors='tf',
                  return_token_type_ids=False,
                  return_attention_mask=True,
                  verbose=True)

y_train = np.array(list(df_train_rand['list_tag']))
y_test = np.array(list(df_test_rand['list_tag']))
```

Figure 20: Splitting train and test dataset into training data and the labels

- The add_special_tokens param is set to True in order to add special BERT tokens like [CLS], [SEP] and [PAD] to the 'tokenized' encodings.
 - The max_length defines the length of the tokenized text.
 - Truncation ensure that the max_length is strictly adhered, whereby longer sentence are truncated to max_length.
 - Padding is added such that, say a length of a description is less than 100 words. After doing tokenizing, it adds 0 at the back such that the entire length is eventually 100.
9. To have a better understanding of what happens after tokenizing, we can have a look at the input_id, how it looks like when the sentence is tokenized, and the attention_mask, which tells the model which tokens should be attended to and which tokens should not.

```

X_train["input_ids"][0]

<tf.Tensor: shape=(100,), dtype=int32, numpy=
array([[ 101, 17712, 2121, 7300, 4021, 3534, 17371, 2078, 1012,
        1038, 14945, 1012, 2001, 2631, 1999, 2384, 1012, 1996,
        2194, 1005, 1055, 2240, 1997, 2449, 2950, 4346, 2968,
        2578, 2006, 1037, 3206, 2030, 7408, 3978, 1012, 102,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0, 0, 0, 0, 0, 0, 0, 0, 0,
         0], dtype=int32)>

# see how the 1st sentence of the training dataset is being encoded
tokenizer.decode(X_train["input_ids"][0])

"[CLS] aker solutions asia pacific sdn. bhd. was founded in 2005. the company's line of business includes providing management services o
n a contract or fee basis. [SEP] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [P
AD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD]
[PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD]
[PAD]"

X_train['attention_mask']

<tf.Tensor: shape=(8423, 100), dtype=int32, numpy=
array([[1, 1, 1, ..., 0, 0, 0],
       [1, 1, 1, ..., 0, 0, 0],
       [1, 1, 1, ..., 1, 1, 1],
       ...,
       [1, 1, 1, ..., 0, 0, 0],
       [1, 1, 1, ..., 0, 0, 0],
       [1, 1, 1, ..., 1, 1, 1]], dtype=int32)>

```

Figure 21: (From top to bottom) The first company description input_id, tokenized form and overall attention mask for X_train

10. Finally, we will one hot encode the labels.

```

label_names = ['sector', 'subsector', 'archetype', 'valuechain']
y_train_multi = {label_names[i] : tf.convert_to_tensor(one_hot(y_train[:, i], tag_counts[i]), dtype=tf.int16) for i in range(4)}
y_test_multi = {label_names[i] : tf.convert_to_tensor(one_hot(y_test[:, i], tag_counts[i]), dtype=tf.int16) for i in range(4)}

```

Figure 22: One-hot encoding

Model Training

1. We will now do a quick check on the Keras version we are using, and then create our multi output branches for namely Sector, Subsector, Archetype and Valuechain.

In our multi_branch output, we will load the first layer from the BERT transformer model as we only need the MainLayer. Then we will define the pooled layers for the respective branches. We set training to False so that we can freeze the pretrained layers. We eventually used a softmax activation function. We will slot in the input_ids and the attention mask. The attention_mask essentially tells the model that we should not focus attention on [PAD] tokens, which is essential for classification problems.

```
import keras

print('--- Version Checking ---')
print("Keras:", keras.__version__)

--- Version Checking ---
Keras: 2.5.0

from tensorflow.keras.regularizers import l2
from tensorflow.keras.layers import Conv1D, Dense, Dropout, MaxPooling1D, GlobalAveragePooling1D, Input, Lambda, Embedding
from tensorflow.keras.layers import ReLU, BatchNormalization, Activation
from tensorflow.keras import Sequential
from tensorflow.keras.models import Model
from tensorflow.keras.initializers import TruncatedNormal

# create multi output model
def multi_branch(inputs, name, output_dim, counter, transformer_model):

    # Load the Transformers BERT model as a layer in a Keras model
    bert = transformer_model.layers[0] # get only the top head
    bert_model = bert(inputs)[1]
    dropout = Dropout(config.hidden_dropout_prob, name="pooled_output_{}".format(counter+1))
    pooled_output = dropout(bert_model, training=False) # freeze pretrained layers

    # build the model output
    output = Dense(output_dim, kernel_initializer=TruncatedNormal(stddev=config.initializer_range), name=name, activation='softmax')(pooled_output)

    return output

# create the multi output branches
def create_model(labels, output_dim):
    # input
    input_ids = Input(shape=(100,), name='input_ids', dtype='int32')
    attention_mask = Input(shape=(100,), name='attention_mask', dtype='int32')
    inputs = {'input_ids': input_ids, 'attention_mask': attention_mask}

    # preload the BERT transformer model
    transformer_model = TFBertModel.from_pretrained('bert-base-uncased', config=config)

    # create the respective branches for the outputs
    branches = []

    for i in range(len(labels)):
        branches.append(multi_branch(inputs, labels[i], output_dim[i], i, transformer_model))

    # fit the entire model into keras Model class
    model = Model(inputs=inputs, outputs=branches, name='company_classification_model')

    return model
```

Figure 23: Creating multi-branches for the categories classification

2. Once we have created the branches, we will create the model, and define categorical_crossentropy as our loss, and Adam as the optimizer. We will be using metrics such as accuracy, precision, recall, area under curve to observe our model performance.

There are 2 sets of codes, one catered for TPU users, and another for GPU users. For those who have earlier ran using GPU as the runtime, the model will directly be created in the else loop. For TPU users, the code will land in the first conditional action, whereby we will create the model all in the strategy_scope, so that we can distribute the workload equally to the TPU devices.

```

from tensorflow.keras.metrics import Precision, Recall, AUC
from tensorflow.keras.optimizers import SGD, Adam
from tensorflow.keras.callbacks import EarlyStopping

# for tpu
if tpu_usage:
    with strategy.scope():
        model = create_model(label_names, tag_counts)
        opt = Adam(learning_rate=5e-05, epsilon=1e-08, decay=0.01, clipnorm=1.0)
        losses = {i: 'categorical_crossentropy' for i in label_names}
        model.compile(loss=losses, optimizer=opt, metrics=['accuracy', Precision(), Recall(), AUC(name='auc_precision_recall', num_thresholds=10000)])
# for gpu
else:
    # create the model
    model = create_model(label_names, tag_counts)
    opt = Adam(learning_rate=5e-05, epsilon=1e-08, decay=0.01, clipnorm=1.0)
    losses = {i: 'categorical_crossentropy' for i in label_names}
    model.compile(loss=losses, optimizer=opt, metrics=['accuracy', Precision(), Recall(), AUC(name='auc_precision_recall', num_thresholds=10000)])

# print model summary
model.summary()

```

Figure 24: Code to create the model

- Once done creating the model, we see the summary of the model, and plot out the model as a diagram.

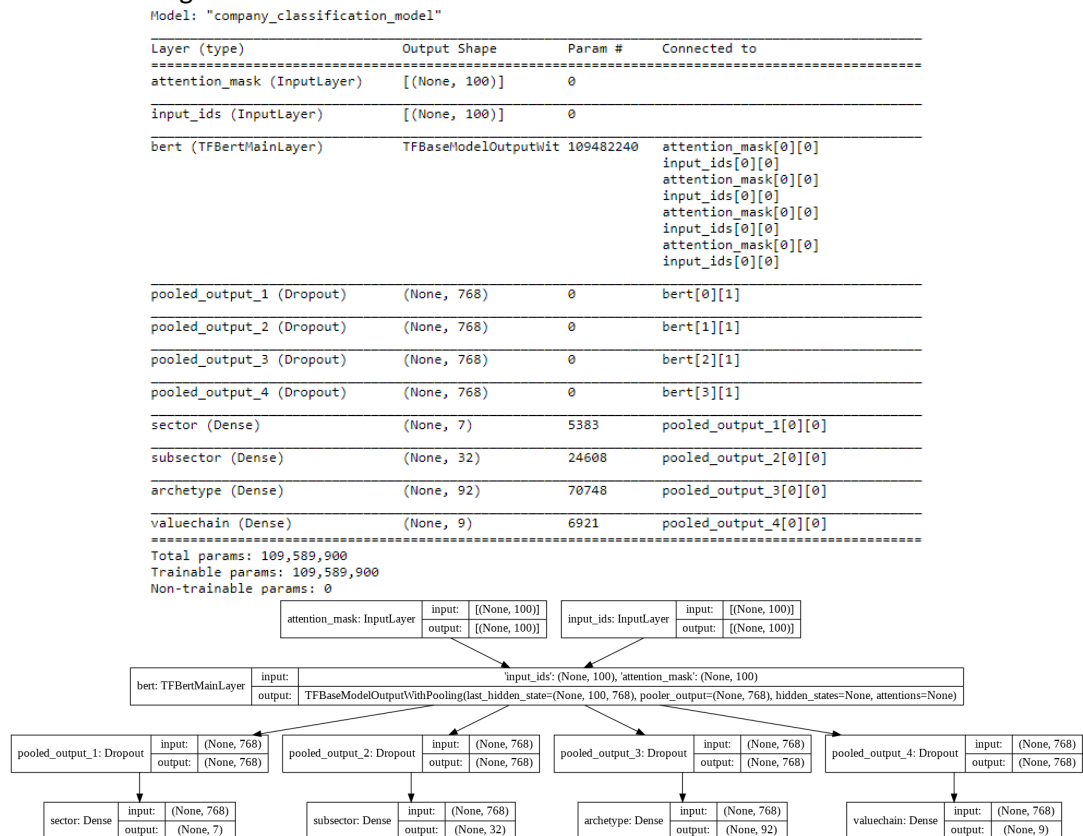


Figure 25: Model summary and diagram

- We then allow the training to proceed. Here, we have a validation split of 0.2, where we partition 20% data from the train dataset for validation. This is done to monitor overfitting.


```

from tensorflow.keras.metrics import Precision, Recall, AUC, CategoricalAccuracy
from tensorflow.keras.optimizers import SGD, Adam
from tensorflow.keras.callbacks import EarlyStopping
from timeit import default_timer as timer

# setting optimizer - based on recommended vals
opt = Adam(learning_rate=5e-05, epsilon=1e-08, decay=0.01, clipnorm=1.0)

# defining loss
losses = {i : 'categorical_crossentropy' for i in label_names}

# train model
start = timer()
history = model.fit({'input_ids': X_train['input_ids'], 'attention_mask': X_train['attention_mask']}, y_train_multi, validation_split=0.2, epochs=15, batch_size=16)
end = timer()
print("Total time taken for execution: ", end-start, "s")

```

Figure 26: Start of model training

Model Evaluation

1. When our model training has been completed, we will use our X_test dataset to evaluate the model. This is a breakdown of the loss, accuracy, precision, recall, PR AUC and F1 Score of the model for the 4 categories.

```
44/44 [=====] - 37s 837ms/step - loss: 4.0287 - sector_loss: 0.3630 - subsector_loss: 0.9883 - archetype_loss: 1.8817 - valuechain_loss: 0.7958 - sector_accuracy: 0.8936 - sector_precision: 0.9037 - sector_recall: 0.8850 - sector_auc_precision_recall: 0.9851 - subsector_accuracy: 0.7336 - subsector_precision: 0.8395 - subsector_recall: 0.6014 - subsector_auc_precision_recall: 0.9782 - archetype_accuracy: 0.5179 - archetype_precision: 0.8336 - archetype_recall: 0.3150 - archetype_auc_precision_recall: 0.9704 - valuechain_accuracy: 0.7429 - valuechain_precision: 0.8201 - valuechain_recall: 0.6250 - valuechain_auc_precision_recall: 0.9592
```

Label	Loss	Accuracy	Precision	Recall	PR AUC	F1 Score
sector	0.363	89.357%	0.734	0.518	0.743	0.304
subsector	0.988	90.372%	0.839	0.834	0.820	0.418
archetype	1.882	88.500%	0.601	0.315	0.625	0.207
valuechain	0.796	98.505%	0.978	0.970	0.959	0.487

Figure 27: Evaluated score of the model

2. We decided to further analyse the details of the result, and created a custom metric analysis based on prediction made using X_test datasets. Here, we create a breakdown of the accuracy of the Sector and Subsector. The correct rows are the number of correctly classified rows. This means that for those classification with 0 correctly classified rows, all the results are wrongly classified.

Label	Accuracy	Correct Rows	% Correct Rows	Precision	Recall

CG	78.0%	137	68.5%	0.872	0.780
> agribusiness	88.89%	8	4.0%	1.000	0.889
> animal protein	78.95%	15	7.5%	1.000	0.789
> fmcg	73.68%	28	14.0%	0.622	0.737
> retail n distribution	66.42%	86	43.0%	0.795	0.664

CNI	95.5%	142	71.0%	0.864	0.955
> building_material	93.48%	86	43.0%	0.835	0.935
> buildings & industrial	91.11%	41	20.5%	0.477	0.911
> cni_equipment suppliers	0.00%	0	0.0%	0.000	0.000
> cni_service providers	11.43%	4	2.0%	0.667	0.114
> infrastructure	0.00%	0	0.0%	0.000	0.000
> utilities	91.67%	11	5.5%	0.550	0.917

IND	88.0%	169	84.5%	0.926	0.880
> auto & mec	90.65%	97	48.5%	0.907	0.907
> diversified	69.57%	16	8.0%	0.842	0.696
> metals and mining	80.00%	56	28.0%	0.889	0.800

ONG	94.5%	159	79.5%	0.917	0.945
> ioc/noc	0.00%	0	0.0%	0.000	0.000
> o&g	87.80%	36	18.0%	0.706	0.878
> o&g_service providers/contractors	80.56%	28	14.0%	0.659	0.806
> ong_traders	0.00%	0	0.0%	0.000	0.000
> petrochemical	93.14%	95	47.5%	0.848	0.931

OOS	87.0%	173	86.5%	0.773	0.870
> others	88.00%	173	86.5%	0.755	0.880

REH	91.5%	118	59.0%	0.953	0.915
> commercial	0.00%	0	0.0%	0.000	0.000
> hotels and accommodation	83.33%	15	7.5%	0.789	0.833
> industrial	0.00%	0	0.0%	0.000	0.000
> mixed	0.00%	0	0.0%	0.000	0.000
> residential	88.24%	60	30.0%	0.550	0.882
> restaurants, catering & services	91.67%	43	21.5%	0.759	0.917

TMT	91.0%	119	59.5%	0.973	0.910
> consumer electronics	96.15%	73	36.5%	0.652	0.962
> digital_business	0.00%	0	0.0%	0.000	0.000
> it_services	96.77%	30	15.0%	0.612	0.968
> media	0.00%	0	0.0%	0.000	0.000
> semiconductor	3.45%	1	0.5%	1.000	0.034
> telecommunication	65.22%	15	7.5%	0.600	0.652

Figure 28: Detailed breakdown of classification results for Sector and Sub-sector

3. We also plotted a heatmap for sector and subsector for a easier identification of the sectors and subsectors mainly classified correctly and wrongly.

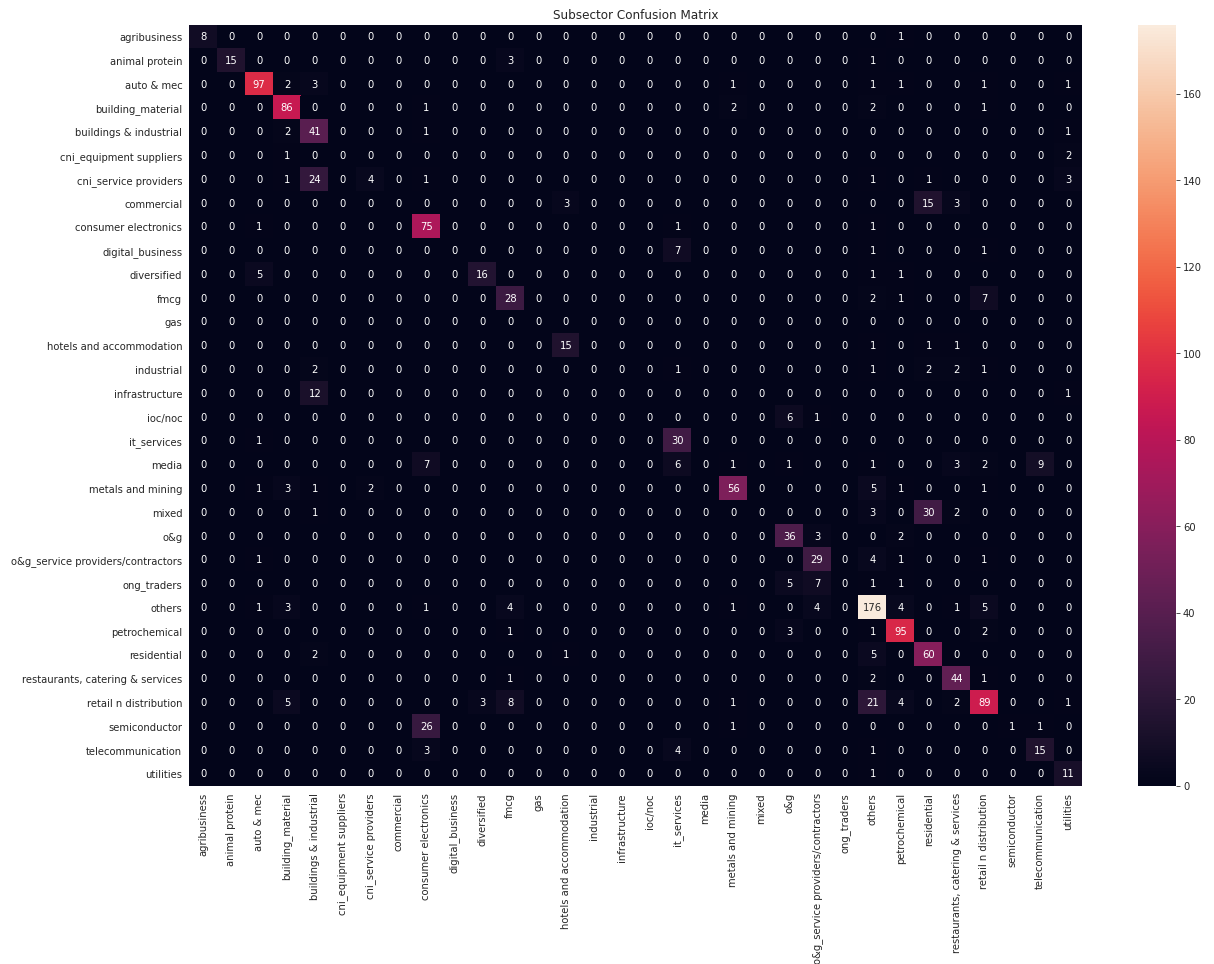
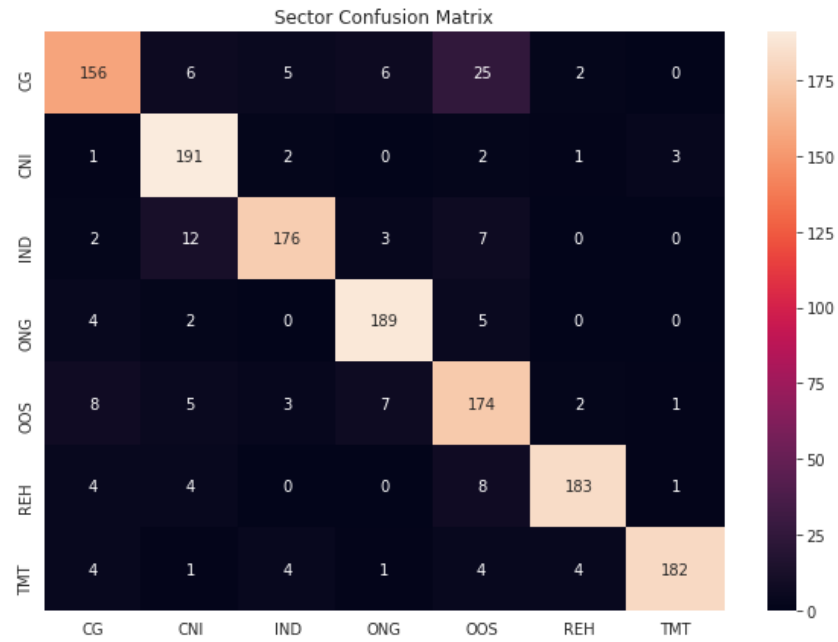
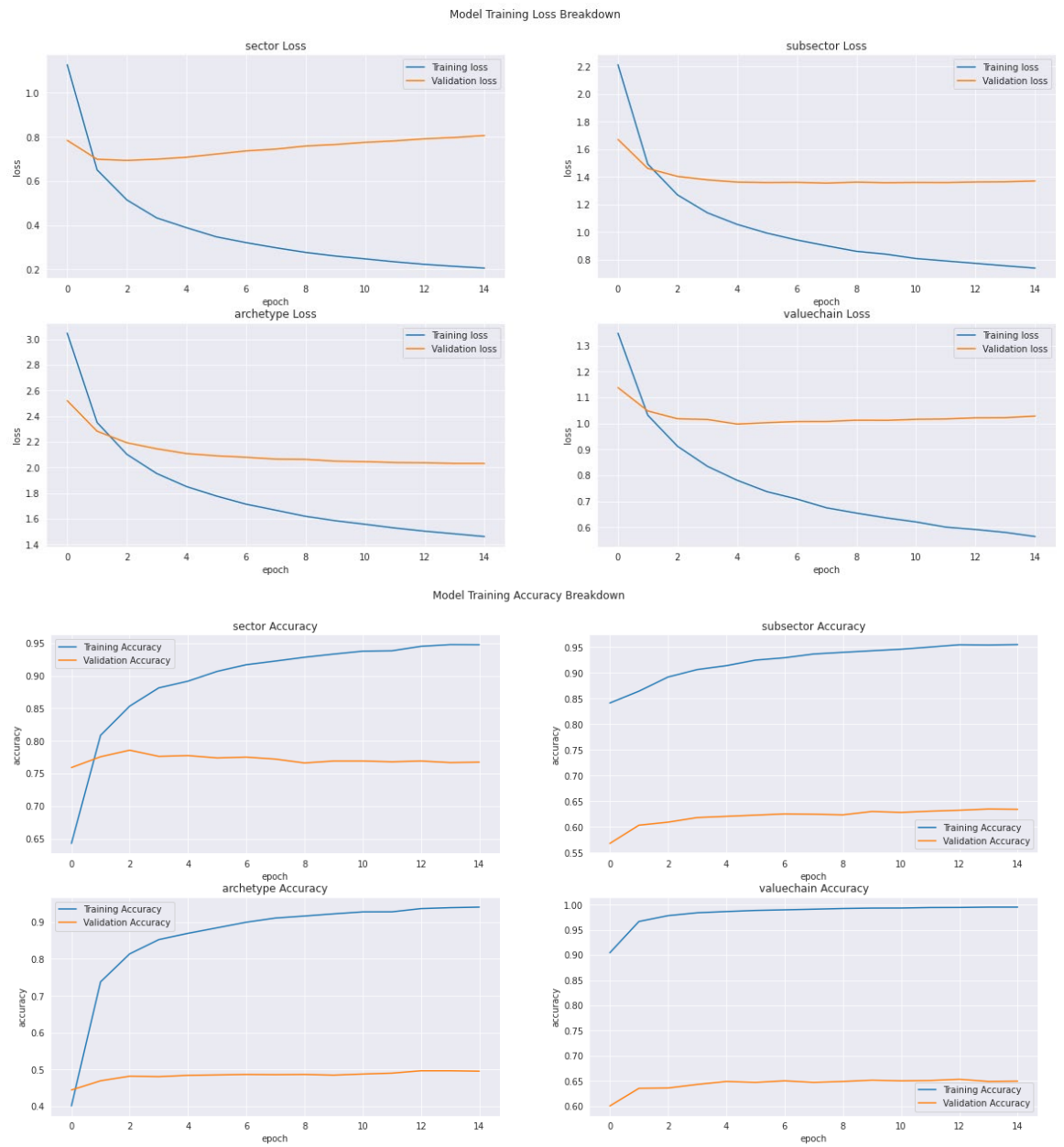


Figure 29: Heatmap for Sector and Subsector classification

4. We also plotted graphs for the sector loss, accuracy and the respective metrics as defined earlier on.



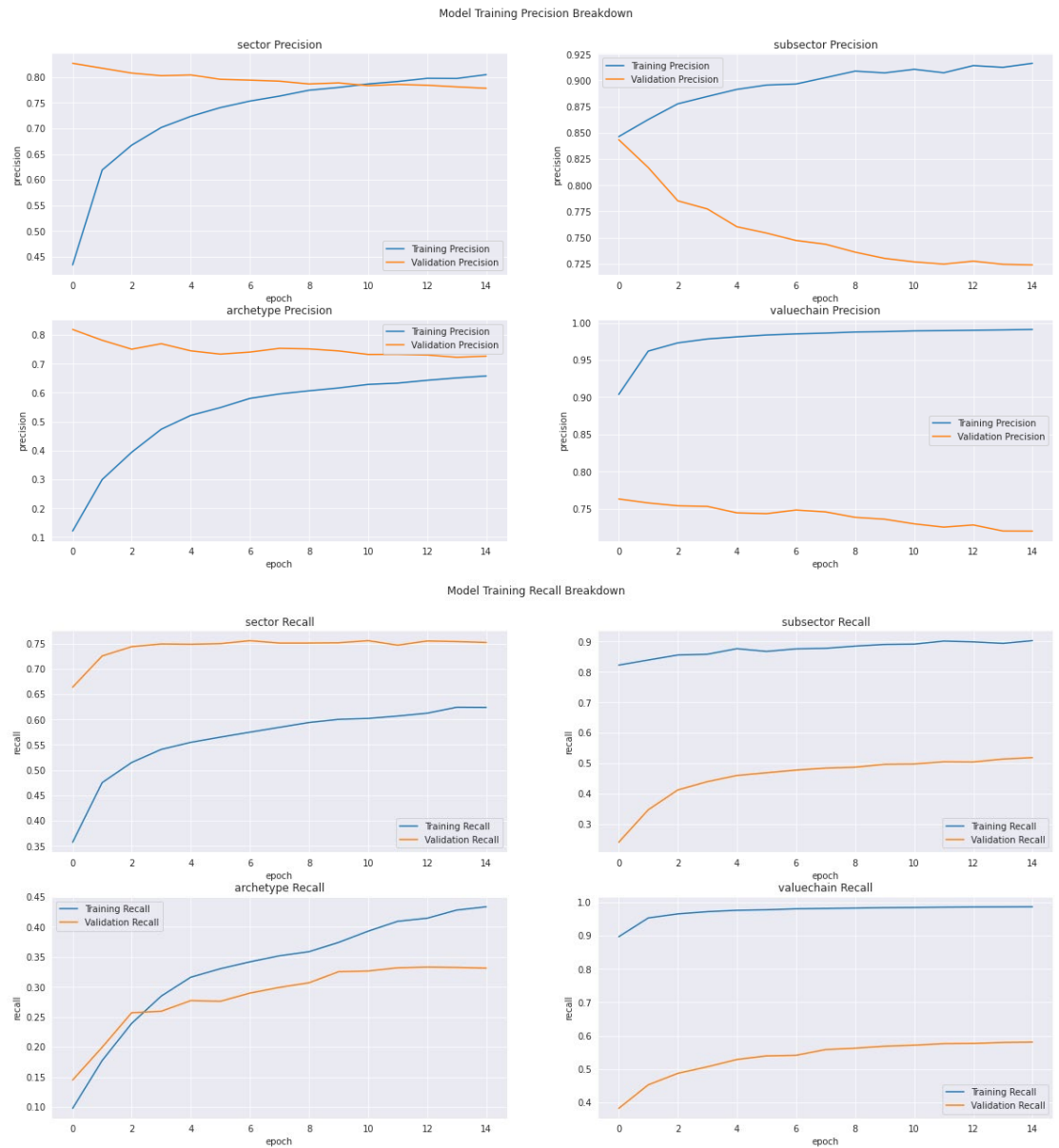


Figure 30: Graphs for the model metrics