# Supplementary Material

## In Theory and in Practice:

## Trade-Offs Between Bounded and Open Data Analysis

## for Large Scale Online Controlled Experiments

### S1: Proof of Results Under Constant Treatment Effect

The goal of this section is to show that both methods produce unbiased estimated of a constant treatment effect. Consider the following additive model for user's daily values for the metric

$$Y_u(t) = R_u(t)(c_u + \varepsilon + Z_u\tau) \tag{1}$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ are iid, $R_u(t)$ is an identity function reflecting whether a user $u$, is present on day $t$ or not. An estimator of the average treatment effect is unbiased if its expected value is $\tau$ for double average class of metrics.

$$\widehat{M}_{u\theta}^z = \frac{1}{N_\theta^z} \sum_{i=1}^a \sum_{u:t_u^0=i,z_u=z} \underbrace{\frac{\sum_{t\in\mathbb{I}_{u\theta}} Y_u(t, z_u=z)}{\sum_{t\in\mathbb{I}_{u\theta}} R_u(t)}}_{m_{u\theta}^z} \tag{2}$$

where $\theta$ is the method indicator, ($\theta = b$ for bounded and $\theta = o$ for open), $a$ is the last day that users are admitted to the experiment, $t_u^0$ is the first day a user is admitted to the experiment.

Let $D_u$, be the random variable showing the number of active days for each user. Using the law of total expectations:

$$\mathbb{E}[m_{u\theta}^z] = \mathbb{E}[\mathbb{E}[m_{u\theta}^z|D_u]] = \mathbb{E}\left[\frac{D_u(C_u+Z_u\tau)}{D_u}\right] = C_u + Z_u\tau \tag{3}$$

therefore, the expected value of the estimator of the treatment effect, $\mathbb{E}[\widehat{\Delta}]$, would be

$$\mathbb{E}[\widehat{\Delta_\theta}] = \mathbb{E}[\widehat{M}_{u\theta}^T] - \mathbb{E}[\widehat{M}_{u\theta}^C] \tag{4}$$

$$\approx \frac{1}{\mathbb{E}[N^T]} \sum_{i=1}^a \mathbb{E}[N_G^i]\, \mathbb{E}[\widehat{M}_{u\theta}^T] - \frac{1}{\mathbb{E}[N^C]} \sum_{i=1}^a \mathbb{E}[N_G^i]\, \mathbb{E}[\widehat{M}_{u\theta}^C]$$

$$= ((C_u + \tau) - C_u) = \tau$$

Where $N_G^i$ is the number of users with start day of $i$.

### S2: Intuition about the root of the bias under model 1

Intuitively, the bias originates from the difference in the ratio of active weekend days to all active days in the analysis sample and population. Let's study a very simple model to develop this intuition further for bounded analysis. Table S1 illustrates the simplest model to study time varying treatment effect for bounded analysis: an experiment with total length of 4 days where the even days have an additional treatment effect $\tau'$. The bounded analysis has an admission window of 2 days ($a = 2$) and observation window of 2 days (d=2). Each user's active days list is shown with 4 characters, 1 for active,

0 for absent and x for 'don't care' which is shown in column 1. The desired value for the ratio of number of even days to all days is 0.5 but different sets of users show values between 0 and 1.

**Table S1:** Illustration of cause of bias in the bounded approach in a simple experiment with even day effect and an admission period and observation period of 2 days.

| Active days of the user | Probability | #Active even days | # Active monitored days | Expected ratio |
|---|---|---|---|---|
| 10xx | $(1-p)p$ | 0 | 1 | 0 |
| 11xx | $p^2$ | 1 | 2 | $0.5p^2$ |
| 010x | $(1-p)^2p$ | 1 | 1 | $(1-p)^2p$ |
| 011x | $p^2(1-p)$ | 1 | 2 | $0.5p^2(1-p)$ |
| 0010 | $(1-p)^2p$ | 0 | 1 | 0 |
| 0011 | $(1-p)^2p^2$ | 1 | 2 | $0.5(1-p)^2p^2$ |

The expected value of the average number of active even days is now $\frac{0.5p^2+p(1-p)^2+0.5p^2(1-p)+ 0.5(p^2)(1-p)^2}{1-(1-p)^4}$ which is different from the theoretical value of 0.5 and varies between 0.25 and 0.5 depending on the value of p.

Using the open method, the desired value for the ratio of number of even active days to all active days is still 0.5. Table S2 shows the breakdown of possible outcomes of user activity and its probability.

**Table S2:** Illustration of expected observation for the open approach in a simple experiment with even day effect and experiment duration of 4 days.

| Active days of the user | Probability | #Active even days | # Active monitored days |
|---|---|---|---|
| 1000, 0010 | $2(1-p)^3p$ | 0 | 1 |
| 0100,0001 | $2(1-p)^3p$ | 1 | 1 |
| 1001, 1100 ,0110,0011 | $4(1-p)^2p^2$ | 1 | 2 |
| 1010 | $(1-p)^2p^2$ | 0 | 2 |
| 0101 | $(1-p)^2p^2$ | 2 | 2 |
| 1011,1110 | $2(1-p)p^3$ | 1 | 3 |
| 1101,0111 | $2(1-p)p^3$ | 2 | 3 |
| 1111 | $p^4$ | 2 | 4 |

Although different users observe different number of even days, the expected value for the average of these users is

$$\left(2(1-p)^3 p\left(\frac{1}{1}\right) + 4(1-p)^2 p^2\left(\frac{1}{2}\right) + (1-p)^2 p^2\left(\frac{2}{2}\right) + 2(1-p)p^3\left(\frac{1}{3}\right) + +2(1-p)p^3\left(\frac{2}{3}\right) + p^4\left(\frac{2}{4}\right)\right)$$
$$/(1-(1-p)^4) = \frac{1}{2}$$

which shows no bias for the open method under this scenario.

## S2: Proof of Results Under Model 1

The goal of this section is to calculate the expected values of the estimate, $\widehat{\Delta}$, and its variance, $\widehat{Var(\Delta)}$.

$$\mathbb{E}[\widehat{\Delta}_\theta] = \mathbb{E}[\widehat{M}^T_{u\theta}] - \mathbb{E}[\widehat{M}^C_{u\theta}]$$

$$\mathbb{E}[Var(\widehat{\Delta}_\theta)] \approx \frac{\mathbb{E}\left[\widehat{s^2_{M^T_{u\theta}}}\right]}{\mathbb{E}[N^T_\theta]} + \frac{\mathbb{E}\left[\widehat{s^2_{M^C_{u\theta}}}\right]}{\mathbb{E}[N^C_\theta]}$$

where $m^z_{u\theta}$ is the user level metrics, $\widehat{M}^z_{u\theta}$ is the sample mean of the metrics and $\widehat{s^2_{M^T_{u\theta}}}$ is the sample variance of the metric under study. Consider the following additive model for user's daily values for the metric

$$Y_u(t) = R_u(t)\left(c_u + \varepsilon + Z_u(\tau + \tau' I(t = Weekend))\right) \quad (5)$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ are iid. Under model 1, an estimator of the average treatment effect is unbiased if its expected value is $\tau + \frac{2}{7}\tau'$ for double average class of metrics.

Let $G_i$ be the set of users with first day as $i$, then the size of each cohort, i.e., the number of users entering the experiment on each day is $N^i_G$ with the expected value of

$$\mathbb{E}[N^i_G] = N(1-p)^{i-1}p, \quad (6)$$

and the total number of admitted users is:

$$\mathbb{E}[N_b] = \sum_{i=1}^{a} \mathbb{E}[N^i_G] = N(1-(1-p)^a). \quad (7)$$

Let's start by the bounded methods. There are several possibilities for $m^z_{ub}$ based on the user's number of active weekdays and weekends, let $w$ be the iterator over these possible outcomes, and $N^w_G$ be the number users with that number of weekdays and weekends ($w$ iterates from 0 weekday, 0 weekend to 5 weekday, 2 weekends):

$$\mathbb{E}[\widehat{M}^z_{ub}] = \frac{1}{\mathbb{E}[N_b]}\sum_{i=1}^{7}\sum_{u:t^0_u=i} m^z_{ub}$$
$$= \frac{1}{N(1-(1-p)^7)}\sum_{i=1}^{7}\sum_{w}\mathbb{E}[N^w]\mathbb{E}[m^{z,w}_{ub}] = \quad (8)$$
$$\frac{1}{1-(1-p)^7}$$

$$\sum_{i=1to\ 5}(1-p)^{i-1}p\left(\sum_{j=0to\ 4}\binom{4}{j}p^j(1-p)^{4-j}\left((1-p)^2\left(C_u + Z(\tau + \frac{\mathbf{0\tau'}}{1+j})\right)\right.\right.$$
$$+2p(1-p)\left(C_u + Z(\tau + \frac{\mathbf{1\tau'}}{2+j})\right)$$
$$+p^2\left(C_u + Z(\tau + \frac{\mathbf{2\tau'}}{3+j})\right)\bigg)\bigg) +$$

$$\sum_{i=6to\ 7}(1-p)^{i-1}p\left(\sum_{j=0to\ 5}\binom{5}{j}p^j(1-p)^{5-j}\left((1-p)\left(C_u + Z(\tau + \frac{\mathbf{1\tau'}}{1+j})\right)\right.\right.$$
$$+p\left(C_u + Z(\tau + \frac{\mathbf{2\tau'}}{2+j})\right)\bigg)\bigg)$$

where $i$ is the first day the user became active, $j$ is the number of active weekdays (other than the first active day of the user). The bold segments are the expected value of the metric for that subgroup of users. The expected value of the treatment effect, $\mathbb{E}[\widehat{\Delta}_b]$, then can easily be calculated.

$$\mathbb{E}[\widehat{\Delta}_b] = \mathbb{E}[\widehat{M}^T_{ub}] - \mathbb{E}[\widehat{M}^C_{ub}] = \tau + \rho_b(p)\,\tau' \quad (9)$$

where $\rho_b$ is a polynomial fraction that is calculated and is shown in Figure S1. The expected bias of this method is the different of this value and $\tau + \frac{2}{7}\tau'$ for double average class of metrics.

Similar breakdown of users based on the number of active weekends and weekdays can be done for the open analysis method. There is no need to use the concept of admission here. Each user can be active for up to 10 weekdays and 4 weekend days, let $N^{i,j}_G$ be the number of users in each group:

$$\mathbb{E}[\widehat{M}^z_{uo}] = \frac{1}{\mathbb{E}[N_o]}\sum_{i=0\ to\ 10}\sum_{j=0\ to\ 4}\mathbb{E}[N^{i,j}]\mathbb{E}[m^{z,i,j}_{uo}]$$
$$= \frac{1}{1-(1-p)^{14}}\sum_{i=0to}\sum_{\substack{j=0to\ 4 \\ j\neq 0\ if\ i=0}}\binom{10}{i}p^i(1 \quad (10)$$
$$-p)^{10-i}\binom{4}{j}p^j(1-p)^{4-j}\left(C_u + Z(\tau + \frac{\mathbf{j\tau'}}{i+j})\right)$$

where $i$ and $j$ count the active week and weekend days. Again, the bold segment is the expected value of the metric for that subgroup of users. This equation can be simplified and it is shown that:

$$\mathbb{E}[\widehat{\Delta}_o] = \mathbb{E}[\widehat{M}^T_{uo}] - \mathbb{E}[\widehat{M}^C_{uo}] = \tau + \frac{2}{7}\tau' \quad (11)$$
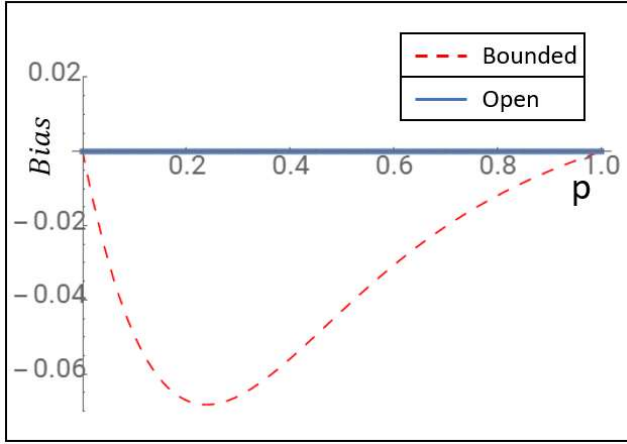
**Figure S1** Bias of open and bounded methods under model 1.

The expected value of the sample variance can similarly be calculated:

$$\mathbb{E}\left[\widehat{s^2_{M^Z_{u\theta}}}\right] = \mathbb{E}\left[\frac{1}{N_\theta}\sum_{i=1}^{a}\sum_{u:t^0_u=i}\left(M^z_{u\theta}-\widehat{M}^z_{ub}\right)^2\right]$$
$$\approx \frac{1}{\mathbb{E}[N_\theta]}\sum_{i=1}^{a}\sum_{u:t^0_u=i}\mathbb{E}\left[\left(M^z_{u\theta}-\widehat{M}^z_{ub}\right)^2\right] \quad (12)$$

For the bounded method, this equation reduces to:

$$\mathbb{E}\left[\widehat{s^2_{M^Z_{ub}}}\right] \approx \frac{1}{\mathbb{E}[N_b]}\sum_{i=1}^{a}\sum_{u:t^0_u=i}\mathbb{E}\left[\left(M^z_{u\theta}-\overline{M^z_{u\theta}}\right)^2\right]$$
$$= \frac{1}{N(1-(1-p)^7)}$$

$$\sum_{i=1\,to\,5}(1-p)^{i-1}p\sum_{j=0\,to\,4}\binom{4}{j}p^j(1-p)^{4-j}\Bigg((1$$

$$-p)^2\,\mathbb{E}\left[\left(\frac{\sum^{1+j}\boldsymbol{\varepsilon}}{1+j}+Z\left(\frac{0\boldsymbol{\tau'}}{1+j}-\boldsymbol{\rho_b}(\boldsymbol{p})\boldsymbol{\tau'}\right)\right)^2\right]$$

$$+2p(1$$

$$-p)\,\mathbb{E}\left[\left(\frac{\sum^{2+j}\boldsymbol{\varepsilon}}{2+j}+Z\left(\frac{1\boldsymbol{\tau'}}{2+j}-\boldsymbol{\rho_b}(\boldsymbol{p})\boldsymbol{\tau'}\right)\right)^2\right]$$

$$+p^2\mathbb{E}\Bigg[\Bigg(\frac{\sum^{3+j}\boldsymbol{\varepsilon}}{3+j}$$

$$+Z\left(\frac{2\boldsymbol{\tau'}}{3+j}-\boldsymbol{\rho_b}(\boldsymbol{p})\boldsymbol{\tau'}\right)\Bigg)^2\Bigg]\Bigg)$$

$$+$$

$$\sum_{i=6t\,7}(1-p)^{i-1}p\sum_{j=0to\,5}\binom{5}{j}p^j(1-p)^{5-j}\Bigg((1$$

$$-p)\mathbb{E}\left[\left(\frac{\sum^{1+j}\boldsymbol{\varepsilon}}{1+j}+Z\left(\frac{1\boldsymbol{\tau'}}{1+j}-\boldsymbol{\rho_b}(\boldsymbol{p})\boldsymbol{\tau'}\right)\right)^2\right]$$

$$+p\mathbb{E}\left[\left(\frac{\sum^{2+j}\boldsymbol{\varepsilon}}{2+j}+Z\left(\frac{2\boldsymbol{\tau'}}{2+j}-\boldsymbol{\rho_b}(\boldsymbol{p})\boldsymbol{\tau'}\right)\right)^2\right]\Bigg)$$

$$\mathbb{E}[\text{Var}(\widehat{\Delta_b})] \approx \frac{\mathbb{E}\left[\widehat{s^2_{M^T_{ub}}}\right]}{\mathbb{E}[N^T_b]}+\frac{\mathbb{E}\left[\widehat{s^2_{M^C_{ub}}}\right]}{\mathbb{E}[N^C_b]}=\eta_b(p)\sigma^2+\zeta_b(p)\tau'^2 \quad (13)$$

where $\eta_b(p)$ and $\zeta_b(p)$ is depicted in Figure S2.

The expected value of the sample variance for the open method can be calculated using the same breakdown we used for the bias calculation:

$$\mathbb{E}\left[\widehat{s^2_{M^Z_{uo}}}\right]$$

$$\approx \frac{1}{\mathbb{E}[N_o]}\sum_{i=0\,to\,10}\sum_{j=0\,to\,4}\mathbb{E}[N^{i,j}_G]\mathbb{E}\left[\left(m^{z,i,j}_{ub}-C_u-Z(\tau\right.\right.$$
$$\left.\left.+\frac{2}{7}\tau')\right)^2\right]=$$

$$=\frac{1}{1-(1-p)^{14}}\sum_{i=0to10}\sum_{j=0to\,4}\binom{10}{i}p^i(1$$

$$-p)^{10-i}\binom{4}{j}p^j(1-p)^{4-j}\,\mathbb{E}\left[\frac{\sum^{i+j}\boldsymbol{\varepsilon}}{i+j}+Z\left(\frac{j\boldsymbol{\tau'}}{i+j}-\frac{2}{7}\right)\right]^2 \quad (14)$$

$$\mathbb{E}[\text{Var}(\widehat{\Delta_o})] \approx \frac{\mathbb{E}\left[\widehat{s^2_{M^T_{uo}}}\right]}{\mathbb{E}[N^T_o]}+\frac{\mathbb{E}\left[\widehat{s^2_{M^C_{uo}}}\right]}{\mathbb{E}[N^C_o]}=\eta_o(p)\sigma^2+\zeta_o(p)\tau'^2 \quad (15)$$

where $\zeta_o(p)$ and $\gamma_o(p)$ is depicted in Figure S2.
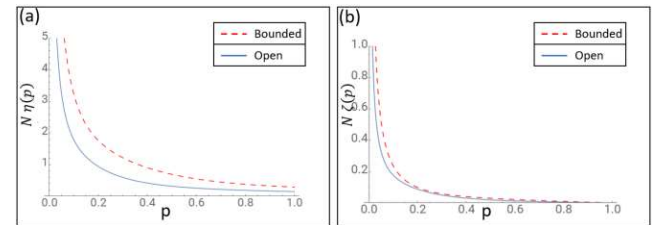


**Figure S2** a) Coefficient of base activity part of the variance. b) Coefficient of the week-day interaction part of the variance for open and bounded methods under model 1.

## A3: Proof of Results Under Model 2

For the metric of interest $Y$ consider the following additive model

$$Y_u(t) = R_u(t)\left(c_u + \varepsilon + Z_u\left(\tau + \tau' I(t = Weekend)\right)\right) \qquad (16)$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ are iid. Under this model, an estimator for the treatment effect in a week is unbiased if its expected value is equal to $\tau + \frac{2}{7}\tau'$ for double average metric. We first show that the bounded observation always provides an unbiased estimator under model 2 if observation window is one week.

$$\widehat{\Delta_b} = \widehat{M}_{ub}^T - \widehat{M}_{ub}^C \qquad (17)$$

$$\widehat{M}_{ub}^z = \frac{1}{N_b^T}\sum_{i=1}^{7}\sum_{u\in G_i}\frac{\sum_{t=i}^{i+d-1} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{i+d-1} R_u(t)}$$
$$= \frac{1}{N_b^T}\sum_{i=1}^{7}\sum_{u\in G_i}\frac{\sum_{t=i}^{i+d-1}\left(C_u + Z_u(\tau+\tau' I(weekend)) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{i+d-1} R_u(t)} \qquad (18)$$

$$\mathbb{E}[\widehat{M}_{ub}^z] = \frac{1}{N_b^T}\sum_{i=1}^{7}\sum_{u\in G_i}\left(C_u + Z_u(\tau + \frac{2}{7}\tau')\right)$$
$$= \frac{1}{N_b^T}\sum_{i=1}^{7}\left(N_s(C_u + Z_u(\tau + \frac{2}{7}\tau'))\right)$$
$$= \frac{1}{7N_s}\left(7N_s\left(C_u + Z_u(\tau + \frac{2}{7}\tau')\right)\right) = C_u + Z_u(\tau + \frac{2}{7}\tau')$$

The variance term can similarly be calculated:

$$\widehat{\sigma_{M_{ub}^z}^2} = \frac{1}{N_b^z}\sum_{i=1}^{7}\sum_{u\in G^i}\left(\frac{\sum_{t=i}^{i+6} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{i+d-1} R_u(t)} - C_u - Z_u(\tau + \frac{2}{7}\tau')\right)^2$$

$$\mathbb{E}\left[\widehat{\sigma_{M_{ub}^z}^2}\right]$$
$$= \frac{1}{\mathbb{E}[N_b^z]}\sum_{i=1}^{7} N_s\,\mathbb{E}\left[\left(\frac{\sum_{t=i}^{i+6} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{i+d-1} R_u(t)}\right.\right. \qquad (19)$$
$$\left.\left. - C_u - Z_u\left(\tau + \frac{2}{7}\tau'\right)\right)^2\right]$$

$$= \frac{1}{\mathbb{E}[N_b^z]}\sum_{i=1}^{7} N_s\,\mathbb{E}\left[\left(\frac{\sum_{t=i}^{i+6} \mathcal{N}(0,\sigma^2)}{7}\right)^2\right] = \frac{\sigma^2}{7}$$

therefore,

$$\mathbb{E}[\widehat{\Delta_b}] = \tau + \frac{2}{7}\tau' \text{ and } \mathbb{E}[Var(\widehat{\Delta_b})] = \frac{2}{49N_s}\sigma^2 \qquad (20)$$

For open analysis method:

$$\mathbb{E}[\widehat{M}_{uo}^z] = \frac{1}{N_o^z}\sum_{i=1}^{14}\mathbb{E}\sum_{u\in G_i}\frac{\sum_{t=i}^{14} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{14} R_u(t)}$$

$$= \frac{1}{N_o^z}\sum_{i=1}^{14}\sum_{u\in G_i}\frac{\sum_{t=i}^{14}\left(C_u + Z_u(\tau + \tau' I(weekend))\right)}{15-i}$$
$$= \frac{1}{N_o^T}\sum_{i=1}^{14}\sum_{u\in G_i} C_u + Z_u\left(\tau + \left(\frac{n_{weekend}(i)}{15-i}\right)\tau'\right)$$
$$= \frac{1}{N_o^z}\sum_{i=1}^{14} N_s\left(C_u + Z_u\left(\tau + \left(\frac{n_{weekend}(i)}{15-i}\right)\tau'\right)\right)$$
$$\cong \frac{1}{14N_s} N_s\left(14 C_u + Z_u(14\tau + 6.7\tau')\right)$$
$$= C_u + Z_u\left(\tau + \frac{6.7}{14}\tau'\right)$$

The variance term can similarly be calculated:

$$\widehat{\sigma_{M_{uo}^z}^2} = \frac{1}{N_o^z}\sum_{i=1}^{14}\sum_{u\in G^i}\left(\frac{\sum_{t=i}^{14} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{14} R_u(t)}\right.$$
$$\left. - C_u - Z_u(\tau + \frac{6.7}{14}\tau')\right)^2$$

$$\mathbb{E}\left[\widehat{\sigma_{M_{uo}^z}^2}\right]$$
$$= \frac{1}{\mathbb{E}[N_o^z]}\sum_{i=1}^{14} N_s\,\mathbb{E}\left[\left(\frac{\sum_{t=i}^{i+6} R_u(t)\left(C_u + Z_u\tau_u(t) + \mathcal{N}(0,\sigma^2)\right)}{\sum_{t=i}^{i+d-1} R_u(t)}\right.\right.$$
$$\left.\left. - C_u - Z_u\left(\tau + \frac{6.7}{14}\tau'\right)\right)^2\right] \qquad (21)$$

$$= \frac{1}{\mathbb{E}[N_o^z]}\sum_{i=1}^{14} N_s\,\mathbb{E}\left[\left(\frac{\sum_{t=i}^{14}\mathcal{N}(0,\sigma^2)}{15-i}\right.\right.$$
$$\left.\left. + \left(\frac{n_{weekend}(i)}{15-i} - \frac{6.7}{14}\right)\tau'\right)^2\right]$$
$$= \frac{1}{\mathbb{E}[N_o^z]}\sum_{i=1}^{14} N_s\,\mathbb{E}\left[\left(\frac{\sum_{t=i}^{14}\mathcal{N}(0,\sigma^2)}{15-i}\right)^2\right]$$
$$+ \mathbb{E}\left[\left(\left(\frac{n_{weekend}(i)}{15-i} - \frac{6.7}{14}\right)\tau'\right)^2\right]$$
$$= \frac{3.25}{14}\sigma^2 + Z_u(\frac{0.76}{14}\tau'^2)$$

therefore,

$$\mathbb{E}[\widehat{\Delta_o}] = \tau + \frac{6.7}{14}\tau' \text{ and } \mathbb{E}[Var(\widehat{\Delta_o})] = \frac{6.5}{14^2 N_s}\sigma^2 + \frac{0.76}{14^2 N_s}\tau'^2 \qquad (22)$$