

# Notes on Maximum Likelihood & Maximum *a Posteriori* Estimation

Paul F. Roysdon, Ph.D.

## I. DERIVATIONS OF THE MLE AND MAP

### A. Estimator Background

When the PDF is viewed as a function of the unknown parameter (with  $\mathbf{x}$  fixed) it is termed the “likelihood function”. Given a generic model

$$\mathbf{x}[0] = \mathbf{A} + \boldsymbol{\eta}[0]$$

where

$$\boldsymbol{\eta}[0] \sim \mathcal{N}(0, \sigma^2)$$

and PDF

$$p_i(\mathbf{x}[0]; \mathbf{A}) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{\left[-\frac{1}{2\sigma_i^2}(\mathbf{x}[0] - \mathbf{A})^2\right]}.$$

The log-likelihood simplifies the equation

$$\ln p(\mathbf{x}[0]; \mathbf{A}) = -\ln \sqrt{2\pi\sigma^2} - \frac{1}{2\sigma^2}(\mathbf{x}[0] - \mathbf{A})^2.$$

The first-derivative is

$$-\frac{\partial \ln p(\mathbf{x}[0]; \mathbf{A})}{\partial \mathbf{A}} = \frac{1}{\sigma^2}(\mathbf{x}[0] - \mathbf{A}).$$

The second-derivative is

$$-\frac{\partial^2 \ln p(\mathbf{x}[0]; \mathbf{A})}{\partial \mathbf{A}^2} = \frac{1}{\sigma^2}.$$

The efficiency of an estimator is measured by the curvature of the estimator. The curvature can be computed by the variance of the estimator

$$\text{var} \langle \hat{\mathbf{A}} \rangle = \frac{1}{-\frac{\partial^2 \ln p(\mathbf{x}[0]; \mathbf{A})}{\partial \mathbf{A}^2}},$$

or more conveniently as the expected value

$$-\mathbb{E} \left\langle \frac{\partial^2 \ln p(\mathbf{x}[0]; \mathbf{A})}{\partial \mathbf{A}^2} \right\rangle.$$

### B. General Model

Assume a linear model

$$\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \boldsymbol{\eta}$$

where  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ ,  $\mathbf{H} \in \mathbb{R}^{n \times p}$  with  $n > p$  and  $\text{rank}(\mathbf{H}) = p$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{p \times 1}$ , and  $\boldsymbol{\eta} \in \mathbb{R}^{n \times 1}$  and  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ .

The PDF of  $\mathbf{x}$  is

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\mathbf{C})}} e^{\left[-\frac{1}{2}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^\top \mathbf{C}^{-1}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta})\right]}$$

### C. Cramer-Rao Lower Bound

The Cramer-Rao Lower Bound (CRLB) is the lower bound for all unbiased estimators, thus it is the Minimum Variance Unbiased Estimator (MVUE). Given the model in Section I-B, the CRLB is defined as

$$\frac{\partial \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$$

where

$$\hat{\boldsymbol{\theta}}_{\text{CRLB}} = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}^{-1} \mathbf{x}$$

Note if  $\hat{\boldsymbol{\theta}}$  is the MVUE, it is an “efficient estimator”, and the minimum variance is the diagonal entries of the covariance matrix

$$\mathbf{C}_{\hat{\boldsymbol{\theta}}} = \mathbf{I}^{-1}(\boldsymbol{\theta}) = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1}$$

### D. Maximum Likelihood Estimator

Given the model in Section I-B, the Maximum Likelihood Estimator (MLE) maximizes the likelihood of the parameter  $\boldsymbol{\theta}$  such that

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{MLE}} &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}; \boldsymbol{\theta}) \\ &= (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}^{-1} \mathbf{x} \end{aligned}$$

If  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  is “efficient”, then  $\hat{\boldsymbol{\theta}}_{\text{MLE}} \mapsto \text{CRLB}$  and thus is the MVUE.

Note, in general the MLE is not optimal. However, as  $n \rightarrow \infty$  then the MLE converges to the MVUE asymptotically.

### E. Maximum *A Posteriori* Estimator

Given the model in Section I-B, the Maximum *a Posteriori* (MAP) Estimator maximizes the posterior PDF, such that

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{x}),$$

where

$$p(\boldsymbol{\theta} | \mathbf{x}) = \frac{p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{x})}, \quad (1)$$

also known as “Bayes’ Rule” or “Bayes’ Risk”. Therefore

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{MAP}} &= \arg \max_{\boldsymbol{\theta}} (p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta})) \\ &= \arg \max_{\boldsymbol{\theta}} (\ln p(\mathbf{x} | \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})). \end{aligned}$$

Note, under normality constraints the MAP is equivalent to the Minimum Mean Square Error (MMSE).

## II. TUTORIAL

Both the MLE and MAP can be used for *estimation* and *prediction*. We will consider both cases, and define how the MLE differs from MAP and Bayesian methods. Throughout this section we will use as example a "fair" coin, meaning that both heads and tails are equally likely.

To begin, let our *evidence*,  $\mathbf{x}$ , be a set of *independent* observations

$$\mathbf{x} = \{x_i\}_{i=1}^{|\mathbf{x}|}, \quad (2)$$

where each  $x_i$  is a realization (a.k.a., an observation or measurement), of a random variable  $x$ , and the notation  $|\mathbf{x}|$  represents the cardinality, or total number, of observations in the set  $\mathbf{x}$  in a multidimensional space.

Let  $\theta$  represent the *set of probability distribution parameters that best explain the evidence  $\mathbf{x}$* .

- **Parameter Estimation:** Given  $\mathbf{x}$ , we can estimate  $\theta$  using Bayes' Rule (repeated here):

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

where  $p(A)$  is the probability of  $A$ , and  $p(A|B)$  is the *conditional* probability of  $A$  given  $B$ .

- **Data Prediction:** If we are given new observations  $\tilde{\mathbf{x}}$ , we can compute the probability of the new observation supported by, or given, the evidence

$$p(\tilde{\mathbf{x}}|\mathbf{x}).$$

### A. Parameter Estimation

Notice that Eqn. 1 can be interpreted as

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}.$$

Therefore, the likelihood function is

$$\text{likelihood} = p(\mathbf{x}|\theta).$$

1) **MLE of  $\theta$ :** To estimate  $\theta$ , we seek the value for  $\theta$  that maximizes the likelihood. Let  $\hat{\theta}_{ML}$  denote the maximum likelihood estimate of  $\theta$ .

From [1] we know that the joint probability of a set of *independent* random variables (rv's) is a *product* of the probabilities associated with the individual rv's in the set. Thus, because  $\mathbf{x}$  is a set of *independent* observations,  $\{x_1, x_2, \dots\}$ , we seek the value  $\theta$  that maximizes

$$\prod_{x_i \in \mathbf{x}} p(x_i|\theta). \quad (3)$$

Because the logarithm is a monotonically increasing function of its argument, we can simplify Eqn. 3 by replacing the product with a sum, and using the notation  $\mathcal{L}$  to represent the logarithm:

$$\mathcal{L} = \sum_{x_i \in \mathbf{x}} \log p(x_i|\theta). \quad (4)$$

Therefore, we seek the values for the parameters in  $\theta$  that maximize  $\mathcal{L}$ , defined as

$$\hat{\theta}_{MLE} = \arg \max_{\theta} (\mathcal{L}). \quad (5)$$

Eqn. 5 is solved by by setting the first derivative to zero, such that

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = 0 \quad \forall \theta_i \in \theta. \quad (6)$$

**Note:** The logarithm, specifically the *natural logarithm*, will further simplify our later examples because most distributions are from exponential families (see Section III), e.g., given  $\gamma^b$ , let

$$a = \gamma^b,$$

then the logarithm of base  $\gamma$  is

$$\log_{\gamma}(a) = b.$$

The natural logarithm for an exponential function  $a = e^b$  is defined as

$$\ln(a) = \log_e(a).$$

Therefore,

$$\ln(a) = b.$$

Notice that if in Eqn. 4,  $p(x_i|\theta)$  is equal to an exponential family, e.g.,  $e^{b(x_i, \theta)}$ , then the expression simplifies to

$$\mathcal{L} = \sum_{x_i \in \mathbf{x}} b(x_i, \theta).$$

An application of this fact is provided in Sections I-B thru I-E, and [2]. For additional reference, see [3], [4].

2) **MAP Estimate of  $\theta$ :** In maximum *a posteriori* estimation, we seek the value for  $\theta$  that maximizes the *posterior* of Eqn. 1, e.g.,  $p(\theta|\mathbf{x})$ , denoted as  $\hat{\theta}_{MAP}$ .

$$\begin{aligned} \hat{\theta}_{MAP} &= \arg \max_{\theta} p(\theta|\mathbf{x}) \\ &= \arg \max_{\theta} \left( \frac{p(\mathbf{x}|\theta) p(\theta)}{p(\mathbf{x})} \right) \end{aligned} \quad (7)$$

$$= \arg \max_{\theta} (p(\mathbf{x}|\theta) p(\theta)) \quad (8)$$

$$= \arg \max_{\theta} \left( \prod_{x_i \in \mathbf{x}} p(x_i|\theta) p(\theta) \right)$$

Note the denominator in Eqn. 7 is dropped in Eqn. 8 because it has no functional dependence on the maximization of  $\theta$ .

Again, because  $\mathbf{x}$  is a set of *independent* observations, we can change the product to a sum, and we can simplify the expression by taking the logarithm of the posteriors:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \sum_{x_i \in \mathbf{x}} \log p(x_i|\theta) + \log p(\theta) \quad (9)$$

3) **Bayesian Estimate of  $\theta$ :** Given the evidence  $\mathbf{x}$ , the MLE holds  $\theta$  constant and seeks a value for  $\theta$  that maximizes the belief in the evidence. Recall that there is no prior in MLE, therefore we cannot use prior belief in likely values for  $\theta$  in the estimation.

On the contrary, the MAP estimate of  $\theta$  provides a method to consider values from a distribution (not the distribution of values) that express our prior beliefs regarding the parameters, i.e., the MAP estimate returns the values of  $\theta$  where the probability  $p(\theta|\mathbf{x})$  is maximized.

By contrast, a Bayesian estimate calculates the full posterior *distribution*  $p(\theta|\mathbf{x})$ . That is, all of the  $\theta$  values that are possible, we select a value for each element of the vector that is, by some definition, “best”. For example, we may choose the expected value of  $\theta$ , assuming its variance is sufficiently small.

Recall the definition of Bayes’ Risk, Eqn. 1. In Bayesian estimation the denominator, the *probability of the evidence*, cannot be ignored. This can be reformulated into an integral

$$p(\mathbf{x}) = \int_{\theta} p(\mathbf{x}|\theta) p(\theta) d\theta \quad (10)$$

Note: for a given likelihood function, if we have a choice regarding how we express our prior beliefs, we must use the form that allows us to carry out the integration of Eqn 10. This is the notion of conjugate priors, see Sections III-F & III-G.

Finally, as with MAP estimation, Bayesian estimation also requires the prior belief in possible values of the parameter vector  $\theta$  to be expressed in the form of a distribution.

4) **Example and Comparison:** The MAP estimate differs from the MLE by allowing us to inject our prior beliefs of the parameter values in  $\theta$ .

Consider a Bernoulli experiment of  $n$  trials flipping a fair coin. Let  $p_h$  be the probability of heads, and for each observation  $x_i$ , a scalar, the value of  $x_i$  is either heads or tails.

**MLE Example:** For the scalar parameter  $\theta_h$ , let’s calculate the MLE,  $\hat{\theta}_{h(MLE)}$ . The evidence,  $\mathbf{x}$ , is the set of *independent* observations

$$\mathbf{x} = \{x_i = \begin{matrix} \text{heads} \\ \text{tails} \end{matrix}, \quad i = 1, \dots, n\}$$

The log likelihood function is

$$\begin{aligned} \log p(\mathbf{x}|\theta_h) &= \sum_{i=1}^n \log p(x_i|\theta_h) \\ &= \sum_i \log p(x_i = \text{heads}) + \sum_i \log p(x_i = \text{tails}) \\ &= n_h \cdot \log \theta_h + (n - n_h) \cdot \log(1 - \theta_h), \end{aligned}$$

where  $n_h$  is the number of trials you predict will be heads. Let

$$\mathcal{L} = \log p(\mathbf{x}|\theta_h).$$

The MLE for  $\theta_h$  is solved by setting

$$\frac{\partial \mathcal{L}}{\partial \theta_h} = 0,$$

which yields the equation

$$\frac{n_h}{\theta_h} - \frac{(n - n_h)}{(1 - \theta_h)} = 0,$$

whose solution is the MLE

$$\hat{\theta}_{h(MLE)} = \frac{n_h}{n}.$$

If the number of trials,  $n = 20$ , and our prediction of heads,  $n_h = 12$ , the *maximum likelihood estimate of our parameter*  $\theta_h$  is  $\hat{\theta}_{h(MLE)} = 12/20 = \mathbf{0.6}$ .

**MAP Example:** For the scalar parameter  $\theta_h$ , let’s calculate the MAP estimate,  $\hat{\theta}_{h(MAP)}$ . The MAP requires a prior belief distribution for the parameter  $\theta_h$ , with the following constraints:

- The prior for  $\theta_h$  must be zero outside the interval  $[0, 1]$ .
- Within the interval  $[0, 1]$ ,  $\theta_h$  may be any real value.
- A distribution that peaks within the interval  $[0, 1]$  is desirable, e.g., the *beta distribution*.

The beta distribution,  $Beta(p|\alpha, \beta)$ , has “shape” parameters  $\alpha$  and  $\beta$  (we also refer to these as *hyper-parameters* to the parameter  $\theta_h$ ):

$$p(\theta_h|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta_h^{\alpha-1} (1 - \theta_h)^{\beta-1}, \quad (11)$$

where  $B = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the beta function. For  $\alpha, \beta > 0$  the mode (maximum value) of the beta distribution is located at

$$\frac{\alpha - 1}{\alpha + \beta - 2}$$

Returning to our coin example, assume that prior experience with the same coin yielded a prior distribution for  $\theta_h$  with a peak at 0.5, i.e., a fair coin. Setting  $\alpha = \beta$  results in a distribution with a peak in the center of the interval  $[0, 1]$ .

The variance of the beta distribution is given by

$$\text{var} \langle Beta \rangle = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

Again, based on prior experience, let  $\alpha = \beta = 5$ . Therefore, the variance,  $\text{var} \langle Beta \rangle \approx 0.025$ , and standard deviation,  $\text{std} \langle Beta \rangle \approx 0.16$ .

Eqn. 9 can be rewritten to be

$$\begin{aligned} \hat{\theta}_{h(MAP)} &= \arg \max_{\theta_h} \left( n_h \cdot \log \theta_h + (n - n_h) \right. \\ &\quad \left. \cdot \log(1 - \theta_h) + \log p(\theta_h) \right). \end{aligned}$$

The MAP for  $\theta_h$  is solved by substituting the beta distribution for  $p(\theta_h)$  and setting the first derivative equal to zero:

$$\frac{n_h}{\theta_h} - \frac{n - n_h}{1 - \theta_h} + \frac{\alpha - 1}{\theta_h} - \frac{\beta - 1}{1 - \theta_h} = 0$$

whose solution is the MAP

$$\begin{aligned}\hat{\theta}_{h(MAP)} &= \frac{n_h + \alpha - 1}{n + \alpha + \beta - 2} \\ &= \frac{n_h + 4}{n + 8}\end{aligned}$$

Based on prior experience,  $\alpha = \beta = 5$ , if the number of trials,  $n = 20$ , and our prediction of heads,  $n_h = 12$ , the *maximum a posteriori estimate of our parameter  $\theta_h$*  is  $\hat{\theta}_{h(MAP)} = (12 + 4)/(20 + 8) = \mathbf{0.571}$ .

**Bayesian Estimation Example:** Again consider the Bernoulli trial example of a fair coin and our measurement, and prediction, of heads and tails. Our prior for that example is given by the beta distribution, Eqn. 11. Note the explicit dependence of the prior on the hyper-parameters  $\alpha$  and  $\beta$ . Inserting Eqn. 11 into Eqn. 10,

$$\begin{aligned}p(\mathbf{x}) &= \int_0^1 p(\mathbf{x}|\theta_h) p(\theta_h) d\theta_h \\ &= \int_0^1 \left( \prod_{i=1}^n p(x_i|\theta_h) \right) p(\theta_h) d\theta_h \\ &= \int_0^1 (\theta_h^{n_h} (1 - \theta_h)^{n - n_h}) p(\theta_h) d\theta_h.\end{aligned}$$

Notice that when we multiply a beta distribution with either a power of  $\theta_h$  or a power of  $(1 - \theta_h)$ , the result is a different beta distribution. Thus, the probability of the evidence  $\mathbf{x}$  is equivalent to a constant,  $z$ , whose value is dependent on the values  $\alpha, \beta$  and measurement  $n_h$ . Reformulating Eqn. 1, we have

$$\begin{aligned}p(\theta_h|\mathbf{x}) &= \frac{p(\mathbf{x}|\theta_h) p(\theta_h)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\theta_h) p(\theta_h)}{z} \\ &= \frac{1}{z} p(\mathbf{x}|\theta_h) p(\theta_h) \\ &= \frac{1}{z} \left( \prod_{i=1}^n p(x_i|\theta_h) \right) p(\theta_h) \\ &= \frac{1}{z} (\theta_h^{n_h} (1 - \theta_h)^{n - n_h}) p(\theta_h) \\ &= \text{Beta}(\theta_h|\alpha + n_h, \beta + n - n_h).\end{aligned}$$

The last step is obtained from the fact of products or powers of distributions, mentioned in Section II-A.3. The result is a closed form expression for the posterior distribution for the parameter vector  $\theta_h$  to be estimated. If instead we seek a single value for  $\theta_h$ , the expected value of the distribution can be calculated.

Recall from the example, based on prior experience, let  $\alpha = \beta = 5$ , then

$$\begin{aligned}\hat{\theta}_{h(Bayesian)} &= E\langle \theta_h|\mathbf{x} \rangle \\ &= \frac{\alpha + n_h}{\alpha + \beta + n} \\ &= \frac{5 + n_h}{10 + n}.\end{aligned}$$

If the number of trials,  $n = 20$ , and our prediction of heads,  $n_h = 12$ , the *Bayesian estimate of our parameter  $\theta_h$*  is  $\hat{\theta}_{h(Bayesian)} = (5 + 12)/(10 + 20) = \mathbf{0.567}$ .

A benefit of the Bayesian estimation is the ability to also calculate the variance of the estimate. Using the standard form for the variance of the beta distribution, the variance of the Bayesian estimate is 0.0079.

**Comparison:** We notice a few things:

- While the MLE only uses our prediction, the MAP and Bayesian estimates use our prediction and “pulls” the estimate toward our *prior* experience. Therefore a larger prior will more strongly influence the MAP and Bayesian estimates. Compare,  $\hat{\theta}_{h(MLE)} = 0.6$  vs.  $\hat{\theta}_{h(MAP)} = 0.571$  vs.  $\hat{\theta}_{h(Bayesian)} = 0.567$ .
- Larger, but equal, values for  $\alpha$  and  $\beta$ , narrow the peak of the beta distribution around  $\theta_h = 0.5$ , also moving the MAP and Bayesian estimates toward the prior.
- The values for  $\alpha$  and  $\beta$  tend to “smooth” the result of the measurement  $n_h$ .

While we used the Beta distribution in this example, we can certainly use, or assume, other distributions, e.g., in [2] both a Gaussian and Laplacian distribution are used.

## B. Data Prediction

Given the evidence, Eq. 2, and new measurements  $\tilde{\mathbf{x}}$ , we want to know how much our new measurements are supported by the evidence, i.e.,

$$p(\tilde{\mathbf{x}}|\mathbf{x}). \quad (12)$$

In other words, given the evidence, can we *predict* the new measurements. Prediction, also called *regression*, can be solved by either the MLE, MAP or Bayesian estimators.

1) **MLE Prediction of  $\tilde{\mathbf{x}}$ :** Expanding Eq. 12,

$$\begin{aligned}p(\tilde{\mathbf{x}}|\mathbf{x}) &= \int_{\theta} p(\tilde{\mathbf{x}}|\theta) p(\theta|\mathbf{x}) d\theta \\ &\approx \int_{\theta} p(\tilde{\mathbf{x}}|\hat{\theta}_{MLE}) p(\theta|\mathbf{x}) d\theta \\ &= p(\tilde{\mathbf{x}}|\hat{\theta}_{MLE}).\end{aligned}$$

Therefore, the probability for the new measurements  $\tilde{\mathbf{x}}$  is the same for all prior measurements, based on the evidence  $\mathbf{x}$  and probability model  $\hat{\theta}_{MLE}$ .

2) **MAP Prediction of  $\tilde{\mathbf{x}}$ :** Expanding Eq. 12,

$$\begin{aligned}p(\tilde{\mathbf{x}}|\mathbf{x}) &= \int_{\theta} p(\tilde{\mathbf{x}}|\theta) p(\theta|\mathbf{x}) d\theta \\ &\approx \int_{\theta} p(\tilde{\mathbf{x}}|\hat{\theta}_{MAP}) p(\theta|\mathbf{x}) d\theta \\ &= p(\tilde{\mathbf{x}}|\hat{\theta}_{MAP}).\end{aligned}$$

Therefore, the probability for the new measurements  $\tilde{\mathbf{x}}$  is the same for all prior measurements, based on the evidence  $\mathbf{x}$  and probability model  $\hat{\theta}_{MAP}$ .

3) **Bayesian Prediction of  $\tilde{x}$ :** As before, we begin with Eq. 12, but now we must use Bayes' Rule for the posterior  $p(\theta|x)$ :

$$\begin{aligned} p(\tilde{x}|x) &= \int_{\theta} p(\tilde{x}|\theta) p(\theta|x) d\theta \\ &= \int_{\theta} p(\tilde{x}|\theta) \frac{p(x|\theta) p(\theta)}{p(x)} d\theta. \end{aligned}$$

### III. EXPONENTIAL FAMILIES

Many *exponential family* distributions are used in both MLE and MAP estimation. This section defines the most common distributions, derives the *natural parameters* and both the *first* and *second derivatives* for use in numerical (software) applications. This section closes with a brief explanation of *conjugate priors*.

#### A. Background

Consider the general exponential function

$$p(x) = h(x)e^{\theta^T T(x) - A(\theta)},$$

where  $\theta$  is a parameter vector,  $T(x)$  is a vector of “sufficient statistics”,  $A(\theta)$  is a cumulant generating function, and a general input function  $h(x)$ .

A normalized distribution for any  $\theta$  is obtained by taking the integral with respect to  $x$

$$\int p(x) dx = e^{-A(\theta)} \int h(x) e^{\theta^T T(x)} d\theta = 1,$$

therefore

$$e^{A(\theta)} = \int h(x) e^{\theta^T T(x)} d\theta.$$

Note: when  $T(x) = x$ , then  $A(\theta)$  is the log of the Laplace transform of  $h(x)$ .

#### B. Examples

Gaussian:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\|x-\mu\|^2/(2\sigma^2)} \quad x \in \mathbb{R}$$

Bernoulli:

$$p(x) = \alpha^x (1-\alpha)^{1-x} \quad x \in \{0, 1\}$$

Binomial:

$$p(x) = \binom{n}{k} \alpha^x (1-\alpha)^{n-x} \quad x \in \{0, 1, 2, \dots, n\}$$

Multinomial:

$$\begin{aligned} p(x) &= \frac{n!}{x_1! x_2! \dots x_n!} \prod_{i=1}^n \alpha_i^{x_i} \\ x &\in \{0, 1, 2, \dots, n\}, \sum_i x_i = n \end{aligned}$$

Exponential:

$$p(x) = \lambda e^{-\lambda x} \quad x \in \mathbb{R}^+$$

Poisson:

$$p(x) = \frac{e^{-\lambda}}{x!} \lambda^x \quad x \in \{0, 1, 2, \dots, n\}$$

Dirichlet:

$$p(x) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i x_i^{\alpha_i-1} \quad x \in [0, 1], \sum_i x_i = n$$

#### C. Natural Parameter Forms

1) **Bernoulli:** Recall the general equation

$$p(x) = h(x) e^{\theta^T T(x) - A(\theta)},$$

then

$$\begin{aligned} p(x) &= \alpha^x (1-\alpha)^{1-x} \\ &= \exp [\log(\alpha^x (1-\alpha)^{1-x})] \\ &= \exp [x \log \alpha + (1-x) \log(1-\alpha)] \\ &= \exp \left[ x \log \frac{\alpha}{1-\alpha} + \log(1-\alpha) \right] \\ &= \exp [x\theta - \log(1+e^\theta)], \end{aligned}$$

where  $T(x) = x$ ,  $\theta = \log \frac{\alpha}{1-\alpha}$ , and  $A(\theta) = \log(1+e^\theta)$ .

2) **Gaussian:**

$$\begin{aligned} p(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\|x-\mu\|^2/(2\sigma^2)} \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\log \sigma - \frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2} \right) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( \theta^T T(x) - \log \sigma - \frac{\mu^2}{2\sigma^2} \right) \end{aligned}$$

where  $h(x) = \frac{1}{\sqrt{2\pi\sigma^2}}$ , and  $A(\theta) = \log \sigma - \frac{\mu^2}{2\sigma^2}$ . Then

$$T(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}, \quad \theta = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix},$$

and

$$\begin{aligned} A(\theta) &= \frac{\mu^2}{2\sigma^2} + \log \sigma \\ &= -\frac{[\theta]_1^2}{4[\theta]_2^2} - \frac{1}{2} \log(-2[\theta]_2). \end{aligned}$$

3) **Multivariate Gaussian:**

$$p(x) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-(x-\mu)\Sigma^{-1}(x-\mu)/2}$$

where  $h(x) = (2\pi)^{D/2}$ , and

$$T(x) = \begin{pmatrix} x \\ x x^T \end{pmatrix}, \quad \theta = \begin{pmatrix} \Sigma^{-1} \mu \\ -\frac{1}{2} \Sigma^{-1} \end{pmatrix},$$

#### D. First & Second Derivatives

1) **First Derivative:** Let,

$$A(\theta) = \log \left[ \int h(x) e^{\theta^T T(x)} dx \right],$$

where  $Q(\theta) = \int h(x) e^{\theta^T T(x)} dx$ , then

$$\begin{aligned} \frac{\partial A(\theta)}{\partial \theta} &= \frac{1}{Q(\theta)} \frac{\partial Q(\theta)}{\partial \theta} = \frac{Q'(\theta)}{Q(\theta)} \\ &= \frac{\int h(x) e^{\theta^T T(x)} T(x) dx}{\int h(x) e^{\theta^T T(x)} dx} \\ &= \frac{\int h(x) e^{\theta^T T(x) - A(\theta)} T(x) dx}{\int h(x) e^{\theta^T T(x) - A(\theta)} dx} \\ &= E_{p_\theta} \langle T(x) \rangle. \end{aligned}$$

2) **Second Derivative:** Using the previous definition for  $A(\theta)$  and  $Q(\theta)$ ,

$$\begin{aligned} \frac{\partial^2 A(\theta)}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \left[ \frac{Q'(\theta)}{Q(\theta)} \right] \\ &= \frac{\partial}{\partial \theta} \left[ Q'(\theta) \frac{1}{Q(\theta)} \right] \\ &= \frac{Q''(\theta)}{Q(\theta)} - \frac{(Q'(\theta))^2}{(Q(\theta))^2}, \end{aligned}$$

therefore

$$\begin{aligned} \frac{\partial^2 A(\theta)}{\partial \theta^2} &= \frac{\int h(x) e^{\theta^T T(x)} T^2(x) dx}{\int h(x) e^{\theta^T T(x)} dx} - (E_{p_\theta} \langle T(x) \rangle)^2 \\ &= \frac{\int h(x) e^{\theta^T T(x) - A(\theta)} T^2(x) dx}{\int h(x) e^{\theta^T T(x) - A(\theta)} dx} - (E_{p_\theta} \langle T(x) \rangle)^2 \\ &= E_{p_\theta} \langle T^2(x) \rangle - (E_{p_\theta} \langle T(x) \rangle)^2 \\ &= \text{Cov}_{p_\theta} \langle T(x) \rangle \succeq 0, \end{aligned}$$

where the symbol  $\succeq$  defines a positive definite matrix, and  $A(\theta)$  is convex.

#### E. Products of Distributions

Products of exponential family distributions are exponential family distributions, e.g.,

$$\begin{aligned} &\left( h(x) e^{\theta_1^T T(x) - A(\theta_1)} \right) \times \left( h(x) e^{\theta_2^T T(x) - A(\theta_2)} \right) \\ &= \tilde{h}(x) e^{(\theta_1 + \theta_2)^T T(x) - \tilde{A}(\theta_1 + \theta_2)}, \end{aligned}$$

however the parametric form is more complex; see [1]. The product of two Gaussians is *always* a Gaussian.

#### F. Conjugate Priors

Given

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

notice the denominator is not a function of  $\theta$ , rather it is a normalizing term. Thus, we have two parametric terms,  $p(\theta)$  and  $p(x|\theta)$ :

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

Conjugacy requires  $p(\theta)$  and  $p(\theta|x)$  to be of the same form, e.g.,

$$p(\theta) \mapsto p(x|\theta)p(\theta) \mapsto p(\theta|x),$$

where  $p(\theta)$  and  $p(\theta|x)$  are Dirichlet, and  $p(x|\theta)$  is Multinomial. Then,  $p(\theta)$  and  $p(\theta|x)$  are called conjugate distributions.

For example: the Dirichlet in  $\theta$ , with  $\Gamma(x) = (x-1)!$  is defined as

$$p(\theta) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta^{\alpha_i - 1},$$

and Multinomial in  $x$  is defined as

$$p(x|\theta) = \frac{(\sum_i x_i)!}{x_1! x_2! \dots x_n!} \prod_i \theta_i^{x_i}$$

Applying the definition above (repeated here)

$$p(\theta|x) \propto p(x|\theta)p(\theta) = \dots \times \prod_i \theta_i^{x_i + \alpha_i - 1},$$

which is again Dirichlet:

$$p(\theta|x) = \frac{\Gamma(\sum_i \alpha_i + x_i)}{\prod_i \Gamma(\alpha_i + x_i)} \prod_i \theta_i^{x_i + \alpha_i - 1}.$$

#### G. Conjugate Pairs

In the prior example we demonstrated conjugate pairs, now we define several common pairs.

Prior	Conditional
Gaussian: $e^{-\ \mu - \mu_0\ ^2 / (2\sigma^2)}$	Gaussian: $e^{-\ \mu - \mu_0\ ^2 / (2\sigma^2)}$
Beta: $\frac{\Gamma(r+s)}{\Gamma(r)\Gamma(s)} \alpha^{r-1} (2-\alpha)^{s-1}$	Bernoulli: $\alpha^x (1-\alpha)^{1-x}$
Dirichlet: $\frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta^{\alpha_i - 1}$	Multinomial: $\frac{(\sum_i x_i)!}{\prod_i x_i!} \prod_i \theta_i^{x_i}$
Inverse Wishart	Gaussian (cov.)

Note that conjugacy is mutual, e.g.,

$$\text{Dirichlet} \mapsto \text{Multinomial} \mapsto \text{Dirichlet}$$

$$\text{Multinomial} \mapsto \text{Dirichlet} \mapsto \text{Multinomial}$$

#### IV. MULTINOMIAL DISTRIBUTIONS FOR MACHINE LEARNING

In this section we will further examine Conjugate Pairs with Multinomial Distributions and explore their connection to modeling likelihoods for images and text data.

add text here...

#### REFERENCES

- [1] S. M. Kay, *Fundamentals of Statistical Signal Processing, Vol. I - Estimation Theory*. Prentice Hall PTR, 2013.
- [2] P. Roysdon and J. Farrell, "Robust GPS-INS Outlier Accommodation using a Sliding Window Filter: A Comparative Study," 2017.
- [3] P. F. Roysdon, *Math Handbook for Engineers and Scientists*. Fibonacci Press; 2nd edition, 2023.
- [4] I. Bronshtein, K. Semendyayev, G. Musiol, and H. Muehlig, *Handbook of Mathematics*. Springer, 2007.