| Classification Accuracy Table | | Sentence Level Backdoor Attack | | |
|---|---|---|---|---|
| | | Trigger Sentence Length | | |
| ML Models | Poison Rate | Short | Medium | Long |
| LR | 0.05 | 0.807 | 0.809 | 0.816 |
| DT | 0.05 | 0.771 | 0.778 | 0.779 |
| NB | 0.05 | 0.781 | 0.779 | 0.789 |
| RF | 0.05 | 0.811 | 0.810 | 0.813 |
| LR | 0.1 | 0.813 | 0.809 | 0.809 |
| DT | 0.1 | 0.781 | 0.785 | 0.773 |
| NB | 0.1 | 0.790 | 0.786 | 0.785 |
| RF | 0.1 | 0.808 | 0.804 | 0.800 |
| LR | 0.3 | 0.803 | 0.804 | 0.807 |
| DT | 0.3 | 0.775 | 0.782 | 0.768 |
| NB | 0.3 | 0.796 | 0.782 | 0.778 |
| RF | 0.3 | 0.801 | 0.884 | 0.797 |

| Backdoor Attack Success Rate Table | | Sentence Level Backdoor Attack | | |
|---|---|---|---|---|
| | | Trigger Sentence Length | | |
| ML Models | Poison Rate | Short | Medium | Long |
| LR | 0.05 | 0.836 | 0.866 | 0.894 |
| DT | 0.05 | 0.813 | 0.828 | 0.863 |
| NB | 0.05 | 0.800 | 0.822 | 0.840 |
| RF | 0.05 | 0.845 | 0.875 | 0.880 |
| LR | 0.1 | 0.856 | 0.890 | 0.901 |
| DT | 0.1 | 0.836 | 0.864 | 0.863 |
| NB | 0.1 | 0.812 | 0.833 | 0.839 |
| RF | 0.1 | 0.859 | 0.884 | 0.890 |
| LR | 0.3 | 0.867 | 0.919 | 0.934 |
| DT | 0.3 | 0.831 | 0.899 | 0.922 |
| NB | 0.3 | 0.825 | 0.873 | 0.932 |
| RF | 0.3 | 0.862 | 0.914 | 0.931 |

| Word Level Backdoor Attack (without defense) | | | Word Level Backdoor Attack (with defense) | | |
|---|---|---|---|---|---|
| Number of Trigger Words | | | Number of Trigger Words | | |
| 1 | 3 | 5 | 1 | 3 | 5 |
| 0.805 | 0.809 | 0.808 | 0.797 | 0.794 | 0.785 |
| 0.780 | 0.800 | 0.776 | 0.718 | 0.726 | 0.707 |
| 0.791 | 0.794 | 0.783 | 0.791 | 0.795 | 0.781 |
| 0.810 | 0.808 | 0.805 | 0.790 | 0.797 | 0.799 |
| 0.804 | 0.818 | 0.805 | 0.755 | 0.771 | 0.799 |
| 0.774 | 0.791 | 0.775 | 0.703 | 0.717 | 0.706 |
| 0.782 | 0.795 | 0.805 | 0.767 | 0.779 | 0.768 |
| 0.805 | 0.813 | 0.803 | 0.768 | 0.772 | 0.787 |
| 0.810 | 0.808 | 0.797 | 0.797 | 0.805 | 0.801 |
| 0.773 | 0.788 | 0.773 | 0.669 | 0.718 | 0.756 |
| 0.792 | 0.793 | 0.781 | 0.753 | 0.760 | 0.754 |
| 0.801 | 0.810 | 0.796 | 0.772 | 0.789 | 0.790 |

| Word Level Backdoor Attack (without defense) | | | Word Level Backdoor Attack (with defense) | | |
|---|---|---|---|---|---|
| Number of Trigger Words | | | Number of Trigger Words | | |
| 1 | 3 | 5 | 1 | 3 | 5 |
| 0.809 | 0.891 | 0.907 | 0.803 | 0.876 | 0.934 |
| 0.792 | 0.862 | 0.885 | 0.702 | 0.792 | 0.857 |
| 0.792 | 0.844 | 0.860 | 0.789 | 0.855 | 0.926 |
| 0.808 | 0.888 | 0.903 | 0.796 | 0.855 | 0.932 |
| 0.807 | 0.869 | 0.905 | 0.854 | 0.865 | 0.914 |
| 0.783 | 0.833 | 0.886 | 0.782 | 0.777 | 0.819 |
| 0.780 | 0.824 | 0.890 | 0.834 | 0.845 | 0.927 |
| 0.799 | 0.863 | 0.909 | 0.861 | 0.859 | 0.916 |
| 0.799 | 0.915 | 0.931 | 0.892 | 0.873 | 0.928 |
| 0.785 | 0.889 | 0.913 | 0.696 | 0.763 | 0.913 |
| 0.790 | 0.884 | 0.935 | 0.852 | 0.844 | 0.949 |
| 0.799 | 0.906 | 0.924 | 0.882 | 0.848 | 0.932 |