| Color coding: Relative difference between test runs and original experiment | | dataset SST | | dataset imdb | | dataset Yelp | | dataset Amazon | dataset 20News | dataset Tweets |
|---|--|---|--|---|--|---|--|--|--|--|
| MODEL | | | EST RU | | EST RU | | TEST RU | PAPER delta TEST RU | PAPER delta TEST RU | PAPER delta TEST RU |
| TABLE 2 vanilla lstm test accurately conicity conicity. | uracy bes evaluate.json mean evaluate.json | 0,818 -2% | 0,803 0,713 0,173 | 0,895 0% 0,690 -13% | VG 0,893 0,602 0,135 | rel% 0,956 -1% 0,530 1% n/avail - | 0,949 0,536 0,121 | rel% AVG 0,937 0% 0,937 0,500 -9% 0,457 n/avail - 0,084 | rel% AVG 0,936 -3% 0,908 0,770 -1% 0,761 n/avail - 0,189 | rel% AVG 0,870 -4% 0,833 0,770 1% 0,776 - 0,159 |
| ortho lstm test accuracy conicity, conicity, | | | 0,776 0,283 0,058 | 0,180 -9% | 0,883 0,163 0,041 | 0,960 -2% 0,180 3% n/avail - | 0,945 0,186 0,044 | 0,930 1% 0,936 0,160 7% 0,171 n/avail - 0,032 | 0,922 0% 0,919 0,230 2% 0,235 n/avail - 0,068 | 0,832 1% 0,839 0,270 -4% 0,260 - 0,062 |
| | uracy bes evaluate.json _mean evaluate.json | 0,800 0% 0,200 -6% n/avail - | 0,800 0,188 0,050 | 0,885 -1% 0,080 11% | 0,878 0,089 0,021 | 0,954 -2% 0,060 478% n/avail - | 0,938 0,347 0,014 | 0,929 | 0,910 0% 0,908 0,150 -10% 0,135 n/avail - 0,046 | 0,870 -2% 0,854 0,240 -4% 0,231 - 0,067 |
| FIGURE 3 Indicato Fraction of hid vanilla lstm - ATI median Box plots 1st quar visual inspection paper 3rd qual RANDOM median 1st quar 3rd qual | Source of info importance_ranking_ importance_ranking_ importance_ranking_ importance_ranking_ importance_ranking_ importance_ranking_ | P - P - P - P - P - P - P - P - P - P - | 0,820 0,480 1,000 0,850 0,850 1,000 | 0,75 -27% 0,18 -33% 1,00 0% 0,98 -1% 0,88 -6% | 0,550 0,120 1,000 0,970 0,830 1,000 | 0,90 -6% 0,90 -6% 1,00 0% 0,93 -3% 0,93 -3% 1,00 0% | 0,85 0,85 1,00 0,90 0,90 1,00 | 0,90 | 0,60 | 1,000 1,000 1,000 0,930 0,930 |
| <u>orth lstm - ATTN</u> median 1st quar 3rd quar <u>RANDOM</u> median 1st quar 3rd quai | rtile importance_ranking_ importance_ranking_ importance_ranking_ | P - P - | 0,220 0,100 0,350 0,780 0,570 1,000 | 0,02 0% 0,12 -17% 0,92 -5% 0,73 -15% | 0,040 0,020 0,100 0,870 0,620 0,970 | 0,25 -28% 0,13 -8% 0,33 -18% 0,93 -1% 0,83 -10% 1,00 -2% | 0,180 0,120 0,270 0,920 0,750 0,980 | 0,30 0,20 - 0,40 - 0,90 0,80 - 0,97 | 0,15 -20% 0,120 0,07 -43% 0,040 0,22 45% 0,320 0,93 -1% 0,920 0,80 -8% 0,740 0,99 -1% 0,980 | 0,530 0,190 0,750 0,930 0,750 1,000 |
| diversity Istm - A median 1st quar 3rd quar RANDOM median 1st quar 3rd quar | rtile importance_ranking_ importance_ranking_ tile importance_ranking_ | p – | 0,180 0,080 0,390 0,770 0,500 1,000 | 0,02 0% 0,14 -29% 0,88 0% 0,62 5% | 0,050 0,020 0,100 0,880 0,650 0,970 | 0,13 246% 0,08 125% 0,18 383% 0,90 6% 0,73 12% 0,95 5% | 0,450 0,180 0,870 0,950 0,820 1,000 | 0,17 | 0,10 -20% 0,080 0,04 -25% 0,030 0,23 13% 0,260 0,93 -3% 0,900 0,75 -1% 0,740 1,00 -3% 0,970 | 0,370 0,200 0,530 0,880 0,680 1,000 |
| FIGURE 4 Comparison Violin plots visual inspecti [0.59-0.75] median [0.75-1.00] median | r Source of info Permutation.png Permutation.png Permutation.png Permutation.png | - - - - | 0,050 0,080 0,230 | | 0,180 0,200 | 0,03 67% 0,05 40% 0,20 -25% n/appl - | 0,05 0,07 0,15 | 0,05 | 0,05 | - 0,030 - 0,050 - 0,100 |
| ortho lstm - [0.0] median [0.25-0.50] median [0.50-0.75] median [0.75-1.00] median | Permutation.png Permutation.png Permutation.png Permutation.png | - - - - | 0,120 0,230 0,330 0,360 | 0,43 12% 0,43 21% | 0,450 0,480 0,520 0,570 | 0,37 16% 0,33 21% 0,35 -6% 0,47 -36% | 0,430 0,400 0,330 0,300 | 0,37 0,33 0,40 n/appl - | 0,49 -12% 0,430 0,49 -6% 0,460 0,48 -10% 0,430 0,50 -6% 0,470 | - 0,180 - 0,200 - 0,250 - 0,300 |
| <u>diversity lstm - [</u> 1 median [0.25-0.50] median [0.50-0.75] median [0.75-1.00] median | Permutation.png Permutation.png Permutation.png Permutation.png | - - - - | 0,100 0,230 0,350 0,420 | 0,44 7% | 0,430 0,470 0,510 | 0,43 -30% 0,47 -38% 0,53 -45% 0,55 -27% | 0,30 0,29 0,29 0,40 | 0,3 0,4 0,5 n/appl - | 0,46 2% 0,470 0,48 -2% 0,470 0,46 4% 0,480 0,49 -4% 0,470 | - 0,200 - 0,250 - 0,300 - 0,350 |
| Rational | e attenti rationale_summary_f e length rationale_summary_f | € 0,240 200% | 0,742 0,719 | 0,217 324% | 0,968 0,920 | 0,44 -2% 0,17 120% | 0,430 0,380 | 0,35 0,16 - | 0,63 -2% 0,617 0,22 173% 0,588 | 0,28 189% 0,820 0,23 260% 0,809 |
| Rational diversity Istm Rational | e attenti rationale_summary_f e length rationale_summary_f e attenti rationale_summary_f | - 0,624 -12% | 0,354 0,102 0,550 | 0,76 20% | 0,921 0,269 0,912 | - - 0,57 <mark>-54%</mark> | 0,535 0,114 0,266 | 0,40 | 0,860 0,240 0,88 6% 0,940 | 0,76 -23% 0,591 |
| TABLE 4 Indicato Comparison to vanilla Istm Overall r | mean Pea Attn_Gradient_X_val son corr Attn_Gradient_X_val | 0,71 -13% 0,21 | 0,180 0,619 0,503 | 0,80 <mark>8%</mark> 0,07 | 0,224 0,863 0,793 | 0,16 19% 0,55 25% 0,16 0,40 15% | 0,191 0,690 0,458 | 0,24 - - 0,43 - 0,19 - 0,43 - | 0,17 53% 0,265 0,72 -29% 0,513 0,28 0,65 -40% 0,390 | 0,31 4% 0,319 0,65 -7% 0,605 0,24 0,56 -14% 0,483 |
| Std Pear Overall r Std JS di | rson corr Attn_Integrated_Grai mean JS c Attn_Gradient_X_val v Attn_Gradient_X_val mean JS c Attn_Integrated_Grai | 0,24 0,10 -10% 0,04 d 0,12 -18% | 0,090 | 0,09 0,09 -11% 0,02 | 0,080 0,114 | 0,19 0,15 -19% 0,04 0,19 -7% 0,05 | 0,121 | 0,19 - 0,17 - 0,04 - 0,21 - 0,06 | 0,32 0,15 21% 0,182 0,07 0,21 6% 0,222 0,06 | 0,25 0,08 -6% 0,075 0,03 0,08 -5% 0,076 0,04 |
| Std Pear Overall r Std Pear Overall r Std JS di | mean JS (Attn_Integrated_Grad | - - - | 0,726 0,689 0,143 0,148 | n/avail - n/avail - | 0,903 0,714 0,126 0,183 | n/avail - n/avail - n/avail - | 0,686 0,668 0,163 0,170 | n/avail n/avail n/avail n/avail | - 0,878 - 0,858 - 0,174 - 0,152 | n/avail - 0,747 n/avail - 0,708 n/avail - 0,184 n/avail - 0,186 |
| Std Pear Overall I Std Pear Overall I Std JS di | mean JS (Attn_Integrated_Gra | 0,19 0,79 3% 0,22 0,08 18% 0,05 0,09 12% | | 0,04 | | 0,79 -18% 0,12 0,79 -29% 0,14 0,13 33% 0,04 0,19 -2% | 0,645 0,559 0,173 0,187 | 0,77 0,14 - 0,77 0,14 - 0,12 - 0,12 - 0,12 - 0,04 | 0,96 -5% 0,914 0,08 0,67 29% 0,867 0,11 0,06 180% 0,168 0,04 0,07 91% 0,134 | 0,80 |
| %attenti %attenti | POS pre quant_pos_attn.png ion PUNC ion ADJ | [NOUN, 0% 20% | NOUN, 0% 17% | | 0% 28% | 28% 17% | [NOUN, 0% 17% | | [NOUN, | [NOUN, 10% 11% 12% 10% |
| %attenti %attenti <u>diversity Istm</u> Order of | FPOS pre quant_pos_attn.png ion PUNC | [NOUN, 0 28% | 0% 25% NOUN, 0% 25% | | 0% 35% 0% 35% | [ADJ,NC 3% 30% | [NOUN, - - [NOUN, 0% 17% | | [NOUN, 0% 10% 0% 10% 10% 8% | [NOUN, 3% 14% [NOUN, 3% 5% 13% 12% |
| | ion PUNC | 40% | 47% | | 25% | 76% | 0% | | -9% 0% | 8% 20% |
| TABLE 5 Indicato Preference given to vanilla vs diversity me | r Source of info | 70/0 | 7//0 | | 23/0 | 70/0 | U70 | | 570 U% | 20% |