Consider a world in which all people have a (hidden, latent) health $z$, and have the option of getting a treatment, with an additive effect $\delta$, which we wish to learn from data. We observe for each person a binary variable for treatment $t$ and observe their total health $y = z + \delta t$.

We need to define two meanings of "expectation".

First consider the distribution from which all people are drawn: $p(z,t)$. From this we could compute $\delta$ in terms of expecations:

$$E(y|t) - E(y|\neg t) = (E(z|t) + \delta) - E(z|\neg t) = \delta + p(z|t) - p(z|\neg t)$$

Again, of course, this distribution is not known. Let's define, just to be suggestive,

$$\hat{\delta} \equiv E(y|t) - E(y|\neg t) = \delta + p(z|t) - p(z|\neg t) \equiv \delta + b$$

where the true selection bias $b = p(z|t) - p(z|\neg t)$. Note that if $z$ (health) and treatment $t$ are indepentent, $b = p(z|t) - p(z|\neg t) = p(z) - p(z) = 0$.

Now consider a sample of $N$ individuals drawn from $p(z,t)$. For these we can define sample averages

$$\bar{\delta} \equiv \overline{y|t} - \overline{y|\neg t} = \qquad \frac{1}{N_t} \sum_{i|t} y_i - \frac{1}{N_{\neg t}} \sum_{i|\neg t} y_i \qquad (1)$$

$$= \qquad \frac{1}{N_t} \sum_{i|t} (z_i + \delta) - \frac{1}{N_{\neg t}} \sum_{i|\neg t} z_i \qquad (2)$$

$$= \qquad \delta + \frac{N_{z=1,t=1}}{N_{t=1}} - \frac{N_{z=1,t=0}}{N_{t=0}} \qquad (3)$$

$$= \qquad \delta + \hat{p}(z|t) - \hat{p}(z|\neg t) \qquad (4)$$

Since $\hat{p}$ refers to a sampled quantity it will in general be nonzero, even if $p(z,t) = p(z)p(t)$ in the true distribution. We can define, to be suggestive, $\bar{b} = \hat{p}(z,t) - \hat{p}(z,\neg t)$