

Maximum Likelihood Estimation

CSE 312 Summer 21
Lecture 21

Important Dates!

- Real World 2 – Wednesday, Aug 11
- Review Summary 3 – Friday, Aug 13
- Problem Set 7 – Monday, Aug 16
- Final Released – Friday, Aug 13
- Final Due & Key Released – Tuesday, Aug 17

Asking The Opposite Question

So far:

Give you rules for an experiment.

Give you the event/outcome we're interested in.

You calculate/estimate/bound what the probability is.

Today:

Give you some of the rules of the experiment.

Tell you what happened.

You estimate what the rest of the rules of the experiment were.

Example

Suppose you flip a coin independently 10 times, and you see

HTTTHHTHHH

What is your estimate of the probability the coin comes up heads?

- a) $2/5$
- b) $1/2$
- c) $3/5$
- d) $55/100$

Fill out the poll everywhere so
Kushal knows how long to explain
Go to pollev.com/cse312su21

Maximum Likelihood Estimation

Idea: we got the results we got.

High probability events happen more often than low probability events.

So, guess the rules that maximize the probability of the events we saw (relative to other choices of the rules).

Since that event happened, might as well guess the set of rules for which that event was most likely.

Maximum Likelihood Estimation

Formally, we are trying to estimate a parameter of the experiment (here: the probability of a coin flip being heads).

The likelihood of an event E given a parameter θ is

$\mathcal{L}(E; \theta)$ is $\mathbb{P}(E)$ when the experiment is run with θ

We'll use the notation $\mathbb{P}(E; \theta)$ for probability when run with parameter θ where the semicolon means "extra rules" rather than conditioning

We will choose $\hat{\theta} = \operatorname{argmax}_{\theta} \mathcal{L}(E; \theta)$

argmax is the argument that produces the maximum so the θ that causes $\mathcal{L}(E; \theta)$ to be maximized.

Notation comparison

$\mathbb{P}(X|Y)$ probability of X , conditioned on the **event** Y having happened (Y is a subset of the sample space)

$\mathbb{P}(X; \theta)$ probability of X , where to properly define our probability space we need to know the extra piece of information θ . Since θ isn't an event, this is not conditioning

$\mathcal{L}(X; \theta)$ the likelihood of event X , given that an experiment was run with parameter θ . Likelihoods don't have all the properties we associate with probabilities (e.g. they don't all sum up to 1) and this isn't conditioning on an event (θ is a parameter/rule of how the event could be generated).

MLE

Maximum Likelihood Estimator

The maximum likelihood estimator of the parameter θ is:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \mathcal{L}(E; \theta)$$

θ is a variable, $\hat{\theta}$ is a number (or formula given the event).

We'll also use the notation $\hat{\theta}_{\text{MLE}}$ if we want to emphasize how we found this estimator.

The Coin Example

$$\mathcal{L}(\text{HTTTTHHTHHH} ; \theta) = \theta^6(1 - \theta)^4$$

Where is θ maximized?

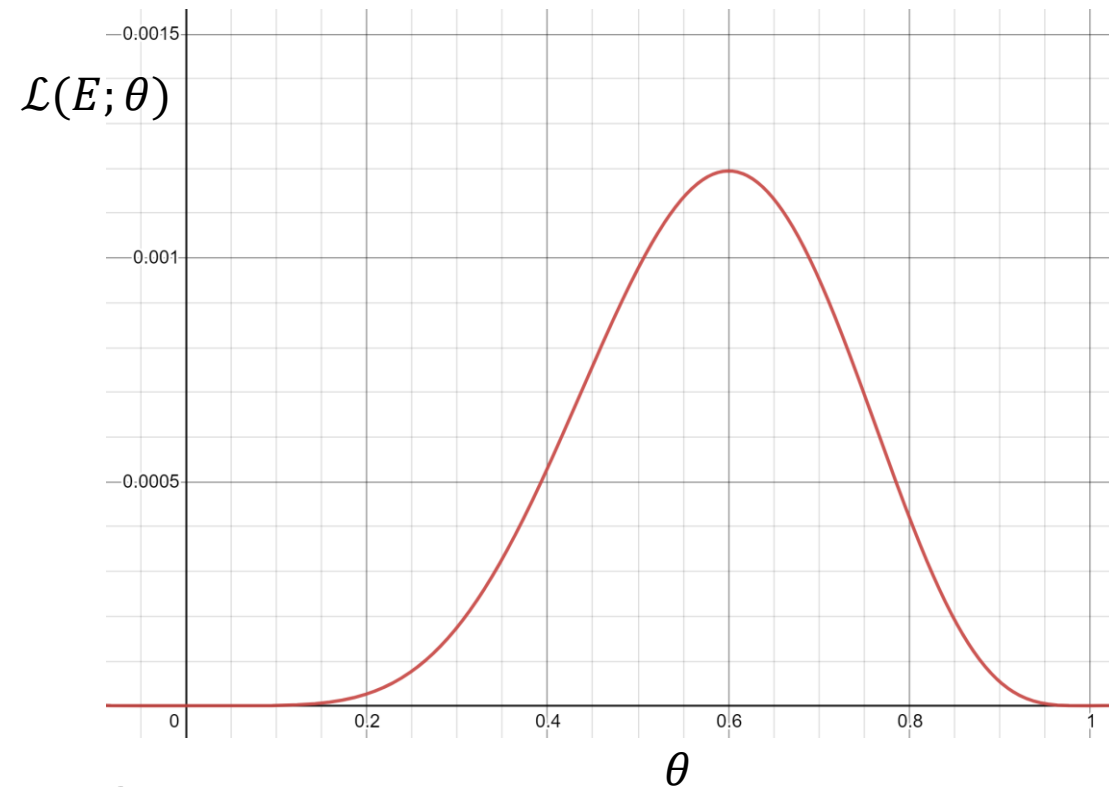
How do we usually find a maximum?

Calculus!!

$$\frac{d}{d\theta} \theta^6(1 - \theta)^4 = 6\theta^5(1 - \theta)^4 - 4\theta^6(1 - \theta)^3$$

Set equal to 0 and solve

$$6\hat{\theta}^5(1 - \hat{\theta})^4 - 4\hat{\theta}^6(1 - \hat{\theta})^3 = 0 \Rightarrow 6(1 - \hat{\theta}) - 4\hat{\theta} = 0 \Rightarrow -10\hat{\theta} = -6 \Rightarrow \hat{\theta} = \frac{3}{5}$$



The Coin Example

For this problem, θ must be in the closed interval $[0,1]$. Since $\mathcal{L}()$ is a continuous function, the maximum must occur at an endpoint or where the derivative is 0.

Evaluate $\mathcal{L}(\cdot; 0) = 0, \mathcal{L}(\cdot; 1) = 0$

at $\theta = 0.6$ we get a positive value,

so $\theta = 0.6$ is the maximizer on the interval $[0,1]$.

Maximizing a Function

CLOSED INTERVALS

Set derivative equal to 0 and solve.

Evaluate likelihood at endpoints and any critical points.

Maximum value must be maximum on that interval.

SECOND DERIVATIVE TEST

Set derivative equal to 0 and solve.

Take the second derivative. If negative everywhere, then the critical point is the maximizer.

A Math Trick

We're going to be taking the derivative of products a lot.

The product rule is not fun. There has to be a better way!

Take the log!

$$\ln(a \cdot b) = \ln(a) + \ln(b)$$

We don't need the product rule if our expression is a sum!

Can we still take the max? $\ln()$ is an increasing function, so

$$\operatorname{argmax}_{\theta} \ln(\mathcal{L}(E; \theta)) = \operatorname{argmax}_{\theta} \mathcal{L}(E; \theta)$$

Coin flips is easier

$$\mathcal{L}(\text{HTTTTHHTHHH}; \theta) = \theta^6(1 - \theta)^4$$

$$\ln(\mathcal{L}(\text{HTTTTHHTHHH}; \theta)) = 6 \ln(\theta) + 4 \ln(1 - \theta)$$

$$\frac{d}{d\theta} \ln(\mathcal{L}(\cdot)) = \frac{6}{\theta} - \frac{4}{1-\theta}$$

Set to 0 and solve:

$$\frac{6}{\hat{\theta}} - \frac{4}{1-\hat{\theta}} = 0 \Rightarrow \frac{6}{\hat{\theta}} = \frac{4}{1-\hat{\theta}} \Rightarrow 6 - 6\hat{\theta} = 4\hat{\theta} \Rightarrow \hat{\theta} = \frac{3}{5}$$

$\frac{d^2}{d\theta^2} = \frac{-6}{\theta^2} - \frac{4}{(1-\theta)^2} < 0$ everywhere, so any critical point must be a maximum.

What about continuous random variables?

Can't use probability, since the probability is going to be 0.

Can use the density!

It's supposed to show relative chances, that's all we're trying to find anyway.

$$\mathcal{L}(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n f_X(x_i; \theta)$$

Continuous Example

Suppose you get values x_1, x_2, \dots, x_n from independent draws of a normal random variable $\mathcal{N}(\mu, 1)$ (for μ unknown)

We'll also call these "realizations" of the random variable.

$$\mathcal{L}(x_i; \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2\right)$$

$$\ln(\mathcal{L}(x_i; \mu)) = \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{1}{2}(x_i - \mu)^2$$

Finding $\hat{\mu}$

$$\ln(\mathcal{L}) = \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{1}{2} (x_i - \mu)^2$$

$$\frac{d}{d\mu} \ln(\mathcal{L}) = \sum_{i=1}^n x_i - \mu$$

Setting $\mu = 0$ and solving:

$$\sum_{i=1}^n x_i - \hat{\mu} = 0 \Rightarrow \sum_{i=1}^n x_i = \hat{\mu} \cdot n \Rightarrow \hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}$$

Check using the second derivative test:

$$\frac{d^2}{d\mu^2} \ln(\mathcal{L}) = -n$$

Second derivative is negative everywhere, so log-likelihood is concave down and average of the x_i is a maximizer.

Summary

Given: an event E (usually n i.i.d. samples from a distribution with unknown parameter θ).

1. Find likelihood $\mathcal{L}(E; \theta)$

Usually $\prod \mathbb{P}(x_i; \theta)$ for discrete and $\prod f(x_i; \theta)$ for continuous

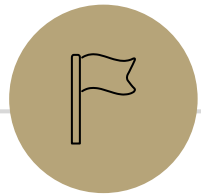
2. Maximize the likelihood. Usually:

A. Take the log (if it will make the math easier)

B. Take the derivative

C. Set the derivative to 0 and solve

3. Use the second derivative test to confirm you have a maximizer



Two Parameter Estimation

Two Parameter Estimation Setup

We just saw that to estimate μ for $\mathcal{N}(\mu, 1)$ we get:

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}$$

Now what happens if we know our data is $\mathcal{N}()$ but nothing else. Both the mean and the variance are unknown.

Log-likelihood

Let θ_μ and θ_{σ^2} be the unknown mean and standard deviation of a normal distribution. Suppose we get independent draws x_1, x_2, \dots, x_n .

$$\mathcal{L}(x_1, \dots, x_n; \theta_\mu, \theta_{\sigma^2}) = \prod_{i=1}^n \frac{1}{\sqrt{\theta_{\sigma^2} 2\pi}} \exp\left(-\frac{1}{2} \cdot \frac{(x_i - \theta_\mu)^2}{\theta_{\sigma^2}}\right)$$

$$\ln\left(\mathcal{L}(x_i; \theta_\mu, \theta_{\sigma^2})\right) = \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{\theta_{\sigma^2} 2\pi}}\right) - \frac{1}{2} \cdot \frac{(x_i - \theta_\mu)^2}{\theta_{\sigma^2}}$$

Expectation

Arithmetic is nearly identical to known variance case.

$$\ln(\mathcal{L}(x_i; \theta_\mu, \theta_{\sigma^2})) = \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{\theta_{\sigma^2} 2\pi}}\right) - \frac{1}{2} \cdot \frac{(x_i - \theta_\mu)^2}{\theta_{\sigma^2}}$$

$$\frac{\partial}{\partial \theta_\mu} \ln(\mathcal{L}) = \sum_{i=1}^n \frac{(x_i - \theta_\mu)}{\theta_{\sigma^2}}$$

Setting equal to 0 and solving

$$\sum_{i=1}^n \frac{(x_i - \widehat{\theta}_\mu)}{\theta_{\sigma^2}} = 0 \Rightarrow \sum_{i=1}^n (x_i - \widehat{\theta}_\mu) = 0 \Rightarrow \sum_{i=1}^n x_i = n \cdot \widehat{\theta}_\mu \Rightarrow \widehat{\theta}_\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$\frac{\partial^2}{\partial \theta_\mu^2} = -\frac{n}{\theta_{\sigma^2}}$$

θ_{σ^2} is an estimate of a variance. It'll never be negative (and as long as the draws aren't identical it won't be 0). So, the second derivative is negative, and we really have a maximizer.

Variance

$$\begin{aligned}\ln(\mathcal{L}(x_i; \theta_\mu, \theta_{\sigma^2})) &= \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{\theta_{\sigma^2} 2\pi}}\right) - \frac{1}{2} \cdot \frac{(x_i - \theta_\mu)^2}{\theta_{\sigma^2}} \\ &= \sum_{i=1}^n -\frac{1}{2} \ln(\theta_{\sigma^2}) - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \cdot \frac{(x_i - \theta_\mu)^2}{\theta_{\sigma^2}} \\ &= -\frac{n}{2} \ln(\theta_{\sigma^2}) - \frac{n \cdot \ln(2\pi)}{2} - \frac{1}{2\theta_{\sigma^2}} \sum_{i=1}^n (x_i - \theta_\mu)^2\end{aligned}$$

$$\frac{\partial}{\partial \theta_{\sigma^2}} \ln(\mathcal{L}) = -\frac{n}{2\theta_{\sigma^2}} + \frac{1}{2(\theta_{\sigma^2})^2} \sum_{i=1}^n (x_i - \theta_\mu)^2$$

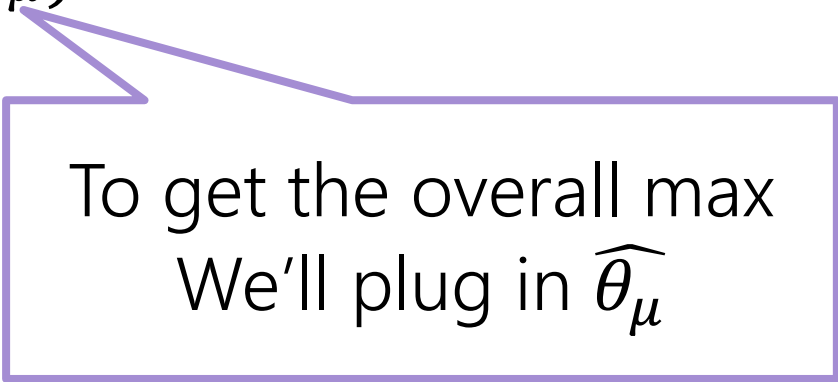
Variance

$$\frac{\partial}{\partial \theta_{\sigma^2}} \ln(\mathcal{L}) = -\frac{n}{2\theta_{\sigma^2}} + \frac{1}{2(\theta_{\sigma^2})^2} \sum_{i=1}^n (x_i - \theta_{\mu})^2$$

$$-\frac{n}{2\widehat{\theta}_{\sigma^2}} + \frac{1}{2(\widehat{\theta}_{\sigma^2})^2} \sum_{i=1}^n (x_i - \theta_{\mu})^2 = 0$$

$$\Rightarrow -\frac{n}{2}\widehat{\theta}_{\sigma^2} + \frac{1}{2} \sum_{i=1}^n (x_i - \theta_{\mu})^2 = 0 \text{ (multiply by } (\widehat{\theta}_{\sigma^2})^2 \text{)}$$

$$\Rightarrow \widehat{\theta}_{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \theta_{\mu})^2$$



To get the overall max
We'll plug in $\widehat{\theta}_{\mu}$

Summary

If you get independent samples x_1, x_2, \dots, x_n from a $\mathcal{N}(\mu, \sigma^2)$ where μ and σ^2 are unknown, the maximum likelihood estimates of the normal is:

$$\widehat{\theta}_\mu = \frac{\sum_{i=1}^n x_i}{n} \quad \text{and} \quad \widehat{\theta}_{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \widehat{\theta}_\mu)^2$$

The maximum likelihood estimator of the mean is the **sample mean** that is the estimate of μ is the average value of all the data points.

The MLE for the variance is: the variance of the experiment "choose one of the x_i at random"