# Web-based Supplementary Materials for "Comparing trends in cancer rates across overlapping regions" by Y. Li and R. Tiwari

## Derivation of Equation (8)

To proceed, we assume that $t_1 \leq t_{s+1} < t_m \leq t_{s+I}$ and note that

$$
\begin{aligned}
Cov(\hat{\beta}_{11}, \hat{\beta}_{21}) &= \frac{1}{\sigma_1^2 \sigma_2^2} Cov \left\{ \sum_{i=1}^{m} (t_i - \bar{t}_1) y_{1i}, \sum_{s+1}^{s+I} (t_i - \bar{t}_2) y_{2i} \right\} \\
&= \frac{1}{\sigma_1^2 \sigma_2^2} \sum_{s+1}^{m} (t_i - \bar{t}_1)(t_i - \bar{t}_2) Cov(y_{1i}, y_{2i}). \quad (14)
\end{aligned}
$$

Recall that we use superscript 'O' to denote the intersection of Regions 1 and 2 and 'NO' the non-overlapping subset. We further introduce the following notation. Let $n_{kji}$, $n_{kji}^{(O)}$ and $n_{kji}^{(NO)}$ be the numbers of underlying population at risk for age group $j$ at time $t_i$ in Region $k(k = 1, 2)$, in the overlapping subregion and in the non-overlapping subregions, respectively. Similarly, define $d_{kji}$, $d_{kji}^{(O)}$ and $d_{kji}^{(NO)}$ the corresponding numbers of events (e.g. deaths or cancer cases). Denote by $n_{ki} = \sum_{j=1}^{J} n_{kji}, n_{ki}^{(O)} = \sum_{j=1}^{J} n_{kji}^{(O)}, n_{ki}^{(NO)} = \sum_{j=1}^{J} n_{kji}^{(NO)}$. Also define $d_{ki}, d_{ki}^{(O)}$ and $d_{kji}^{(NO)}$ in the similar fashion. In fact, $d_{kji}^{(O)}$ and $n_{kj}^{(O)}$ are independent of index $k$ (for region) as they correspond to the same common subregion for $k = 1, 2$.

Let $y_i^{(O)} = \log(r_i^{(O)}) = \log \left( \sum_{j=1}^{J} w_j \frac{d_{ji}^{(O)} + 1/J}{n_{ji}^{(O)}} \right)$ be the logarithm of the (zero corrected) age-adjusted rate $r_i^{(O)}$ at time $t_i$ for the overlapping region, and let $y_{1i}^{(NO)}$ and $y_{2i}^{(NO)}$ be defined similarly based on $r_{1i}^{(NO)}$ and $r_{2i}^{(NO)}$, respectively, for the non-overlapping regions/intervals for the two groups.

Dropping the subscript $i$ (for time), we assume the age groups have the same distribution across the overlapping and non-overlapping regions, that is,

$$
\frac{n_{k1}^{(O)}}{n_{k1}} = \frac{n_{k2}^{(O)}}{n_{k2}} = \ldots = \frac{n_{kJ}^{(O)}}{n_{kJ}} = p_k^{(O)}, \text{ and } \frac{n_{k1}^{(NO)}}{n_{k1}} = \frac{n_{k2}^{(NO)}}{n_{k2}} = \ldots = \frac{n_{kJ}^{(NO)}}{n_{kJ}} = p_k^{(NO)}, \quad (15)
$$

for k=1,2. This assumption is common in comparing the age-adjusted rates across different geographical areas (see, e.g., Pickle and White, 1995), under which, we have

$$
r_k = \sum_{j=1}^{J} w_j \frac{d_{kj}}{n_{kj}} = \sum_{j=1}^{J} w_j \frac{d_{kj}^{(O)} + d_{kj}^{(NO)}}{n_{kj}}
$$

$$= \sum_{j=1}^{J} w_j \frac{n_{kj}^{(O)}}{n_{kj}} \frac{d_{kj}^{(O)} + 1/J}{n_{kj}^{(O)}} + \sum_{j=1}^{J} w_j \frac{n_{kj}^{(NO)}}{n_{kj}} \frac{d_{kj}^{(NO)} + 1/J}{n_{kj}^{(NO)}} + c_k$$

$$= p_k^{(O)} r_k^{(O)} + p_k^{(NO)} r_k^{(NO)} + c_k,$$

where $c_k = -\frac{1}{J} \sum_{j=1}^{J} \frac{w_j}{n_{kj}}$, a negligible constant. Again, since $r_1^{(O)} = r_2^{(O)}$, let $r^{(O)}$ denote this common value, and let $y^{(O)} = \log(r^{(O)})$. Now, since $Cov(r_1^{(NO)}, r_2^{(NO)}) = 0$ and $Cov(r^{(O)}, r_k^{(NO)}) = 0, k = 1, 2$, using the delta method, we have,

$$
\begin{aligned}
Cov(y_1, y_2) &= Cov(\log(r_1), \log(r_2)) \\
&\approx \frac{1}{E(r_1)E(r_2)} Cov(r_1, r_2) \\
&= \frac{1}{E(r_1)E(r_2)} p_1^{(O)} p_2^{(O)} Var(r^{(O)}) \\
&= \frac{1}{E(r_1)E(r_2)} p_1^{(O)} p_2^{(O)} Var(e^{y^{(O)}}).
\end{aligned}
$$

Let $y^{(O)}$ satisfy the regression model (3), and let $\mu^{(O)} = E(y^{(O)})$. Since $y^{(O)} \sim N(\mu^{(O)}, \sigma^2)$, using the properties of that log normal distribution, we have that

$$E(r^{(O)}) = E(e^{y^{(O)}}) = e^{\mu^{(O)}} e^{\sigma^2/2},$$
$$Var(e^{y^{(O)}}) = e^{2\mu^{(O)}} e^{\sigma^2} \left(e^{\sigma^2} - 1\right).$$

Furthermore, the null hypothesis implies that $E(y_1) = E(y_2) = E(y^{(O)})$. Hence, adding back the time index $i$, we will have

$$
\begin{aligned}
Cov(y_{1i}, y_{2i}) &= \left(e^{\sigma^2} - 1\right) p_{1i}^{(O)} p_{2i}^{(O)} \\
&\approx \sigma^2 p_{1i}^{(O)} p_{2i}^{(O)},
\end{aligned}
$$

when $\sigma^2$ is small. For the US population, $p_{1i}^{(O)}$ and $p_{2i}^{(O)}$ were found to be constant over years (as confirmed by the SEER population data base). We then write $p_{ki}^{(O)} \equiv p_k^{(O)}$ for $i = s + 1, \ldots, m$, an estimate of which is given by $\hat{p}_k^{(O)} = \frac{n_k^{(O)}}{n_k}$, where $n_k = \sum_{i=s+1}^{m} \sum_{j=1}^{J} n_{kji}$ and $n^{(O)} = \sum_{i=s+1}^{m} \sum_{j=1}^{J} n_{ji}^{(O)}$. Hence, $Cov(y_{1i}, y_{2i}) \approx \sigma^2 \frac{(n^{(O)})^2}{n_1 n_2}$ for $i = s + 1, \ldots, m$. Inserting it back to (14) yields (8).

# Tables for Data Analysis

Table A.1: Comparison of Changes in Age-adjusted cancer mortality rates between California (1990-2004) and the US (1988-2002) for males. $APC_{us}$ and $APC_{ca}$ are the annual percent changes for the US and California respectively. $\sigma^2$ is the common (residual) variance in the Cancer Rate Regression Models (6) and (7).

| sites | $APC_{us}$ (SE) | $APC_{ca}$ (SE) | $\sigma^2$ | p-value (Z-test) | p-value (t-test) |
|---|---|---|---|---|---|
| 1 All Malignant Cancers | -1.14529 ( 0.08562 ) | -1.69304 ( 0.05924 ) | 0.01323 | 0.00000 | 0.00000 |
| 2 Oral Cavity and Pharynx | -2.54187 ( 0.14183 ) | -2.36549 ( 0.30387 ) | 0.04262 | 0.61088 | 0.62439 |
| 3 Lip | -5.04715 ( 0.93868 ) | -2.89026 ( 2.28470 ) | 0.31393 | 0.39824 | 0.41626 |
| 4 Tongue | -2.30819 ( 0.17872 ) | -1.55945 ( 0.48771 ) | 0.06602 | 0.16317 | 0.17961 |
| 5 Salivary Gland | -1.21958 ( 0.39018 ) | -2.69678 ( 0.88342 ) | 0.12274 | 0.13895 | 0.15445 |
| 6 Floor of Mouth | -8.74256 ( 0.54619 ) | -4.60891 ( 1.60742 ) | 0.21577 | 0.01850 | 0.02341 |
| 7 Gum and Other Mouth | -3.59892 ( 0.31909 ) | -4.57622 ( 0.45483 ) | 0.07062 | 0.08882 | 0.10152 |
| 8 Nasopharynx | -2.60014 ( 0.26453 ) | -2.52804 ( 0.67281 ) | 0.09188 | 0.92313 | 0.92602 |
| 9 Tonsil | -1.62000 ( 0.39344 ) | -0.51149 ( 0.97421 ) | 0.13354 | 0.30741 | 0.32600 |
| 10 Oropharynx | -1.02095 ( 0.33416 ) | -0.57017 ( 0.97536 ) | 0.13104 | 0.67232 | 0.68399 |
| 11 Hypopharynx | -5.54162 ( 0.33743 ) | -3.08101 ( 0.99334 ) | 0.13334 | 0.02327 | 0.02900 |
| 12 Other Oral Cavity and Pha | -1.57741 ( 0.31890 ) | -2.39696 ( 0.72259 ) | 0.10039 | 0.31547 | 0.33405 |
| 13 Digestive System | -1.00699 ( 0.03420 ) | -1.03053 ( 0.06669 ) | 0.00953 | 0.76122 | 0.76996 |
| 14 Esophagus | 0.74199 ( 0.04981 ) | 0.03682 ( 0.26946 ) | 0.03483 | 0.01279 | 0.01659 |
| 15 Stomach | -3.18066 ( 0.14032 ) | -2.74636 ( 0.15274 ) | 0.02636 | 0.04280 | 0.05126 |
| 16 Small Intestine | -0.71617 ( 0.35911 ) | -2.28962 ( 0.82727 ) | 0.11462 | 0.09145 | 0.10433 |
| 17 Colon and Rectum | -1.98524 ( 0.04743 ) | -2.30563 ( 0.10036 ) | 0.01411 | 0.00523 | 0.00721 |
| 18 Colon excluding Rectum | -2.18078 ( 0.05359 ) | -2.36600 ( 0.11226 ) | 0.01581 | 0.14974 | 0.16569 |
| 19 Rectum and Rectosigmoid J | -0.88380 ( 0.07602 ) | -2.00738 ( 0.31554 ) | 0.04125 | 0.00081 | 0.00127 |
| 20 Liver and Intrahepatic Bi | 2.70060 ( 0.16143 ) | 2.85857 ( 0.19285 ) | 0.03196 | 0.54342 | 0.55871 |
| 21 Liver | 2.31905 ( 0.13813 ) | 2.71266 ( 0.22187 ) | 0.03322 | 0.14513 | 0.16090 |
| 22 Intrahepatic Bile Duct | 4.89523 ( 0.40718 ) | 3.73589 ( 0.51502 ) | 0.08344 | 0.08758 | 0.10019 |
| 23 Gallbladder | -2.10738 ( 0.22285 ) | -1.53995 ( 0.48949 ) | 0.06836 | 0.30742 | 0.32600 |
| 24 Other Biliary | -3.39084 ( 0.23329 ) | -3.65275 ( 0.53714 ) | 0.07443 | 0.66526 | 0.67715 |
| 25 Pancreas | -0.29431 ( 0.05214 ) | -0.39426 ( 0.19401 ) | 0.02553 | 0.63031 | 0.64325 |
| 26 Retroperitoneum | -4.12158 ( 0.58644 ) | -2.68420 ( 1.11387 ) | 0.15999 | 0.26932 | 0.28778 |
| 27 Peritoneum, Omentum and M | -0.90453 ( 0.66323 ) | 4.60791 ( 1.42707 ) | 0.20001 | 0.00070 | 0.00111 |
| 28 Other Digestive Organs | 3.32072 ( 1.36339 ) | 3.54776 ( 2.10751 ) | 0.31902 | 0.93027 | 0.93289 |
| 29 Respiratory System | -1.60133 ( 0.08997 ) | -2.54261 ( 0.08343 ) | 0.01559 | 0.00000 | 0.00000 |
| 30 Nose, Nasal Cavity and Mi | -2.45896 ( 0.36930 ) | -3.71470 ( 1.15055 ) | 0.15358 | 0.31473 | 0.33332 |
| 31 Larynx | -1.89041 ( 0.15676 ) | -2.21785 ( 0.41215 ) | 0.05604 | 0.47254 | 0.48938 |
| 32 Lung and Bronchus | -1.57318 ( 0.08800 ) | -2.53538 ( 0.08350 ) | 0.01542 | 0.00000 | 0.00000 |
| 33 Pleura | -4.76359 ( 0.69374 ) | -5.63373 ( 0.83073 ) | 0.13756 | 0.43671 | 0.45418 |
| 34 Trachea, Mediastinum and | -4.38175 ( 0.56495 ) | -4.21490 ( 1.43052 ) | 0.19548 | 0.91642 | 0.91956 |
| 35 Bones and Joints | -0.63050 ( 0.22877 ) | 0.50710 ( 0.56841 ) | 0.07787 | 0.07247 | 0.08391 |
| 36 Soft Tissue including Hea | -0.19405 ( 0.32414 ) | -1.86636 ( 0.56160 ) | 0.08241 | 0.01260 | 0.01635 |
| 37 Skin excluding Basal and | -0.20217 ( 0.08928 ) | -0.92676 ( 0.29656 ) | 0.03936 | 0.02362 | 0.02940 |
| 38 Melanoma of the Skin | 0.22384 ( 0.14864 ) | -1.14511 ( 0.35312 ) | 0.04870 | 0.00055 | 0.00088 |
| 39 Other Non-Epithelial Skin | -1.28166 ( 0.33594 ) | -0.34737 ( 0.61161 ) | 0.08869 | 0.19523 | 0.21260 |
| 40 Breast | 0.67233 ( 0.44829 ) | -0.35521 ( 1.08577 ) | 0.14930 | 0.39743 | 0.41545 |
| 41 Male Genital System | -2.17218 ( 0.34136 ) | -3.35652 ( 0.22018 ) | 0.05163 | 0.00479 | 0.00664 |
| 42 Prostate | -2.18583 ( 0.34814 ) | -3.40124 ( 0.21988 ) | 0.05233 | 0.00430 | 0.00600 |
| 43 Testis | -1.32447 ( 0.42559 ) | -1.31581 ( 0.90241 ) | 0.12681 | 0.99330 | 0.99355 |
| 44 Penis | -1.45769 ( 0.29383 ) | 0.44102 ( 1.74277 ) | 0.22463 | 0.29867 | 0.31724 |
| 45 Urinary System | -0.15872 ( 0.06568 ) | -0.46994 ( 0.10662 ) | 0.01592 | 0.01621 | 0.02069 |
| 46 Urinary Bladder | -0.39827 ( 0.07156 ) | -0.66431 ( 0.19502 ) | 0.02640 | 0.21539 | 0.23318 |
| 47 Kidney and Renal Pelvis | 0.09973 ( 0.10326 ) | -0.20732 ( 0.16686 ) | 0.02494 | 0.13010 | 0.14520 |
| 48 Ureter | -1.81235 ( 0.36063 ) | -3.05155 ( 1.85441 ) | 0.24011 | 0.52571 | 0.54142 |
| 49 Other Urinary Organs | 4.67508 ( 1.32642 ) | 4.13840 ( 2.59886 ) | 0.37084 | 0.85877 | 0.86404 |
| 50 Eye and Orbit | -2.59686 ( 0.51488 ) | 2.06275 ( 1.37227 ) | 0.18629 | 0.00210 | 0.00308 |
| 51 Brain and Other Nervous S | -0.59408 ( 0.07447 ) | -0.76199 ( 0.27160 ) | 0.03579 | 0.56409 | 0.57886 |
| 52 Endocrine System | 0.46159 ( 0.22153 ) | 0.94815 ( 0.49324 ) | 0.06872 | 0.38402 | 0.40219 |
| 53 Thyroid | 1.31809 ( 0.32622 ) | 2.55398 ( 0.66252 ) | 0.09386 | 0.10545 | 0.11923 |
| 54 Other Endocrine including | -0.51444 ( 0.26401 ) | -0.90111 ( 0.73065 ) | 0.09874 | 0.63017 | 0.64311 |
| 55 Lymphoma | 0.04612 ( 0.25417 ) | -0.90577 ( 0.31011 ) | 0.05096 | 0.02164 | 0.02710 |
| 56 Hodgkin Lymphoma | -3.77024 ( 0.29593 ) | -2.91687 ( 0.50752 ) | 0.07467 | 0.15996 | 0.17630 |
| 57 Non-Hodgkin Lymphoma | 0.32123 ( 0.28472 ) | -0.76903 ( 0.33308 ) | 0.05569 | 0.01608 | 0.02054 |
| 58 Myeloma | 0.00492 ( 0.15560 ) | -0.63384 ( 0.28698 ) | 0.04149 | 0.05837 | 0.06852 |
| 59 Leukemia | -0.41114 ( 0.07777 ) | -1.16978 ( 0.18611 ) | 0.02564 | 0.00027 | 0.00046 |
| 60 Lymphocytic Leukemia | -0.80381 ( 0.18393 ) | -1.43414 ( 0.38148 ) | 0.05383 | 0.14991 | 0.16587 |
| 61 Acute Lymphocytic Leukemi | -1.88621 ( 0.15635 ) | -0.61542 ( 0.70555 ) | 0.09185 | 0.08891 | 0.10162 |
| 62 Chronic Lymphocytic Leuke | -0.22980 ( 0.25973 ) | -1.58381 ( 0.39517 ) | 0.06010 | 0.00561 | 0.00769 |
| 63 Other Lymphocytic Leukemi | -3.11597 ( 0.26524 ) | -2.84172 ( 1.08602 ) | 0.14209 | 0.81241 | 0.81935 |
| 64 Myeloid and Monocytic Leu | 0.38178 ( 0.10812 ) | -0.34155 ( 0.26334 ) | 0.03618 | 0.01397 | 0.01801 |
| 65 Acute Myeloid Leukemia | 1.85239 ( 0.13780 ) | 1.27926 ( 0.25779 ) | 0.03715 | 0.05786 | 0.06795 |
| 66 Acute Monocytic Leukemia | -5.87966 ( 0.33270 ) | -5.81060 ( 1.38431 ) | 0.18095 | 0.96258 | 0.96398 |
| 67 Chronic Myeloid Leukemia | -4.54161 ( 0.69213 ) | -7.48499 ( 0.97162 ) | 0.15162 | 0.01699 | 0.02162 |
| 68 Other Myeloid/Monocytic L | 3.25551 ( 1.82360 ) | 4.63431 ( 2.02953 ) | 0.34678 | 0.62494 | 0.63804 |
| 69 Other Leukemia | -1.26579 ( 0.15121 ) | -2.35890 ( 0.35968 ) | 0.04959 | 0.00672 | 0.00910 |
| 70 Other Acute Leukemia | -2.69076 ( 0.21664 ) | -4.26489 ( 0.44006 ) | 0.06234 | 0.00191 | 0.00281 |
| 71 Aleukemic, Subleukemic an | 0.19461 ( 0.25052 ) | -0.13660 ( 0.44085 ) | 0.06445 | 0.52745 | 0.54313 |
| 72 Miscellaneous Malignant C | -0.06793 ( 0.38400 ) | -0.04692 ( 0.33004 ) | 0.06435 | 0.96798 | 0.96919 |

Table A.2: Comparison of Changes in Age-adjusted cancer mortality rates between California (1990-2004) and the US (1988-2002) for females.

| sites | $APC_{us}$ (SE) | $APC_{ca}$ (SE) | $\sigma^2$ | p-value (Z-test) | p-value (t-test) |
|---|---|---|---|---|---|
| 1 All Malignant Cancers | -0.4967 ( 0.06367 ) | -1.1995 ( 0.07756 ) | 0.01275 | 0.00000 | 0.00000 |
| 2 Oral Cavity and Pharynx | -2.3100 ( 0.09949 ) | -2.6478 ( 0.32027 ) | 0.04262 | 0.33065 | 0.34845 |
| 3 Tongue | -1.7552 ( 0.20731 ) | -1.9406 ( 0.56414 ) | 0.07639 | 0.76566 | 0.77390 |
| 4 Salivary Gland | -1.5547 ( 0.28108 ) | -1.6580 ( 1.53080 ) | 0.19781 | 0.94885 | 0.95070 |
| 5 Floor of Mouth | -8.5738 ( 0.59746 ) | -10.3092 ( 1.46769 ) | 0.20140 | 0.29016 | 0.30797 |
| 6 Gum and Other Mouth | -2.4059 ( 0.33005 ) | -3.5480 ( 0.89911 ) | 0.12173 | 0.24941 | 0.26696 |
| 7 Nasopharynx | -1.8134 ( 0.28784 ) | -1.7402 ( 1.07235 ) | 0.14112 | 0.94922 | 0.95105 |
| 8 Tonsil | -2.7779 ( 0.45408 ) | -3.0276 ( 1.09826 ) | 0.15105 | 0.83921 | 0.84495 |
| 9 Oropharynx | -0.5663 ( 0.55131 ) | -1.4113 ( 1.74661 ) | 0.23279 | 0.65585 | 0.66754 |
| 10 Hypopharynx | -5.0032 ( 0.69969 ) | -2.0945 ( 2.10828 ) | 0.28233 | 0.20595 | 0.22284 |
| 11 Other Oral Cavity and Pha | -2.4089 ( 0.33467 ) | -2.6857 ( 0.96831 ) | 0.13021 | 0.79414 | 0.80143 |
| 12 Digestive System | -0.9999 ( 0.03521 ) | -1.0174 ( 0.08568 ) | 0.01177 | 0.85554 | 0.86071 |
| 13 Esophagus | -0.1298 ( 0.09611 ) | -0.4936 ( 0.44040 ) | 0.05729 | 0.43558 | 0.45237 |
| 14 Stomach | -2.4744 ( 0.08503 ) | -2.2361 ( 0.26341 ) | 0.03518 | 0.40560 | 0.42281 |
| 15 Small Intestine | -0.5498 ( 0.28495 ) | -1.1258 ( 0.59587 ) | 0.08395 | 0.39957 | 0.41685 |
| 16 Colon and Rectum | -1.7837 ( 0.03854 ) | -2.1068 ( 0.12624 ) | 0.01678 | 0.01804 | 0.02266 |
| 17 Colon excluding Rectum | -1.9450 ( 0.04393 ) | -2.0724 ( 0.12469 ) | 0.01680 | 0.35203 | 0.36973 |
| 18 Rectum and Rectosigmoid J | -0.6662 ( 0.09959 ) | -2.3023 ( 0.35657 ) | 0.04705 | 0.00002 | 0.00004 |
| 19 Anus, Anal Canal and Anor | 0.9835 ( 0.36923 ) | 1.4404 ( 0.50560 ) | 0.07957 | 0.48091 | 0.49693 |
| 20 Liver and Intrahepatic Bi | 2.1121 ( 0.22646 ) | 2.8416 ( 0.25130 ) | 0.04299 | 0.03727 | 0.04471 |
| 21 Liver | 1.0356 ( 0.23798 ) | 2.2415 ( 0.30421 ) | 0.04909 | 0.00256 | 0.00366 |
| 22 Intrahepatic Bile Duct | 5.3993 ( 0.32153 ) | 4.6627 ( 0.45037 ) | 0.07033 | 0.19849 | 0.21523 |
| 23 Gallbladder | -2.3486 ( 0.12983 ) | -1.7201 ( 0.32906 ) | 0.04496 | 0.08616 | 0.09815 |
| 24 Other Biliary | -3.3533 ( 0.27122 ) | -3.2242 ( 0.98554 ) | 0.12992 | 0.90295 | 0.90645 |
| 25 Pancreas | 0.0459 ( 0.06003 ) | -0.3244 ( 0.14192 ) | 0.01958 | 0.02029 | 0.02529 |
| 26 Retroperitoneum | -3.4767 ( 0.42605 ) | -2.4084 ( 2.15884 ) | 0.27968 | 0.63910 | 0.65128 |
| 27 Peritoneum, Omentum and M | 10.6773 ( 0.50603 ) | 11.5268 ( 1.02266 ) | 0.14502 | 0.47208 | 0.48827 |
| 28 Other Digestive Organs | 2.9486 ( 1.25248 ) | 4.0027 ( 1.31273 ) | 0.23060 | 0.57471 | 0.58863 |
| 29 Respiratory System | 1.1074 ( 0.14308 ) | -0.8987 ( 0.13886 ) | 0.02534 | 0.00000 | 0.00000 |
| 30 Nose, Nasal Cavity and Mi | -2.7924 ( 0.50534 ) | -1.6394 ( 0.87092 ) | 0.12798 | 0.26870 | 0.28641 |
| 31 Larynx | -0.9160 ( 0.31894 ) | -3.1695 ( 0.99551 ) | 0.13286 | 0.03732 | 0.04476 |
| 32 Lung and Bronchus | 1.1684 ( 0.14263 ) | -0.8594 ( 0.13851 ) | 0.02527 | 0.00000 | 0.00000 |
| 33 Trachea, Mediastinum and | -4.1474 ( 0.51934 ) | -3.7029 ( 1.76282 ) | 0.23357 | 0.81526 | 0.82183 |
| 34 Bones and Joints | -0.3413 ( 0.23101 ) | -0.0312 ( 0.59852 ) | 0.08154 | 0.64064 | 0.65278 |
| 35 Soft Tissue including Hea | -0.3137 ( 0.50341 ) | -2.3136 ( 0.59262 ) | 0.09883 | 0.01298 | 0.01665 |
| 36 Skin excluding Basal and | -0.6894 ( 0.10824 ) | -1.7968 ( 0.25599 ) | 0.03532 | 0.00012 | 0.00021 |
| 37 Melanoma of the Skin | -0.6677 ( 0.13053 ) | -2.2282 ( 0.27928 ) | 0.03918 | 0.00000 | 0.00000 |
| 38 Other Non-Epithelial Skin | -0.7827 ( 0.28864 ) | 0.1701 ( 0.86128 ) | 0.11545 | 0.31098 | 0.32882 |
| 39 Breast | -2.1080 ( 0.09911 ) | -2.4052 ( 0.12810 ) | 0.02059 | 0.07627 | 0.08752 |
| 40 Female Genital System | -0.7757 ( 0.06950 ) | -0.8266 ( 0.13887 ) | 0.01974 | 0.75116 | 0.75988 |
| 41 Cervix Uteri | -2.5485 ( 0.16808 ) | -2.7864 ( 0.30711 ) | 0.04450 | 0.51164 | 0.52704 |
| 42 Corpus and Uterus, NOS | -0.3500 ( 0.08836 ) | -0.4197 ( 0.22008 ) | 0.03014 | 0.77630 | 0.78419 |
| 43 Corpus Uteri | -0.9847 ( 0.10858 ) | -1.6515 ( 0.29045 ) | 0.03941 | 0.03779 | 0.04529 |
| 44 Uterus, NOS | 0.2709 ( 0.20775 ) | 0.6246 ( 0.30305 ) | 0.04670 | 0.35242 | 0.37012 |
| 45 Ovary | -0.4497 ( 0.10827 ) | -0.3897 ( 0.19006 ) | 0.02780 | 0.79108 | 0.79847 |
| 46 Vagina | -1.4893 ( 0.30881 ) | -0.0450 ( 1.12439 ) | 0.14820 | 0.23152 | 0.24885 |
| 47 Vulva | 0.3715 ( 0.21886 ) | -0.4401 ( 0.74207 ) | 0.09833 | 0.31091 | 0.32875 |
| 48 Other Female Genital Orga | 1.0285 ( 0.89011 ) | -2.8798 ( 1.05902 ) | 0.17583 | 0.00636 | 0.00854 |
| 49 Urinary System | -0.1898 ( 0.10186 ) | -0.3737 ( 0.25074 ) | 0.03440 | 0.51161 | 0.52701 |
| 50 Urinary Bladder | -0.3297 ( 0.12862 ) | -0.2421 ( 0.36919 ) | 0.04969 | 0.82869 | 0.83480 |
| 51 Kidney and Renal Pelvis | -0.0759 ( 0.17167 ) | -0.5114 ( 0.23849 ) | 0.03735 | 0.15231 | 0.16770 |
| 52 Ureter | -1.1593 ( 0.53854 ) | -0.4319 ( 1.29441 ) | 0.17819 | 0.61627 | 0.62909 |
| 53 Other Urinary Organs | 0.97551 ( 0.98382 ) | 0.51241 ( 1.65580 ) | 0.24479 | 0.81634 | 0.82288 |
| 54 Eye and Orbit | -2.29495 ( 0.51385 ) | -1.09212 ( 1.83285 ) | 0.24193 | 0.54163 | 0.55636 |
| 55 Brain and Other Nervous S | -0.60288 ( 0.15861 ) | -0.78950 ( 0.24186 ) | 0.03676 | 0.53314 | 0.54807 |
| 56 Endocrine System | -0.03612 ( 0.17458 ) | -0.82845 ( 0.57794 ) | 0.07673 | 0.20493 | 0.22180 |
| 57 Thyroid | 0.19503 ( 0.22603 ) | -0.11015 ( 0.52083 ) | 0.07216 | 0.60363 | 0.61680 |
| 58 Other Endocrine including | -0.40304 ( 0.20725 ) | -2.21315 ( 0.99215 ) | 0.12882 | 0.08453 | 0.09640 |
| 59 Lymphoma | 0.01028 ( 0.29109 ) | -1.31900 ( 0.40750 ) | 0.06365 | 0.01035 | 0.01347 |
| 60 Hodgkin Lymphoma | -2.55202 ( 0.28254 ) | -1.53544 ( 0.79565 ) | 0.10731 | 0.24484 | 0.26234 |
| 61 Non-Hodgkin Lymphoma | 0.18531 ( 0.30940 ) | -1.29605 ( 0.41882 ) | 0.06618 | 0.00600 | 0.00809 |
| 62 Myeloma | 0.19767 ( 0.15648 ) | -0.87526 ( 0.33092 ) | 0.04652 | 0.00464 | 0.00636 |
| 63 Leukemia | -0.47456 ( 0.09093 ) | -1.19624 ( 0.15836 ) | 0.02321 | 0.00013 | 0.00023 |
| 64 Lymphocytic Leukemia | -0.71969 ( 0.22034 ) | -1.55322 ( 0.35934 ) | 0.05357 | 0.05613 | 0.06563 |
| 65 Acute Lymphocytic Leukemi | -1.33321 ( 0.27870 ) | -0.40742 ( 0.54511 ) | 0.07781 | 0.14412 | 0.15920 |
| 66 Chronic Lymphocytic Leuke | -0.11289 ( 0.28829 ) | -1.78207 ( 0.46202 ) | 0.06922 | 0.00307 | 0.00433 |
| 67 Other Lymphocytic Leukemi | -4.18929 ( 0.22281 ) | -5.55778 ( 1.60359 ) | 0.20577 | 0.41424 | 0.43134 |
| 68 Myeloid and Monocytic Leu | 0.19180 ( 0.09133 ) | -0.37999 ( 0.29947 ) | 0.03979 | 0.07772 | 0.08909 |
| 69 Acute Myeloid Leukemia | 1.60735 ( 0.11928 ) | 1.22486 ( 0.30773 ) | 0.04195 | 0.26296 | 0.28063 |
| 70 Acute Monocytic Leukemia | -7.75354 ( 0.82563 ) | -9.06558 ( 3.75654 ) | 0.48884 | 0.74178 | 0.75081 |
| 71 Chronic Myeloid Leukemia | -4.29532 ( 0.63553 ) | -7.22987 ( 0.99391 ) | 0.14994 | 0.01628 | 0.02057 |
| 72 Other Myeloid/Monocytic L | 1.78219 ( 1.33592 ) | 5.14116 ( 1.49323 ) | 0.25465 | 0.10538 | 0.11859 |
| 73 Other Leukemia | -1.37872 ( 0.11477 ) | -2.39493 ( 0.25479 ) | 0.03552 | 0.00044 | 0.00071 |
| 74 Other Acute Leukemia | -3.25827 ( 0.25236 ) | -4.90300 ( 0.39909 ) | 0.06001 | 0.00077 | 0.00118 |
| 75 Aleukemic, Subleukemic an | 0.67118 ( 0.13108 ) | 0.46868 ( 0.40711 ) | 0.05436 | 0.64745 | 0.65939 |
| 76 Miscellaneous Malignant C | -0.11415 ( 0.29840 ) | -0.54259 ( 0.28608 ) | 0.05254 | 0.31678 | 0.33461 |

Table B.1: Comparison of Changes in Age-adjusted cancer mortality rates between California (1990-2004) and the US (1980-1994) for males. $APC_{us}$ and $APC_{ca}$ are the annual percent changes for the US and California respectively. $\sigma^2$ is the common (residual) variance in the Cancer Rate Regression Models (6) and (7).

| sites | $APC_{us}$ (SE) | $APC_{ca}$ (SE) | $\sigma^2$ | p-value (Z-test) | p-value (t-test) |
|---|---|---|---|---|---|
| 1 All Malignant Cancers | 0.13395 ( 0.05004 ) | -1.69304 ( 0.05924 ) | 0.00986 | 0.00000 | 0.00000 |
| 2 Oral Cavity and Pharynx | -2.09739 ( 0.12845 ) | -2.36549 ( 0.30387 ) | 0.04193 | 0.46045 | 0.44932 |
| 3 Lip | -6.33929 ( 0.87302 ) | -2.89026 ( 2.28470 ) | 0.31086 | 0.20026 | 0.18925 |
| 4 Tongue | -2.22245 ( 0.15640 ) | -1.55945 ( 0.48771 ) | 0.06510 | 0.23971 | 0.22816 |
| 5 Salivary Gland | -0.45950 ( 0.39389 ) | -2.69678 ( 0.88342 ) | 0.12294 | 0.03566 | 0.03129 |
| 6 Floor of Mouth | -6.78522 ( 0.25603 ) | -4.60891 ( 1.60742 ) | 0.20688 | 0.22460 | 0.21323 |
| 7 Gum and Other Mouth | -3.05074 ( 0.23408 ) | -4.57622 ( 0.45483 ) | 0.06501 | 0.00676 | 0.00550 |
| 8 Nasopharynx | -1.06262 ( 0.31636 ) | -2.52804 ( 0.67281 ) | 0.09449 | 0.07342 | 0.06652 |
| 9 Tonsil | -2.82023 ( 0.38469 ) | -0.51149 ( 0.97421 ) | 0.13312 | 0.04528 | 0.04017 |
| 10 Oropharynx | 0.19700 ( 0.37195 ) | -0.57017 ( 0.97536 ) | 0.13267 | 0.50446 | 0.49387 |
| 11 Hypopharynx | -4.45059 ( 0.48011 ) | -3.08101 ( 0.99334 ) | 0.14022 | 0.25954 | 0.24782 |
| 12 Other Oral Cavity and Pha | -0.22622 ( 0.32481 ) | -2.39696 ( 0.72259 ) | 0.10069 | 0.01282 | 0.01075 |
| 13 Digestive System | -0.68619 ( 0.03274 ) | -1.03053 ( 0.06669 ) | 0.00944 | 0.00003 | 0.00002 |
| 14 Esophagus | 1.02825 ( 0.06092 ) | 0.03682 ( 0.26946 ) | 0.03511 | 0.00112 | 0.00083 |
| 15 Stomach | -2.03063 ( 0.11464 ) | -2.74636 ( 0.15274 ) | 0.02427 | 0.00066 | 0.00048 |
| 16 Small Intestine | 0.73969 ( 0.33065 ) | -2.28962 ( 0.82727 ) | 0.11323 | 0.00201 | 0.00155 |
| 17 Colon and Rectum | -1.24173 ( 0.08356 ) | -2.30563 ( 0.10036 ) | 0.01660 | 0.00000 | 0.00000 |
| 18 Colon excluding Rectum | -0.98791 ( 0.12656 ) | -2.36600 ( 0.11226 ) | 0.02150 | 0.00000 | 0.00000 |
| 19 Rectum and Rectosigmoid J | -2.60904 ( 0.19370 ) | -2.00738 ( 0.31554 ) | 0.04706 | 0.13996 | 0.13033 |
| 20 Liver and Intrahepatic Bi | 3.05617 ( 0.12756 ) | 2.85857 ( 0.19285 ) | 0.02939 | 0.43764 | 0.42628 |
| 21 Liver | 2.31517 ( 0.14534 ) | 2.71266 ( 0.22187 ) | 0.03371 | 0.17348 | 0.16298 |
| 22 Intrahepatic Bile Duct | 9.01465 ( 0.30636 ) | 3.73589 ( 0.51502 ) | 0.07616 | 0.00000 | 0.00000 |
| 23 Gallbladder | -2.64656 ( 0.24536 ) | -1.53995 ( 0.48949 ) | 0.06959 | 0.06641 | 0.05990 |
| 24 Other Biliary | -2.81449 ( 0.18682 ) | -3.65275 ( 0.53714 ) | 0.07228 | 0.18065 | 0.17000 |
| 25 Pancreas | -0.40317 ( 0.05874 ) | -0.39426 ( 0.19401 ) | 0.02576 | 0.96816 | 0.96737 |
| 26 Retroperitoneum | -5.63701 ( 0.49823 ) | -2.68420 ( 1.11387 ) | 0.15509 | 0.02796 | 0.02427 |
| 27 Peritoneum, Omentum and M | -0.31829 ( 0.89814 ) | 4.60791 ( 1.42707 ) | 0.21431 | 0.00797 | 0.00653 |
| 28 Other Digestive Organs | -4.00785 ( 0.30807 ) | 3.54776 ( 2.10751 ) | 0.27071 | 0.00127 | 0.00096 |
| 29 Respiratory System | 0.15618 ( 0.10760 ) | -2.54261 ( 0.08343 ) | 0.01730 | 0.00000 | 0.00000 |
| 30 Nose, Nasal Cavity and Mi | -2.16152 ( 0.27623 ) | -3.71470 ( 1.15055 ) | 0.15039 | 0.23318 | 0.22170 |
| 31 Larynx | -0.63743 ( 0.09716 ) | -2.21785 ( 0.41215 ) | 0.05382 | 0.00070 | 0.00051 |
| 32 Lung and Bronchus | 0.20164 ( 0.11152 ) | -2.53538 ( 0.08350 ) | 0.01771 | 0.00000 | 0.00000 |
| 33 Pleura | 0.49245 ( 0.40844 ) | -5.63373 ( 0.83073 ) | 0.11766 | 0.00000 | 0.00000 |
| 34 Trachea, Mediastinum and | -3.84059 ( 0.35496 ) | -4.21490 ( 1.43052 ) | 0.18733 | 0.81758 | 0.81310 |
| 35 Bones and Joints | -1.03432 ( 0.41720 ) | 0.50710 ( 0.56841 ) | 0.08961 | 0.04708 | 0.04183 |
| 36 Soft Tissue including Hea | 1.03457 ( 0.10806 ) | -1.86636 ( 0.56160 ) | 0.07269 | 0.00000 | 0.00000 |
| 37 Skin excluding Basal and | 1.48499 ( 0.21683 ) | -0.92676 ( 0.29656 ) | 0.04669 | 0.00000 | 0.00000 |
| 38 Melanoma of the Skin | 1.65189 ( 0.13590 ) | -1.14511 ( 0.35312 ) | 0.04809 | 0.00000 | 0.00000 |
| 39 Other Non-Epithelial Skin | 1.10131 ( 0.47089 ) | -0.34737 ( 0.61161 ) | 0.09811 | 0.08827 | 0.08060 |
| 40 Breast | 0.40317 ( 0.47453 ) | -0.35521 ( 1.08577 ) | 0.15060 | 0.56104 | 0.55129 |
| 41 Male Genital System | 1.38340 ( 0.11120 ) | -3.35652 ( 0.22018 ) | 0.03135 | 0.00000 | 0.00000 |
| 42 Prostate | 1.45612 ( 0.11206 ) | -3.40124 ( 0.21988 ) | 0.03137 | 0.00000 | 0.00000 |
| 43 Testis | -3.17493 ( 0.29347 ) | -1.31581 ( 0.90241 ) | 0.12061 | 0.07517 | 0.06817 |
| 44 Penis | -2.15477 ( 0.41309 ) | 0.44102 ( 1.74277 ) | 0.22764 | 0.18806 | 0.17726 |
| 45 Urinary System | -0.28045 ( 0.08727 ) | -0.46994 ( 0.10662 ) | 0.01751 | 0.21163 | 0.20044 |
| 46 Urinary Bladder | -1.12675 ( 0.16025 ) | -0.66431 ( 0.19502 ) | 0.03208 | 0.09612 | 0.08809 |
| 47 Kidney and Renal Pelvis | 1.00260 ( 0.09281 ) | -0.20732 ( 0.16686 ) | 0.02427 | 0.00000 | 0.00000 |
| 48 Ureter | -1.06126 ( 0.43471 ) | -3.05155 ( 1.85441 ) | 0.24208 | 0.34258 | 0.33066 |
| 49 Other Urinary Organs | -2.59350 ( 0.62897 ) | 4.13840 ( 2.59886 ) | 0.33984 | 0.02222 | 0.01909 |
| 50 Eye and Orbit | -1.74425 ( 0.43342 ) | 2.06275 ( 1.37227 ) | 0.18291 | 0.01627 | 0.01379 |
| 51 Brain and Other Nervous S | 0.85214 ( 0.10548 ) | -0.76199 ( 0.27160 ) | 0.03703 | 0.00000 | 0.00000 |
| 52 Endocrine System | -0.08767 ( 0.18546 ) | 0.94815 ( 0.49324 ) | 0.06698 | 0.07421 | 0.06726 |
| 53 Thyroid | 0.06835 ( 0.26616 ) | 2.55398 ( 0.66252 ) | 0.09075 | 0.00157 | 0.00119 |
| 54 Other Endocrine including | -0.24578 ( 0.25539 ) | -0.90111 ( 0.73065 ) | 0.09837 | 0.44189 | 0.43057 |
| 55 Lymphoma | 2.05381 ( 0.10812 ) | -0.90577 ( 0.31011 ) | 0.04174 | 0.00000 | 0.00000 |
| 56 Hodgkin Lymphoma | -3.66815 ( 0.28420 ) | -2.91687 ( 0.50752 ) | 0.07393 | 0.24076 | 0.22921 |
| 57 Non-Hodgkin Lymphoma | 2.66716 ( 0.12986 ) | -0.76903 ( 0.33308 ) | 0.04544 | 0.00000 | 0.00000 |
| 58 Myeloma | 1.44945 ( 0.08579 ) | -0.63384 ( 0.28698 ) | 0.03807 | 0.00000 | 0.00000 |
| 59 Leukemia | -0.31509 ( 0.08696 ) | -1.16978 ( 0.18611 ) | 0.02611 | 0.00016 | 0.00011 |
| 60 Lymphocytic Leukemia | -0.06745 ( 0.15905 ) | -1.43414 ( 0.38148 ) | 0.05253 | 0.00267 | 0.00208 |
| 61 Acute Lymphocytic Leukemi | -1.07408 ( 0.23533 ) | -0.61542 ( 0.70555 ) | 0.09453 | 0.57541 | 0.56590 |
| 62 Chronic Lymphocytic Leuke | 0.88433 ( 0.25397 ) | -1.58381 ( 0.39517 ) | 0.05970 | 0.00000 | 0.00000 |
| 63 Other Lymphocytic Leukemi | -4.22759 ( 0.26089 ) | -2.84172 ( 1.08602 ) | 0.14196 | 0.25976 | 0.24804 |
| 64 Myeloid and Monocytic Leu | -1.23576 ( 0.19512 ) | -0.34155 ( 0.26334 ) | 0.04166 | 0.01321 | 0.01109 |
| 65 Acute Myeloid Leukemia | -0.78771 ( 0.25898 ) | 1.27926 ( 0.25779 ) | 0.04644 | 0.00000 | 0.00000 |
| 66 Acute Monocytic Leukemia | -5.35300 ( 0.52077 ) | -5.81060 ( 1.38431 ) | 0.18798 | 0.77870 | 0.77332 |
| 67 Chronic Myeloid Leukemia | -0.58622 ( 0.16262 ) | -7.48499 ( 0.97162 ) | 0.12521 | 0.00000 | 0.00000 |
| 68 Other Myeloid/Monocytic L | -7.31063 ( 0.43678 ) | 4.63431 ( 2.02953 ) | 0.26385 | 0.00000 | 0.00000 |
| 69 Other Leukemia | 1.00458 ( 0.22000 ) | -2.35890 ( 0.35968 ) | 0.05359 | 0.00000 | 0.00000 |
| 70 Other Acute Leukemia | 1.57263 ( 0.29680 ) | -4.26489 ( 0.44006 ) | 0.06746 | 0.00000 | 0.00000 |
| 71 Aleukemic, Subleukemic an | 0.37945 ( 0.25853 ) | -0.13660 ( 0.44085 ) | 0.06496 | 0.35908 | 0.34720 |
| 72 Miscellaneous Malignant C | -0.10363 ( 0.23525 ) | -0.04692 ( 0.33004 ) | 0.05151 | 0.89888 | 0.89637 |

Table B.2: Comparison of Changes in Age-adjusted cancer mortality rates between California (1990-2004) and the US (1980-1994) for females.

| sites | $APC_{us}$ (SE) | $APC_{ca}$ (SE) | $\sigma^2$ | p-value (Z-test) | p-value (t-test) |
|---|---|---|---|---|---|
| 1 All Malignant Cancers | 0.40400 ( 0.03737 ) | -1.1995 ( 0.07756 ) | 0.01094 | 0.00000 | 0.00000 |
| 2 Oral Cavity and Pharynx | -1.35357 ( 0.09378 ) | -2.6478 ( 0.32027 ) | 0.04241 | 0.00042 | 0.00031 |
| 3 Tongue | -1.33078 ( 0.19952 ) | -1.9406 ( 0.56414 ) | 0.07605 | 0.35408 | 0.34273 |
| 4 Salivary Gland | -0.79551 ( 0.27112 ) | -1.6580 ( 1.53080 ) | 0.19759 | 0.61392 | 0.60550 |
| 5 Floor of Mouth | -4.07290 ( 0.32488 ) | -10.3092 ( 1.46769 ) | 0.19106 | 0.00016 | 0.00011 |
| 6 Gum and Other Mouth | -1.43953 ( 0.25762 ) | -3.5480 ( 0.89911 ) | 0.11887 | 0.04038 | 0.03585 |
| 7 Nasopharynx | -0.88013 ( 0.29470 ) | -1.7402 ( 1.07235 ) | 0.14135 | 0.48192 | 0.47154 |
| 8 Tonsil | -3.36049 ( 0.37876 ) | -3.0276 ( 1.09826 ) | 0.14765 | 0.79443 | 0.78964 |
| 9 Oropharynx | 0.53419 ( 0.70546 ) | -1.4113 ( 1.74661 ) | 0.23941 | 0.34767 | 0.33631 |
| 10 Hypopharynx | -3.84815 ( 0.54731 ) | -2.0945 ( 2.10828 ) | 0.27684 | 0.46413 | 0.45355 |
| 11 Other Oral Cavity and Pha | -0.08469 ( 0.37118 ) | -2.6857 ( 0.96831 ) | 0.13180 | 0.02257 | 0.01954 |
| 12 Digestive System | -1.05530 ( 0.04113 ) | -1.0174 ( 0.08568 ) | 0.01208 | 0.71694 | 0.71049 |
| 13 Esophagus | -0.02695 ( 0.10507 ) | -0.4936 ( 0.44040 ) | 0.05755 | 0.34862 | 0.33726 |
| 14 Stomach | -2.19679 ( 0.11168 ) | -2.2361 ( 0.26341 ) | 0.03636 | 0.90055 | 0.89819 |
| 15 Small Intestine | 0.49738 ( 0.23347 ) | -1.1258 ( 0.59587 ) | 0.08134 | 0.02110 | 0.01822 |
| 16 Colon and Rectum | -1.71438 ( 0.06537 ) | -2.1068 ( 0.12624 ) | 0.01807 | 0.01207 | 0.01017 |
| 17 Colon excluding Rectum | -1.56494 ( 0.09316 ) | -2.0724 ( 0.12469 ) | 0.01978 | 0.00303 | 0.00240 |
| 18 Rectum and Rectosigmoid J | -2.70471 ( 0.26231 ) | -2.3023 ( 0.35657 ) | 0.05626 | 0.40840 | 0.39733 |
| 19 Anus, Anal Canal and Anor | 1.54210 ( 0.37044 ) | 1.4404 ( 0.50560 ) | 0.07966 | 0.88270 | 0.87993 |
| 20 Liver and Intrahepatic Bi | 2.41870 ( 0.16904 ) | 2.8416 ( 0.25130 ) | 0.03849 | 0.20425 | 0.19367 |
| 21 Liver | 1.17094 ( 0.20476 ) | 2.2415 ( 0.30421 ) | 0.04661 | 0.00794 | 0.00657 |
| 22 Intrahepatic Bile Duct | 8.52092 ( 0.33496 ) | 4.6627 ( 0.45037 ) | 0.07134 | 0.00000 | 0.00000 |
| 23 Gallbladder | -2.94543 ( 0.15996 ) | -1.7201 ( 0.32906 ) | 0.04650 | 0.00233 | 0.00182 |
| 24 Other Biliary | -3.13665 ( 0.17143 ) | -3.2242 ( 0.98554 ) | 0.12714 | 0.93656 | 0.93505 |
| 25 Pancreas | 0.34744 ( 0.06179 ) | -0.3244 ( 0.14192 ) | 0.01967 | 0.00008 | 0.00005 |
| 26 Retroperitoneum | -4.32978 ( 0.40012 ) | -2.4084 ( 2.15884 ) | 0.27906 | 0.42619 | 0.41525 |
| 27 Peritoneum, Omentum and M | 4.51348 ( 0.89926 ) | 11.5268 ( 1.02266 ) | 0.17308 | 0.00000 | 0.00000 |
| 28 Other Digestive Organs | -4.25304 ( 0.33605 ) | 4.0027 ( 1.31273 ) | 0.17223 | 0.00000 | 0.00000 |
| 29 Respiratory System | 3.61472 ( 0.13458 ) | -0.8987 ( 0.13886 ) | 0.02458 | 0.00000 | 0.00000 |
| 30 Nose, Nasal Cavity and Mi | -0.70758 ( 0.43469 ) | -1.6394 ( 0.87092 ) | 0.12371 | 0.38404 | 0.37282 |
| 31 Larynx | 1.21371 ( 0.25728 ) | -3.1695 ( 0.99551 ) | 0.13068 | 0.00011 | 0.00007 |
| 32 Lung and Bronchus | 3.71166 ( 0.13993 ) | -0.8594 ( 0.13851 ) | 0.02502 | 0.00000 | 0.00000 |
| 33 Trachea, Mediastinum and | -2.31106 ( 0.38885 ) | -3.7029 ( 1.76282 ) | 0.22944 | 0.48327 | 0.47290 |
| 34 Bones and Joints | -0.48883 ( 0.37556 ) | -0.0312 ( 0.59852 ) | 0.08981 | 0.55596 | 0.54658 |
| 35 Soft Tissue including Hea | 1.54851 ( 0.22162 ) | -2.3136 ( 0.59262 ) | 0.08042 | 0.00000 | 0.00000 |
| 36 Skin excluding Basal and | 0.13866 ( 0.14570 ) | -1.7968 ( 0.25599 ) | 0.03744 | 0.00000 | 0.00000 |
| 37 Melanoma of the Skin | 0.17192 ( 0.15360 ) | -2.2282 ( 0.27928 ) | 0.04051 | 0.00000 | 0.00000 |
| 38 Other Non-Epithelial Skin | 0.00164 ( 0.23994 ) | 0.1701 ( 0.86128 ) | 0.11363 | 0.86394 | 0.86073 |
| 39 Breast | -0.08345 ( 0.12945 ) | -2.4052 ( 0.12810 ) | 0.02315 | 0.00000 | 0.00000 |
| 40 Female Genital System | -0.70389 ( 0.07288 ) | -0.8266 ( 0.13887 ) | 0.01993 | 0.47665 | 0.46621 |
| 41 Cervix Uteri | -1.81997 ( 0.12578 ) | -2.7864 ( 0.30711 ) | 0.04218 | 0.00810 | 0.00671 |
| 42 Corpus and Uterus, NOS | -1.45753 ( 0.09517 ) | -0.4197 ( 0.22008 ) | 0.03047 | 0.00008 | 0.00006 |
| 43 Corpus Uteri | -1.25283 ( 0.14610 ) | -1.6515 ( 0.29045 ) | 0.04132 | 0.26488 | 0.25367 |
| 44 Uterus, NOS | -1.66602 ( 0.18009 ) | 0.6246 ( 0.30305 ) | 0.04480 | 0.00000 | 0.00000 |
| 45 Ovary | 0.12034 ( 0.07301 ) | -0.3897 ( 0.19006 ) | 0.02588 | 0.02274 | 0.01969 |
| 46 Vagina | -1.52270 ( 0.33520 ) | -0.0450 ( 1.12439 ) | 0.14912 | 0.25213 | 0.24101 |
| 47 Vulva | -0.19144 ( 0.26520 ) | -0.4401 ( 0.74207 ) | 0.10016 | 0.77419 | 0.76896 |
| 48 Other Female Genital Orga | -0.51159 ( 0.45107 ) | -2.8798 ( 1.05902 ) | 0.14630 | 0.06138 | 0.05545 |
| 49 Urinary System | 0.12841 ( 0.07580 ) | -0.3737 ( 0.25074 ) | 0.03329 | 0.08136 | 0.07437 |
| 50 Urinary Bladder | -0.88939 ( 0.15915 ) | -0.2421 ( 0.36919 ) | 0.05110 | 0.14319 | 0.13390 |
| 51 Kidney and Renal Pelvis | 1.17505 ( 0.15240 ) | -0.5114 ( 0.23849 ) | 0.03597 | 0.00000 | 0.00000 |
| 52 Ureter | -0.86113 ( 0.44480 ) | -0.4319 ( 1.29441 ) | 0.17396 | 0.77551 | 0.77031 |
| 53 Other Urinary Organs | -1.59386 ( 0.32699 ) | 0.51241 ( 1.65580 ) | 0.21451 | 0.25648 | 0.24533 |
| 54 Eye and Orbit | -2.35049 ( 0.44521 ) | -1.09212 ( 1.83285 ) | 0.23973 | 0.54409 | 0.53454 |
| 55 Brain and Other Nervous S | 0.96806 ( 0.12100 ) | -0.78950 ( 0.24186 ) | 0.03437 | 0.00000 | 0.00000 |
| 56 Endocrine System | -0.64446 ( 0.20063 ) | -0.82845 ( 0.57794 ) | 0.07776 | 0.78450 | 0.77950 |
| 57 Thyroid | -1.08132 ( 0.26928 ) | -0.11015 ( 0.52083 ) | 0.07452 | 0.13204 | 0.12308 |
| 58 Other Endocrine including | 0.07706 ( 0.22568 ) | -2.21315 ( 0.99215 ) | 0.12932 | 0.04069 | 0.03613 |
| 59 Lymphoma | 1.61631 ( 0.06812 ) | -1.31900 ( 0.40750 ) | 0.05251 | 0.00000 | 0.00000 |
| 60 Hodgkin Lymphoma | -3.46456 ( 0.24364 ) | -1.53544 ( 0.79565 ) | 0.10576 | 0.03503 | 0.03091 |
| 61 Non-Hodgkin Lymphoma | 2.11389 ( 0.06817 ) | -1.29605 ( 0.41882 ) | 0.05393 | 0.00000 | 0.00000 |
| 62 Myeloma | 1.31123 ( 0.08985 ) | -0.87526 ( 0.33092 ) | 0.04358 | 0.00000 | 0.00000 |
| 63 Leukemia | -0.32315 ( 0.08684 ) | -1.19624 ( 0.15836 ) | 0.02295 | 0.00001 | 0.00001 |
| 64 Lymphocytic Leukemia | -0.11952 ( 0.19751 ) | -1.55322 ( 0.35934 ) | 0.05212 | 0.00148 | 0.00113 |
| 65 Acute Lymphocytic Leukemi | -0.90026 ( 0.28043 ) | -0.40742 ( 0.54511 ) | 0.07791 | 0.46476 | 0.45419 |
| 66 Chronic Lymphocytic Leuke | 0.91209 ( 0.23062 ) | -1.78207 ( 0.46202 ) | 0.06563 | 0.00000 | 0.00000 |
| 67 Other Lymphocytic Leukemi | -4.62562 ( 0.39316 ) | -5.55778 ( 1.60359 ) | 0.20985 | 0.60770 | 0.59917 |
| 68 Myeloid and Monocytic Leu | -1.23508 ( 0.15510 ) | -0.37999 ( 0.29947 ) | 0.04286 | 0.02114 | 0.01826 |
| 69 Acute Myeloid Leukemia | -0.75770 ( 0.17158 ) | 1.22486 ( 0.30773 ) | 0.04478 | 0.00000 | 0.00000 |
| 70 Acute Monocytic Leukemia | -4.54293 ( 0.54879 ) | -9.06558 ( 3.75654 ) | 0.48252 | 0.27871 | 0.26741 |
| 71 Chronic Myeloid Leukemia | -0.77315 ( 0.24368 ) | -7.22987 ( 0.99391 ) | 0.13007 | 0.00000 | 0.00000 |
| 72 Other Myeloid/Monocytic L | -7.61620 ( 0.54563 ) | 5.14116 ( 1.49323 ) | 0.20206 | 0.00000 | 0.00000 |
| 73 Other Leukemia | 1.14041 ( 0.20278 ) | -2.39493 ( 0.25479 ) | 0.04139 | 0.00000 | 0.00000 |
| 74 Other Acute Leukemia | 1.45185 ( 0.30917 ) | -4.90300 ( 0.39909 ) | 0.06416 | 0.00000 | 0.00000 |
| 75 Aleukemic, Subleukemic an | 0.75916 ( 0.16949 ) | 0.46868 ( 0.40711 ) | 0.05605 | 0.54922 | 0.53974 |
| 76 Miscellaneous Malignant C | -0.36323 ( 0.16305 ) | -0.54259 ( 0.28608 ) | 0.04185 | 0.62040 | 0.61209 |