

• Example

$X_1, X_2, \dots, X_6 \stackrel{iid}{\sim} \text{Bern}(\theta)$

$\vec{X} = (X_1, X_2, \dots, X_6)$

If $\vec{X} = (0, 0, 1, 0, 1, 0)$,

$$p(\vec{X} = (0, 0, 1, 0, 1, 0); \theta) = \prod_{i=1}^6 \theta^{x_i} (1-\theta)^{1-x_i} = \theta^2 (1-\theta)^4$$

+ Statistical Inference

• When θ is unknown but we want to know it, we'd use Statistical Inference

• There are 3 goals:

① We want the best guess for θ ;
We'll call this $\hat{\theta}$.

② Confidence set - range of likely values of θ .

③ Test theories about θ .

• Example: Consider $X_1, \dots, X_6 \sim \text{Bern}(0.5)$;

$$p(\vec{X} = (0, 0, 1, 0, 1, 0); \theta = 0.5) = 0.5^6 = 0.0156$$

• How about $\theta = 0.25$? $\Rightarrow p(\dots) = (0.25)^2 (0.75)^4 =$

• Note that $p(\vec{X}; 0.25) > p(\vec{X}; 0.5)$ so $\theta = 0.25$ is

a better estimator than $\hat{\theta} = 0.5$.

• Pick \mathcal{T} , a class of parametric models then

use Inference to get the best $\hat{\theta}$.

+ Likelihood Function $L(\theta, \vec{x})$

• The Likelihood function $L(\theta, \vec{x}) = p(\vec{x}, \theta)$

is a function of θ , which is just the joint pdf in which \vec{x} is fixed.

It gives the "likelihood" (as a probability)

of some θ occurring.

• So, it would make sense to take $\hat{\theta}$ to be the one with the biggest likelihood - i.e., the one that maximizes $L(\theta)$.

• $L(\theta)$ is not really necessarily a pdf - e.g. $\sum_{\theta \in \Theta} L(\theta, \vec{x}) \neq 1$ always.

+ Maximum Likelihood Estimator (MLE)

• $\hat{\theta}_{MLE} := \arg \max_{\theta \in \Theta} \{L(\theta; x)\}$, $\theta \in \Theta$

is called the **MLE** & is the maximum point of $L(\theta; x)$, if it exists.

• Note that for any fn g , $\arg \max(g) = \arg \max(\ln g)$, for e.g. (fact from calc.).

- This fact will be useful when looking for the MLE, since as we know, the maximum point for a function is found by taking the derivative, and it is often easier to take the derivative of $\ln g$ than just g .

• So if $\ell(\theta; x) := \ln(L(\theta; x))$ (called the log likelihood), then $\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \{\ell(\theta; x)\}$.

+ Side Note

With likelihood, we really want to know how probable the value of θ is when given data is fixed

— this is the «inverse question»!

1/31/2018

+ More on Maximum Likelihood Estimator (MLE)

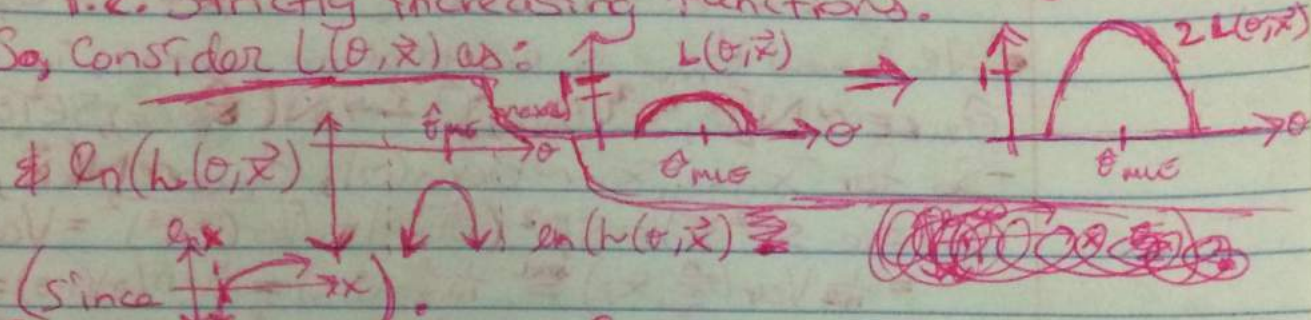
• $L(\theta; x_1, x_2, \dots, x_n) = P(X_1, \dots, X_n; \theta)$ joint density / prob function

• Goal of inferencing: looking for most likely θ .

• $\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{L(\theta; \vec{x})\} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{g(\theta)\} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{\ln(L(\theta; \vec{x}))\}$

• Note: A "positive monotonic transformation" g is one that preserves: $\forall x, y, x > y \Rightarrow g(x) > g(y)$; i.e. Strictly increasing functions.

• So, consider $L(\theta; \vec{x})$ as:



• Ex: What is $\hat{\theta}_{MLE}$ for $X_1, X_2, \dots, X_6 \sim \text{Bern}(\theta)$?

$$L(\theta; \vec{x}) = P(\vec{x}; \theta) = \prod_{i=1}^6 P(X_i; \theta) = \prod_{i=1}^6 \theta^{x_i} (1-\theta)^{1-x_i};$$

To find MLE, derive wrt θ & set equal to 0:

$$\begin{aligned} l(\theta; \vec{x}) &= \ln\left(\prod_{i=1}^6 \theta^{x_i} (1-\theta)^{1-x_i}\right) = \sum_{i=1}^6 \ln[\theta^{x_i} (1-\theta)^{1-x_i}] \\ &= \sum_{i=1}^6 [x_i \ln(\theta) + (1-x_i) \ln(1-\theta)] = 6\bar{x} \ln \theta + \ln(1-\theta) \cdot (6-6\bar{x}) \end{aligned}$$

$$l'(\theta; \vec{x}) = \frac{6\bar{x}}{\theta} - \frac{6-6\bar{x}}{1-\theta} \Rightarrow 0 \stackrel{\text{set}}{=} l'(\theta; \vec{x}) \Rightarrow \frac{\bar{x}}{\theta} - \frac{1-\bar{x}}{1-\theta} = 0 \Rightarrow$$

$$\frac{\bar{x}}{\theta} = \frac{1-\bar{x}}{1-\theta} \Rightarrow (1-\theta)\bar{x} = \theta(1-\bar{x}) \Rightarrow \theta = \bar{x} \Rightarrow \hat{\theta}_{MLE} = \bar{x} = \hat{p}.$$

So here, if $\vec{x} = \{0, 0, 1, 0, 1, 0\}$, $\hat{\theta}_{MLE} = \bar{x} = 2/6 = 1/3$.

• $\hat{\theta}_{MLE} = \bar{x}$ is an estimator & $\hat{\theta}_{MLE} = \bar{x}$ is an estimate.

\bar{x} is r.v. of average of \vec{x} , & \bar{x} is a realization of \bar{x} .

• Why do we care about MLE's?

① $\hat{\theta}_{MLE} \xrightarrow{P} \theta$ (converges in probability)

As n gets large, $\hat{\theta}_{MLE}$ looks more & more like θ . - i.e. $\forall \epsilon > 0, \lim_{n \rightarrow \infty} P(|\hat{\theta}_{MLE} - \theta| < \epsilon) = 1$.

* ② Asymptotic Normality $\hat{\theta}_{MLE} \xrightarrow{d} N(\theta, SE[\hat{\theta}_{MLE}])$.

③ Hypothesis testing via:

$$R_{\text{ref}} \text{ or } R_{\text{reg}} = \left[\theta_0 \pm \frac{2\alpha}{2} SE[\hat{\theta}_{MLE}] | \theta_0 \right]$$

③ Efficiency: $SE(\hat{\theta}_{MLE})$ is theoretically the lowest SE for all consistent estimators (proposition 1).

Ex: What is $\hat{\theta}_{MLE}$ for $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Geom}(\theta)$, where $X \sim \text{Geom}(p) \Rightarrow P(X=x) = (1-p)^{x-1} p$? $\rightarrow \theta = (0, 1)$

$$l(\theta; \vec{x}) = \prod_{i=1}^n P(X_i, \theta) = \prod_{i=1}^n (1-\theta)^{x_i-1} \theta \Rightarrow$$

$$l(\theta; \vec{x}) = \ln \left(\prod_{i=1}^n (1-\theta)^{x_i-1} \theta \right) = \sum_{i=1}^n \ln[(1-\theta)^{x_i-1} \theta] =$$

$$\sum_{i=1}^n [\ln(1-\theta)^{x_i-1} + \ln \theta] = \sum_{i=1}^n [x_i \ln(1-\theta) + \ln \theta] =$$

$$\sum_{i=1}^n x_i \ln(1-\theta) + \sum_{i=1}^n \ln \theta = n \bar{x} \ln(1-\theta) + n \ln \theta \Rightarrow$$

$$l'(\theta; \vec{x}) = \frac{n}{\theta} - \frac{n \bar{x}}{1-\theta} \Rightarrow 0 \stackrel{\text{set}}{=} n \left(\frac{1}{\theta} - \frac{\bar{x}}{1-\theta} \right) \Rightarrow \frac{1}{\theta} = \frac{\bar{x}}{1-\theta} \Rightarrow$$

$$1-\theta = \theta \bar{x} \Rightarrow 1 = \theta(\bar{x}+1) \Rightarrow \theta = \frac{1}{\bar{x}+1} \Rightarrow \hat{\theta}_{MLE} = \frac{1}{\bar{x}+1}.$$

We know $\hat{\theta}_{MLE}$ is asymptotic normal; i.e.

$$\hat{\theta}_{MLE} \sim N[\theta, SE^2(\hat{\theta}_{MLE})] \xrightarrow{P} N[\hat{\theta}_{MLE}, SE^2(\hat{\theta}_{MLE})]$$

- So for $X_1, \dots, X_n \sim \text{Bern}(\theta)$, $\hat{\theta}_{MLE} \sim N(\bar{x}, SE^2(\bar{x}))$,

$$\text{where } SE^2(\bar{x}) = SE^2\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \text{Var}\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \text{Var}\left(\frac{\sum_{i=1}^n \frac{X_i}{n}\right)$$

$$= \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \stackrel{iid}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \cdot n \cdot \theta(1-\theta) = \frac{\theta(1-\theta)}{n}$$

- for $X_1, \dots, X_n \sim \text{Geom}(\theta)$, $\hat{\theta}_{MLE} \sim N\left(\frac{1}{\bar{x}+1}, SE^2\left(\frac{1}{\bar{x}+1}\right)\right)$, Where SE is not so easy to calculate.

Goals of Inference Based on MLE (Frequentist Model)

① Point of estimation for θ .

② Confidence sets for θ ; (region of likely values).

③ Testing theories about θ .

Inference via MLEs

$\hat{\theta}_{MLE}$ is a ~~continuous~~ number.

(1) Confidence Interval.

$$CI_{\alpha, 1-\alpha} := [\hat{\theta}_{MLE} \pm z_{\alpha/2} SE[\hat{\theta}_{MLE}]].$$

(3) Hypothesis Testing

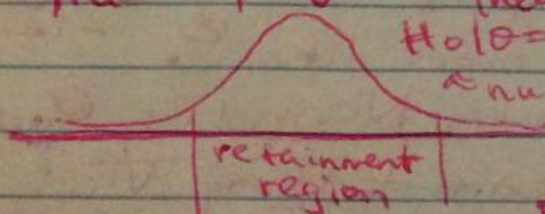
Two-Sided hypothesis test for θ :

$H_0: \theta = \theta_0$ - default theory

$H_a: \theta \neq \theta_0$ - theory you wish to prove.

$H_0: \theta = \theta_0$

A null distribution



Justified by CLT

$$\text{For } \text{Bern}(\theta), \text{ Ret. Region } \alpha := \left[\theta_0 \pm z_{\alpha/2} \sqrt{\frac{\theta_0(1-\theta_0)}{n}} \right].$$

level of confidence best guess

Recipe for Inference

Observe data \rightarrow Pick θ_0 (Parametric model) \rightarrow Do Inference via MLE.

- What's wrong with this?