

Data Analysis and Uncertainty

Part 2: Estimation

Instructor: Sargur N. Srihari

University at Buffalo
The State University of New York
srihari@cedar.buffalo.edu

Topics in Estimation

1. Estimation
2. Desirable Properties of Estimators
3. Maximum Likelihood Estimation
Examples: Binomial, Normal
4. Bayesian Estimation
Examples: Binomial, Normal
 1. Jeffrey's Prior

Estimation

- In inference we want to make statements about entire population from which sample is drawn
- Two most important methods for estimating parameters of a model:
 1. Maximum Likelihood Estimation
 2. Bayesian Estimation

Desirable Properties of Estimators

- Let $\hat{\theta}$ be an estimate of Parameter θ
- Two measures of estimator quality $\hat{\theta}$
 1. Expected Value of Estimate (Bias) Expectation over all possible data sets of size n
 - Difference between expected and true value

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

– Measures Systematic departure from true value

2. Variance of Estimate

$$\text{Var}(\hat{\theta}) = E[\hat{\theta} - E[\hat{\theta}]]^2$$

– Data driven component of error in estimation procedure

– E.g., Always saying $\hat{\theta} = 1$ has a variance of zero but high bias

- Mean Squared Error $E[(\hat{\theta} - \theta)^2]$ can be partitioned as sum of bias² and variance

Bias-Variance in Point Estimate

True height of Chinese emperor: 200cm, about 6'6"
Poll a random American: ask "How tall is the emperor?"
We want to determine how wrong they are, on average



Each scenario has expected value of 180 (or bias error = 20), but increasing variance in estimate

200 — Bias
180 — No variance

- **Scenario 1**
- Everyone believes it is 180 (variance=0)
- Answer is always 180
- The error is always -20
- Ave squared error is 400
- Average bias error is 20
- 400=400+0

Bias 200 —
Some variance 180 —

- **Scenario 2**
- Normally distributed beliefs with mean 180 and std dev 10 (variance 100)
- Poll two: One says 190, other 170
- Bias Errors are -10 and -30
 - Average bias error is -20
- Squared errors: 100 and 900
 - Ave squared error: 500
- 500 = 400 + 100

Bias 200 —
More variance 180 —

- **Scenario 3**
- Normally distributed beliefs with mean 180 and std dev 20 (variance=400)
- Poll two: One says 200 and other 160
- Errors: 0 and -40
 - Ave error is -20
- Sq. errors: 0 and 1600
 - Ave squared error: 800
- 800 = 400 + 400

Squared error = Square of bias error + Variance
As variance increases, error increases

Mean Squared Error as a Criterion for $\hat{\theta}$

- Natural decomposition as sum of squared bias and its variance

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta)^2] \\ &= (E[\hat{\theta}] - \theta)^2 + E[(\hat{\theta} - E[\hat{\theta}])^2] \\ &= (\text{Bias}(\hat{\theta}))^2 + \text{Var}(\hat{\theta}) \end{aligned}$$

- Mean squared error (over data sets) is a useful criterion since it incorporates both bias and variance

Maximum Likelihood Estimation

- Most widely used method for parameter estimation
- Likelihood Function is probability that data D would have arisen for a given value of θ

$$\begin{aligned}L(\theta \mid D) &= L(\theta \mid x(1), \dots, x(n)) \\ &= p(x(1), \dots, x(n) \mid \theta) \\ &= \prod_{i=1}^n p(x(i) \mid \theta)\end{aligned}$$

- A scalar function of θ
- Value of θ for which the data has the highest probability is the MLE

Example of MLE for Binomial

- Customers either purchase/ not purchase milk
 - We want estimate of proportion purchasing
- Binomial with unknown parameter θ
- Binomial is a generalization of Bernoulli
 - Bernoulli: probability of a binary variable $x=1$ or 0

Denoting $p(x=1)=\theta$ probability mass function

$$\text{Bern}(x|\theta) = \theta^x (1-\theta)^{1-x}$$

- Mean= θ , variance = $\theta(1-\theta)$
- Binomial: probability of r successes in n trials

$$\text{Bin}(r|n,\theta) = \binom{n}{r} \theta^r (1-\theta)^{n-r}$$

- Mean= $n\theta$, variance = $n\theta(1-\theta)$

Likelihood Function: Binomial/Bernoulli

- Samples $x(1), \dots, x(1000)$ where r purchase milk
- Assuming conditional independence, likelihood function is

$$L(\theta | x(1), \dots, x(1000)) = \prod_i \theta^{x(i)} (1 - \theta)^{n-x(i)} = \theta^r (1 - \theta)^{1000-r}$$

- Binomial pdf includes every possible way of getting r successes so it has ${}^n C_r$ additive terms
- Log-likelihood Function

$$l(\theta) = \log L(\theta) = r \log \theta + (1000 - r) \log(1 - \theta)$$

- Differentiating and setting equal to zero

9

$$\hat{\theta}_{ML} = r/1000$$

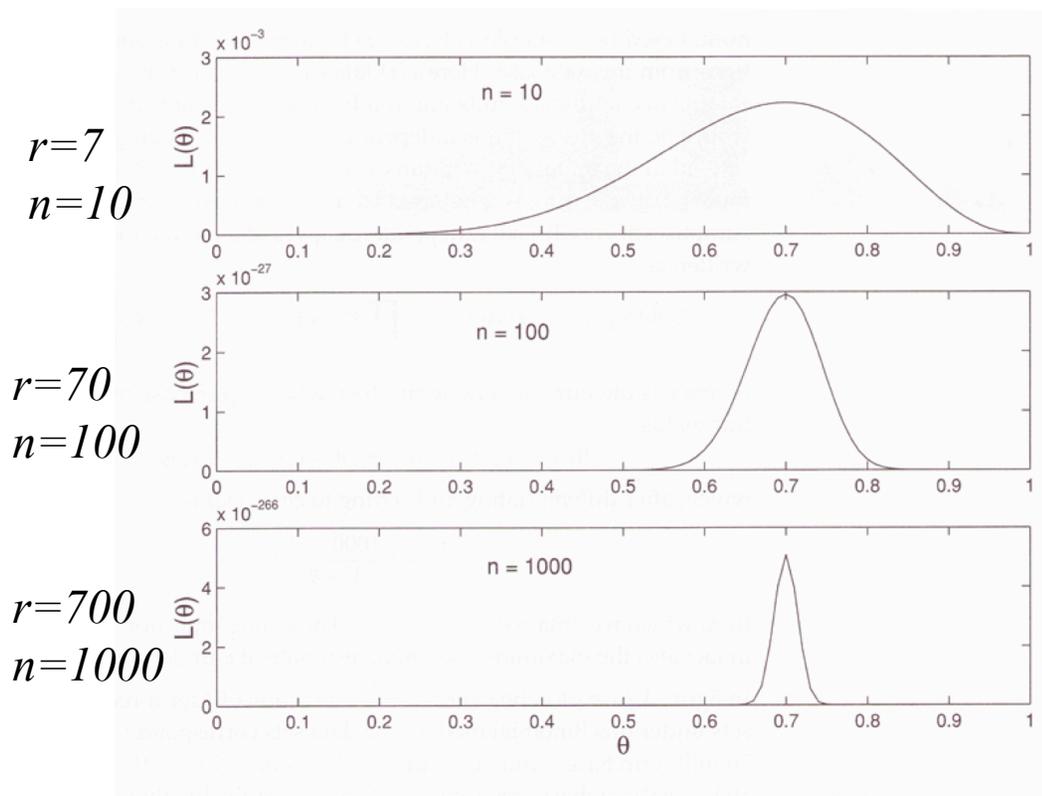
Binomial: Likelihood Functions

Likelihood function for three data sets

Binomial distribution

r milk purchases out of n customers

θ is the probability that milk is purchased by random customer



Uncertainty becomes smaller as n increases

Likelihood under Normal Distribution

- Unit variance, Unknown mean
- Likelihood function

$$\begin{aligned}L(\theta | x(1), \dots, x(n)) &= \prod_{i=1}^n (2\pi)^{-1/2} \exp\left(-\frac{1}{2}(x(i) - \theta)^2\right) \\ &= (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x(i) - \theta)^2\right)\end{aligned}$$

- Log-likelihood function

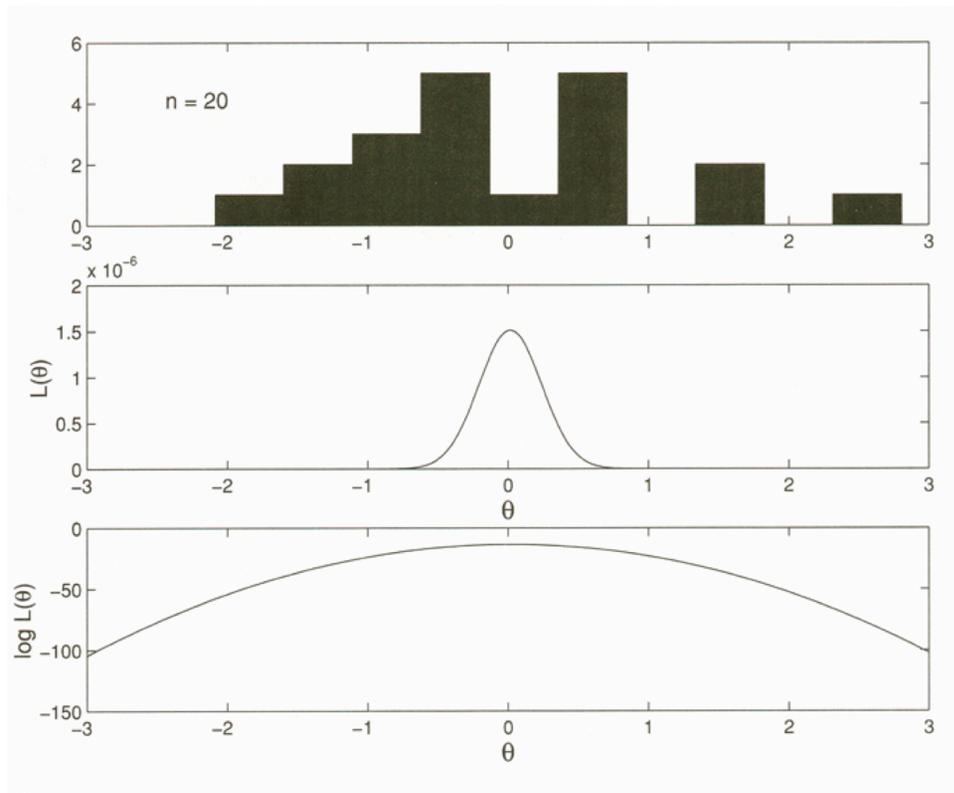
$$l(\theta | x(1), \dots, x(n)) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^n (x(i) - \theta)^2$$

- To find mle set derivative $d/d\theta$ to zero

$$\hat{\theta}_{ML} = \sum_i x(i) / n$$

Normal: Histogram, Likelihood Log-Likelihood

Estimate unknown mean θ

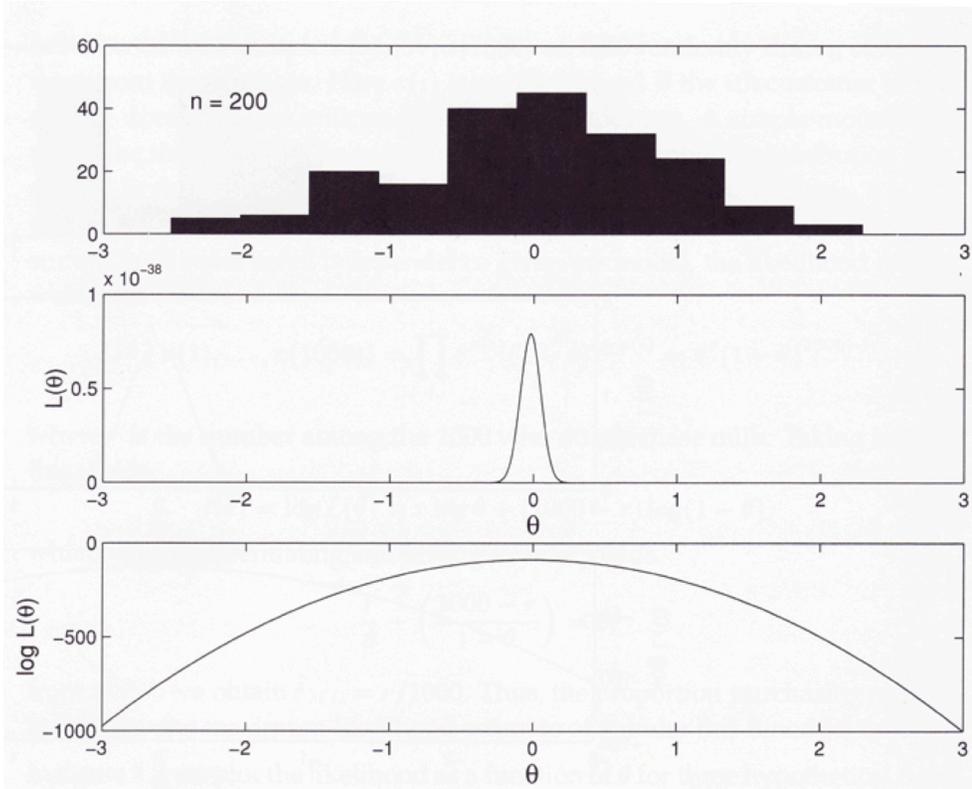


Histogram of 20 data points drawn from zero mean, unit variance

Likelihood function

Log-Likelihood function

Normal: More data points



Histogram of 200 data points drawn from zero mean, unit variance

Likelihood function

Log-Likelihood function

Sufficient Statistics

- Useful general concept in statistical estimation
- Quantity $s(D)$ is a sufficient statistic for θ if the likelihood $l(\theta)$ only depends on the data through $s(D)$
- Examples
 - For binomial parameter θ number of successes r is sufficient
 - For the mean of normal distribution sum of the observations $\sum x(i)$ is sufficient for likelihood function of mean (which is a function of θ)

Interval Estimate

- Point estimate does not convey uncertainty associated with it
- Interval estimates provide a confidence interval
- Example:
 - 100 observations from $N(\text{unknown } \mu, \text{known } \sigma^2)$
 - we want 95% confidence interval for estimate of μ
 - Since distribution of sample mean is $N(\mu, \sigma^2/100)$
 - 95% of normal distribution lies within 1.96 std dev of mean

$$P(\mu - 1.96\sigma/10 \leq \bar{x} \leq \mu + 1.96\sigma/10) = 0.95$$

Rewritten as

$$P(\bar{x} - 1.96\sigma/10 \leq \mu \leq \bar{x} + 1.96\sigma/10) = 0.95$$

$l(x) = \bar{x} - 1.96\sigma/10$ and $u(x) = \bar{x} + 1.96\sigma/10$ is a 95% interval

Bayesian Approach

- Frequentist Approach
 - Parameters are fixed but unknown
 - Data is a random sample
 - Intrinsic variability lies in data $D=\{x(1),\dots,x(n)\}$
- Bayesian Statistics
 - Data are known
 - Parameters θ are random variables
 - θ has a distribution of values
 - $p(\theta)$ reflects degree of belief on where true parameters θ may be

Bayesian Estimation

- Distribution of probabilities for θ is the *prior* $p(\theta)$
 - Analysis of data leads to modified distribution, called the *posterior* $p(\theta/D)$

- Modification done by Bayes rule

$$p(\theta | D) = \frac{p(D | \theta)p(\theta)}{p(D)} = \frac{p(D | \theta)p(\theta)}{\int_{\psi} p(D | \psi)p(\psi)\psi}$$

- Leads to a distribution rather than single value
 - Single value obtainable: mean or mode (latter known as maximum a posteriori (MAP) method)
 - MAP and MLE of θ may well coincide
 - Since prior is flat preferring no single value
 - MLE can be viewed as a special case of MAP procedure, which in turn is a restricted form of Bayesian estimation

Summary of Bayesian Approach

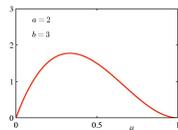
- For a given data set D and a particular model (model = distributions for prior and likelihood)

$$p(\theta | D) \propto p(D | \theta)p(\theta)$$

- In words: Posterior distribution given D (distribution conditioned on having observed the data) is
 - Proportional to product: prior $p(\theta)$ & likelihood $p(D|\theta)$
- If we have a weak belief about parameter before collecting data, choose a wide prior (normal with large variance)
- Larger the data, more dominant is the likelihood

Bayesian Binomial

- Single Binary variable X : wish to estimate $\theta = p(X=1)$
- Prior for parameter in $[0, 1]$ is the Beta distribtn

$$p(\theta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$$


Where $\alpha > 0, \beta > 0$ are two parameters of this model

$$\text{Beta}(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

$$E[\theta] = \frac{\alpha}{\alpha + \beta} \quad \text{mode}[\theta] = \frac{\alpha - 1}{\alpha + \beta - 2}$$

$$\text{var}[\theta] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

- Likelihood (same as for MLE): $L(\theta | D) = \theta^r (1 - \theta)^{n-r}$
- Combining likelihood and prior

$$p(\theta | D) \propto p(D | \theta)p(\theta) = \theta^r (1 - \theta)^{n-r} \theta^{\alpha-1} (1 - \theta)^{\beta-1} = \theta^{r+\alpha-1} (1 - \theta)^{n-r+\beta-1}$$

- We get another Beta distribution

- With parameters $r + \alpha$ and $n - r + \beta$ and mean = $E[\theta] = \frac{r + \alpha}{n + \alpha + \beta}$
- If $\alpha = \beta = 0$ we get the standard MLE of r/n

Advantages of Bayesian Approach

- Retain full knowledge of all problem uncertainty
 - Eg, calculating full posterior distribution on θ
 - E.g., Prediction of new point $x(n+1)$ not in training set D is done by averaging over all possible θ

$$\begin{aligned} p(x(n+1) | D) &= \int p(x(n+1), \theta | D) d\theta \\ &= \int p(x(n+1) | \theta) p(\theta | D) d\theta \end{aligned}$$

Since $x(n+1)$ is conditionally independent of training data D

- Can average over all possible models
 - Considerably more computation than max likelihood
- Natural updating of distribution sequentially

$$p(\theta | D_1, D_2) \propto p(D_2 | \theta) p(D_1 | \theta) p(\theta)$$

Predictive Distribution

- In equation to modify prior to posterior

$$p(\theta | D) = \frac{p(D | \theta)p(\theta)}{p(D)} = \frac{p(D | \theta)p(\theta)}{\int_{\psi} p(D | \psi)p(\psi)\psi}$$

- Denominator is called predictive distribution of D
 - Represents predictions about value of D
 - Includes uncertainty about θ via $p(\theta)$ and uncertainty about D when θ is known
- Useful for model checking
 - If observed data have only small probability then it is unlikely to be correct

Bayesian: Normal Distribution

- Suppose x comes from a Normal Distribution with unknown mean θ and known variance α

$$x \sim N(\theta, \alpha)$$

- Prior distribution for θ is $\theta \sim N(\theta_0, \alpha_0)$

$$p(\theta | x) \propto p(x | \theta)p(\theta)$$

$$= \frac{1}{\sqrt{2\pi\alpha}} \exp\left(-\frac{1}{2\alpha}(x - \theta)^2\right) \frac{1}{\sqrt{2\pi\alpha_0}} \exp\left(-\frac{1}{2\alpha_0}(x - \theta_0)^2\right)$$

- After some algebra

$$p(\theta | x) = \frac{1}{\sqrt{2\pi\alpha_1}} \exp\left(-\frac{1}{2}(\theta - \theta_1)^2 / \alpha_1\right)$$

where $\alpha_1 = (\alpha_0^{-1} + \alpha^{-1})^{-1}$ is the sum of precisions and $\theta_1 = \alpha_1(\theta_0/\alpha_0 + x/\alpha)$ weighted sum of mean and datum

- Normal prior altered to yield normal posterior

Improper Priors

- If we have no idea of the normal mean then we could give a uniform distribution for prior: but would not be a density function
- But could adopt an improper prior which is uniform over regions where parameter may occur

Jeffrey's Prior

- Results for one prior may differ from another
- Use a reference prior
- Fisher Information

$$I(\theta | x) = -E\left[\frac{\partial^2 \log L(\theta | x)}{\partial \theta^2}\right]$$

- Negative of expectation second derivative of log-likelihood
- Measures curvature or flatness of likelihood function
- Jeffrey's Prior $p(\theta) \propto \sqrt{I(\theta | x)}$
 - Consistent no matter how parameter is transformed

Conjugate Priors

- Beta to Beta
- Normal to Normal
- Advantage of using conjugate families is that updating process replaced by simple updating of parameters

Credibility Interval

- Can obtain point estimate or interval estimate from posterior distribution
- When there is a single parameter the interval containing a given probability (say 90%) is a credibility interval
- Interpretation is more straightforward than frequentist confidence intervals

Stochastic Estimation

- Bayesian methods involve complicated joint distributions
- Drawing random samples from estimated distributions enable properties of the distributions of parameters to be estimated
 - Called Markov Chain Monte Carlo (MCMC)