

FISH 604

Module 2 (part 1) :

Basic statistical concepts

Instructor: Franz Mueter

Lena Point, Rm 315

796-5448

fmueter@alaska.edu



Today's goals

You should understand...

- ... difference between arithmetic, geometric, and harmonic mean
- ... standard error, standard deviation, and coefficient of variation

You should be able to...

- ... estimate location and spread for any random variable
- ... quickly summarize data statistically & graphically
- ... estimate variance of the sum, the product, and the ratio of random variables

Basic statistical concepts

- Data summaries
 - Location: Mean, median, quantiles
 - Spread: Variance, Standard deviation, MAD, IQR
 - Graphical summaries (Histograms, dotplots, boxplots)
- Probability and probability distributions
- Distributions
 - Discrete (Binomial, Multinomial, Poisson)
 - Continuous (Uniform, Exponential, Normal or Gaussian, Log-normal, Gamma)

Data summaries

MSL

FISH

604

- Data: x_1, x_2, \dots, x_n
- Measures of location / central tendency:
 - Mean (sample mean)

Arithmetic mean: $AM_x = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

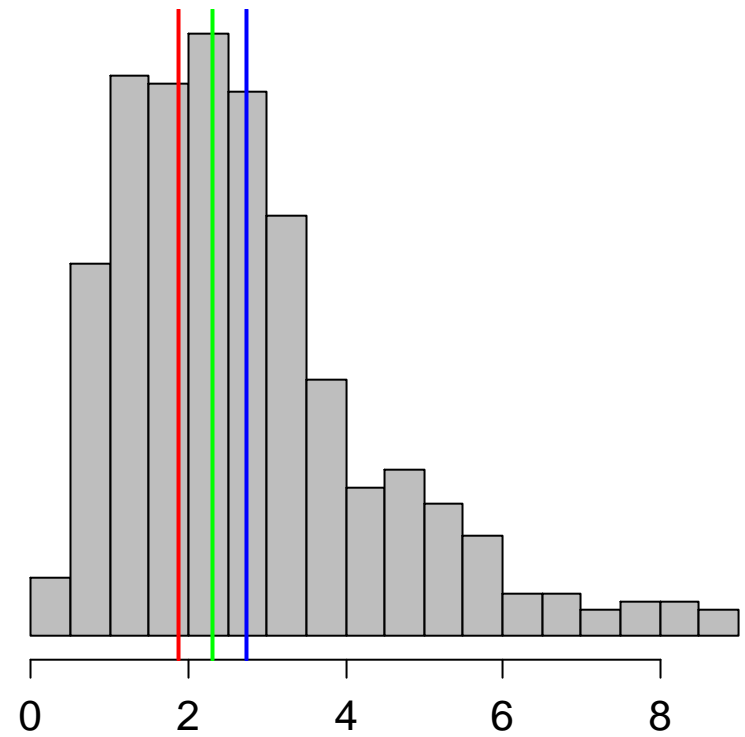
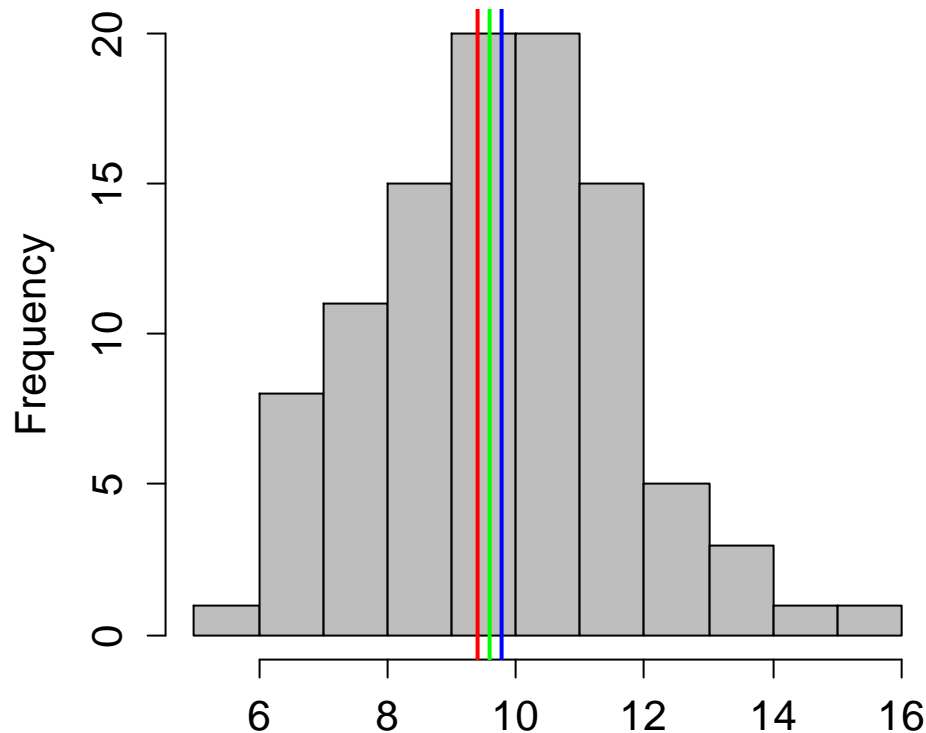
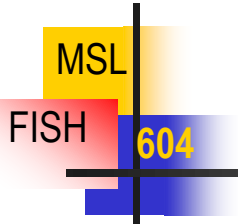
Geometric mean: $GM_x = \exp\left(\frac{1}{n} \sum_{i=1}^n \log(x_i)\right) = \sqrt[n]{x_1 x_2 \dots x_n}$

Harmonic mean: $HM_x = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}}$

$$HM \leq GM \leq AM$$

(always!)

Examples: Arithmetic, geometric, and harmonic means



→ See script *Mod 2(1).R* to explore other examples

Data summaries

MSL

FISH

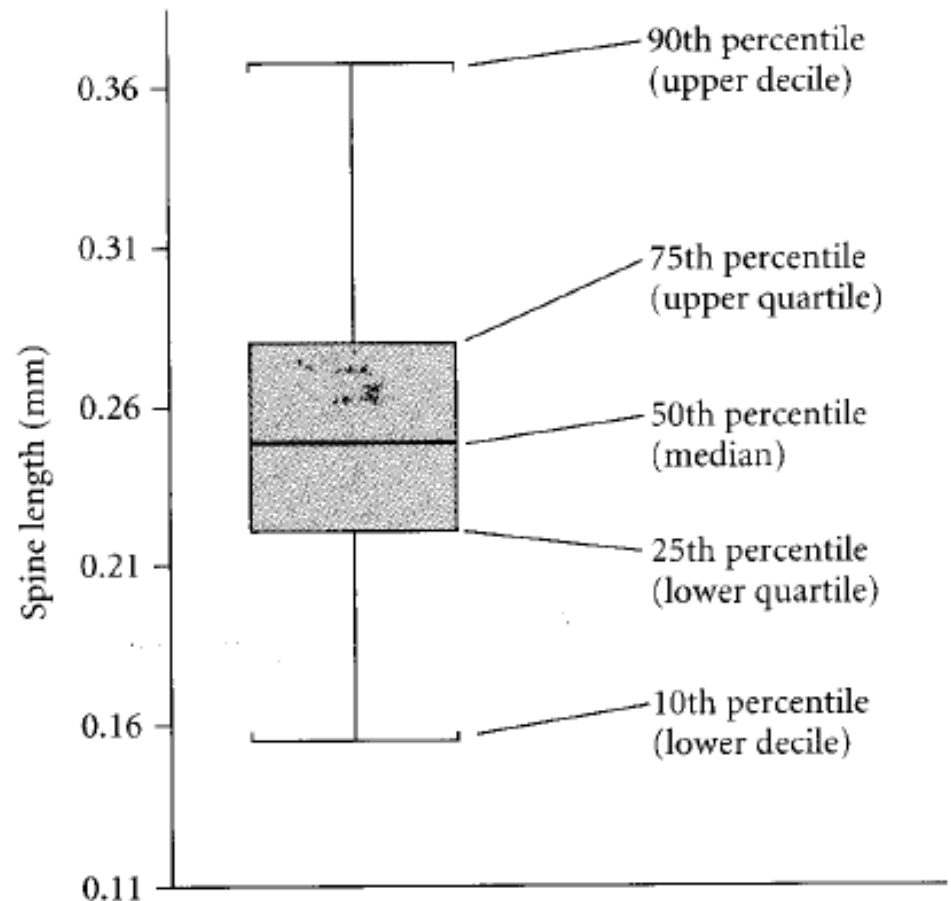
604

- Data: x_1, x_2, \dots, x_n
- Other measures of location:
 - Median (robust measure)
 - $\Pr(x_i < \text{Median}) = 0.5$
 - 50% of data points < Median < 50% of data points
 - Quantiles
 - Quartiles
 - Quantiles corresponding to any probability

→ See R script '*Mod 2(1).R*'

Boxplots

Figure 3.6 Box plot illustrating quantiles of data from Table 3.1 ($n = 50$). The line indicates the 50th percentile (median), and the box encompasses 50% of the data, from the 25th to the 75th percentile. The vertical lines extend from the 10th to the 90th percentile.





MSL
FISH
604

Data summaries

Measures of spread / dispersion

- Variance / Standard deviation
- Coefficient of Variation:

Data summaries

Measures of spread / dispersion

- Variance / Standard deviation

$$Var(x) = \sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)$$

$$SD(x) = s = \sqrt{Var(x)}$$

- Coefficient of Variation: $CV(x) = s / \bar{x}$

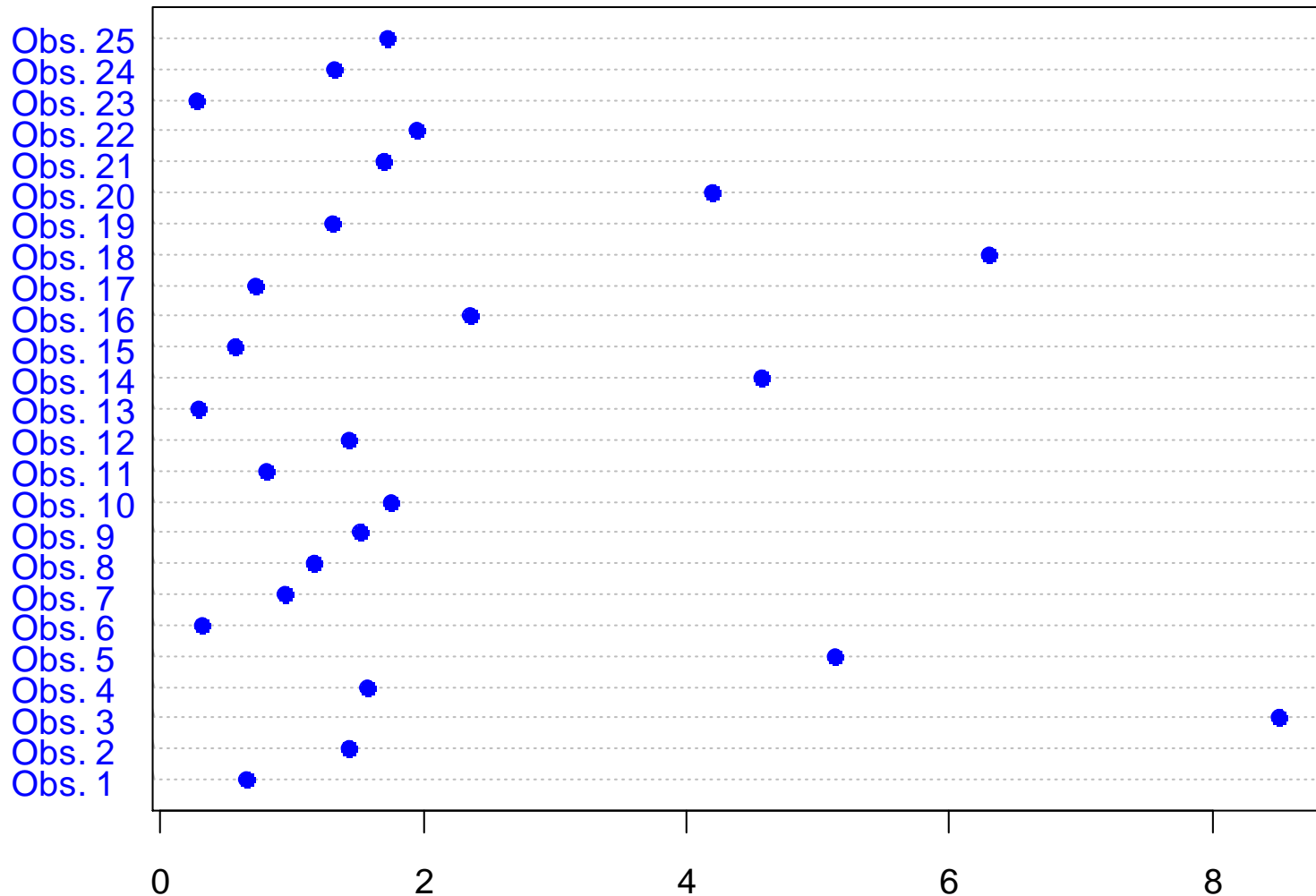
- Median Absolute Deviation (robust measure)

$$MAD(x) = 1.4826 \cdot Median|(\mathbf{x} - \bar{x})|$$

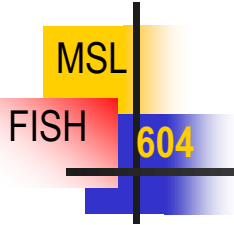
- Range, interquartile range

Graphical data summaries

Distribution of log-normal variable, $n = 25$

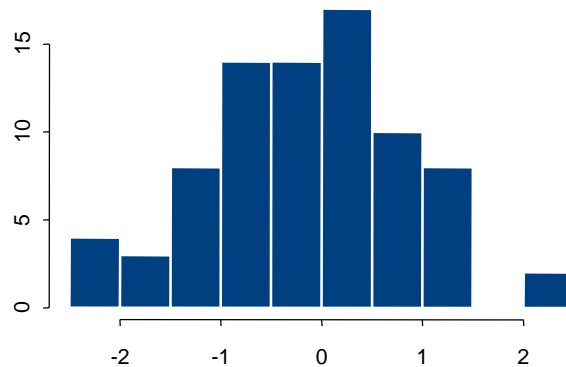


Graphical data summaries

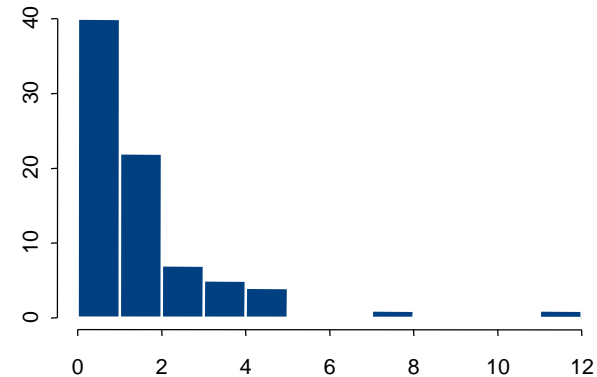


Histograms: Simulated data

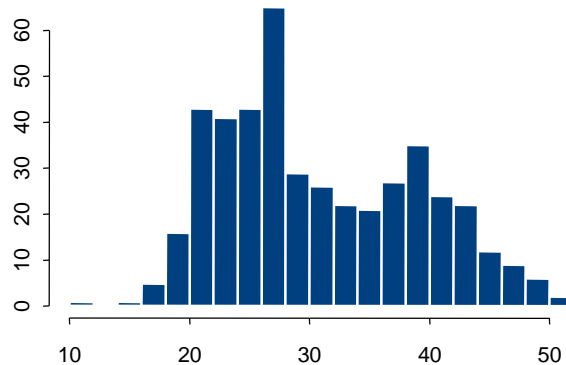
Standardized normal



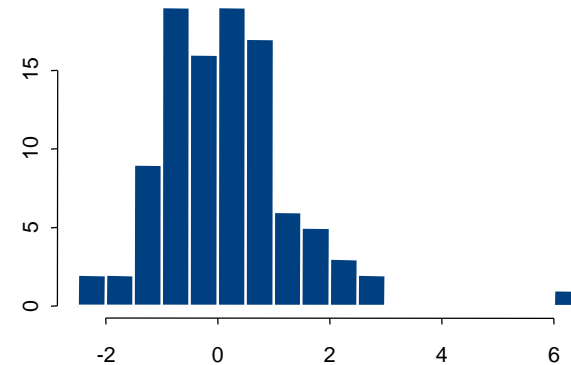
Log-normal



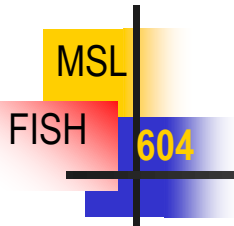
Mixture of two Normals



Normal with outlier

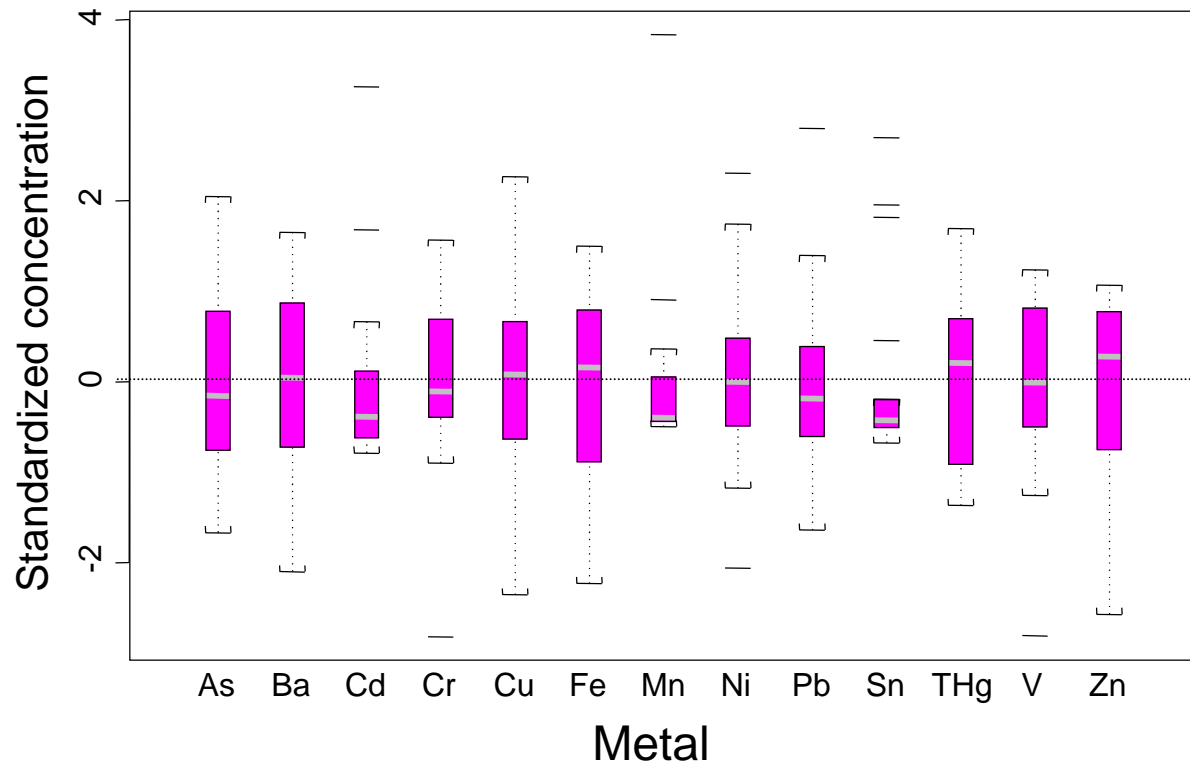


Graphical data summaries



Boxplots

Standardized concentrations of 13 metals in Beaufort Lagoon sediments





Expectation and Variance

Expectation or expected value (mean)

- Formal definition:

(discrete)

(continuous)

- Some useful theorems:

(assuming independence)

Expectation and Variance



Expectation or expected value (mean)

■ Formal definition:

$$E(x) = x_1P(X = x_1) + \dots + x_nP(X = x_n) = \sum_{j=1}^n x_jP(X = x_j) \quad \textbf{(discrete)}$$

$$E(x) = \int_{-\infty}^{\infty} xf(x)dx \quad \textbf{(continuous)}$$

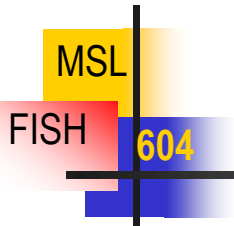
■ Some useful theorems:

$$E(cX) = cE(X)$$

$$E(X + Y) = E(X) + E(Y)$$

$$E(XY) = E(X)E(Y) \quad \textbf{(assuming independence)}$$

Example (part 1: Expected value)



What was the total catch C (in pounds) of age-3 pollock caught in the Eastern Bering Sea in 2011?

- Estimated catch at age-3 in numbers from stock assessment (Ianelli et al. 2012):

$$E(N_3) = 193.2 \text{ million fish}$$

with an assumed error of 5% ($CV = 0.05$)

- Mean weight of age-3 fish from random sampling was $E(w_3) = 0.290 \text{ kg}$ with a sampling error of 2% ($CV = 0.02$)
- $1 \text{ kg} = 2.21 \text{ pounds}$ (conversion factor $f = 2.21 \text{ lb/kg}$)

→ Compute expected value of catch:

$$E(C) = E(f * w_3 * N_3)$$

Expected value of catch:

$$\begin{aligned} E(C) &= E(f * w_3 * N_3) \\ &= ?? \end{aligned}$$

'f' is a constant!

'w₃', 'N₃' are independent!

Variance

- Quantifies "spread" or uncertainty
- Definition:
- Some useful theorems:

Variance

- Quantifies “spread” or uncertainty

- Definition: $\text{Var}(X) = \sigma^2 = E[(X - \mu)^2]$

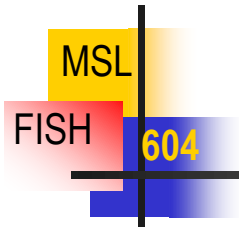
- Some useful theorems:

$$\text{Var}(X) = E[(X - \mu)^2] = E(X^2) - \mu^2$$

$$\text{Var}(cX) = c^2 \cdot \text{Var}(X)$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + \underbrace{2\text{Cov}(X, Y)}$$

=0 if X,Y uncorrelated



Covariance

- Definition:
- Useful theorems:

Covariance

MSL

FISH

604

- Definition:

$$\begin{aligned}\text{Cov}(X, Y) &= E[(X - \mu_X) \cdot (Y - \mu_Y)] \\ &= E(XY) - E(X)E(Y)\end{aligned}$$

- Useful theorems:

$$\begin{aligned}\text{Cov}(aX, bY) &= ab \cdot \text{Cov}(X, Y) \\ \text{Cov}(X, Y + Z) &= \text{Cov}(X, Y) + \text{Cov}(X, Z)\end{aligned}$$



Correlation

- Definition:

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

(dimensionless)

$$-1 < \rho < 1$$

Variance of a product

- General law for variance of a product:

$$\text{var}(x_1 x_2) \approx \mu_2^2 \text{var}(x_1) + \mu_1^2 \text{var}(x_2) + 2 \text{cov}(x_1 x_2) \mu_1 \mu_2$$

where: $\mu_i = E(x_i)$

- Two independent random variables:

$$\text{var}(x_1 x_2) = \mu_2^2 \text{var}(x_1) + \mu_1^2 \text{var}(x_2) + \text{var}(x_1) \text{var}(x_2)$$

Example (part 2: Variance)

MSL

FISH

604

What is the variance of the total catch of age-3 pollock caught in the Eastern Bering Sea in 2011?

- Estimated catch at age-3 in numbers from stock assessment (Ianelli et al. 2012):

$$E(N_3) = 193.2 \text{ million fish}$$

with an assumed error of 5% ($CV = 0.05$)

- Mean weight of age-3 fish from random sampling was $E(w_3) = 0.290 \text{ kg}$ with a sampling error of 2% ($CV = 0.02$)
- 1 pound = 0.453 kg (conversion factor $c = 0.453$)

→ Compute variance of catch:

$$\text{var}(C) = \text{var}(c * w_3 * N_3)$$

Variance, catch of age-3 pollock:

$$\begin{aligned}\text{var}(C) &= \text{var}(f * w_3 * N_3) \\ &= ??\end{aligned}$$

Variance, catch of age-3 pollock:

$$\begin{aligned}\text{var}(C) &= \text{var}(f * w_3 * N_3) \\ &= ?? \\ &= f^2 * \text{var}(w_3 * N_3) \\ &= f^2 * \{E(N_3)^2 * \text{var}(w_3) + E(w_3)^2 * \text{var}(N_3) \\ &\quad + \text{var}(w_3) * \text{var}(N_3)\}\end{aligned}$$

Given: $CV(N_3) = \text{sd}(N_3) / E(N_3) = 0.05$

hence: $\text{sd}(N_3) = 0.05 * E(N_3)$

and: $\text{var}(N_3) = \text{sd}(N_3)^2 = 0.05^2 * E(N_3)^2$

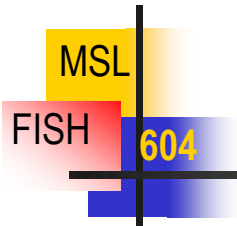
similarly: $\text{var}(w_3) = 0.02^2 * E(w_3)^2$

Substituting these expressions for $\text{var}(N_3)$ and $\text{var}(w_3)$ above yields: $\text{var}(C) = 44.28$

$$\text{sd}(C) = 6.65$$

and a 95% CI is given by: $\{180.2, 206.2\}$

Variance estimation



- If no analytical solution can be obtained for the variance of a random variable, we can use:
 - Delta method (Approximation based on Taylor series expansion)
 - Jackknife (largely superceded by bootstrap)
 - Bootstrap (computer intensive method based on re-sampling)

Delta method

Approximate method to estimate expected value and variance of a function of X , if $E(X)=\mu_x$ and $\text{var}(X)=\sigma_x^2$ are known

- Let $Y = g(X)$, then:

$$E(Y) = \mu_y \approx g(\mu_x) \quad (\text{rough approximation})$$

$$E(Y) = \mu_y \approx g(\mu_x) + \frac{1}{2} \sigma_x^2 \cdot g''(\mu_x) \quad (\text{better approximation})$$

$$\text{Var}(Y) = \sigma_y^2 \approx \sigma_x^2 [g'(\mu_x)]^2$$

- Example: Approximate mean and variance of the inverse of a variable: $Y = g(X) = 1/X$

Bootstrap estimate of variance

- A method to provide an estimate the variance of a summary statistic (T) from a sample alone, without making any distributional assumptions
- The sample is taken to “represent” the distribution of the whole population
- Many random “bootstrap samples” are drawn from this “population” (with replacement).
- The statistic, T , is calculated from each bootstrap sample and tallied (=bootstrap estimates of T , or T_B)
- The variance of the bootstrap estimates (T_b) provides an estimate of the variance of T (in fact, the full distribution of T_B estimates the full distribution of T)



Exercise (in class)

- Delta method example to compute variance of inverse of a variable
- Simple Bootstrap examples to estimate variance of the median

Review / preview

Basic statistical concepts

- Data summaries
 - Location: Mean, median, quantiles
 - Spread: Variance, Standard deviation, MAD, IQR
 - Graphical summaries
 - Expectation & Variance; Variance estimation
- Probability and probability distributions
- Distributions
 - Discrete (Binomial, Multinomial, Poisson)
 - Continuous (Uniform, Exponential, Normal or Gaussian, Log-normal, Gamma)

Next
time



Reading

Reading assignment

- Gotelli, N.J., and Ellison, A.M. 2004. A Primer of Ecological Statistics. **Chapter 3**

Further background (theoretical)

- Casella, G., and Berger, R. 2002. Statistical Inference. Duxbury, Pacific Grove, CA.
- Rice, J.A., 1995. Mathematical Statistics and Data Analysis. Duxbury Press, Belmont, CA.