

## LECTURE 16 - BIASED/UNBIASED ESTIMATORS AND CONJUGATE PRIORS

### 1. BIASED AND UNBIASED ESTIMATORS

- Last class we introduced the concept of Maximum Likelihood Estimation
- Maximum Likelihood Estimation (MLE) is a method of estimating parameters by maximizing the data likelihood with respect to those parameters
- Suppose we have assume a Gaussian data likelihood. What is the MLE solution for the mean?
- The MLE solution for the variance is
- A point estimate (e.g., the MLE estimate of the mean or variance) is said to be *unbiased* if the expected value of that estimator equals the “true” underlying value.
- In other words, if you have data generated from a Gaussian distribution and you used that data to compute the MLE of the mean of the Gaussian distribution. If the expected value of the MLE of the mean equals the true mean of the Gaussian, then the MLE solution is *unbiased*. If the expected value does not equal the true value, then the point estimate is *biased*.

#### 1.1. MLE of Mean of Gaussian.

- Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$  be samples from a multi-variance Normal distribution with known covariance matrix and an unknown mean. Given this data, obtain the ML estimate of the mean vector.

$$(1) \quad p(\mathbf{x}_k|\mu) = \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}_k - \mu)^T \Sigma^{-1} (\mathbf{x}_k - \mu) \right)$$

- We can define our likelihood given the  $N$  data points. We are assuming these data points are drawn independently but from an identical distribution (i.i.d.):

$$(2) \quad \prod_{n=1}^N p(\mathbf{x}_n|\mu) = \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right)$$

– We can apply our “trick” to simplify

$$(3) \quad \mathcal{L} = \ln \prod_{n=1}^N p(\mathbf{x}_n | \mu) = \ln \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right)$$

$$(4) \quad = \sum_{n=1}^N \ln \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right)$$

$$(5) \quad = \sum_{n=1}^N \left( \ln \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} + \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right) \right)$$

$$(6) \quad = -N \ln(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}} + \sum_{n=1}^N \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right)$$

– Now, lets maximize:

$$(7) \quad \frac{\partial \mathcal{L}}{\partial \mu} = \frac{\partial}{\partial \mu} \left[ -N \ln(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}} + \sum_{n=1}^N \left( -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right) \right] = 0$$

$$(8) \quad \rightarrow \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}_n - \mu) = 0$$

$$(9) \quad \rightarrow \sum_{n=1}^N \Sigma^{-1} \mathbf{x}_n = \sum_{n=1}^N \Sigma^{-1} \mu$$

$$(10) \quad \rightarrow \Sigma^{-1} \sum_{n=1}^N \mathbf{x}_n = \Sigma^{-1} \mu N$$

$$(11) \quad \rightarrow \sum_{n=1}^N \mathbf{x}_n = \mu N$$

$$(12) \quad \rightarrow \frac{\sum_{n=1}^N \mathbf{x}_n}{N} = \mu$$

$$(13)$$

– So, the ML estimate of  $\mu$  is the sample mean!

### 1.2. Is MLE of Mean of Gaussian a Biased or Unbiased Estimator?

– An “unbiased” estimator means that the expected value of the estimator equals the true value.

$$(14) \quad \mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N x_n \right] = \frac{1}{N} N \mathbb{E}[x] = \mathbb{E}[x] = \mu$$

– First equality makes use of the fact that the expected value of any (i.i.d.) sample is the same.

– We showed the last equality in the last class.

### 1.3. Is MLE of Variance of Gaussian a Biased or Unbiased Estimator?

– The ML estimate of the variance of a Gaussian given  $N$  samples and a known mean is:

$$(15) \quad \sigma_{MLE}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

– The expected value of this estimator is:

$$(16) \quad \mathbb{E}[\sigma_{MLE}^2] = \mathbb{E}\left[\frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2\right] = \frac{N-1}{N} \sigma^2$$

– So, the  $\mathbb{E}[\sigma_{MLE}^2] \neq \sigma^2$ . Thus, the MLE of the variance of a Gaussian is a *biased* estimate.

– The *unbiased* estimate of the variance turns out to be:

$$(17) \quad \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu)^2$$

## 2. METHOD OF MAXIMUM A POSTERIORI (MAP)

- Bayes Rule:  $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$
- Lets look at this in terms of binary variables, e.g., Flipping a coin:  $X = 1$  is “heads”,  $X = 2$  is “tails”
- Let  $\mu$  be the probability of heads. If we know  $\mu$ , then:  $P(x = 1|\mu) = \mu$  and  $P(x = 0|\mu) = 1 - \mu$

$$(18) \quad P(x|\mu) = \mu^x (1 - \mu)^{1-x} = \begin{cases} \mu & \text{if } x = 1 \\ 1 - \mu & \text{if } x = 0 \end{cases}$$

- This is called the *Bernoulli* distribution

$$(19) \quad E[x] = \mu$$

$$(20) \quad E[(x - \mu)^2] = \mu(1 - \mu)$$

- So, suppose we conducted many Bernoulli trials (e.g., flip a coin) and we want to estimate  $\mu$
- **Method: Maximum Likelihood**

$$(21) \quad p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu)$$

$$(22) \quad = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}$$

- Maximize : (*What trick should we use?*)

$$(23) \quad \mathcal{L} = \sum_{n=1}^N x_n \ln \mu + (1 - x_n) \ln(1 - \mu)$$

$$(24) \quad \frac{\partial \mathcal{L}}{\partial \mu} = 0 = \frac{1}{\mu} \sum_{n=1}^N x_n - \frac{1}{1 - \mu} \sum_{n=1}^N (1 - x_n)$$

$$(25) \quad 0 = \frac{(1 - \mu) \sum_{n=1}^N x_n - \mu \sum_{n=1}^N (1 - x_n)}{\mu(1 - \mu)}$$

$$(26) \quad 0 = \sum_{n=1}^N x_n - \mu \sum_{n=1}^N x_n - \mu \sum_{n=1}^N 1 + \mu \sum_{n=1}^N x_n$$

$$(27) \quad 0 = \sum_{n=1}^N x_n - \mu N$$

$$(28) \quad \mu = \frac{1}{N} \sum_{n=1}^N x_n = \frac{m}{N}$$

where  $m$  is the number of successful trials.

- So, if we flip a coin 1 time and get heads, then  $\mu = 1$  and probability of getting tails is 0. *Would you believe that? We need a prior!*
- Look at several independent trials. Consider  $N = 3$  and  $m = 2$  ( $N$  is number of trials,  $m$  is number of successes) and look at all ways to get 2 H and 1 T:
  - H H T  $\rightarrow \mu\mu(1 - \mu) = \mu^2(1 - \mu)$
  - H T H  $\rightarrow \mu(1 - \mu)\mu = \mu^2(1 - \mu)$
  - T H H  $\rightarrow (1 - \mu)\mu\mu = \mu^2(1 - \mu)$
- $\binom{3}{2} \mu^2(1 - \mu) \rightarrow \binom{N}{m} \mu^m(1 - \mu)^{N-m} = \frac{N!}{(N-m)!m!} \mu^m(1 - \mu)^{N-m}$
- This is the Binomial Distribution, gives the probability of  $m$  observations of  $x = 1$  out of  $N$  independent trials
- So, what we saw is that we need a prior. We want to incorporate our prior belief. Let us place a prior on  $\mu$

$$(29) \quad \text{Beta}(\mu|a, b) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1 - \mu)^{b-1}$$

$$(30) \quad E[\mu] = \frac{a}{a + b}$$

$$(31) \quad \text{Var}[\mu] = \frac{ab}{(a + b)^2(a + b + 1)}$$

- Note:  $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$  and when  $x$  is an integer, then it simplifies to  $x!$

- Calculation of the posterior, Take  $N = m + l$  observations:

$$(32) \quad p(\mu|m, l, a, b) \propto \text{Bin}(m, l|\mu) \text{Beta}(\mu|a, b)$$

$$(33) \quad \propto \mu^m (1 - \mu)^l \mu^{a-1} (1 - \mu)^{b-1}$$

$$(34) \quad = \mu^{m+a-1} (1 - \mu)^{l+b-1}$$

- What does this look like? Beta:  $a \leftarrow m + a$ ,  $b \leftarrow l + b$
- So, what's the posterior?

$$(35) \quad p(\mu|m, l, a, b) = \frac{\Gamma(m + a + l + b)}{\Gamma(m + a) \Gamma(l + b)} \mu^{m+a-1} (1 - \mu)^{l+b-1}$$

- **Conjugate Prior Relationship:** When the posterior is the same form as the prior

- Now we can maximize the (log of the) posterior:

$$(36) \quad \max_{\mu} (m + a - 1) \ln \mu + (l + b - 1) \ln(1 - \mu)$$

$$(37) \quad \frac{\partial \mathcal{L}}{\partial \mu} = 0 = \frac{m + a - 1}{\mu} - \frac{l + b - 1}{1 - \mu}$$

$$(38) \quad = (1 - \mu)(m + a - 1) - \mu(l + b - 1)$$

$$(39) \quad = (m + a - 1) - \mu(m + a - 1) - \mu(l + b - 1)$$

$$(40) \quad \mu = \frac{m + a - 1}{m + a + l + b - 2}$$

- This is the MAP solution. *So, what happens now when you flip one heads, two heads, etc.?*
- Discuss online updating of the prior. Eventually the data takes over the prior.