

Appendix A

Formation of the Statistical Model Underlying the Simulations

The statistical model forming the basis of our computer simulations was constructed under a set of assumptions and calculations. All of the model parameters were estimated from the data on the women participating in the EATS study. The details of the estimation procedures are contained here. A brief description of the model is presented in the main text.

Some food groups are not consumed every day by all individuals. We refer to days on which a given food group is consumed by a given individual as that individual's "consumption days," the remaining days being the individual's "non-consumption days."

Power transformation

Our first assumption was as follows:

Normality Assumption: Distributions of intake on consumption days, both between individuals and within individuals, are assumed normal after a suitable power transformation. The deciles of the distribution of intake on consumption days were calculated (10th percentile to 90th percentile). Then the values were raised to various powers (0.1 to 0.9 in steps of 0.1), and the transformation that best achieved symmetry (10th percentile and 90th percentile equidistant from median) was chosen. The power transformation for each food/nutrient chosen is shown in the second column of Supplemental Table 1.

Estimation of the parameters of the intake distributions on consumption days under the assumption of equal probability of consumption

We next estimated the distribution of intakes on consumption days under the following assumption.

Assumption (i): *The probability of consumption is the same for all individuals.*

Online Supporting Material

This probability and the mean, standard deviations and correlations (between a food/nutrient and energy intake) of the between-person and within-person distributions on consumption days were obtained by transforming to the chosen power scale and then performing the calculations given below. The method of estimation is known as the method of moments, which is a little simpler than the method of maximum likelihood, and for normal distributions gives very similar results. (If all individuals had the same number of consumption days, then the two methods would give identical results.) The parameter values that were computed by this method and then used in the computer simulations are shown in the remaining columns of Supplemental Table 1.

Take Total Fruit as the example. We use the four 24HR reports in EATS. Let the number of persons in the EATS sample be N . Then the following nine parameter estimates were calculated.

For food group or nutrient (F):

1. Estimated probability of consumption = total number of consumption days reported by the 738 women/total number of days reported by the 738 women.
2. Suppose person i consumes (untransformed) amount f_{ij} on consumption day j .
Suppose person i consumes the food on n_i days. Then the overall mean

$$\text{consumption on the power scale} = \bar{f} = \frac{\sum_{i=1}^{N_1} n_i \bar{f}_i}{\sum_{i=1}^{N_1} n_i}, \text{ where } \bar{f}_i = \sum_{j=1}^{n_i} f_{ij}^\lambda / n_i, N_1 \text{ is}$$

the number of persons with at least one consumption day, and λ is the power used for the transformation of food F.

3. The within-person standard deviation on consumption days was computed as follows. For the N_2 persons who consumed the food on at least 2 occasions

Online Supporting Material

$$\hat{\sigma}_{Fw}^2 = \frac{\sum_{i=1}^{N_2} \sum_{j=1}^{n_i} (f_{ij}^\lambda - \bar{f}_i)^2}{\sum_{i=1}^{N_2} (n_i - 1)}$$

was calculated, which estimates the within-person

variance. The within-person standard deviation is the square root of this estimate.

4. The between-person standard deviation was computed as follows. First the

between-person mean square was calculated as follows: $s_{Fb}^2 = \frac{\sum_{i=1}^{N_1} (\bar{f}_i - \bar{f})^2}{N_1 - 1}$.

Then the between-person variance was estimated as $\hat{\sigma}_{Fb}^2 = s_{Fb}^2 - \left\{ \hat{\sigma}_{Fw}^2 \frac{\sum_{i=1}^{N_1} 1/n_i}{N_1} \right\}$.

The between-person standard deviation is the square root of this estimate.

For energy (e):

5. Suppose person i consumes energy e_{ij} on reporting day j . Suppose person i reports on m_i days. The mean energy consumption on the power scale was

computed as $\bar{e} = \frac{\sum_{i=1}^N m_i \bar{e}_i}{\sum_{i=1}^N m_i}$, where $\bar{e}_i = \sum_{j=1}^{m_i} e_{ij}^\kappa / m_i$, N is the number of persons

in the sample, and $\kappa=0.5$ is the power used for transforming energy.

6. The within-person variance for energy was estimated: $\hat{\sigma}_{Ew}^2 = \frac{\sum_{i=1}^N \sum_{j=1}^{m_i} (e_{ij}^\kappa - \bar{e}_i)^2}{\sum_{i=1}^N (m_i - 1)}$.

The within-person standard deviation is the square root of this estimate.

7. The between-person mean square was computed as: $s_{Eb}^2 = \frac{\sum_{i=1}^N (\bar{e}_i - \bar{e})^2}{N-1}$. Then the between-person variance was estimated as:

$$\hat{\sigma}_{Eb}^2 = s_{Eb}^2 - \left\{ \hat{\sigma}_{Ew}^2 \frac{\sum_{i=1}^{N_1} 1/m_i}{N} \right\}.$$

The between-person standard deviation is the square root of this estimate.

For energy and food jointly:

8. The within-person covariance of food and energy on food consumption days among the individuals who report consuming the food on at least two days, was computed as follows:

$$\hat{\sigma}_{FEw} = \frac{\sum_{i=1}^{N_2} \sum_{j=1}^{n_i} (f_{ij}^\lambda - \bar{f}_i)(e_{ij}^\kappa - \bar{e}_i^*)}{\sum_{i=1}^{N_2} (n_i - 1)}, \text{ where } \bar{e}_i^* = \sum_{j=1}^{n_i} e_{ij}^\kappa / n_i \text{ is the mean energy}$$

intake on the power scale for individual i on his/her consumption days. The within-person correlation is this estimate divided by the product of the within-person standard deviations of food and energy.

9. The between-person mean cross-product of food and energy consumption among persons who report consuming the food was computed as follows:

$$s_{FEb} = \frac{\sum_{i=1}^{N_1} (\bar{f}_i - \bar{f})(\bar{e}_i - \bar{e})}{N_1 - 1}. \text{ Then the between-person covariance was estimated}$$

Online Supporting Material

$$\text{as } \hat{\sigma}_{FEb} = s_{FEb} - \left\{ \hat{\sigma}_{FEw} \frac{\sum_{i=1}^{N_1} 1/n_i}{N_1} \right\}. \text{ The between-person correlation is this}$$

estimate divided by the product of the between-person standard deviations of food and energy.

Estimation of the parameters of the intake distributions on consumption days under the assumption of different probabilities of consumption, and no dependence between probability of consumption and amount consumed on consumption days

As mentioned, the above estimation procedure was based on assumption (i) that the probability of consumption is the same for all individuals. The estimated probabilities of consumption under this assumption are shown in the second column of Supplemental Table 2.

Unfortunately, assumption (i) is not supported by the EATS data, where too many individuals report consuming a particular food group on either no days or on all 4 days.

Therefore, we examined a more complex assumption as follows:

Assumption (ii): There are five subclasses of individuals consuming the food on 0%, 25%, 50%, 75% or 100% of days, respectively.

Implicitly we also assumed that the distribution of intakes on consumption days was independent of the probability to consume.

The proportions of persons in each class were estimated from the EATS data. Among the 650 women participating in EATS who reported on all four days, we calculated the observed proportions of those consuming the food/nutrient on 0, 1, 2, 3 and all 4 days. We then postulated 5 subclasses of individuals who over the long-term consumed on 0%, 25%, 50%, 75% and 100% with equal probability on each day, with stochastic independence between days. In estimating the proportion in each subclass, maximum likelihood can be difficult to execute because each of the five probabilities is bounded

Online Supporting Material

between 0 and 1 and their sum is likewise bounded. In addition, for such problems with limited data, as in the EATS study, maximum likelihood can produce solutions that seem unrealistic (for example, distributions with individuals consuming either 25% of the time or 100% of the time, with no other classes represented). Therefore, a combination of the method of moments and careful visual inspection was used, that is, we found by inspection a realistic set of proportions for each probability subclass that would lead to expected values of the five proportions equaling (approximately) those observed in the data. The resulting estimated proportions are shown in Supplemental Table 2 for each food group.

As mentioned, assumption (ii) included the assumption that the distribution of intakes on consumption days was independent of the individual's probability of consumption. For this reason we did not need to re-estimate the parameters of these distributions, as they would be unaffected by the change in assumption about the probability of consumption.

Estimation of the parameters of the intake distributions on consumption days under the assumption of different probabilities of consumption, with dependence between probability of consumption and the mean amount consumed on consumption days

Unfortunately, the assumption that the distribution of intakes on consumption days was independent of the individual's probability of consumption was also not supported by the EATS data. In fact, reported intakes on consumption days have previously been reported to correlate positively with the probability to consume [See reference 6 of main text]. We therefore explored a third model that was based on an assumption as follows.

Assumption (iii): There are the same five subclasses of probability of consumption as in assumption (ii), but each class has its own mean intake of the food/nutrient on consumption days.

The proportions of individuals in each subclass were estimated as under assumption (ii), and the mean intakes on consumption days for each subclass were estimated from the

Online Supporting Material

EATS data, again using a combination of the method of moments and careful visual inspection. That is, we found a set of mean intakes for each subclass that, together with the estimated proportions in each subclass, would lead to an expected mean intake on consumption days equal (approximately) to that observed in the data, paying attention to requirement that the estimated mean intake on consumption days should increase monotonically with the probability of consumption. The subclass mean intakes are shown in Supplemental Table 2.

Within each subclass the between-person and within-person standard deviations were chosen to be proportional to the mean in that subclass, with the constant of proportionality chosen so that the overall variances were equal to those estimated under assumption (i).

Appendix B

The Simulations

Simulation programs were written in S-Plus that (a) generated data from the food/nutrient and energy distributions under the three different assumptions regarding the probability to consume and (b) computed the three estimates of population mean HEI-2005 score (mean score, score of the mean ratio, and score of the population ratio). Each simulation generated a population of 10,000 persons and a single day of intake for that individual to be used in computing the three estimates. Whether or not the day was a consumption day was generated as a Bernoulli variable with the appropriate probability, under assumption (i). Under assumptions (ii) and (iii), a random multinomial variable was first generated to determine the probability subclass, and then a Bernoulli variable with the corresponding probability was generated. The food/nutrient and energy intakes on consumption days, on the transformed scale, were generated as bivariate normal distributions with the appropriate means, variances and correlations, depending on the assumption regarding the probability of consumption.

In addition, the S-Plus programs computed a true usual intake for each person. The “true” usual intake values for each simulated person were calculated by generating 1000 values from the within-person distribution for that individual on the suitably transformed scale, applying the inverse transformation, and finally taking the mean of the 1000 values. For each individual, the “true” usual intake of the component was then divided by the true usual energy intake, and from this ratio the individual’s “true” HEI-2005 component score was calculated. These HEI-2005 component scores were averaged over 10,000 individuals in the simulation to obtain the “true” population mean usual HEI-2005 component score, against which the three estimates could be compared.

Even after the power transformations, it was found that that the left-hand tails of some of the distributions were “clumped.” The main reason for this was that for several food groups, some individuals consumed only a tiny amount on a substantial proportion of days. An example might be the fruit consumed in fruit-flavored yogurt or the lemon juice in mayonnaise. These clumps of small values led to standard errors that were quite large

Online Supporting Material

relative to the means on the transformed scale. The most extreme example was for total fruit consumption, where the standard error was 74% of its mean. This in turn led to problems in the simulations where a substantial proportion of negative values (on the transformed scale) were generated. To overcome this problem, we set all simulated values less than 0.05 to the threshold value of 0.05 on the transformed scale. When the 0.05 is transformed back to the original scale, it still represents a very low intake; the adjustment, therefore, does not materially affect the mean HEI-2005 estimates.

Appendix C

Sensitivity Analysis to Examine Robustness of our Conclusions

We conducted a sensitivity analysis to check whether our conclusion that the score of the population ratio is the optimal score for currently available US population data is robust to the sampling errors involved when estimating the parameters from the sample of 738 women participating in EATS. To check this we performed a sensitivity analysis under the assumption of constant probability of consumption in which we perturbed these parameter estimates by a random normal deviate with standard deviation equal to the standard error of the estimate. We generated 5 realizations for each component so that there were 60 realizations in total. Among these, 7 out of the 60 realizations showed a change in the optimal method from that shown in Table 3 of the main text. In three realizations, the score of the population ratio had been the optimal method and was surpassed by another method. In another three realizations, the score of the mean ratio had been the optimal method (SoFAAS component) and was surpassed by the score of the population ratio. In the other realization (also the SoFAAS component), the score of the mean ratio had been optimal and was surpassed by the mean score. We then repeated this exercise making random perturbations with double the standard deviation, and again there were no substantial changes in the optimal method.

Appendix D

Confidence interval for the population mean total HEI-2005 score

The following algorithm is suggested for computing the confidence interval for the population mean total HEI-2005 score, although the resulting interval may be conservative, i.e., it may be somewhat wider than really required.

Determine the 12 “pre-truncated” component scores for each individual in the sample. For example, suppose that an individual reports consuming 56% of energy from solid fat, alcohol, and sugar (SoFAAS). Since 20% scores 20 and 50% scores zero, then without truncation, 56% would score -4 (by linear extrapolation). We call this the pre-truncated score for SoFAAS. Thus, an individual’s pre-truncated score can in some circumstances fall above the maximum for that component or below the minimum, as in the example above. Pre-truncated scores for sodium and for saturated fat are a bit more complicated since the scoring for them is not strictly linear between zero and 10. As a consequence, all individuals with intakes below the threshold for zero or above the threshold for 8 are assigned a pre-truncated score by linearly extrapolating the scores between zero (15% energy from saturated fat or 2.0 grams sodium per 1000 kcal) and 8 (10% energy from saturated fat and 1.1 grams sodium per 1000 kcal).

Returning to the algorithm, sum the 12 pre-truncated scores for each individual and then multiply by the individual’s energy intake. Using a standard survey package, estimate the ratio of the population total for this value to the population total for energy intake and, more importantly, estimate the standard error of that ratio. The ratio can be viewed as the estimated “raw” (or untruncated) total HEI-2005 score for the population. Construct a confidence interval around the estimated true total population HEI-2005 score (the sum of the possibly-truncated 12 component scores) using the estimated standard error for the raw total score. Truncate the end point(s) of the confidence interval if necessary.

This interval may be conservative because the algorithm can overestimate the contributions to the standard error from components with truncated or near-truncated scores.

Appendix E

Biases, Standard Errors and Mean Squared Errors of the Three Estimators

Our main comparison of the three estimators was based on their biases and not on their standard errors. We considered the standard error of the estimators to be of secondary importance to the bias, because in the relatively large samples that we envisage the bias will dominate the error of the estimate, especially in this case where the biases are often large. To check this further, we computed from our simulation (under the assumption of a varying probability of consumption that is correlated with amount of intake on consumption days) the standard error of the three estimates that would be expected from a sample of 1000 individuals. Average standard errors over the 12 components were 0.09 for the mean score, 0.18 for the score of the mean ratio, and 0.14 for the score of the population ratio, compared to average absolute biases of 0.73, 0.66 and 0.37 respectively. The standard error was smaller than the bias in 11 out of 12 components for the mean score and for the score of the mean ratio, and in 9 out of 12 components for the score of the population ratio. Moreover, when we compared the estimators according to their mean squared error (bias squared plus standard error squared), a measure that takes both bias and standard error into account, the same ordering was obtained as when using the bias.

Supplemental Table 1: A: Model parameters of the distribution of power-transformed intake on consumption days for HEI components and energy* (1 kcal = 4.184 kJ; 1 cup equivalent = 236.6 mL; 1 oz. equivalent = 28.35g)

Dietary component (original scale units)	Power transform	Median (on original scale)	Within- person SD (power scale)	Within- person correlation with energy	Between- person SD (power scale)	Between- person correlation with energy
Energy (kcal)	0.5	1673	6.45	-	5.92	-
Total fruit (cup equiv)	0.6	0.95	0.593	0.10	0.406	0.22
Whole fruit (cup equiv)	0.6	0.84	0.518	-0.04	0.302	0.19
Total vegetables (cup equiv)	0.5	1.42	0.415	0.30	0.241	0.31
Dark green and orange vegetables and legumes (cup equiv)	0.4	0.35	0.299	0.11	0.153	0.11
Total grains (oz. equiv)	0.6	5.27	0.802	0.58	0.508	0.72
Whole grains (oz. equiv)	0.4	0.98	0.430	0.03	0.217	0.27
Milk (cup equiv)	0.4	1.05	0.330	0.28	0.247	0.54
Meat and beans (oz. equiv)	0.4	3.94	0.491	0.44	0.239	0.58
Oils (g)	0.4	9.98	1.140	0.34	0.521	0.57
Saturated fat (g)	0.3	17.7	0.382	0.78	0.307	0.86
Sodium (mg)	0.3	2784	1.410	0.71	1.010	0.83
Calories from SoFAAS (kcal)	0.5	557	5.900	0.79	5.390	0.82

* Estimated from data reported by 738 women participating in EATS

Supplemental Table 2: Probabilities of consumption on a given day, with mean intakes (on power scale) within each probability subclass, for HEI components and energy*

Dietary component	Overall probability of consumption	Proportions of individuals (and their mean intake on consumption days on power scale) found in each consumption sub-class				
		Pr=0**	Pr=0.25	Pr=0.5	Pr=0.75	Pr=1
Total fruit	0.78	.00	.05	.10	.50	.35
		-	(.20)	(.20)	(.85)	(1.23)
Whole fruit	0.55	.05	.24	.31	.26	.14
		-	(.65)	(.85)	(.85)	(1.15)
Total vegetables	0.96	.00	.005	.005	.15	.84
		-	(.05)	(.05)	(1.07)	(1.21)
Dark green and orange vegetables and legumes	0.51	.03	.05	.80	.07	.05
		-	(.10)	(.65)	(.70)	(.83)
Total grains	0.99	.00	.00	.005	.045	.95
		-	-	(.60)	(2.50)	(2.74)
Whole grains	0.53	.05	.13	.41	.38	.03
		-	(.70)	(.89)	(1.05)	(1.44)
Milk	0.89	.00	.01	.03	.37	.59
		-	(.10)	(.20)	(.90)	(1.10)
Meat and beans	0.96	.00	.00	.005	.15	.845
		-	-	(.10)	(1.74)	(1.74)
Oils	0.88	.00	.005	.02	.445	.53
		-	(1.50)	(1.50)	(2.39)	(2.63)

* Estimated from data reported by 738 women participating in EATS

** Pr=0 represents the sub-class of individuals whose probability of consumption is zero, i.e., non-consumers; likewise, Pr=0.25 represents the subclass of individuals who consume the food on 25% of days, etc.