

Algorithms and Inference

Statistical Learning

From Efron and Hastie (2016), chapter 1

Algorithms and Inference

Statistics is the science of *learning from experience*, particularly experience that arrives a little bit at a time

- the successes and failures of a new experimental drug
- the uncertain measurements of an asteroid's path toward Earth
- etc.

There are two main aspects of statistical analysis:

- the *algorithmic* aspect
- the *inferential* aspect

The distinction begins with the most basic, and most popular, statistical method, averaging.

Suppose we have observed y_1, \dots, y_n , realizations of the random variables Y_1, \dots, Y_n i.i.d. Y , applying to some phenomenon of interest, with the *parameter of interest* being $\mu = \mathbb{E}(Y)$

Averaging is the *algorithm*

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

How accurate is that number?

The standard error provides an inference on the algorithm's accuracy

$$\widehat{\text{se}} = \sqrt{\frac{1}{n} \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

It is a surprising, and crucial, aspect of statistical theory that the same data that supplies an estimate can also assess its accuracy.

Of course, $\widehat{\text{se}}$ is itself an algorithm, which could be (and is) subject to further inferential analysis concerning its accuracy.

The point is that the algorithm comes *first* and the inference follows at a *second* level of statistical consideration.

In practice this means that algorithmic invention is a more free-wheeling and adventurous enterprise, with inference playing catch-up as it strives to assess the accuracy, good or bad, of some hot new algorithmic methodology.

Basic Gaussian model

In the basic Gaussian model with Y_1, \dots, Y_n i.i.d. $Y \sim \mathcal{N}(\mu, \sigma^2)$, where

- the unknown parameter μ is the *parameter of interest*, and we wish to assess the relation of the data to that value
- σ^2 is the *nuisance parameter*, which can be either known or unknown.

One *estimator* for μ is

$$\bar{Y} \sim N(\mu, \sigma^2/n)$$

with standard error

$$\text{se}(\bar{Y}) = \sqrt{\text{Var}(\bar{Y})} = \sigma \sqrt{1/n}$$

If σ^2 is unknown, we use the estimator of the standard error

$$\hat{\text{se}}(\bar{Y}) = \hat{\sigma} \sqrt{1/n}$$

where the estimator for σ^2 is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n-1} \sim \sigma^2 \chi_{n-1}^2 / (n-1)$$

A pivot for inference about μ is

$$T = \frac{\bar{Y} - \mu}{\sigma \sqrt{1/n}} \cdot \frac{\sigma}{\hat{\sigma}} \sim \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2 / (n-1)}} \sim t_{n-1}$$

with $\Pr(-t_{n-1}^{1-\alpha/2} \leq T \leq t_{n-1}^{1-\alpha/2}) = 1 - \alpha$ where $t_{n-1}^{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the Student t distribution with $n - 1$ degrees of freedom.

A $1 - \alpha$ *confidence interval* for μ is

$$[\underline{\mu}, \bar{\mu}] = \bar{Y} \pm t_{n-1}^{1-\alpha/2} \cdot \hat{\text{se}}(\bar{Y})$$

The following simulation shows that the *coverage* probability of the confidence interval is

$$\Pr([\underline{\mu}, \bar{\mu}] \ni \mu) = 1 - \alpha$$

```
sim = function(alpha=0.05, n=5, mu=0, sigma=1){  
  y = rnorm(n, mean=mu, sd=sigma)  
  bary = mean(y)  
  hatse = sqrt( var(y) / n )  
  k = qt(alpha/2, df = n-1, lower.tail = F)  
  ci = bary + c(-1,1) * k * hatse  
  cover = (mu >= ci[1] & mu <= ci[2])  
  return(cover)  
}  
set.seed(123)  
mean( replicate(100, sim(alpha=0.05) ) )
```

```
## [1] 0.96
```

When the target not μ , but the future realization y of Y .

The point prediction is still \bar{y} , but the $1 - \alpha$ *prediction interval* is

$$[\underline{Y}, \bar{Y}] = \bar{Y} \pm t_{n-1}^{1-\alpha/2} \hat{\sigma} \cdot \sqrt{1 + 1/n}$$

This can be obtained from the pivot

$$T = \frac{\bar{Y} - Y}{\sigma \sqrt{1 + 1/n}} \cdot \frac{\sigma}{\hat{\sigma}} \sim \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2 / (n-1)}} \sim t_{n-1}$$

and guarantees that

$$\Pr([\underline{Y}, \bar{Y}] \ni Y) = 1 - \alpha$$

```
sim = function(alpha=0.05, n=5, mu=0, sigma=1){  
  y = rnorm(n, mean=mu, sd=sigma)  
  bary = mean(y)  
  k = qt(alpha/2, df = n-1, lower.tail = F)  
  pi = bary + c(-1,1) * k * sqrt( var(y) * (1 + (1/n)) )  
  y = rnorm(1, mean=mu, sd=sigma)  
  cover = (y >= pi[1] & y <= pi[2])  
  return(cover)  
}  
set.seed(123)  
mean( replicate(100, sim(alpha=0.05) ))
```

```
## [1] 0.93
```

References

- Efron and Hastie (2016) Computer-Age Statistical Inference: Algorithms, Evidence, and Data Science, Cambridge University Press