# Coder reliability

## Analysis of coder agreement

We consider a few possible statistical indicators of agreement between coders.

### Cohen's $\kappa$

Cohen's $\kappa$ [1] is a very popular indicator to compare the agreement between two coders, based on the equation

$$\kappa = \frac{p - p_r}{1 - p_r},\tag{1}$$

where $p$ stands for the agreement rate between coders and $p_r$ for the probability of random agreement. The agreement between pairs of coders according to Cohen's $\kappa$ is shown in table 1.

Table 1: Agreement between pairs of coders according to Cohen's $\kappa$ statistics. $C_i - C_j$ stands for agreement between coder $i$ and $j$.

| Pair | Purpose | Gender | Relation | Min Age | Avg Age | Max Age |
|---|---|---|---|---|---|---|
| $C_1 - C_2$ | 0.815 | 0.961 | 0.636 | 0.476 | 0.582 | 0.555 |
| $C_1 - C_3$ | 0.923 | 0.978 | 0.728 | 0.808 | 0.839 | 0.866 |
| $C_2 - C_3$ | 0.810 | 0.944 | 0.647 | 0.449 | 0.508 | 0.526 |

These results show that in general the agreement is higher for gender, followed by purpose and relation. The agreement between coders 1 and 3 is similar also concerning age, while the agreement with coder 2 is quite poor in these categories. Although there is no real sound mathematical way to evaluate the absolute value of these numbers, according to popular benchmarks, an agreement between 0.8 and 1 is considered as "almost perfect", an agreement between 0.6 and 0.8 as "substantial", while an agreement between 0.4 and 0.6 is only "moderate"[2].

### Fleiss' $\kappa$

It generalises eq. 1 to deal with multiple coders and categories[3]. The corresponding values are shown in table 2.

We see that, in relative terms, agreement is higher for gender, followed by purpose and relation, and lowest for age. In absolute terms, according to the benchmarks, we have almost perfect agreement in gender and purpose, substantial in relation and "fair"

---

[1]Cohen, Jacob *A coefficient of agreement for nominal scales.* Educational and Psychological Measurement 20 (1): 3746 (1960), citation [50] in the main text.

[2]Landis, J.R.; Koch, G.G. (1977). *The measurement of observer agreement for categorical data* Biometrics. 33 (1): 159174, citation [51] in the main text.

[3]Fleiss, J. L. (1971) *Measuring nominal scale agreement among many raters*, Psychological Bulletin, Vol. 76, No. 5 pp. 378–382, citation [52] in the main text.

Table 2: Agreement between coders according to Fleiss' $\kappa$ statistics.

| Purpose | Gender | Relation | Min Age | Avg Age | Max Age |
|---------|--------|----------|---------|---------|---------|
| 0.849   | 0.961  | 0.669    | 0.289   | 0.332   | 0.300   |

(i.e., worst than "moderate") for age indicators, due to the effect of the different coding by coder 2.

Anyway, if we try to plot the age difference between coders, as in figure 1, we see that although disagreement with coder 2 is substantial, it is almost completely limited to a tendency of coder 2 to put pedestrians in a slightly younger category, i.e. the difference in age between the codings is limited. Nevertheless the Fleiss indicator does not take in account the magnitude of difference, and is thus not completely adequate to deal with ordered data.
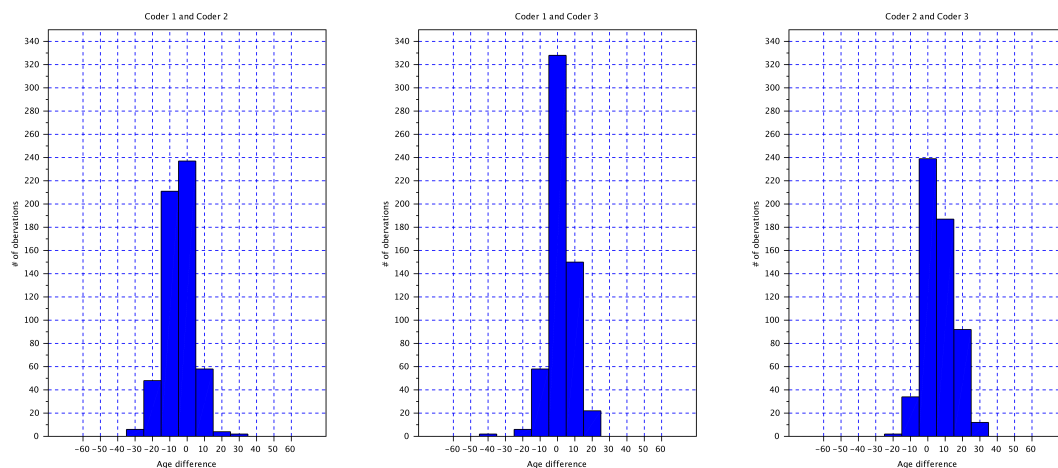


Figure 1: Histograms of age differences between coders.

## Krippendorff's $\alpha$

The Krippendorff $\alpha$ statistics[4], that allows for consideration of quantitative differences between coding results, gives the results shown in table 3.

Krippendorff does not provide any "magic number" but suggests to use data with at least $\alpha > 0.667$ (satisfied by all our categories) and require $\alpha > 0.8$, satisfied by purpose and gender, for reliable results ($\alpha$ between 0.667 and 0.8 could be used for "tentative

---

[4]Krippendorff, Klaus. *Reliability in content analysis*, Human communication research 30.3 (2004): 411-433, citation [53] in the main text.

Table 3: Agreement between coders according to Krippendorff's $\alpha$ statistics. Purpose, gender and relation are "nominal" data, age is on an "interval", according to the definition of $\alpha$ statistics.

| Purpose | Gender | Relation | Min Age | Avg Age | Max Age |
|---------|--------|----------|---------|---------|---------|
| 0.849   | 0.961  | 0.669    | 0.709   | 0.730   | 0.729   |

conclusions").

## Discussion

Using popular indicators of coder reliability, we have found that, in relative terms, the most reliable coding regards gender, followed by purpose. In absolute terms, according to the Krippendorff $\alpha$ statistics that can better cope with the nature of our data, we may see that the purpose and gender codings may be considered as enough reliable to provide sound findings, while the relation and age codings are reliable enough for reporting tentative findings.

The analysis based on these indicators provides an estimate on the reliability of coding of pedestrians in different categories. We may nevertheless use another approach to test the reliability of our findings when based on different coding processes. Since for each category we analyse the values of the observables $V$, $r$, $x$ and $y$, we may compare these quantitative results between different coders.

This comparison, which has also the advantage of being based on more mathematically sound statistical indicators (standard errors, ANOVA analysis) is performed in the following section, and shows again that for purpose and gender we have an almost perfect quantitative agreement, while for relation and age, although the agreement is less good, the major patterns of behaviour are qualitatively observed regardless of coders.

## Quantitative comparison of results

### Purpose

The results (on the common subset of data) for the purpose dependence of all observables between the main coder (coder 1) and the secondary coders are compared in tables 4, 5 and 6.

The differences between coders are thus always of one standard error or smaller, and the extremely significant statistical differences in the $x$ and $V$ distribution (along with the less significant $y$ and $r$ ones) are reported by all coders.

### Relation

The results (on the common subset of data) for the relation dependence of all observables between the main coder (coder 1) and the secondary coders are compared in tables 7,

Table 4: Observable dependence on purpose for dyads according to coder 1 (common data set only). Lengths in millimetres, times in seconds.

| Purpose | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Leisure | 136 | $1085 \pm 19$ ($\sigma=220$) | $796 \pm 21$ ($\sigma=248$) | $636 \pm 13$ ($\sigma=151$) | $351 \pm 28$ ($\sigma=327$) |
| Work | 132 | $1257 \pm 14$ ($\sigma=157$) | $829 \pm 17$ ($\sigma=196$) | $723 \pm 12$ ($\sigma=143$) | $303 \pm 21$ ($\sigma=241$) |
| $F_{1,266}$ | | 53.1 | 1.41 | 23.5 | 1.88 |
| $p$ | | $< 10^{-8}$ | 0.236 | $2.14 \cdot 10^{-6}$ | 0.171 |
| $R^2$ | | 0.166 | 0.00529 | 0.0811 | 0.00703 |
| $\delta$ | | 0.893 | 0.146 | 0.594 | 0.168 |

Table 5: Observable dependence on purpose for dyads according to coder 2 (common data set only). Lengths in millimetres, times in seconds.

| Purpose | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Leisure | 151 | $1093 \pm 17$ ($\sigma=212$) | $793 \pm 20$ ($\sigma=243$) | $641 \pm 12$ ($\sigma=147$) | $344 \pm 26$ ($\sigma=318$) |
| Work | 117 | $1269 \pm 15$ ($\sigma=159$) | $837 \pm 18$ ($\sigma=196$) | $728 \pm 13$ ($\sigma=146$) | $306 \pm 22$ ($\sigma=243$) |
| $F_{1,56}$ | | 56.2 | 2.56 | 23.4 | 1.13 |
| $p$ | | $< 10^{-8}$ | 0.111 | $2.18 \cdot 10^{-6}$ | 0.289 |
| $R^2$ | | 0.175 | 0.00954 | 0.081 | 0.00422 |
| $\delta$ | | 0.927 | 0.198 | 0.599 | 0.131 |

Table 6: Observable dependence on purpose for dyads according to coder 3 (common data set only). Lengths in millimetres, times in seconds.

| Purpose | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Leisure | 133 | $1077 \pm 19$ ($\sigma=217$) | $789 \pm 22$ ($\sigma=250$) | $626 \pm 13$ ($\sigma=145$) | $354 \pm 29$ ($\sigma=330$) |
| Work | 133 | $1262 \pm 14$ ($\sigma=156$) | $836 \pm 17$ ($\sigma=195$) | $732 \pm 12$ ($\sigma=144$) | $302 \pm 21$ ($\sigma=239$) |
| $F_{1,264}$ | | 63.6 | 2.93 | 35.6 | 2.13 |
| $p$ | | $< 10^{-8}$ | 0.0881 | $< 10^{-8}$ | 0.145 |
| $R^2$ | | 0.194 | 0.011 | 0.119 | 0.00802 |
| $\delta$ | | 0.982 | 0.211 | 0.734 | 0.18 |

8 and 9. While all the major trends exposed in the main text are confirmed, quantitative results between coders may sometimes be different (we refer in particular to the $y$ distribution for couples, extremely narrow according to coder 3).

## Gender

The results (on the common subset of data) for the gender dependence of all observables between the main coder (coder 1) and the secondary coders are compared in tables 10, 11 and 12, showing that there is basically no difference in the coding of gender.

Table 7: Observable dependence on relation for dyads according to coder 1 (common data set only). Lengths in millimetres, times in seconds.

| Relation | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Colleagues | 125 | $1256 \pm 14$ ($\sigma$=154) | $829 \pm 18$ ($\sigma$=196) | $725 \pm 13$ ($\sigma$=142) | $301 \pm 21$ ($\sigma$=239) |
| Couples | 28 | $1087 \pm 37$ ($\sigma$=194) | $690 \pm 33$ ($\sigma$=174) | $611 \pm 21$ ($\sigma$=112) | $248 \pm 37$ ($\sigma$=198) |
| Families | 40 | $1051 \pm 24$ ($\sigma$=153) | $864 \pm 54$ ($\sigma$=341) | $594 \pm 21$ ($\sigma$=134) | $492 \pm 69$ ($\sigma$=438) |
| Friends | 56 | $1121 \pm 36$ ($\sigma$=271) | $777 \pm 24$ ($\sigma$=182) | $669 \pm 19$ ($\sigma$=145) | $286 \pm 32$ ($\sigma$=243) |
| $F_{3,245}$ | | 16.4 | 4.19 | 11.8 | 6.12 |
| $p$ | | $< 10^{-8}$ | 0.00651 | $3.06 \cdot 10^{-7}$ | 0.0005 |
| $R^2$ | | 0.167 | 0.0488 | 0.126 | 0.0697 |
| $\delta$ | | 1.33 | 0.612 | 0.934 | 0.678 |

Table 8: Observable dependence on relation for dyads according to coder 2 (common data set only). Lengths in millimetres, times in seconds.

| Relation | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Colleagues | 116 | $1267 \pm 14$ ($\sigma$=156) | $839 \pm 18$ ($\sigma$=197) | $729 \pm 14$ ($\sigma$=147) | $308 \pm 23$ ($\sigma$=244) |
| Couples | 44 | $1082 \pm 28$ ($\sigma$=184) | $703 \pm 21$ ($\sigma$=140) | $582 \pm 19$ ($\sigma$=125) | $296 \pm 33$ ($\sigma$=221) |
| Families | 42 | $1054 \pm 25$ ($\sigma$=164) | $894 \pm 53$ ($\sigma$=341) | $651 \pm 25$ ($\sigma$=163) | $451 \pm 70$ ($\sigma$=457) |
| Friends | 66 | $1131 \pm 31$ ($\sigma$=254) | $786 \pm 23$ ($\sigma$=188) | $673 \pm 17$ ($\sigma$=136) | $304 \pm 29$ ($\sigma$=238) |
| $F_{3,264}$ | | 19 | 6.55 | 11.9 | 3.13 |
| $p$ | | $< 10^{-8}$ | 0.000276 | $2.54 \cdot 10^{-7}$ | 0.0262 |
| $R^2$ | | 0.178 | 0.0692 | 0.119 | 0.0344 |
| $\delta$ | | 1.35 | 0.74 | 1.04 | 0.437 |

Table 9: Observable dependence on relation for dyads according to coder 3 (common data set only) Lengths in millimetres, times in seconds..

| Relation | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Colleagues | 136 | $1259 \pm 14$ ($\sigma$=158) | $834 \pm 17$ ($\sigma$=194) | $727 \pm 13$ ($\sigma$=147) | $304 \pm 21$ ($\sigma$=242) |
| Couples | 23 | $1070 \pm 42$ ($\sigma$=204) | $624 \pm 20$ ($\sigma$=96.4) | $578 \pm 20$ ($\sigma$=95.1) | $182 \pm 17$ ($\sigma$=81.2) |
| Families | 50 | $1053 \pm 24$ ($\sigma$=172) | $867 \pm 44$ ($\sigma$=312) | $612 \pm 20$ ($\sigma$=140) | $478 \pm 59$ ($\sigma$=416) |
| Friends | 54 | $1084 \pm 32$ ($\sigma$=235) | $780 \pm 27$ ($\sigma$=196) | $663 \pm 22$ ($\sigma$=159) | $298 \pm 33$ ($\sigma$=245) |
| $F_{3,259}$ | | 23.4 | 7.57 | 12.4 | 7.52 |
| $p$ | | $< 10^{-8}$ | $7.11 \cdot 10^{-5}$ | $1.36 \cdot 10^{-7}$ | $7.61 \cdot 10^{-5}$ |
| $R^2$ | | 0.213 | 0.0807 | 0.125 | 0.0801 |
| $\delta$ | | 1.27 | 0.915 | 1.06 | 0.849 |

**Age**

The results (on the common subset of data) for the minimum age dependence of all observables between the main coder (coder 1) and the secondary coders are compared in tables 13, 14 and 15. Sadly, almost no groups with children are present in the common set. The drop in velocity with age is, on the other hand, confirmed in a statistically significant way by all coders.

Table 10: Observable dependence on gender for dyads according to coder 1 (common data set only). Lengths in millimetres, times in seconds.

| Gender | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Two females | 55 | 1076 ± 32 ($\sigma$=240) | 745 ± 21 ($\sigma$=155) | 629 ± 15 ($\sigma$=112) | 290 ± 33 ($\sigma$=242) |
| Mixed | 86 | 1095 ± 19 ($\sigma$=173) | 820 ± 31 ($\sigma$=287) | 641 ± 17 ($\sigma$=159) | 384 ± 39 ($\sigma$=360) |
| Two males | 127 | 1261 ± 16 ($\sigma$=178) | 836 ± 17 ($\sigma$=195) | 727 ± 13 ($\sigma$=150) | 305 ± 22 ($\sigma$=243) |
| $F_{2,265}$ | | 27.2 | 3.25 | 12.8 | 2.48 |
| $p$ | | $< 10^{-8}$ | 0.0404 | $5.09 \cdot 10^{-6}$ | 0.0855 |
| $R^2$ | | 0.171 | 0.0239 | 0.0879 | 0.0184 |
| $\delta$ | | 0.93 | 0.494 | 0.699 | 0.292 |

Table 11: Observable dependence on gender for dyads according to coder 2 (common data set only). Lengths in millimetres, times in seconds.

| Gender | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Two females | 53 | 1078 ± 33 ($\sigma$=241) | 747 ± 22 ($\sigma$=158) | 637 ± 15 ($\sigma$=106) | 286 ± 32 ($\sigma$=233) |
| Mixed | 89 | 1093 ± 18 ($\sigma$=173) | 814 ± 30 ($\sigma$=283) | 635 ± 17 ($\sigma$=159) | 382 ± 38 ($\sigma$=360) |
| Two males | 126 | 1263 ± 16 ($\sigma$=177) | 838 ± 17 ($\sigma$=194) | 728 ± 13 ($\sigma$=150) | 306 ± 22 ($\sigma$=244) |
| $F_{2,265}$ | | 28.2 | 3.12 | 13.3 | 2.48 |
| $p$ | | $< 10^{-8}$ | 0.0459 | $3.22 \cdot 10^{-6}$ | 0.0853 |
| $R^2$ | | 0.176 | 0.023 | 0.091 | 0.0184 |
| $\delta$ | | 0.935 | 0.494 | 0.604 | 0.3 |

Table 12: Observable dependence on gender for dyads according to coder 3 (common data set only). Lengths in millimetres, times in seconds.

| Gender | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| Two females | 55 | 1074 ± 32 ($\sigma$=239) | 742 ± 21 ($\sigma$=153) | 636 ± 14 ($\sigma$=103) | 281 ± 31 ($\sigma$=230) |
| Mixed | 89 | 1093 ± 19 ($\sigma$=175) | 824 ± 31 ($\sigma$=288) | 634 ± 17 ($\sigma$=161) | 397 ± 39 ($\sigma$=368) |
| Two males | 124 | 1267 ± 16 ($\sigma$=173) | 834 ± 17 ($\sigma$=190) | 730 ± 13 ($\sigma$=150) | 298 ± 21 ($\sigma$=232) |
| $F_{2,265}$ | | 30.4 | 3.44 | 14.2 | 3.99 |
| $p$ | | $< 10^{-8}$ | 0.0336 | $1.36 \cdot 10^{-6}$ | 0.0196 |
| $R^2$ | | 0.187 | 0.0253 | 0.0969 | 0.0293 |
| $\delta$ | | 0.987 | 0.511 | 0.622 | 0.359 |

Table 13: Observable dependence on minimum age for dyads according to coder 1 (common data set only). Lengths in millimetres, times in seconds.

| Minimum age | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| 10-19 years | 16 | 1157 ± 86 ($\sigma$=343) | 715 ± 31 ($\sigma$=123) | 653 ± 23 ($\sigma$=92.3) | 223 ± 38 ($\sigma$=151) |
| 20-29 years | 58 | 1183 ± 28 ($\sigma$=215) | 765 ± 28 ($\sigma$=211) | 666 ± 20 ($\sigma$=149) | 268 ± 33 ($\sigma$=252) |
| 30-39 years | 96 | 1186 ± 21 ($\sigma$=203) | 817 ± 21 ($\sigma$=211) | 689 ± 17 ($\sigma$=166) | 327 ± 27 ($\sigma$=262) |
| 40-49 years | 41 | 1193 ± 25 ($\sigma$=161) | 811 ± 27 ($\sigma$=173) | 684 ± 22 ($\sigma$=143) | 327 ± 38 ($\sigma$=245) |
| 50-59 years | 31 | 1210 ± 29 ($\sigma$=160) | 880 ± 46 ($\sigma$=254) | 696 ± 29 ($\sigma$=160) | 407 ± 65 ($\sigma$=360) |
| 60-69 years | 21 | 1017 ± 35 ($\sigma$=160) | 869 ± 66 ($\sigma$=304) | 671 ± 34 ($\sigma$=156) | 401 ± 85 ($\sigma$=388) |
| $\geq$ 70 years | 5 | 949 ± 15 ($\sigma$=34.2) | 913 ± 170 ($\sigma$=379) | 608 ± 28 ($\sigma$=61.7) | 551 ± 210 ($\sigma$=470) |
| $F_{6,261}$ | | 3.35 | 1.81 | 0.462 | 1.91 |
| $p$ | | 0.00337 | 0.0974 | 0.836 | 0.0789 |
| $R^2$ | | 0.0715 | 0.0399 | 0.0105 | 0.0421 |
| $\delta$ | | 1.73 | 0.964 | 0.578 | 1.29 |

Table 14: Observable dependence on minimum age for dyads according to coder 2 (common data set only). Lengths in millimetres, times in seconds.

| Minimum age | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| 0-9 years | 2 | 1190 ± 220 ($\sigma$=312) | 749 ± 110 ($\sigma$=152) | 700 ± 87 ($\sigma$=123) | 202 ± 55 ($\sigma$=77.8) |
| 10-19 years | 16 | 1169 ± 84 ($\sigma$=334) | 682 ± 20 ($\sigma$=80.1) | 646 ± 20 ($\sigma$=78.4) | 172 ± 18 ($\sigma$=73.3) |
| 20-29 years | 107 | 1163 ± 19 ($\sigma$=196) | 765 ± 16 ($\sigma$=165) | 655 ± 13 ($\sigma$=138) | 288 ± 22 ($\sigma$=231) |
| 30-39 years | 78 | 1217 ± 24 ($\sigma$=209) | 869 ± 28 ($\sigma$=244) | 727 ± 21 ($\sigma$=185) | 362 ± 33 ($\sigma$=290) |
| 40-49 years | 32 | 1181 ± 31 ($\sigma$=176) | 855 ± 44 ($\sigma$=249) | 645 ± 22 ($\sigma$=124) | 418 ± 66 ($\sigma$=373) |
| 50-59 years | 24 | 1074 ± 32 ($\sigma$=158) | 853 ± 60 ($\sigma$=293) | 706 ± 29 ($\sigma$=143) | 343 ± 74 ($\sigma$=363) |
| 60-69 years | 9 | 1047 ± 42 ($\sigma$=127) | 861 ± 100 ($\sigma$=311) | 655 ± 45 ($\sigma$=135) | 432 ± 120 ($\sigma$=375) |
| $F_{6,261}$ | | 2.1 | 3.09 | 2.3 | 2.13 |
| $p$ | | 0.0533 | 0.0061 | 0.0349 | 0.0505 |
| $R^2$ | | 0.0461 | 0.0663 | 0.0503 | 0.0467 |
| $\delta$ | | 0.842 | 0.833 | 0.482 | 1.14 |

Table 15: Observable dependence on minimum age for dyads according to coder 3 (common data set only). Lengths in millimetres, times in seconds.

| Minimum age | $N_g^k$ | $V$ | $r$ | $x$ | $y$ |
|---|---|---|---|---|---|
| 10-19 years | 14 | 1163 ± 98 ($\sigma$=367) | 701 ± 31 ($\sigma$=117) | 623 ± 35 ($\sigma$=130) | 218 ± 48 ($\sigma$=181) |
| 20-29 years | 64 | 1157 ± 29 ($\sigma$=236) | 758 ± 24 ($\sigma$=194) | 658 ± 20 ($\sigma$=158) | 274 ± 27 ($\sigma$=220) |
| 30-39 years | 50 | 1197 ± 27 ($\sigma$=193) | 830 ± 32 ($\sigma$=227) | 685 ± 23 ($\sigma$=162) | 349 ± 43 ($\sigma$=302) |
| 40-49 years | 77 | 1205 ± 19 ($\sigma$=163) | 832 ± 23 ($\sigma$=205) | 684 ± 16 ($\sigma$=141) | 351 ± 33 ($\sigma$=293) |
| 50-59 years | 36 | 1205 ± 25 ($\sigma$=152) | 820 ± 35 ($\sigma$=207) | 722 ± 28 ($\sigma$=168) | 300 ± 39 ($\sigma$=233) |
| 60-69 years | 20 | 1025 ± 40 ($\sigma$=179) | 903 ± 74 ($\sigma$=332) | 699 ± 29 ($\sigma$=129) | 418 ± 97 ($\sigma$=436) |
| $\geq$ 70 years | 7 | 956 ± 26 ($\sigma$=69.1) | 881 ± 120 ($\sigma$=326) | 605 ± 36 ($\sigma$=96.4) | 503 ± 140 ($\sigma$=382) |
| $F_{6,261}$ | | 3.69 | 2.04 | 1.33 | 1.67 |
| $p$ | | 0.00153 | 0.0604 | 0.245 | 0.129 |
| $R^2$ | | 0.0783 | 0.0449 | 0.0296 | 0.0369 |
| $\delta$ | | 1.74 | 0.755 | 0.732 | 1.09 |