ECE 830 Fall 2010 Statistical Signal Processing

instructor: R. Nowak , scribe: K. Surender

Lecture 18: Bayesian Signal Processing

Most signal processing algorithms are designed and based on prior knowledge of signal and noise characteristics. It is therefore natural (and useful) to view them as Bayesian inference strategies. Let's first review the Bayesian set-up:

Prior: $p(\theta)$ - Prior probability distribution on signal parameters

Likelihood: $p(x|\theta)$ - Conditional Distribution of observation x given θ

Posterior: $p(\theta|x) = \frac{p(x,\theta)}{p(x)} = \frac{p(x|\theta)p(\theta)}{\int p(x|\theta)p(\theta) d\theta}$ - Posterior probability distribution of θ given x

Example 1 Let $S(\theta) \in \mathbb{R}^n, \theta \in \mathbb{R}^k$. S is a n-dimensional signal vector determined by $k \leq n$ parameters θ . We observe $X = S(\theta) + W, W \sim \mathcal{N}(0, \sigma^2 I)$

Prior: $p(\theta)$

Likelihood: $X|\theta \sim \mathcal{N}(S(\theta), \sigma^2 I)$

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

$$p(x|\theta)p(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} exp\left(-\frac{1}{2\sigma^2}\left(X - S(\theta)\right)^T \left(X - S(\theta)\right)\right) p(\theta)$$

A loss function is formed by taking the negative log of the posterior probability. Note that this loss function is comprised of two θ dependent terms. One term is based on the observed data and the other term is based on the prior probability.

$$-\log(p(\theta|x)) = \frac{\|X - S(\theta)\|_2^2}{2\sigma^2} - \log(p(\theta)) + constants$$

1 Point Estimators

Usually we are interested in obtaining an estimator of θ given x. Here are the two most common Bayesian estimators.

1.1 Maximum A Posteriori (MAP) Estimator

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|x)$$
$$= \arg \max_{\theta} p(x|\theta) p(\theta)$$

Note if $p(\theta) = \text{constant}$ then $\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$, where $\hat{\theta}_{MLE}$ is the maximum likelihood estimate of θ .

1.2 Posterior Mean Estimator

The posterior mean estimator is defined as:

$$\hat{\theta}_{PM} = \mathbb{E}[\theta|x] = \int \theta \, p(\theta|x) \, d\theta$$

Note that $p(\theta|x)$ requires knowledge of the normalizing term p(x). Specifically,

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\int p(x|\theta)p(\theta) \, d\theta}$$

where the denominator represents the normalization.

1.3 Posterior Mean and Bayesian MSE

The Bayesian MSE is $\mathbb{E}\left[\left\|\hat{\theta} - \theta\right\|_{2}^{2}\right]$. The estimator that minimizes this MSE is the posterior mean:

$$\begin{split} \hat{\theta} &= \arg\min_{\theta} \mathbb{E} \Big[\parallel \hat{\theta} - \theta \parallel_{2}^{2} \Big] \\ &= \min_{\tilde{\theta}} \iint \parallel \tilde{\theta} - \theta \parallel_{2}^{2} p(x,\theta) \, dx \, d\theta \\ &= \arg\min_{\tilde{\theta}} \int \parallel \tilde{\theta} - \theta \parallel_{2}^{2} p(\theta|x) \, d\theta \end{split}$$

Now differentiate with respect to $\tilde{\theta}$ and set equal to zero to find the minimizer $\hat{\theta}$

$$\frac{\partial}{\partial \tilde{\theta}} \int \| \tilde{\theta} - \theta \|_{2}^{2} p(\theta|x) d\theta = \int 2 \left(\hat{\theta} - \theta\right) p(\theta|x) d\theta = 0$$

$$\implies \hat{\theta} \int p(\theta|x) d\theta = \int \theta p(\theta|x) d\theta$$

$$\int p(\theta|x) d\theta = 1 \text{ therefore: } \hat{\theta} = \int \theta p(\theta|x) d\theta$$

$$= \hat{\theta}_{PM}$$

Example 2 Coin Tossing. Let $\theta \sim Uniform[0,1]$ be a parameter describing the probability that a coin toss lands heads. Additionally assume we observe x = 10 heads in 10 flips. Then the posterior probability is given as $p(\theta|x) = 11\theta^{10}$ and the point estimators of θ are:

$$\begin{aligned} \theta_{MAP} &= 1 \\ \hat{\theta}_{PM} &= \int_0^1 \theta \, p(\theta | x) \, d\theta = 11 \int_0^1 \theta^{11} \, d\theta = 11 \frac{\theta^{12}}{12} |_0^1 = \frac{11}{12} \end{aligned}$$

In the above example it is seen that the MAP estimator is more aggresive than the posterior mean estimator.

2 Linear Minimum MSE (LMMSE) Estimators

Suppose $p(\theta)$ is such that $\int \theta p(\theta) d\theta = 0$. Also, assume that $\mathbb{E}[x] = 0$. Consider an estimator of the form $\hat{\theta} = A^T X$ where A is a (kxn) matrix. Let's find the A to minimize the Bayesian MSE:



Figure 1: Posterior Probability Function from Example 2.

$$MSE(A) = \mathbb{E}\left[\left\| \theta - A^{T}X \right\|_{2}^{2} \right]$$

$$= \mathbb{E}\left[tr\left((\theta - A^{T}X)(\theta - A^{T}X)^{T} \right) \right]$$

$$= tr\left(\mathbb{E}\left[(\theta - A^{T}X)(\theta - A^{T}X)^{T} \right] \right)$$

$$= tr\left(\mathbb{E}[\theta\theta^{T}] - A^{T}\mathbb{E}[X\theta^{T}] - \mathbb{E}[\theta X^{T}]A + A^{T}\mathbb{E}[XX^{T}]A \right)$$

$$= tr\left(\Sigma_{\theta\theta} - A^{T}\Sigma_{x\theta} - \Sigma_{\theta x}A + A^{T}\Sigma_{xx}A \right)$$

Now differentiate with respect to A and set equal to zero to find the minimizer \hat{A}

$$\frac{\partial}{\partial A}MSE(A) = -2\Sigma_{x\theta} + 2\Sigma_{xx}\hat{A} = 0$$

The following equation is called the Wiener Hopf Equation and provides a solution for \hat{A} .

$$\begin{split} \Sigma_{x\theta} &= \Sigma_{xx} \hat{A} \\ \implies \hat{A} &= \Sigma_{xx}^{-1} \Sigma_{x\theta} \end{split}$$

Therefore the LMMSE Estimator is given by the Wiener Filter and the observed data X

-

$$\hat{\theta}_{LMMSE} = \Sigma_{x\theta} \Sigma_{xx}^{-1} X$$

 $\Sigma_{x\theta} \Sigma_{xx}^{-1}$ is the Wiener Filter

Example 3 Let $X = H\theta + W$ where $\theta \sim \mathcal{N}(0, \sigma_{\theta}^2 I)$ and $W \sim \mathcal{N}(0, \sigma_W^2 I)$ are independent. Therefore $X \sim \mathcal{N}(0, \sigma_{\theta}^2 H H^T + \sigma_W^2 I)$. To find the LMMSE define the covariance matrices:

$$\Sigma_{xx} = \sigma_{\theta}^2 H H^T + \sigma_W^2 I$$

$$\Sigma_{x\theta} = \mathbb{E}[X\theta^T] = \mathbb{E}[(H\theta + W)\theta^T] = H\Sigma_{\theta\theta} = \sigma_{\theta}^2 H$$

Now the LMMSE can be found as $\hat{A}^T X$.

$$\hat{A} = (\sigma_{\theta}^2 H H^T + \sigma_W^2 I)^{-1} \sigma_{\theta}^2 H$$

Lecture 18: Bayesian Signal Processing

$$\hat{\theta}_{LMMSE} = \sigma_{\theta}^2 H^T (\sigma_{\theta}^2 H H^T + \sigma_W^2 I)^{-1} X = H^T (H H^T + \frac{\sigma_W^2}{\sigma_{\theta}^2} I)^{-1} X$$

Note that as SNR increases the LMMSE goes to the MLE. That is:

$$As \; \frac{\sigma_{\theta}^2}{\sigma_W^2} \to \infty \quad \hat{\theta}_{LMMSE} \to (H^T H)^{-1} H^T X = \hat{\theta}_{MLE}$$

We can rewrite the LMMSE by letting H be a matrix of k orthonormal column vectors $H = [h_1 | h_2 | \cdots | h_k]$ and U an orthonormal basis for $\mathbb{R}^N U = [h_1 | h_2 | \cdots | h_k | \tilde{h}_{k+1} | \cdots | \tilde{h}_n]$.

$$\hat{\theta}_{LMMSE} = \sigma_{\theta}^{2} H^{T} (\sigma_{\theta}^{2} U \begin{bmatrix} I_{k \times k} & 0_{k \times n-k} \\ 0_{n-k \times k} & 0_{n-k \times n-k} \end{bmatrix} U^{T} + \sigma_{W}^{2} U U^{T})^{-1} X$$

$$= \sigma_{\theta}^{2} H^{T} \left(U (\sigma_{\theta}^{2} \begin{bmatrix} I_{k \times k} & 0_{k \times n-k} \\ 0_{n-k \times k} & 0_{n-k \times n-k} \end{bmatrix} + \sigma_{W}^{2}) U^{T} \right)^{-1} X = \sigma_{\theta}^{2} H^{T} U D^{-1} U^{T} X$$

Note that:

$$H^{T}U = \begin{bmatrix} I_{k \times k} & 0_{k \times (n-k)} \end{bmatrix} \quad D^{-1} = \begin{bmatrix} \frac{1}{(\sigma_{\theta}^{2} + \sigma_{W}^{2})} & & & \\ & \ddots & & \\ & & \frac{1}{(\sigma_{\theta}^{2} + \sigma_{W}^{2})} & & \\ & & & \frac{1}{\sigma_{W}^{2}} & \\ & & & & \ddots & \\ & & & & & \frac{1}{\sigma_{W}^{2}} \end{bmatrix}$$

Then:

$$H^T U D^{-1} U^T = \frac{1}{\sigma_\theta^2 + \sigma_W^2} \begin{bmatrix} I_{k \times k} & 0_{k \times n} \end{bmatrix} U^T = \frac{1}{\sigma_\theta^2 + \sigma_W^2} H^T$$

Now the LMMSE can be written as the MLE multiplied by a shrinkage term $\frac{\sigma_{\theta}^2}{\sigma_{\theta}^2 + \sigma_W^2}$ by noting that in this case $H^T H = I$.

$$\hat{\theta}_{LMMSE} = \frac{\sigma_{\theta}^2}{\sigma_{\theta}^2 + \sigma_W^2} (H^T H)^{-1} H^T X$$