Surveillance Research Program, NCI, Technical Report #2018-01

## Detecting outliers with sample size adjustment in the SEER Box Plot Outlier Tool

Huann-Sheng Chen[1], Benmei Liu[1] and James Cucinelli[2]

1. Statistical Methodology and Applications Branch, Division of Cancer Control and Population Sciences, National Cancer Institute
2. Information Management Services Inc.

## Introduction

The SEER Box Plot Outlier Tool (BOT) is developed to detect outliers in the Surveillance, Epidemiology, and End Results (SEER) database for 118 variables (either at person level or tumor case level) across all of the SEER cancer registries and multiple diagnosis years from 1988 to the most recent year (2015). Table 1 provides the full list of variables that are currently implemented in the SEER BOT.

For each variable included in the SEER BOT, for each specific value or group of values of interest, three analyses are performed: 1) Compute percentage of a specific value (or group of values) out of all records in each registry with all years combined, and then compare the percentages across the SEER registries and detect outlier registries if any; 2) Compute percentage of a specific value (or group of values) out of all records in each registry using the most recent (2015) data only, and then compare the percentages across the SEER registries and detect outlier registries if any; 3) Compute percentage of a specific value (or group of values) out of all records in each diagnosis year within an individual cancer registry, and then compare the percentages across the diagnosis years and detect outlier years if any.

The SEER Box Plot Outlier Tool (BOT) is integrated in the SEER Quality Assurance Program. Each year, new data submissions are examined using the BOT and any outliers are investigated.

This brief technical report describes the statistical procedure used in SEER BOT to detect outlier registries or diagnosis years in the SEER database.

## Statistical Methods

For analysis 1, assume y is a variable of interest (e.g., Race 1), $y_{ikt}$ denotes the individual value of y for person or case $i$ ($i = 1, \ldots, n_{kt}$) in cancer registry $k$ ($k = 1, \ldots, K$) at diagnosis year $t$ ($t = 1988, \ldots, 2015$), where $K$ is the total number of cancer registries included in the SEER database. Let $\tilde{y}_{ikt}$ denote the binary recode for a specific value (or groups of values) of interest. Without loss of generosity, we use the specific value of interest as those "unknowns" for illustration below. That is, $\tilde{y}_{ikt} = \begin{cases} 1, & if\ y_{ijt} = "Unknown" \\ 0, & otherwise \end{cases}$.

Assume each $\tilde{y}_{ikt}$ follows a Bernoulli distribution with probability of $p$ to get value of 1 and $1 - p$ to get value of 0. Let $\tilde{y}_k = \sum_{t=1988}^{2015} \sum_{i=1}^{n_{kt}} \tilde{y}_{ikt}$, $n_k = \sum_{t=1988}^{2015} n_{kt}$.

Then, for each $k$,

$$\tilde{y}_k \sim Binomial(n_k, p).$$

Let $\hat{p}_k = \frac{\tilde{y}_k}{n_k}$, then $\hat{p}_k$'s are estimates of $p$, and

$$E(\hat{p}_k) = p, \quad Var(\hat{p}_k) = \frac{p(1-p)}{n_k}.$$

The confidence interval for $\hat{p}_k$ is:

$$\hat{p}_k \pm 1.96 \sqrt{\frac{\hat{p}_k(1-\hat{p}_k)}{n_k}}.$$

It is noted that with the same $\hat{p}$, if the sample size $n$ is smaller, the standard deviation is larger. Therefore, there is more variation for a smaller registry, and it is more likely become outliers just by chance. To take this factor into account, we use box-plot of the point estimates and the confidence intervals associated with each point estimates together to identify real outliers.

We use the estimate $\hat{p}_k, k = 1, ..., K$, i.e., the registry level percentage of a specific value (or groups of values) for a specific variable of interest, to draw box-plot (Figure 1).

To determine whether the outside values are really outliers, the confidence interval of each $\hat{p}_k$ is compared to the box boundaries calculated from the inter-quartile range (IQRs), the width of the box and whisker plot. $IQR = Q_3 - Q_1$.

The lower boundary is defined as Lower=max(0, $Q_1 - 1.5 * IQR$), and the upper boundary is defined as $Q_3 + 1.5 * IQR$. So the box boundary range is: (max(0, $Q_1 - 1.5 * IQR$), $Q_3 + 1.5 * IQR$).

For any $\hat{p}_k$, if its confidence interval falls completely outside of the boundary range, then it is considered as an outlier. Values with confidence intervals overlap or partially overlap with the box boundary range cannot be considered as outliers. If the confidence interval falls completely outside of the boundary range formed using (max(0, $Q_1 - 3*$IQR), $Q_3 + 3*$IQR), then the $\hat{p}_k$ is considered as an extreme outlier.

Same outlier detection procedure is used for all the three analyses, all the variables and all the specific values (or group of values) of interest built in the SEER BOT. The most recent years of data will be updated when new years of SEER data is available.

Table 1: SEER Variables Currently included in the SEER BOT

| | |
|---|---|
| Age at Dx | Grade |
| Age Recode with < 1 Year Olds | Histologic Type ICD-O-3 |
| Ambiguous Terminology DX | ICCC Site Recode ICD-O-3 / WHO 2008 |
| AYA Site Recode / WHO 2008 | ICD Revision Number |
| Behavior ICD-O-3 | ICD-O-3 Conversion Flag |
| Behavior Recode for Analysis | IHS Link |
| Cause of Death | Laterality |
| | Lymphoma Subtype Recode / WHO |
| Census Tract 2000 | 2008 |
| Census Tract Certainty 2000 | Marital Status at DX |
| Computed Ethnicity | Month of Birth |
| County at DX | Month of Conclusive DX |
| CS Extension | Month of Dx |
| CS Lymph Nodes | Month of Multiple Tumors |
| CS Lymph Nodes Eval | Mult Tum Rpt as One Prim |
| CS Mets at Dx | Multiplicity Counter |
| CS Mets Eval | NAPIIA Derived Pacific Islander Origin |
| CS Site-Specific Factor 1 | NHIA Derived Hispanic Origin |
| | Origin Recode NHIA (Hispanic, Non- |
| CS Site-Specific Factor 2 | Hisp) |
| CS Site-Specific Factor 3 | Override Site / Type |
| CS Site-Specific Factor 4 | Place of Birth |
| CS Site-Specific Factor 5 | Primary Payer at DX |
| CS Site-Specific Factor 6 | Primary Site |
| CS Site-Specific Factor 7 | Race 1 |
| CS Site-Specific Factor 8 | Race 2 |
| CS Site-Specific Factor 9 | Race 3 |
| CS Site-Specific Factor 10 | Race 4 |
| CS Site-Specific Factor 11 | Race 5 |
| CS Site-Specific Factor 12 | Race Recode (W, B, AI, API) |
| CS Site-Specific Factor 13 | Race Recode (White, Black, Other) |
| CS Site-Specific Factor 14 | Reason For No Surgery |
| CS Site-Specific Factor 15 | Regional Nodes Examined |
| CS Site-Specific Factor 16 | Regional Nodes Positive |
| CS Site-Specific Factor 17 | RX Summ -- BRM |
| CS Site-Specific Factor 18 | RX Summ -- Chemo |
| CS Site-Specific Factor 19 | RX Summ -- Hormone |
| CS Site-Specific Factor 20 | RX Summ -- Other |
| CS Site-Specific Factor 21 | RX Summ -- Radiation |
| CS Site-Specific Factor 22 | RX Summ -- Scope Reg LN Sur |

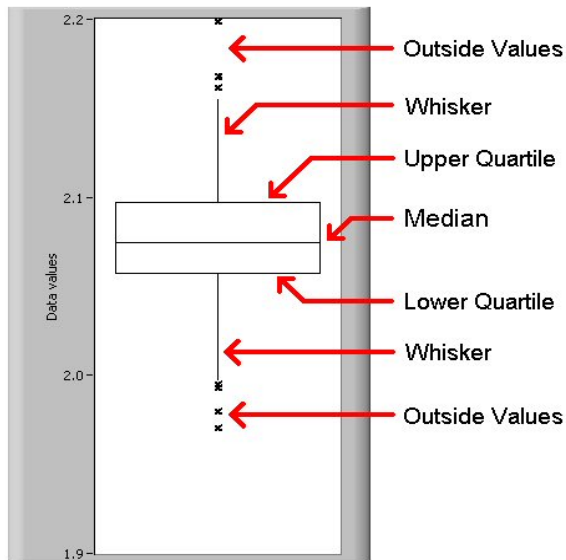| | |
|---|---|
| CS Site-Specific Factor 23 | RX Summ -- Surg / Rad Seq |
| CS Site-Specific Factor 24 | RX Summ -- Surg Oth Reg/Dis |
| CS Site-Specific Factor 25 | RX Summ -- Surg Prim Site |
| CS Size/Ext Eval | RX Summ -- Transplant / Endocr |
| CS Tumor Size | RX summ--Systemic Surg Seq |
| CS Version 1st | SEER Record Number |
| CS Version Latest | SEER Summary Stage 1977 |
| Date of Initial RX -- SEER - Month | SEER Summary Stage 2000 |
| Date of Last Contact - Month | SEER Type of Follow-up |
| Date of Last Contact - Year | Sequence Number - Central |
| Day of Conclusive DX | Sex |
| Day of Multiple Tumors | Site Recode ICD-O-3 / WHO 2008 |
| Derived AJCC 6th Stage | Spanish / Hispanic Origin |
| Derived AJCC 7th Stage | Tumor Marker 1 |
| Diagnostic Confirmation | Tumor Marker 2 |
| DSS2000S | Tumor Marker 3 |
| EOD10 Extension | Type of Reporting Source |
| EOD10 Extension Prostate Path | Vital Status |
| EOD10 Lymph Node Involvement | Vital Status Recode |
| EOD10 Tumor Size | Year of Birth |
| Ethnicity Source | Year of Conclusive DX |



*Figure 1: A boxplot illustration. The ends of the vertical lines or "whiskers" indicate the minimum and maximum data values, unless outliers are present in which case the whiskers extend to a maximum of 1.5 times the inter-quartile range (IQR). Extreme outliers are the ones that falls 3 times IQR outside of the box boundary. Outliers are the ones that falls 1.5 times IQR outside of the box boundary.*