# SDSC3006 FUNDAMENTALS OF MACHINE LEARNING I PROJECT

**Tittle： Comparing Logistic Regression Models with and without Principal Component**

DATE: 2023-11-25
TSUI NGA CHING 56621075
LAM LOK HIN     56474716
FUXIAN ZHAO 40149034

# TEAM MEMBER DISTRIBUTION

Tsui Nga Ching: Info, Problem Formulation, Strategies and Methods, Conclusion & discussion

Lam Lok Hin :Logistic Regression- justification, data analysis, results

Fuxian Zhao:PCA, try LDA,SVM-Justification, results

# 1. INFO

- Two datasets: **Training.xlsx (n = 519)** and **test.xlsx (n = 50)**
- Variables: 30 predictors (X1, X2,...,X30), all numerical
- 1 response variable (Y) with two classes (0 and 1)
- Test data has missing values for the response variable

**TRAINING**

| x1 | x2 | x3 | x4 | x5 | x6 | x7 | x8 | x9 | x10 | x11 | x12 | x13 | x14 | x15 | x16 | x17 | x18 | x19 | x20 | x21 | x22 | x23 | x24 | x25 | x26 | x27 | x28 | x29 | x30 | y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 17.99 | 10.38 | 122.8 | 1001 | 0.1184 | 0.2776 | 0.3001 | 0.1471 | 0.2419 | 0.07871 | 1.095 | 0.9053 | 8.589 | 153.4 | 0.006399 | 0.04904 | 0.05373 | 0.01587 | 0.03003 | 0.006193 | 25.38 | 17.33 | 184.6 | 2019 | 0.1622 | 0.6656 | 0.7119 | 0.2654 | 0.4601 | 0.1189 | 0 |
| 19.69 | 21.25 | 130 | 1203 | 0.1096 | 0.1599 | 0.1974 | 0.1279 | 0.2069 | 0.05999 | 0.7456 | 0.7869 | 4.585 | 94.03 | 0.00615 | 0.04006 | 0.03832 | 0.02058 | 0.0225 | 0.004571 | 23.57 | 25.53 | 152.5 | 1709 | 0.1444 | 0.4245 | 0.4504 | 0.243 | 0.3613 | 0.08758 | 0 |
| 11.42 | 20.38 | 77.58 | 386.1 | 0.1425 | 0.2830 | 0.2414 | 0.1052 | 0.2597 | 0.09744 | 0.4956 | 1.156 | 3.445 | 27.23 | 0.00911 | 0.07458 | 0.05661 | 0.01867 | 0.05963 | 0.009208 | 14.91 | 26.5 | 98.87 | 567.7 | 0.2098 | 0.8663 | 0.6869 | 0.2575 | 0.6638 | 0.173 | 0 |
| 20.29 | 14.34 | 135.1 | 1297 | 0.1003 | 0.1328 | 0.198 | 0.1043 | 0.1809 | 0.05883 | 0.7572 | 0.7813 | 5.438 | 94.44 | 0.01149 | 0.02461 | 0.05688 | 0.01885 | 0.01756 | 0.005115 | 22.54 | 16.67 | 152.2 | 1575 | 0.1374 | 0.205 | 0.4 | 0.1625 | 0.2364 | 0.07678 | 0 |
| 12.45 | 15.7 | 82.57 | 477.1 | 0.1278 | 0.17 | 0.1578 | 0.08089 | 0.2087 | 0.07613 | 0.3345 | 0.8902 | 2.217 | 27.19 | 0.00751 | 0.03345 | 0.03672 | 0.01137 | 0.02165 | 0.005082 | 15.47 | 23.75 | 103.4 | 741.6 | 0.1791 | 0.5249 | 0.5355 | 0.1741 | 0.3985 | 0.1244 | 0 |
| 18.25 | 19.98 | 119.6 | 1040 | 0.09463 | 0.109 | 0.1127 | 0.074 | 0.1794 | 0.05742 | 0.4467 | 0.7732 | 3.18 | 53.91 | 0.004314 | 0.01382 | 0.02254 | 0.01039 | 0.01369 | 0.002179 | 22.88 | 27.66 | 153.2 | 1606 | 0.1442 | 0.2576 | 0.3784 | 0.1932 | 0.3063 | 0.08368 | 0 |
| 13.71 | 20.83 | 90.2 | 577.9 | 0.1189 | 0.1645 | 0.09366 | 0.05985 | 0.2196 | 0.07451 | 0.5835 | 1.377 | 3.856 | 50.96 | 0.008805 | 0.03029 | 0.02488 | 0.01448 | 0.01486 | 0.005412 | 17.06 | 28.14 | 110.6 | 897 | 0.1654 | 0.3682 | 0.2678 | 0.1556 | 0.3196 | 0.1151 | 0 |
| 13 | 21.82 | 87.5 | 519.8 | 0.1273 | 0.1932 | 0.1859 | 0.09353 | 0.235 | 0.07389 | 0.3063 | 1.002 | 2.406 | 24.32 | 0.005731 | 0.03502 | 0.03553 | 0.01226 | 0.02143 | 0.003749 | 15.49 | 30.73 | 106.2 | 739.3 | 0.1703 | 0.5401 | 0.539 | 0.206 | 0.4378 | 0.1072 | 0 |
| 12.46 | 24.04 | 83.97 | 475.9 | 0.1186 | 0.2396 | 0.2273 | 0.08543 | 0.203 | 0.08243 | 0.2976 | 1.599 | 2.039 | 23.94 | 0.007149 | 0.07217 | 0.07743 | 0.01432 | 0.01789 | 0.01008 | 15.09 | 40.68 | 97.65 | 711.4 | 0.1853 | 1.058 | 1.105 | 0.221 | 0.4366 | 0.2075 | 0 |
| 15.78 | 17.89 | 103.6 | 781 | 0.0971 | 0.1292 | 0.09954 | 0.06606 | 0.1842 | 0.06082 | 0.5058 | 0.9849 | 3.564 | 54.16 | 0.005771 | 0.04061 | 0.02791 | 0.01282 | 0.02008 | 0.004144 | 20.42 | 27.28 | 136.5 | 1299 | 0.1396 | 0.5609 | 0.3965 | 0.181 | 0.3792 | 0.1048 | 0 |
| 19.17 | 24.8 | 132.4 | 1123 | 0.0974 | 0.2458 | 0.2065 | 0.1118 | 0.2397 | 0.078 | 0.9555 | 3.568 | 11.07 | 116.2 | 0.003139 | 0.08297 | 0.0889 | 0.0409 | 0.04484 | 0.01284 | 20.96 | 29.94 | 151.7 | 1332 | 0.1037 | 0.3903 | 0.3639 | 0.1767 | 0.3176 | 0.1023 | 0 |
| 15.85 | 23.95 | 103.7 | 782.7 | 0.08401 | 0.1002 | 0.09938 | 0.05364 | 0.1847 | 0.05338 | 0.4033 | 1.078 | 2.903 | 36.58 | 0.009789 | 0.03126 | 0.05051 | 0.01992 | 0.02981 | 0.003002 | 16.84 | 27.66 | 112 | 876.5 | 0.1131 | 0.1924 | 0.2322 | 0.1119 | 0.2809 | 0.06287 | 0 |
| 14.68 | 20.13 | 94.74 | 684.5 | 0.09867 | 0.072 | 0.07395 | 0.05259 | 0.1586 | 0.05922 | 0.4727 | 1.24 | 3.195 | 45.4 | 0.005718 | 0.01162 | 0.01998 | 0.01109 | 0.0141 | 0.002085 | 19.07 | 30.88 | 123.4 | 1138 | 0.1464 | 0.1871 | 0.2914 | 0.1609 | 0.3029 | 0.08216 | 0 |
| 16.13 | 20.68 | 108.1 | 798.8 | 0.117 | 0.2022 | 0.1722 | 0.1028 | 0.2164 | 0.07356 | 0.5692 | 1.073 | 3.854 | 54.18 | 0.007026 | 0.02501 | 0.03188 | 0.01297 | 0.01689 | 0.004142 | 20.96 | 31.48 | 136.8 | 1315 | 0.1789 | 0.4233 | 0.4784 | 0.2073 | 0.3706 | 0.1142 | 0 |
| 19.81 | 22.15 | 130 | 1260 | 0.09831 | 0.1027 | 0.1479 | 0.09498 | 0.1582 | 0.05395 | 0.7582 | 1.017 | 5.865 | 112.4 | 0.006494 | 0.01893 | 0.03391 | 0.01521 | 0.01356 | 0.001997 | 27.32 | 30.88 | 186.8 | 2398 | 0.1512 | 0.315 | 0.5372 | 0.2388 | 0.2768 | 0.07615 | 0 |
| 13.54 | 14.36 | 87.46 | 566.3 | 0.09779 | 0.08129 | 0.06664 | 0.04781 | 0.1885 | 0.05766 | 0.2699 | 0.7886 | 2.058 | 23.56 | 0.008462 | 0.0146 | 0.02387 | 0.01315 | 0.0198 | 0.0023 | 15.11 | 19.26 | 99.7 | 711.2 | 0.144 | 0.1773 | 0.239 | 0.1288 | 0.2977 | 0.07259 | 1 |
| 13.08 | 15.71 | 85.63 | 520 | 0.1075 | 0.127 | 0.04568 | 0.0311 | 0.1967 | 0.06811 | 0.1852 | 0.7477 | 1.383 | 14.67 | 0.004007 | 0.01808 | 0.01608 | 0.00649 | 0.01678 | 0.002425 | 14.5 | 20.49 | 96.09 | 630.5 | 0.1312 | 0.2776 | 0.189 | 0.07283 | 0.3184 | 0.08183 | 1 |
| 9.504 | 12.44 | 60.34 | 273.9 | 0.1024 | 0.06492 | 0.02956 | 0.02076 | 0.1815 | 0.06905 | 0.2773 | 0.9768 | 1.909 | 15.7 | 0.009606 | 0.01432 | 0.01985 | 0.01421 | 0.02027 | 0.002968 | 10.23 | 15.66 | 65.13 | 314.9 | 0.1324 | 0.1148 | 0.08867 | 0.06227 | 0.245 | 0.07773 | 1 |
| 15.34 | 14.26 | 102.5 | 704.4 | 0.1073 | 0.2135 | 0.2077 | 0.09756 | 0.2521 | 0.07032 | 0.4388 | 0.7096 | 3.384 | 44.91 | 0.006789 | 0.05328 | 0.06446 | 0.02252 | 0.03672 | 0.004394 | 18.07 | 19.08 | 125.1 | 980.9 | 0.139 | 0.5954 | 0.6305 | 0.2393 | 0.4667 | 0.09946 | 0 |
| 21.16 | 23.04 | 137.2 | 1404 | 0.09428 | 0.1022 | 0.1097 | 0.08632 | 0.1769 | 0.05278 | 0.6917 | 1.127 | 4.303 | 93.99 | 0.004728 | 0.01259 | 0.01715 | 0.01038 | 0.01083 | 0.001987 | 29.17 | 35.59 | 188 | 2615 | 0.1401 | 0.26 | 0.3155 | 0.2009 | 0.2822 | 0.07526 | 0 |
| 16.65 | 21.38 | 110 | 904.6 | 0.1121 | 0.1457 | 0.1525 | 0.0917 | 0.1995 | 0.0633 | 0.8068 | 0.9017 | 5.455 | 102.6 | 0.006048 | 0.01882 | 0.02741 | 0.0113 | 0.01468 | 0.002801 | 26.46 | 31.56 | 177 | 2215 | 0.1805 | 0.3578 | 0.4695 | 0.2095 | 0.3613 | 0.09564 | 0 |
| 17.14 | 16.4 | 116 | 912.7 | 0.1186 | 0.2276 | 0.2229 | 0.1401 | 0.304 | 0.07413 | 1.046 | 0.976 | 7.276 | 111.4 | 0.008029 | 0.03799 | 0.03732 | 0.02397 | 0.02308 | 0.007444 | 22.25 | 21.4 | 152.4 | 1461 | 0.1545 | 0.3949 | 0.3853 | 0.255 | 0.4066 | 0.1059 | 0 |
| 14.58 | 21.53 | 97.41 | 644.8 | 0.1054 | 0.1868 | 0.1425 | 0.08783 | 0.2252 | 0.06924 | 0.2545 | 0.9832 | 2.11 | 21.05 | 0.004452 | 0.03055 | 0.02681 | 0.01352 | 0.01454 | 0.003711 | 17.62 | 33.21 | 122.4 | 896.9 | 0.1525 | 0.6643 | 0.5539 | 0.2701 | 0.4264 | 0.1275 | 0 |
| 18.61 | 20.25 | 122.1 | 1094 | 0.0944 | 0.149 | 0.1731 | 0.1697 | 0.2676 | 0.05994 | 0.8529 | 1.849 | 5.632 | 93.54 | 0.01075 | 0.02722 | 0.05081 | 0.01911 | 0.02293 | 0.004217 | 21.31 | 27.26 | 139.9 | 1403 | 0.1338 | 0.2117 | 0.3446 | 0.149 | 0.2341 | 0.07421 | 0 |
| 15.3 | 25.27 | 102.4 | 732.4 | 0.1082 | 0.1697 | 0.1683 | 0.08751 | 0.1926 | 0.0654 | 0.439 | 1.012 | 3.498 | 43.5 | 0.005233 | 0.03057 | 0.03576 | 0.01083 | 0.01768 | 0.002967 | 20.27 | 36.71 | 149.3 | 1269 | 0.1641 | 0.611 | 0.6335 | 0.2024 | 0.4027 | 0.09876 | 0 |
| 17.57 | 15.05 | 115 | 955.1 | 0.09847 | 0.1157 | 0.09875 | 0.07953 | 0.1739 | 0.06149 | 0.6003 | 0.8225 | 4.655 | 61.1 | 0.005627 | 0.03033 | 0.03407 | 0.01354 | 0.01925 | 0.003742 | 20.01 | 19.52 | 134.9 | 1227 | 0.1255 | 0.2812 | 0.2489 | 0.1456 | 0.2756 | 0.07919 | 0 |
| 18.63 | 25.11 | 124.8 | 1088 | 0.1064 | 0.1887 | 0.2319 | 0.1244 | 0.2183 | 0.06197 | 0.8307 | 1.466 | 5.574 | 105 | 0.006248 | 0.03374 | 0.05196 | 0.01158 | 0.02007 | 0.00456 | 23.15 | 34.01 | 160.5 | 1670 | 0.1491 | 0.4257 | 0.6133 | 0.1848 | 0.3444 | 0.09782 | 0 |
| 17.02 | 23.98 | 112.8 | 899.3 | 0.1197 | 0.1496 | 0.2417 | 0.1203 | 0.2248 | 0.06382 | 0.6009 | 1.398 | 3.999 | 67.78 | 0.008268 | 0.03082 | 0.05042 | 0.01112 | 0.02102 | 0.003854 | 20.88 | 32.09 | 136.1 | 1344 | 0.1634 | 0.3559 | 0.5588 | 0.1847 | 0.353 | 0.08842 | 0 |

**TESTING**

| x1 | x2 | x3 | x4 | x5 | x6 | x7 | x8 | x9 | x10 | x11 | x12 | x13 | x14 | x15 | x16 | x17 | x18 | x19 | x20 | x21 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 13.4 | 20.52 | 88.64 | 556.7 | 0.1106 | 0.1469 | 0.1445 | 0.08172 | 0.2116 | 0.07325 | 0.3906 | 0.9306 | 3.093 | 33.67 | 0.005414 | 0.02265 | 0.03452 | 0.01334 | 0.01705 | 0.004005 | 16.41 |
| 13.21 | 25.25 | 84.1 | 537.9 | 0.08791 | 0.05205 | 0.02772 | 0.02068 | 0.1619 | 0.05584 | 0.2084 | 1.35 | 1.314 | 17.58 | 0.005768 | 0.008082 | 0.0151 | 0.006451 | 0.01347 | 0.001828 | 14.35 |
| 14.02 | 15.66 | 89.59 | 606.5 | 0.07966 | 0.05581 | 0.02087 | 0.02652 | 0.1589 | 0.05586 | 0.2142 | 0.6549 | 1.606 | 19.25 | 0.004837 | 0.009238 | 0.009213 | 0.01076 | 0.01171 | 0.002104 | 14.91 |
| 14.26 | 18.17 | 91.22 | 633.1 | 0.06576 | 0.0522 | 0.02475 | 0.01374 | 0.1635 | 0.05586 | 0.23 | 0.669 | 1.661 | 20.56 | 0.003169 | 0.01377 | 0.01079 | 0.01079 | 0.01103 | 0.001957 | 16.22 |
| 13.03 | 18.42 | 82.61 | 523.8 | 0.08983 | 0.03766 | 0.02562 | 0.02923 | 0.1467 | 0.05863 | 0.1839 | 2.342 | 1.17 | 14.16 | 0.004352 | 0.004899 | 0.01343 | 0.01164 | 0.02671 | 0.001777 | 13.3 |
| 11.34 | 18.61 | 72.76 | 391.2 | 0.1049 | 0.08499 | 0.04302 | 0.02594 | 0.1927 | 0.06211 | 0.243 | 1.01 | 1.491 | 18.19 | 0.008577 | 0.01441 | 0.04304 | 0.01107 | 0.02434 | 0.00217 | 12.47 |
| 12.05 | 22.72 | 78.75 | 447.8 | 0.06935 | 0.1073 | 0.07943 | 0.02978 | 0.1203 | 0.06659 | 0.1194 | 1.434 | 1.778 | 9.549 | 0.005042 | 0.0456 | 0.04305 | 0.01667 | 0.0247 | 0.007358 | 12.57 |
| 11.7 | 19.11 | 74.33 | 418.7 | 0.08814 | 0.05253 | 0.01583 | 0.01148 | 0.1936 | 0.06128 | 0.1601 | 1.43 | 1.109 | 11.28 | 0.006064 | 0.00911 | 0.01042 | 0.007638 | 0.02349 | 0.006161 | 12.61 |
| 7.729 | 25.49 | 47.98 | 178.8 | 0.08098 | 0.04878 | 0 | 0 | 0.187 | 0.07285 | 0.3777 | 1.462 | 2.492 | 19.14 | 0.01266 | 0.009692 | 0 | 0 | 0.02882 | 0.006872 | 9.077 |
| 10.26 | 14.71 | 66.2 | 321.6 | 0.09882 | 0.09159 | 0.03581 | 0.02037 | 0.1633 | 0.07005 | 0.338 | 2.509 | 2.394 | 19.33 | 0.01736 | 0.04671 | 0.02611 | 0.01296 | 0.03675 | 0.006758 | 10.88 |
| 14.69 | 13.98 | 98.22 | 656.1 | 0.1031 | 0.1836 | 0.145 | 0.063 | 0.2086 | 0.07406 | 0.5462 | 1.511 | 4.795 | 49.45 | 0.009976 | 0.05244 | 0.05278 | 0.0158 | 0.02653 | 0.005444 | 16.46 |
| 14.62 | 24.02 | 94.57 | 662.7 | 0.08974 | 0.08606 | 0.03102 | 0.02957 | 0.1685 | 0.05866 | 0.3721 | 1.111 | 2.279 | 33.76 | 0.004868 | 0.01818 | 0.01121 | 0.008606 | 0.02085 | 0.002893 | 16.11 |
| 9.397 | 21.68 | 59.75 | 268.8 | 0.07969 | 0.06053 | 0.03735 | 0.005128 | 0.1274 | 0.06724 | 0.1186 | 1.182 | 1.174 | 6.802 | 0.005515 | 0.02674 | 0.03735 | 0.005128 | 0.01951 | 0.004583 | 9.965 |
| 16.84 | 19.46 | 108.4 | 880.2 | 0.07445 | 0.07223 | 0.0515 | 0.02771 | 0.1844 | 0.05268 | 0.4789 | 2.06 | 3.479 | 46.61 | 0.003443 | 0.02674 | 0.03056 | 0.0111 | 0.0152 | 0.001519 | 18.22 |
| 14.64 | 15.24 | 95.77 | 651.9 | 0.1132 | 0.1339 | 0.09966 | 0.07064 | 0.2116 | 0.06346 | 0.5115 | 0.7372 | 3.814 | 42.76 | 0.005508 | 0.04412 | 0.04436 | 0.01623 | 0.02427 | 0.004841 | 16.34 |
| 15.46 | 11.89 | 102.5 | 736.9 | 0.1257 | 0.1555 | 0.2032 | 0.1097 | 0.1966 | 0.07069 | 0.4209 | 0.6583 | 2.805 | 44.64 | 0.005393 | 0.02321 | 0.04303 | 0.0132 | 0.01792 | 0.004168 | 18.79 |
| 9.042 | 18.9 | 60.07 | 244.5 | 0.09968 | 0.1972 | 0.1975 | 0.04908 | 0.233 | 0.08743 | 0.4653 | 1.911 | 3.769 | 24.2 | 0.009845 | 0.0659 | 0.1027 | 0.02527 | 0.03491 | 0.007877 | 10.06 |
| 20.51 | 27.81 | 134.4 | 1319 | 0.09159 | 0.1074 | 0.1554 | 0.0834 | 0.1448 | 0.05592 | 0.524 | 1.189 | 3.767 | 70.01 | 0.00502 | 0.02062 | 0.03457 | 0.01091 | 0.01298 | 0.002887 | 24.47 |
| 19.55 | 23.21 | 128.9 | 1174 | 0.101 | 0.1318 | 0.1856 | 0.1021 | 0.1989 | 0.05884 | 0.6107 | 2.836 | 5.383 | 70.1 | 0.01124 | 0.04097 | 0.07469 | 0.03441 | 0.02768 | 0.00624 | 20.82 |
| 20.94 | 23.56 | 138.9 | 1364 | 0.1007 | 0.1606 | 0.2712 | 0.131 | 0.2205 | 0.05898 | 1.004 | 0.8208 | 6.372 | 137.9 | 0.005283 | 0.03908 | 0.09518 | 0.01864 | 0.02401 | 0.005002 | 25.58 |
| 11.84 | 18.7 | 77.93 | 440.6 | 0.1109 | 0.1516 | 0.1218 | 0.05182 | 0.2301 | 0.07799 | 0.4825 | 1.03 | 3.475 | | 0.005551 | 0.03414 | 0.04205 | 0.01044 | 0.02273 | 0.005667 | |

# 2.PROBLEM FORMULATION

- The problem at hand is to apply the knowledge and skills acquired in data analysis to analyze real data.
- The objective is to interpret the results of the data analysis and present the findings.
- Specifically, we apply **LDA,SVM,Logistic Regression** to the training dataset for prediction and decide **Logistic Regression** as
- our final model because of its smallest overfitting problem.
- To solve overfitting problem further ,we use **PCA** to reduce dimension of data.
- Programming Language used: **Python**

# 3.STRATEGIES AND METHODS

## Strategy:
The strategy is to compare **accuracy** and **sensitivity** (Recall) between different models and whether the same model uses PCA accoring to cross validation and find the best model.

## Methods:

**Logistic regression**: a statistical model used to predict binary outcomes by fitting a logistic function to the observed data.

**Linear discriminant analysis:** a supervised dimensionality reduction technique that optimally transforms input features to maximize the separation between classes while minimizing within-class variance

**SVM(linear kernel):** a supervised machine learning algorithm that separates classes in a dataset by finding the hyperplane that maximally separates the support vectors, representing instances near the class boundaries.

# 4.JUSTIFICATION

## 1) Data Processing

- import data and relevant package
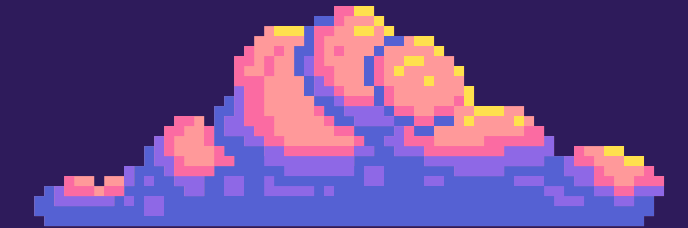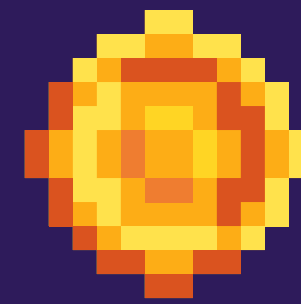- check null values and outliers

## 2)Principal Component Analysis

- Standardizes the training data
- Fits a PCA model to the standardized data
- visualization

# 5.DATA ANALYSIS

After the data preparation, the training data will be fitted in the models. The accuracy and sensitivity of cross validation are used to measure model performance and the accuracy on training data is used to measure overfitting level.For each model ,we use original data and data with top 8 principal variables.

# 6.RESULTS

| MODEL | FITTING ON DATA | ACCURACY | RECALL | ACCURACY ON TRAINING DATA |
|---|---|---|---|---|
| LDA | ORIGINAL DATA | 0.9262574257425742 | 0.9814903846153846 | 0.940 |
| | PRINCIPAL DATA | 0.9268297236743838 | 0.9723543123543124 | 0.933 |
| SVM(LINEAR KERNEL) | ORIGINAL DATA(C=25) | 0.9560990099009901 | 0.9722596153846155 | 0.964 |
| | PRINCIPAL DATA(C=15) | 0.9162178217821783 | 0.953798076923077 | 0.924 |
| LOGISTIC REGRESSION | ORIGINAL DATA | 0.9422516803584763 | 0.9633100233100234 | 0.946 |
| | PRINCIPAL DATA | 0.9042376237623764 | 0.947644230769231 | 0.915 |

Considering accuracy and recall, SVM fitted on the original data performs best, but there is overfitting problem.So,Logistic regression fitted on data with top8 principal variables perform best which has lowest accuracy on training data.

# 7.CONCLUSION& DISCUSSION

## Conclusion:

- Considering overfitting level ,the logistic regression with top 8 principal variables is selected as the final model.According to cross validation, the estimated accuracy and sensitivity are 90.42% and 94.76%,which means most positive data are predicted correctly with high accuracy.

# 7. CONCLUSION& DISCUSSION

## limitation

- Although we try different models and use PCA to reduce dimensions, we don't solve overfitting problem because the accuracy on training data still high.
- However，logistic regression are relevantly simple model，which doesn't cause overfitting problem easily. So, we think the effective way to avoid this problem is to increase amount of data.
- For overfitting model, the estimation of cross validation is not reliable enough.So, we don't have enough confidence to measure the model performance on test data.