

# *Introduction to Bayesian Statistics: Sufficiency and Likelihood Principles, Prior, Posterior and Posterior Predictive Distributions, Gaussian Examples, Exponential Family of Distributions*

*Prof. Nicholas Zabaras*

*Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame*

*Notre Dame, Indiana, USA*

*Email: [nzabaras@gmail.com](mailto:nzabaras@gmail.com)*

*URL: <https://www.zabaras.com/>*

*September 3, 2018*



# References

---

- C P Robert, [The Bayesian Choice: From Decision-Theoretic Motivations to Computational Implementation](#), Springer-Verlag, NY, 2001 ([online resource](#))
- A Gelman, JB Carlin, HS Stern and DB Rubin, [Bayesian Data Analysis](#), Chapman and Hall CRC Press, 2<sup>nd</sup> Edition, 2003.
- J M Marin and C P Robert, [The Bayesian Core](#), Spring Verlag, 2007 ([online resource](#))
- D. Sivia and J Skilling, [Data Analysis: A Bayesian Tutorial](#), Oxford University Press, 2006.
- Bayesian Statistics for Engineering, [Online Course at Georgia Tech](#), B. Vidakovic.
- [Chris Bishop's PRML book](#), Chapter 2
- M. Jordan, An introduction to Probabilistic Graphical Models, Chapter 8 (pre-print)
- Kevin Murphy's, [Machine Learning: A probabilistic perspective](#), Chapters 2 and 4



# Contents

---

- Parametric modeling, Sufficiency principle, MLE and the Likelihood principles
- MLE for a Univariate Gaussian, MLE for the Multivariate Gaussian,
- Bayesian Statistics, Bayes rule, Prior, Likelihood and Posterior, Posterior Point Estimates, Predictive Distribution, Univariate Gaussian with Unknown Mean, Inference of Precision with Known Mean, Scaled Inverse Chi Squared Prior for the Variance, Inference on Mean and Precision
- Inference for the Multivariate Gaussian, Unknown Mean, Unknown Precision, MAP Estimation and Shrinkage, Unknown Mean and Precision, Marginal Likelihood, Non-informative Prior, Wishart Distribution
- Exponential Family, Bernoulli, Beta, Gamma, Gaussian, Conjugate Priors, Posterior Predictive



# Parametric Modeling

---

- Statistical theory derives from observations of a random phenomenon an inference about the probability distribution underlying this phenomenon.<sup>a</sup>
  - We consider **parametric modeling**: The observations  $x$  are the realizations of a random variable  $X$  of known probability density function  $f(x|\theta)$  where
    - $\theta$  is unknown and belongs to a space  $\Theta$  of finite dimension.
    - The function  $f(x|\theta)$  considered as a function of  $\theta$  for a fixed realization of the observation  $X = x$  is called **the likelihood function**.
- $$\ell(\theta | x) = f(x | \theta)$$

<sup>a</sup> Here we follow closely:

- C. P. Robert, [The Bayesian Choice](#), Springer, 2<sup>nd</sup> edition, [chapter 1](#) (full text available)
- Brani Vidakovic, [Bayesian Statistics for Engineers](#), online course.



# Example of Parametric Modeling

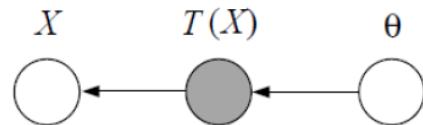
- Consider the problem of forest fires. Ecological and meteorological factors influence their eruption. Determining the probability  $p$  of fire as a function of these factors can be useful in the prevention of forest fires.
- We assume a parametrized shape for the function  $p$ .
- Denoting by  $h$  the humidity rate,  $t$  the average temperature,  $x$  the degree of management of the forest, a logistic model (Bernoulli random variable of parameter  $p$ ) could be proposed as:

$$p = \exp(\alpha_1 h + \alpha_2 t + \alpha_3 x) / [1 + \exp(\alpha_1 h + \alpha_2 t + \alpha_3 x)]$$

- The statistical step is now dealing with the evaluation of  $\alpha_1, \alpha_2, \alpha_3$ .

# Sufficiency Principle

- Consider  $X \sim f(x | \theta)$ . A function  $T$  of  $X$  (called a statistic of  $X$ ) is said to be sufficient if the distribution of  $X$  conditional upon  $T(X)$  is independent of  $\theta$ .



$$f(x | \theta) = h(T(x) | \theta) g(x | T(x))^*$$

- A sufficient statistic  $T(x)$  contains the whole information brought by  $x$  about  $\theta$ .
- Let us consider a simple example. Let  $X = (X_1, X_2, \dots, X_n)$  be i.i.d. from  $\mathcal{N}(\mu, \sigma^2)$  with  $\theta = (\mu, \sigma^2)$ . Then we can write:

$$\begin{aligned} f(x | \theta) &= \prod_{j=1}^N \mathcal{N}(x_j | \theta) = \prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_j - \mu)^2\right) = \frac{1}{(2\pi\theta_2)^{N/2}} \exp\left(-\frac{1}{2\theta_2} \sum_{j=1}^N (x_j - \theta_1)^2\right) = \\ &= \frac{1}{(2\pi\theta_2)^{N/2}} \exp\left(-\frac{1}{2\theta_2} \sum_{j=1}^N x_j^2 + \frac{\theta_1}{\theta_2} \sum_{j=1}^N x_j - \frac{N\theta_1^2}{2\theta_2}\right) \end{aligned}$$

\*  $f(x | \theta) = f(x, T(x) | \theta) = h(T(x) | \theta) g(x | T(x), \theta) = h(T(x) | \theta) g(x | T(x))$



# Sufficiency Principle

$$f(x|\theta) = \frac{1}{(2\pi\theta_2)^{N/2}} \exp\left(-\frac{1}{2\theta_2} \sum_{j=1}^N x_j^2 + \frac{\theta_1}{\theta_2} \sum_{j=1}^N x_j - \frac{N\theta_1^2}{2\theta_2}\right)$$

- We can see that  $f(x|\theta)$  depends only on  $T(x) = \left(\sum_{j=1}^N x_j, \sum_{j=1}^N x_j^2\right)$  which is our set of sufficient statistics.
- Introducing  $\bar{x} = \frac{\sum_{j=1}^N x_j}{N}$ ,  $s^2 = \sum_{j=1}^N (x_j - \bar{x})^2$ , we can also re-write the above equation:

$$f(x|\theta) = \frac{1}{(2\pi\theta_2)^{N/2}} \exp\left(-\frac{s^2 + N\bar{x}^2}{2\theta_2} + \frac{\theta_1 N \bar{x}}{\theta_2} - \frac{N\theta_1^2}{2\theta_2}\right)$$

So  $\bar{x}, s^2$  is an alternative set of sufficient statistics.



# Sufficiency Principle

- **Sufficiency principle:** Two observations  $x$  and  $y$  have the same values of  $T(x) = T(y)$  of statistics sufficient for  $f(\cdot | \theta)$ . Then the inferences about  $\theta$  based on  $x$  and  $y$  should be the same.
- Consider the model  $X_i \sim \mathcal{N}(\mu, 1)$  and we want to estimate  $\mu$  (our  $\theta$ ) based on  $N$  data. In this case, the sufficient statistic is

$$T(x_{1:N}) = \sum_{j=1}^N x_j$$

- Consider the (MLE) estimate of  $\theta$ :  $\hat{\mu}_1 = \sum_{j=1}^N x_j / N$ . It satisfies the sufficiency principle because if we have another dataset  $x'_{1:N}$  such that

$$T(x_{1:N}) = T(x'_{1:N}) \text{ then } \hat{\mu}_2 = \frac{1}{N} \sum_{j=1}^N x'_j = \frac{1}{N} \sum_{j=1}^N x_j = \hat{\mu}_1$$

- On the other hand, the estimate  $\hat{\mu}_1 = x_1$  does not satisfy the sufficiency principle for  $n > 1$  because if we have another dataset  $x'_{1:N}$  such that  $T(x_{1:N}) = T(x'_{1:N})$ , then  $\hat{\mu}_2 = x'_1 \neq \hat{\mu}_1$ , if  $x'_1 \neq x_1$ .

# The Likelihood Principle

---

- **Likelihood Principle.** In the inference about  $\theta$ , the information brought by an observation is entirely contained in the likelihood function
$$\ell(\theta | x) = f(x | \theta).$$
- Also, two likelihood functions contain the same information about  $\theta$  if they are proportional to each other; i.e.

$$\ell_1(\theta | x) = c(x) \ell_2(\theta | x)$$

- It is straight forward to show that the MLE (maximum likelihood procedure) satisfies the likelihood principle

$$\arg \max_{\theta} \ell_1(\theta | x) = \arg \max_{\theta} \ell_2(\theta | x)$$

- Classical approaches do not necessarily satisfy the likelihood principle.

# MLE: Summary

---

- The likelihood principle is fairly vague since it does not lead to the selection of a particular procedure.
- Maximum likelihood estimation is one way to implement the sufficiency and likelihood principles

$$\hat{\theta} = \operatorname{argsup}_{\theta} \ell(\theta | x)$$

- Indeed:

$$\arg \sup_{\theta} \ell(\theta | x) = \arg \sup_{\theta} h(x)g(T(x) | \theta) = \arg \sup_{\theta} g(T(x) | \theta)$$

$$\ell_1(\theta | x) = c(x)\ell_2(\theta | x) \Rightarrow \arg \sup_{\theta} \ell_1(\theta | x) = \arg \sup_{\theta} \ell_2(\theta | x)$$



# Maximum Likelihood for a Gaussian

---

- Suppose that we have a data set of observations  $\mathcal{D} = (x_1, \dots, x_N)^T$ , representing  $N$  observations of the **scalar random** variable  $X$ . The observations are drawn independently from a Gaussian distribution whose mean  $\mu$  and variance  $\sigma^2$  are unknown.
- We would like to determine these parameters from the data set.
- *i.i.d. data: Data points that are drawn independently from the same distribution are said to be **independent and identically distributed**, which is often abbreviated to i.i.d.*



# Maximum Likelihood for a Gaussian

- Because our data set  $\mathcal{D}$  is i.i.d., we can write the probability of the data set, given  $\mu$  and  $\sigma^2$ , *in the form*

Likelihood function:  $p(\mathbf{x} | \mu, \sigma^2) = \prod_{i=1}^N \mathcal{N}(x_i | \mu, \sigma^2)$

*This is seen as a function of  $\mu, \sigma^2$*



# Max Likelihood for a Gaussian Distribution

$$\text{Likelihood function: } p(\mathbf{x} | \mu, \sigma^2) = \prod_{i=1}^N \mathcal{N}(x_i | \mu, \sigma^2)$$

- One common criterion for determining the parameters in a probability distribution using an observed data set is to find the parameter values that *maximize the likelihood function*, i.e. *maximizing the probability of the data given the parameters* (contrast this with maximizing the probability of the parameters given the data).
- We can equivalently maximize the log-likelihood:

$$\max_{\mu, \sigma^2} \ln p(\mathbf{x} | \mu, \sigma^2) = \max_{\mu, \sigma^2} \left( -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \right) \Rightarrow$$

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i, \sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2$$



# Maximum Likelihood for a Gaussian Distribution

$$\mu_{ML} = \underbrace{\frac{1}{N} \sum_{i=1}^N x_i}_{\text{Sample mean}}, \sigma_{ML}^2 = \underbrace{\frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2}_{\text{Sample variance wrt MLE mean (not the exact mean)}}$$

- The MLE underestimates the variance (*bias* due to overfitting) because  $\mu_{ML}$  fitted some of the noise in the data.
- The maximum likelihood solutions  $\mu_{ML}, \sigma_{ML}^2$  are functions of the data set values  $x_1, \dots, x_N$ . Consider the expectations of these quantities with respect to the data set values, which come from a Gaussian.
- Using the equations above you can show that :

$$\mathbb{E}[\mu_{ML}] = \mu, \quad \mathbb{E}[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$$

In this derivation  
you need to use :  
 $\mathbb{E}[x_i x_j] = \mu^2$  for  $i \neq j$   
 $\mathbb{E}[x_i^2] = \sigma^2 + \mu^2$



# Maximum Likelihood for a Gaussian Distribution

---

$$\mathbb{E}[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$$

We use:

$$\mathbb{E}[x_i x_j] = \mu^2 \text{ for } i \neq j$$

$$\mathbb{E}[x_i^2] = \sigma^2 + \mu^2$$

$$\begin{aligned}\mathbb{E}[\sigma_{ML}^2] &= \mathbb{E}\left[\frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2\right] = \frac{1}{N} \mathbb{E}\left[\sum_{n=1}^N \left(x_n - \frac{1}{N} \sum_{m=1}^N x_m\right)^2\right] \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}\left[x_n^2 - \frac{2}{N} x_n \sum_{m=1}^N x_m + \frac{1}{N^2} \sum_{m=1}^N \sum_{l=1}^N x_m x_l\right] \\ &= \frac{1}{N} \left\{ N(\mu^2 + \sigma^2) - N \frac{2}{N} \left( (N-1)\mu^2 + (\mu^2 + \sigma^2) \right) + N \frac{1}{N^2} N \left( (N-1)\mu^2 + (\mu^2 + \sigma^2) \right) \right\} \\ &= \frac{1}{N} \left\{ N(\mu^2 + \sigma^2) - (N\mu^2 + \sigma^2) \right\} \\ &= \frac{(N-1)}{N} \sigma^2\end{aligned}$$



# Maximum Likelihood for a Gaussian Distribution

$$\mathbb{E}[\mu_{ML}] = \mu, \sigma_{ML}^2 = \frac{N-1}{N} \sigma^2, \mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i, \sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2$$

□ On average the MLE estimate obtains the correct mean but will *underestimate the true variance by a factor  $(N - 1)/N$ .*

□ An unbiased estimate of the variance is given as:

$$\bar{\sigma}^2 = \frac{N}{N-1} \sigma_{ML}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_{ML})^2$$

For large  $N$ ,  
the bias is not  
a problem

□ This result can be obtained from “a Bayesian treatment” in which we *marginalize over the unknown mean*.

□ The  $N - 1$  factor takes account the fact that “*1 degree of freedom has been used in fitting the mean*” and removes the bias of the MLE.



# MLE for the Multivariate Gaussian

- We can easily generalize the earlier MLE results for a multivariate Gaussian. The log-likelihood takes the form:

$$\ln p(X | \mathcal{D}, \mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

- Setting the derivatives wrt  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  equal to zero gives the following:

$$\boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T$$

- We provide a proof of the calculation of  $\boldsymbol{\Sigma}_{ML}$  next.

# MLE for the Multivariate Gaussian

$$\ln p(X | \mathcal{D}, \mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

- We differentiate the log likelihood wrt  $\Sigma^{-1}$ . Each contributing term is:

$$-\frac{N}{2} \frac{\partial}{\partial \Sigma^{-1}} \ln |\Sigma| = \frac{N}{2} \frac{\partial}{\partial \Sigma^{-1}} \ln |\Sigma^{-1}| = \frac{N}{2} \Sigma^T = \boxed{\frac{N}{2} \Sigma}$$

A useful trick!

$$-\frac{1}{2} \frac{\partial}{\partial \Sigma^{-1}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = -\frac{1}{2} N \frac{\partial}{\partial \Sigma^{-1}} \text{Tr} \left( \Sigma^{-1} \sum_{n=1}^N \frac{1}{N} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T \right)$$

$$= -\frac{1}{2} N \frac{\partial}{\partial \Sigma^{-1}} \text{Tr}(\Sigma^{-1} S) \quad \text{S symmetric}$$

$$= \boxed{-\frac{1}{2} NS}, \text{ where } S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T$$

- Setting the derivative equal to zero leads to:  $\Sigma_{ML} = S$

- Here we used:  
 $\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{AB}) = \mathbf{B}^T, \frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T,$   
 $|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}, \text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$



# Appendix: Some Useful Matrix Operations

□ Show that

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{AB}) = \mathbf{B}^T \text{ and } \frac{\partial}{\partial \mathbf{B}} \text{Tr}(\mathbf{AB}) = \mathbf{A}^T$$

Indeed

$$\frac{\partial}{\partial A_{mn}} \text{Tr}(\mathbf{AB}) = \frac{\partial}{\partial A_{mn}} (A_{ik} B_{ki}) = B_{nm} \Rightarrow \frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{AB}) = \mathbf{B}^T$$

□ Show that

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T$$

Using the cofactor expansion of the det:

$$\frac{\partial}{\partial A_{mn}} \ln |\mathbf{A}| = \frac{1}{|\mathbf{A}|} \frac{\partial}{\partial A_{mn}} |\mathbf{A}| = \frac{1}{|\mathbf{A}|} \frac{\partial}{\partial A_{mn}} \sum_j (-1)^{i+j} A_{ij} M_{ij} = \frac{1}{|\mathbf{A}|} (-1)^{m+n} M_{mn} = (\mathbf{A}^{-1})_{nm}$$

where in the last step we used Cramer's rule.



# MLE for a Multivariate Gaussian

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \equiv \bar{x}, \quad \Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T = \frac{1}{N} \sum_{n=1}^N x_n x_n^T - \bar{x} \bar{x}^T$$

- Note that *the unconstrained maximization of the log-likelihood gives a symmetric  $\Sigma$ .*
- As for the univariate case, we can define an **unbiased covariance** as:

$$\bar{\Sigma}_{ML} = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T, \quad \mathbb{E}[\bar{\Sigma}_{ML}] = \Sigma$$

- To prove this, you will need to use that:

$$\mathbb{E}[x_n x_m^T] = \mu \mu^T + \delta_{mn} \Sigma$$



# **Bayesian Statistics**

---

- A Bayesian model is made of a parametric statistical model ( $\mathcal{X}, f(x|\theta)$ ) and a prior distribution on the parameters ( $\Theta, \pi(\theta)$ ).
- The unknown parameters are now considered as random.
  - Some statisticians question this approach but most accept the probabilistic modeling on the observations.
- Example: Assume you want to measure the speed of light given some observations. Why should you put a prior on a physical constant?
  - Due to the limited accuracy of the measurement, this constant will never be known exactly.
  - It is thus justified to put a (e.g. uniform) prior on this parameter reflecting this uncertainty.



# *Recall Bayes' rule*

---

- Coming back from a trip, you feel sick and your doctor thinks you might have contracted a rare disease (0.01% of the population has the disease).
- A test is available but not perfect.
  - If a tested patient has the disease, 100% of the time the test will be positive.
  - If a tested patient does not have the disease, 95% of the time the test will be negative (5% false positive).
- Your test is positive, should you really care?



# **Medical Diagnosis Example**

---

- Let A=‘the patient has the disease’
- Let B=‘the test returns a positive result’

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})} = \frac{1 \times 0.0001}{1 \times 0.0001 + 0.05 \times 0.9999} \approx 0.002$$

- Such a test would be a complete waste of money



# Prior, Likelihood and Posterior

---

In the previous example, we can identify the following:

- Data  $x$  (e.g. `the test is positive')
- Hypothesis  $h$  (e.g. `do you have the disease?'). We want to make inferences about  $h$ .

In Bayesian settings all variables are random and all inferences are probabilistic.

We identify three key ingredients of a Bayesian inference approach:

- Prior  $\pi(h)$ : How likely is hypothesis  $h$  before looking at the data
- Likelihood  $f(x | h)$ : How likely is to observe  $x$  assuming  $h$  is true.
- Posterior  $\pi(h | x)$ : How likely is  $h$  after data  $x$  have been observed.

$$\pi(h | x) = \frac{f(x | h)\pi(h)}{m(x)}$$



# Prior $\pi(\theta)$

---

- We use the prior to introduce quantitatively some insights on the parameters of interest.
- This can be as subjective or as objective as you want it to be – and that's why frequentists do not like Bayesian approaches!
- There is no such a thing as a true prior!
- Even when prior information is heavily subjective, the Bayesian inference model is honest.



# Likelihood $f(x|\theta)$

---

- The likelihood encapsulates the mathematical model of the physical phenomena you are investigating.
- If you know the input  $X = x$  to your problem, the likelihood can represent the computed output  $y = f(x)$ .
- It is the most computational expensive part of Bayesian approaches to inference problems (inverse problems).



# **Posterior $\pi(\theta|x)$ : Inference and Prediction**

---

- It combines the prior and likelihood.
- It weights the data and the prior information in making probabilistic inferences
- The posterior distribution is also useful in estimating the probability of observing a future outcome (prediction)

$$\pi(\theta | x) = \frac{f(x | \theta)\pi(\theta)}{m(x)}$$

- The normalizing factor is given as:

$$m(x) = \int f(x|\theta)\pi(\theta)d\theta$$



# Posterior Inference: Point Estimates

$$\pi(\theta | x) = \frac{f(x | \theta)\pi(\theta)}{m(x)}$$

## Maximum A Posteriori estimate (MAP)

$$\theta^* = \arg \max_{\theta} \log(\pi(\theta | x)) = \arg \max_{\theta} (\log \pi(x | \theta) + \log \pi(\theta))$$

## Posterior Mean

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

## Posterior Quantiles

$$\Pr[\theta > a] = \int_a^{\infty} \pi(\theta | x) d\theta$$



# Prediction

Suppose we have observed  $x$  and we want to make a prediction about (future) unknown observables: What is the probability of observing data  $\hat{x}$ ? If we already have observed data  $x$ ?

This means finding  $g(\hat{x}|x)$

We have:

$$\begin{aligned} g(\hat{x}|x) &= \int g(\hat{x}, \theta|x)d\theta = \int \frac{\pi(\hat{x}, \theta, x)}{m(x)} d\theta = \int \frac{\pi(\hat{x}, \theta, x)}{\phi(\theta, x)} \frac{\phi(\theta, x)}{m(x)} d\theta = \\ &= \int f(\hat{x}|\theta, x)\pi(\theta|x)d\theta = \int f(\hat{x}|\theta)\pi(\theta|x)d\theta \end{aligned}$$

Compare this with the normalizing factor in Bayes' rule:

$$m(\hat{x}) = \int f(\hat{x}|\theta)\pi(\theta)d\theta$$



# A Gaussian Example

- Consider  $X_1 | \theta \sim \mathcal{N}(\theta, \sigma^2)$ , with prior  $\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$ .
- Then we can derive the following:

$$\pi(\theta | x_1) \propto f(x_1 | \theta) \pi(\theta) \propto \exp\left(-\frac{(x_1 - \theta)^2}{2\sigma^2} - \frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right) \Rightarrow$$
$$\pi(\theta | x_1) \propto \exp\left(-\frac{\theta^2}{2}\left(\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}\right) + \theta\left(\frac{x_1}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)\right) \propto \exp\left(-\frac{1}{2\sigma_1^2}(\theta - \mu_1)^2\right) \Rightarrow$$

With prior  $\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$  and observation  $x_1$   
 $\theta | x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  with

$$\frac{1}{\sigma_1^2} = \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2} \Rightarrow \sigma_1^2 = \frac{\sigma_0^2 \sigma^2}{\sigma_0^2 + \sigma^2}, \text{ and}$$

$$\mu_1 = \sigma_1^2 \left( \frac{x_1}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)$$



# A Gaussian Example: Continued

- To predict the distribution of a new observation  $X | \theta \sim \mathcal{N}(\theta, \sigma^2)$  in light of  $x_1$ , we use **the predictive distribution** as follows:

$$f(x | x_1) = \int \underbrace{f(x | \theta)}_{Likelihood} \underbrace{\pi(\theta | x_1)}_{Posterior} d\theta \propto \int e^{-\frac{(x-\theta)^2}{2\sigma^2}} e^{-\frac{(\theta-\mu_1)^2}{2\sigma_1^2}} d\theta = \int e^{-\frac{1}{2}\left(\frac{(x-\theta)^2}{\sigma^2} + \frac{(\theta-\mu_1)^2}{\sigma_1^2}\right)} d\theta$$

- You can verify with direct substitution the following:

$$\begin{aligned} -\frac{1}{2} \left( \frac{(x-\theta)^2}{\sigma^2} + \frac{(\theta-\mu_1)^2}{\sigma_1^2} \right) &= -\frac{1}{2} (x - \mu_1 \quad \theta - m_1) \frac{1}{\sigma_1^2 \sigma^2} \begin{bmatrix} \sigma_1^2 & -\sigma_1^2 \\ -\sigma_1^2 & \sigma_1^2 + \sigma^2 \end{bmatrix} \begin{pmatrix} x - \mu_1 \\ \theta - \mu_1 \end{pmatrix} + \dots \\ &= -\frac{1}{2} (x - \mu_1 \quad \theta - \mu_1) \begin{bmatrix} \sigma_1^2 + \sigma^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma_1^2 \end{bmatrix}^{-1} \begin{pmatrix} x - \mu_1 \\ \theta - \mu_1 \end{pmatrix} + \dots \end{aligned}$$

- This is a bivariate Gaussian and thus  $f(x|x_1)$  is the marginal in  $x$ , i.e.

$$X | x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2 + \sigma^2)$$

# Bayesian Inference for the Gaussian

- Consider  $\mathbf{X} = \{x_1, x_2, \dots, x_N\} \sim \mathcal{N}(\mu, \sigma^2)$ , with prior  $\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$ .
- The likelihood takes the form:

$$p(\mathbf{X} | \mu) = \prod_{n=1}^N f(x_n | \mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{\sum_{n=1}^N (x_n - \mu)^2}{2\sigma^2}\right)$$

- Note that in terms of  $\mu$  this is not a probability density and is not normalized. Introducing the conjugate (Gaussian) prior on  $\mu$  leads to:

$$\begin{aligned} \pi(\mu | \mathbf{X}) &= \prod_{n=1}^N f(x_n | \mu) \pi(\mu) \propto \exp\left(-\frac{\sum_{n=1}^N (x_n - \mu)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \Rightarrow \\ \pi(\mu | \mathbf{X}) &\propto \exp\left(-\frac{\mu^2}{2} \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right) + \mu \left(\frac{\sum_{n=1}^N x_n}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)\right) \propto \exp\left(-\frac{1}{2\sigma_N^2} (\mu - \mu_N)^2\right) \end{aligned}$$



# Bayesian Inference for the Gaussian

$$\pi(\mu | X) \propto \exp\left(-\frac{\mu^2}{2}\left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right) + \mu\left(\frac{\sum_{n=1}^N x_n}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)\right) \propto \exp\left(-\frac{1}{2\sigma_N^2}(\mu - \mu_N)^2\right)$$

➤ So the posterior is a Gaussian as before with

$\mu | X \sim \mathcal{N}(\mu_N, \sigma_N^2)$  with

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \Rightarrow \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2}, \text{ and}$$

$$\mu_N = \sigma_N^2 \left( \frac{\sum_{n=1}^N x_n}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) = \sigma_N^2 \left( \frac{N\mu_{ML}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$



# Bayesian Inference for the Gaussian

$\mu | X \sim \mathcal{N}(\mu_N, \sigma_N^2)$  with

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \Rightarrow \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2}, \text{ and } \mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$

- *The posterior precision is the sum of the precision of the prior plus one contribution of the data precision for each observed data point.*
- Observe the posterior mean for  $N \rightarrow \infty$  and  $N \rightarrow 0$ .
- For  $N \rightarrow \infty$  the posterior peaks around the  $\mu_{ML}$  and the posterior variance goes to zero, i.e. the point MLE estimate is recovered within the Bayesian paradigm for infinite data.
- How about when  $\sigma_0^2 \rightarrow \infty$ ? In this case note that  $\sigma_N^2 \rightarrow \frac{\sigma^2}{N}$  and  $\mu_N \rightarrow \mu_{ML}$



# Gaussian Example: Bayes' versus MLE

- We have seen that the ML estimate of  $\theta$  at data point  $N$  is simply:

$$\theta_{ML} = \arg \sup_{\theta} \prod_{i=1}^N f(x_i | \theta) = \frac{1}{N} \sum_{i=1}^N x_i$$

- The posterior of  $\theta$  at time  $N$  is (simply generalizing the earlier result):

$$\theta | x_1, \dots, x_N \sim \mathcal{N}(\mu_N, \sigma_N^2)$$

where

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \Rightarrow \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N \sigma_0^2 + \sigma^2} \underset{N \rightarrow \infty}{\sim} \frac{\sigma^2}{N}$$

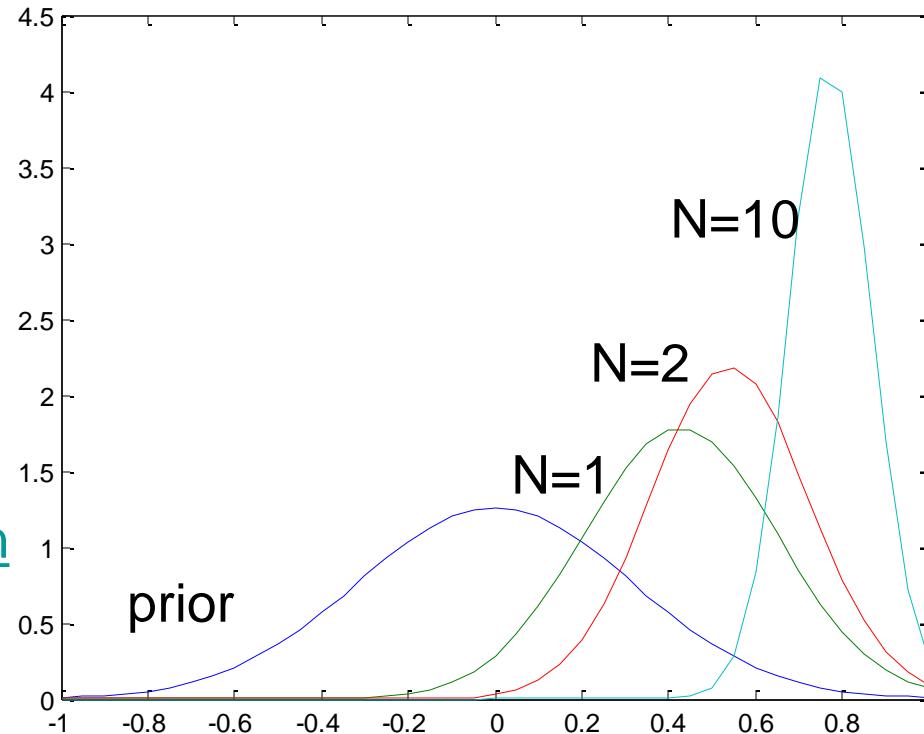
$$\mu_N = \sigma_N^2 \left( \frac{\sum_{i=1}^N x_i}{\sigma^2} + \frac{m_0}{\sigma_0^2} \right) \underset{N \rightarrow \infty}{\sim} \frac{\sum_{i=1}^N x_i}{N}$$

- As  $N \rightarrow \infty$ , the prior is washed out by the data and the posterior mean is the MLE estimate:  $\mathbb{E}[\theta | x_1, \dots, x_N] = \mu_N \simeq \theta_{ML}$



# Bayesian Inference for the Gaussian

$$\mu | X \sim \mathcal{N}(\mu_N, \sigma_N^2) \text{ with } \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2}, \text{ and } \mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$



MatLab  
implementation

$X = \{x_1, x_2, \dots, x_N\} \sim \mathcal{N}(0.8, 0.1)$ , with prior  $\mu \sim \mathcal{N}(0, 0.1)$ .

# Sequential Bayesian Inference

$\mu | X \sim \mathcal{N}(\mu_N, \sigma_N^2)$  with

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \Rightarrow \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2}, \text{ and } \mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$

- We can derive sequential estimates of the posterior variance and mean. They are as follows:

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}, \text{ and } \mu_N = \frac{\sigma_N^2}{\sigma_{N-1}^2} \mu_{N-1} + \frac{\sigma_N^2}{\sigma^2} x_N$$

- Show this is by recognizing the sequential nature of Bayesian inference (as we collect one data at a time, the posterior at the previous step becomes the new prior) and recalling:

With prior  $\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$  and observation  $x_1 \Rightarrow \theta | x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  with

$$\frac{1}{\sigma_1^2} = \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2} \Rightarrow \sigma_1^2 = \frac{\sigma_0^2 \sigma^2}{\sigma_0^2 + \sigma^2}, \text{ and } \mu_1 = \sigma_1^2 \left( \frac{x_1}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right)$$



# Example of Linear Gaussian Systems: Inferring the Mean

- Revisit the Bayesian inference problem for the Gaussian. Consider  $\mathbf{y} = \{y_1, y_2, \dots, y_N\} \sim \mathcal{N}(y | x, \sigma^2 = \lambda_y^{-1})$ , with prior  $x \sim \mathcal{N}(x | \mu_0, \sigma_0^2 = \lambda_0^{-1})$ .
- Likelihood for the  $N$  –dataset in the form of a linear Gaussian system:  
 $p(\mathbf{y} | x) = \mathcal{N}(\mathbf{y} | Ax + \mathbf{b}, \Sigma_y)$ ,  $A = \mathbf{1}_N$  (column vector of 1's),  $\mathbf{b} = \mathbf{0}$ ,  $\Sigma_y^{-1} = \text{diag}(\lambda_y \mathbf{I})$
- Applying conditional Gaussian results (to be reviewed later on):  
$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$
  
$$p(y | x) = \mathcal{N}(y | Ax + b, L^{-1})$$
  
$$p(x | y) = \mathcal{N}\left(x | \left(\Lambda + A^T L A\right)^{-1} \left(\Lambda \mu + A^T L (y - b)\right), \left(\Lambda + A^T L A\right)^{-1}\right)$$
  
$$p(x | y) = \mathcal{N}\left(x | \left(\lambda_0 + \mathbf{1}_N^T \lambda_y \mathbf{I} \mathbf{1}_N\right)^{-1} \left(\lambda_0 \mu_0 + \mathbf{1}_N^T \lambda_y \mathbf{I} (y - \mathbf{0})\right), \left(\lambda_0 + \mathbf{1}_N^T \lambda_y \mathbf{I} \mathbf{1}_N\right)^{-1}\right)$$

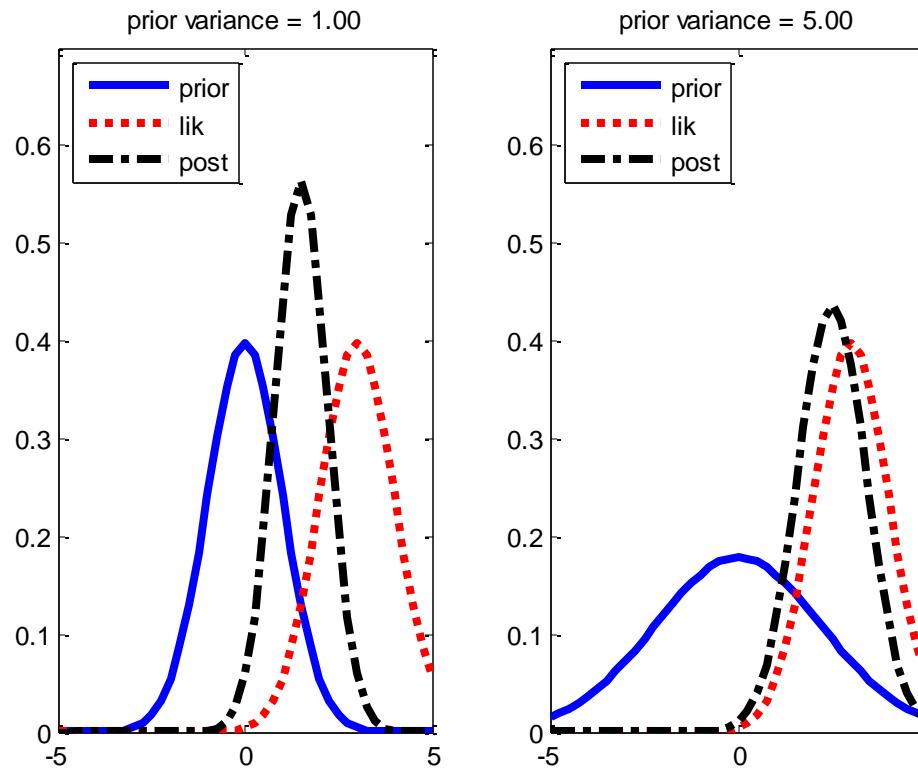
➤ This can be simplified as  $p(x|y) = \mathcal{N}(x|\mu_N, \lambda_N^{-1})$  where:

$$p(x|y) = \mathcal{N}\left(x | \frac{N\lambda_y}{\lambda_0 + N\lambda_y} \bar{y} + \frac{\lambda_0}{\lambda_0 + N\lambda_y} \mu_0, (\lambda_0 + N\lambda_y)^{-1}\right)$$

➤ The precision is the prior precision +  $N$  measurement precisions. The mean is the weighted average of the MLE and prior mean. These are identical results to those obtained earlier.

# Inferring the Mean of a Gayssian

[gaussInferParamsMean1d](#)  
from [PMTK](#)



- Inference about  $x$  given a single noisy observation  $y = 3$ .
  - (a) Strong prior  $\mathcal{N}(0, 1)$ . The posterior mean is “shrunk” towards the prior mean, which is 0.
  - (b) Weak prior  $\mathcal{N}(0, 5)$ . The posterior mean is similar to the MLE

# Shrinkage and Signal-To-Noise Ratio

$$\sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2}, \text{ and } \mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$

- The posterior precision is the sum of the precision of the prior plus one contribution of the data precision for each observed data point. For  $N \rightarrow \infty$  the posterior peaks around the  $\mu_{ML}$  and the posterior variance goes to zero, i.e. MLE estimate is recovered within the Bayesian paradigm.
- If we apply the data sequentially, we can write for the posterior mean *after the collection of one data point ( $N = 1$ )*, i.e.  $\mu_{ML} = y$  the following:

$$\mu_1 = y - (y - \mu_0) \frac{\sigma^2}{\sigma_0^2 + \sigma^2} \quad (\text{shrinkage of the data } y \text{ towards the prior mean } \mu_0)$$

- Shrinkage is often measured also with the *signal-to-noise ratio*:

$$SNR = \frac{\mathbb{E}[X^2]}{\mathbb{E}[\varepsilon^2]} = \frac{\sigma_0^2 + \mu_0^2}{\sigma^2}, \text{ for } y = x + \varepsilon \text{ (observed signal), } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- How about when  $\sigma_0^2 \rightarrow \infty$ ? In this case note that  $\sigma_N^2 \rightarrow \frac{\sigma^2}{N}$  and  $\mu_N \rightarrow \mu_{ML}$

# Inference of Precision with Known Mean

- Consider  $x_n \sim \mathcal{N}(x_n | \mu, \lambda^{-1})$ ,  $n = 1, \dots, N$ . We want to infer the precision  $\lambda = 1/\sigma^2$  with the mean  $\mu$  taken as known.
- The likelihood takes the form:

$$p(X | \lambda) = \prod_{n=1}^N f(x_n | \mu) \propto \lambda^{N/2} \exp\left(-\frac{1}{2} \lambda \sum_{n=1}^N (x_n - \mu)^2\right)$$

- The corresponding “conjugate prior” (a prior that results in a posterior of the same family as the prior) should be proportional to the product of a power of  $\lambda$  and the exponential of a linear function of  $\lambda$ . This corresponds to the gamma distribution.

$$\text{Gamma}(\lambda | a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}, x \in [0, \infty]$$

$$\mathbb{E}[\lambda] = \frac{a}{b}, \quad \text{var}[\lambda] = \frac{a}{b^2}$$

- The gamma distribution has a finite integral if  $a > 0$ , and the distribution itself is finite if  $a \geq 1$ .



# Inference of Precision with Known Mean

- The posterior takes the form:

$$p(\lambda | X, \mu) = \prod_{n=1}^N f(x_n | \mu) \mathcal{Gamma}(\lambda | a_0, b_0) \propto \lambda^{N/2+a_0-1} \exp\left(-b_0\lambda - \frac{1}{2}\lambda \sum_{n=1}^N (x_n - \mu)^2\right)$$

- We can immediately see that the posterior is also a Gamma distribution:

$$p(\lambda | X, \mu) = \mathcal{Gamma}(\lambda | a_N, b_N), a_N = N/2 + a_0, b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2$$

- Here  $\sigma_{ML}^2$  is the MLE of the variance.



# Inference of Precision with Known Mean

- The effect of observing  $N$  data points is to increase the value of  $a$  by  $N/2$  (i.e.  $\frac{1}{2}$  for each data point). Thus we interpret the parameter  $a_0$  as  $2a_0$  ‘effective’ prior observations.

$$p(\lambda | X, \mu) = \text{Gamma}(\lambda | a_N, b_N), a_N = N/2 + a_0, b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2$$

- Each measurement contributes to the parameter  $b$  a variance  $\sigma_{ML}^2 / 2$ . Since we have  $2a_0$  effective prior measurements, each of them contributes to  $b$  an effective prior variance

$$b_0 = \frac{2a_0}{2} \sigma^2 \Rightarrow \sigma^2 = \frac{b_0}{a_0}$$

- The interpretation of a conjugate prior in terms of effective data points is typical for the exponential family of distributions.
- The results above are identical with inference directly of the variance  $\sigma^2$  using as prior  $\text{InvGamma}(\sigma^2 | a_0, b_0)$  resulting in a posterior  $\text{InvGamma}(\sigma^2 | a_N, b_N)$



# Gamma and Inverse Gamma

Gamma

$$\theta \sim \text{Gamma}(\alpha, \beta)$$
$$p(\theta) = \text{Gamma}(\theta | \alpha, \beta)$$

shape  $\alpha > 0$   
inverse scale  $\beta > 0$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad \theta > 0$$

$$E(\theta) = \frac{\alpha}{\beta}$$
$$\text{var}(\theta) = \frac{\alpha}{\beta^2}$$
$$\text{mode}(\theta) = \frac{\alpha-1}{\beta}, \text{ for } \alpha \geq 1$$

Inverse-gamma

$$\theta \sim \text{Inv-gamma}(\alpha, \beta)$$
$$p(\theta) = \text{Inv-gamma}(\theta | \alpha, \beta)$$

shape  $\alpha > 0$   
scale  $\beta > 0$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$$

$$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$$
$$\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$$
$$\text{mode}(\theta) = \frac{\beta}{\alpha+1}$$

If  $\theta \sim \text{Gamma}(\theta | a, b)$  then  $\theta^{-1} \sim \text{InvGamma}(\theta^{-1} | a, b)$

Here :  $\lambda \sim \text{Gamma}(\lambda | a, b)$  then  $\sigma^2 \sim \text{InvGamma}(\sigma^2 | a, b)$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



# Univariate Posterior - Inverse Chi Squared Prior

- Alternative prior for  $\sigma^2$  is the Scaled Inverse Chi-Squared Distribution\*

$$\chi^{-2}(\sigma^2 | v_0, \sigma_0^2) \equiv \text{InvGamma}\left(\sigma^2 | \frac{v_0}{2}, \frac{v_0\sigma_0^2}{2}\right) \propto (\sigma^2)^{-v_0/2-1} \exp\left(-\frac{v_0\sigma_0^2}{2\sigma^2}\right)$$

- Here  $v_0$  represents the strength of the prior and  $\sigma_0^2$  encodes the value of the prior. With this, the posterior takes the form:

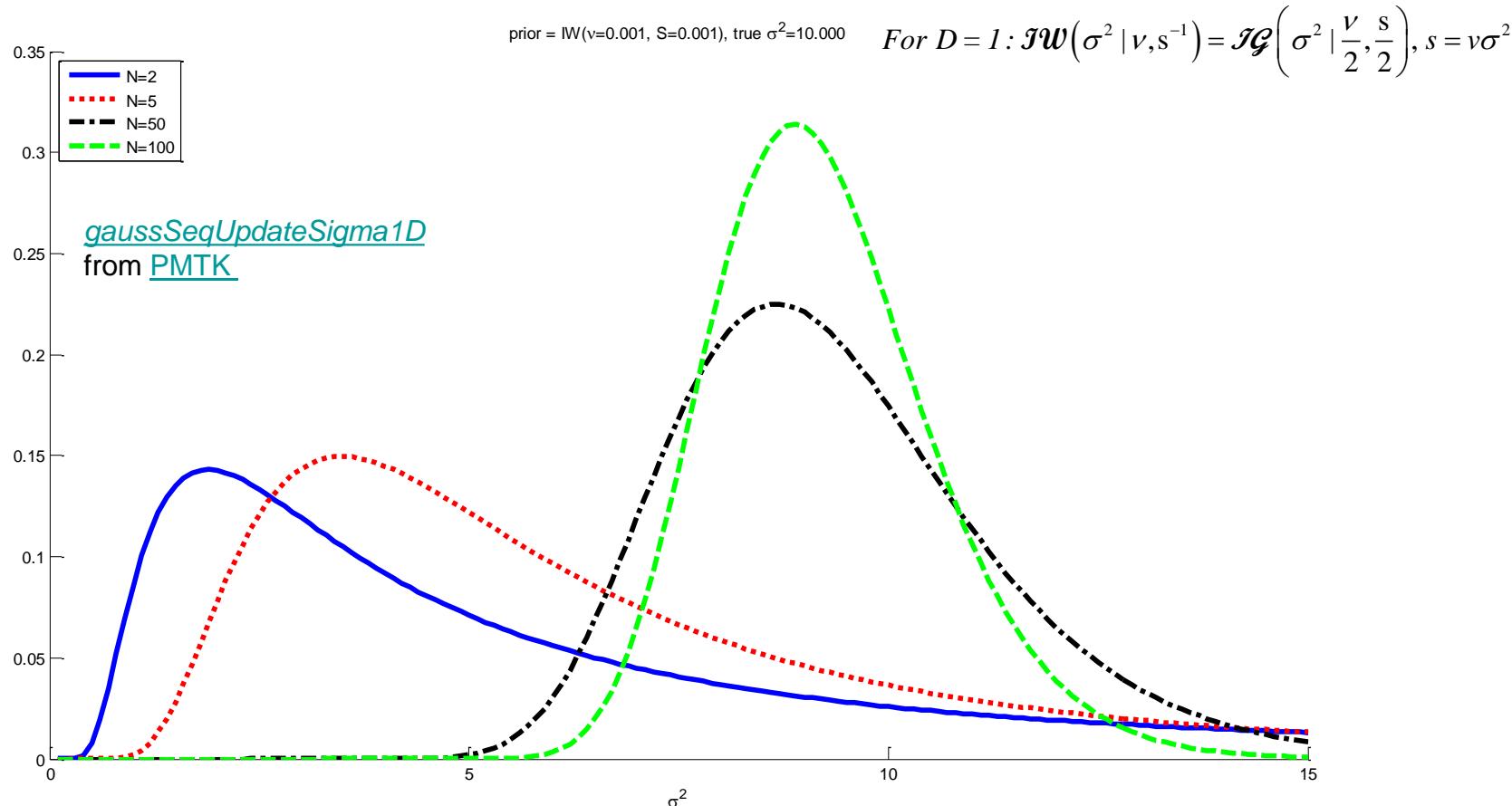
$$\chi^{-2}(\sigma^2 | \mathcal{D}, \mu) = \chi^{-2}(\sigma^2 | v_N, \sigma_N^2), v_N = v_0 + N, \sigma_N^2 = \frac{v_0\sigma_0^2 + \sum_{i=1}^N (x_i - \mu)^2}{v_N}$$

- The posterior dof  $v_N$  is the prior dof plus  $N$ . The posterior sum of squares  $\sigma_N^2 v_N = v_0\sigma_0^2 + \sum_{i=1}^N (x_i - \mu)^2$  is the sum of the prior sum of squares plus the data sum of squares. An **uninformative prior corresponds to zero virtual sample size,  $v_0 = 0$** . This is the prior  $p(\sigma^2) \propto \sigma^{-2}$
- This approach is certainly more appealing.

\* Often denoted as  $\text{Scale - Inv - } \chi^2(\sigma^2 | v_0, \sigma_0^2)$ ,  $\text{mean} = \frac{v_0\sigma_0^2}{v_0 - 2}$  for  $v_0 > 2$   
 $\text{var} = \frac{2v_0^2\sigma_0^4}{(v_0 - 2)^2(v_0 - 4)}$  for  $v_0 > 4$ ,  $\text{mode} = \frac{v_0\sigma_0^2}{v_0 + 2}$



# Sequential Update of the Posterior for $\sigma^2$



- Sequential update of the posterior for  $\sigma^2$  starting from an uninformative prior  $\mathcal{IW}(\sigma^2 | v_0 = 0.001, s_0 = v_0\sigma_0^2 = 0.001)$ . The scalar inverse-Wishart here is the same as the inverse Gamma. The data were generated from  $\mathcal{N}(5, 10)$ .

- Gelman 2006. [Prior distributions for variance parameters in hierarchical models](#). Bayesian Analysis 1(3):515–533



# Bayesian Inference: Unknown Mean and Precision

- Consider  $x_n \sim \mathcal{N}(x_n | \mu, \lambda^{-1})$ ,  $n = 1, \dots, N$ . We want to infer both the precision  $\lambda = 1/\sigma^2$  and the mean  $\mu$ .
- The likelihood takes the form:

$$\begin{aligned} p(\mathbf{X} | \mu, \lambda) &= \prod_{n=1}^N f(x_n | \mu) \propto \lambda^{N/2} \exp\left(-\frac{1}{2}\lambda \sum_{n=1}^N (x_n - \mu)^2\right) \\ &= \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left(\lambda\mu \sum_{n=1}^N x_n - \frac{1}{2}\lambda \sum_{n=1}^N x_n^2\right) \end{aligned}$$

- We need a prior that has a similar functional form in terms of  $\lambda$  and  $\mu$ .

$$\begin{aligned} p(\mu, \lambda) &\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^\beta \exp(\lambda\mu c - \lambda d) \\ &= \underbrace{\exp\left(-\frac{\beta\lambda}{2}\left(\mu - \frac{c}{\beta}\right)^2\right)}_{p(\mu|\lambda)} \underbrace{\lambda^{\beta/2} \exp\left(-\left(d - \frac{c^2}{2\beta}\right)\lambda\right)}_{p(\lambda)} \end{aligned}$$



# Bayesian Inference: Unknown Mean and Precision

$$\begin{aligned} p(\mu, \lambda) &\propto \left[ \lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right) \right]^\beta \exp(\lambda\mu c - \lambda d) \\ &= \underbrace{\left( \beta\lambda \right)^{1/2} \exp\left(-\frac{\beta\lambda}{2} \left( \mu - \frac{c}{\beta} \right)^2 \right)}_{p(\mu|\lambda)} \underbrace{\lambda^{(\beta-1)/2} \exp\left(-\left(d - \frac{c^2}{2\beta}\right)\lambda\right)}_{p(\lambda)} \end{aligned}$$

- We can easily identify that the prior is of the form (**Normal-Gamma**):

$$p(\mu, \lambda) = \mathcal{N}\left(\mu \mid \mu_0 = \frac{c}{b}, (\beta\lambda)^{-1}\right) \mathcal{Gamma}\left(\lambda \mid a = \frac{1+\beta}{2}, b = d - \frac{c^2}{2\beta}\right)$$

- Recall the form of the Gamma distribution:

$$\mathcal{Gamma}(\lambda \mid a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$$



# Bayesian Inference: Unknown Mean and Precision

- Combining the likelihood and prior, we can re-arrange and write:

$$p(\mu, \lambda | X) \propto \lambda^{N/2} \lambda^{a-1} \exp\left(-\left(b + \frac{1}{2} \sum_{n=1}^N x_n^2 + \frac{\beta}{2} \mu_0^2\right)\lambda\right) \times \\ (\lambda(N+\beta))^{1/2} \exp\left(-\frac{\lambda(N+\beta)}{2}\left(\mu^2 - \frac{2}{N+\beta}\left(\beta\mu_0 + \sum_{n=1}^N x_n\right)\mu\right)\right)$$

- Completing the square on the 2<sup>nd</sup> argument gives:

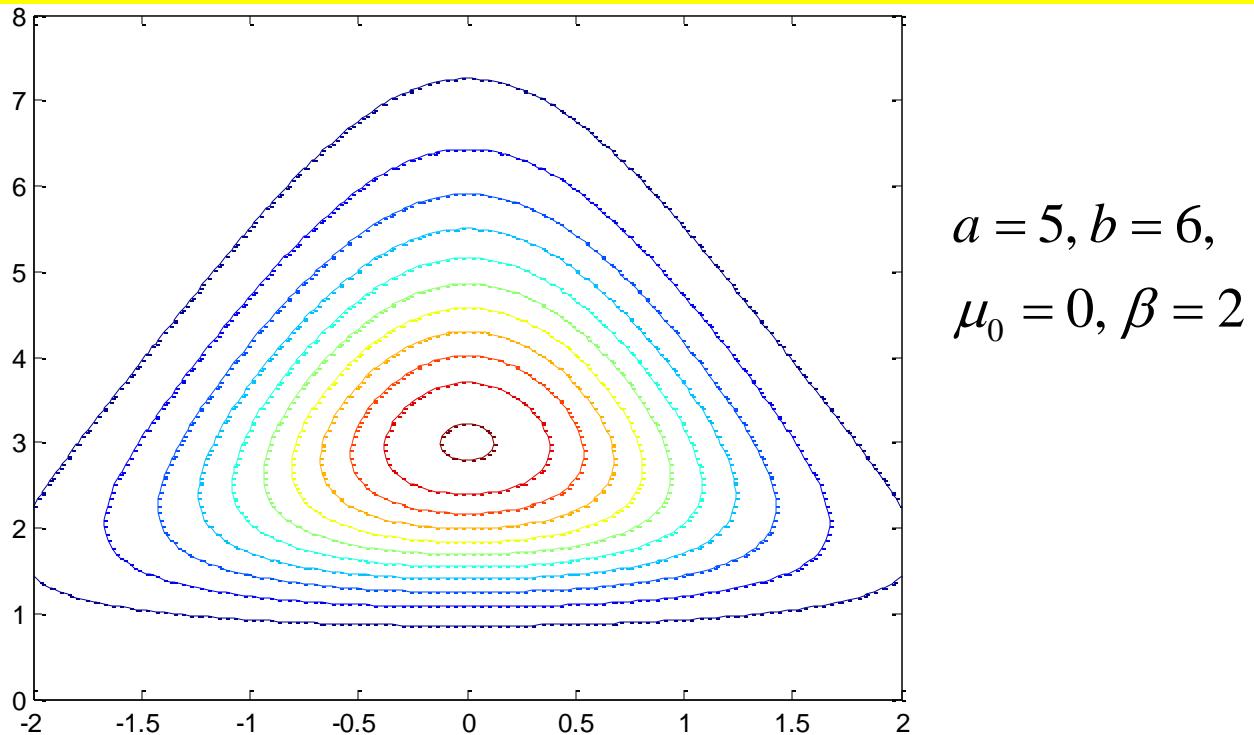
$$p(\mu, \lambda | X) \propto \lambda^{N/2} \lambda^{a-1} \exp\left(-\left(b + \frac{1}{2} \sum_{n=1}^N x_n^2 + \frac{\beta}{2} \mu_0^2 - \underbrace{\frac{\left(\beta\mu_0 + \sum_{n=1}^N x_n\right)^2}{2(N+\beta)}}_{\frac{(N+\beta)}{2}\mu_N^2}\right)\lambda\right) \Leftrightarrow \text{Gamma}\left(\lambda | a_N = \frac{N}{2} + a, b_N\right)$$
  

$$(\lambda(N+\beta))^{1/2} \exp\left(-\frac{\lambda(N+\beta)}{2}\left(\mu - \frac{\beta\mu_0 + \sum_{n=1}^N x_n}{N+\beta}\right)^2\right) \Leftrightarrow \mathcal{N}\left(\mu | \mu_N = \frac{\beta\mu_0 + \sum_{n=1}^N x_n}{N+\beta}, \left(\frac{(N+\beta)}{\beta_N}\lambda\right)^{-1}\right)$$



# The Normal-Gamma Distribution

$$p(\mu, \lambda) = \mathcal{N}\left(\mu \mid \mu_0 = \frac{c}{b}, (\beta\lambda)^{-1}\right) \text{Gamma}\left(\lambda \mid a = 1 + \frac{\beta}{2}, b = d - \frac{c^2}{2\beta}\right)$$



MatLab Code

- K. Murphy, [Conjugate Bayesian Analysis of the Gaussian Distribution](#), 2007 (provides additional results for posterior marginals, posterior predictive, and reference results for an uninformative prior)

# Posterior for $\mu$ and $\sigma$ for Scalar Data

- We can also work directly with  $\sigma^2$ . We use the normal inverse chi-squared distribution (NIX)

$$\begin{aligned} \mathcal{NI}\chi^2(\mu, \sigma^2 | m_0, \kappa_0, v_0, \sigma_0^2) &= \mathcal{N}(\mu | m_0, \sigma^2 / \kappa_0) \chi^{-2}(\sigma^2 | v_0, \sigma_0^2) \\ &\propto \left( \frac{1}{\sigma^2} \right)^{(v_0+3)/2} \exp \left( -\frac{v_0 \sigma_0^2 + \kappa_0 (\mu - m_0)^2}{2\sigma^2} \right) \end{aligned}$$

$$\begin{aligned} \chi^{-2}(\sigma^2 | v_0, \sigma_0^2) &= \mathcal{IG}\left(\sigma^2 | \frac{v_0}{2}, \frac{v_0 \sigma_0^2}{2}\right) \\ &\propto (\sigma^2)^{-v_0/2-1} \exp\left(-\frac{v_0 \sigma_0^2}{2\sigma^2}\right) \end{aligned}$$

- Similarly to our earlier calculations, the posterior is given as:

$$\begin{aligned} p(\mu, \sigma^2 | \mathcal{D}) &= \mathcal{NI}\chi^2(\mu, \sigma^2 | m_N, \kappa_N, v_N, \sigma_N^2), m_N = \frac{\kappa_0 m_0 + N \bar{x}}{\kappa_N} = \frac{\kappa_0}{\kappa_0 + N} m_0 + \frac{N}{\kappa_0 + N} \bar{x} \\ \kappa_N &= \kappa_0 + N, v_N = v_0 + N, v_N \sigma_N^2 = v_0 \sigma_0^2 + \sum_{i=1}^N (x_i - \bar{x})^2 + \frac{N \kappa_0}{\kappa_0 + N} (m_0 - \bar{x})^2 \end{aligned}$$

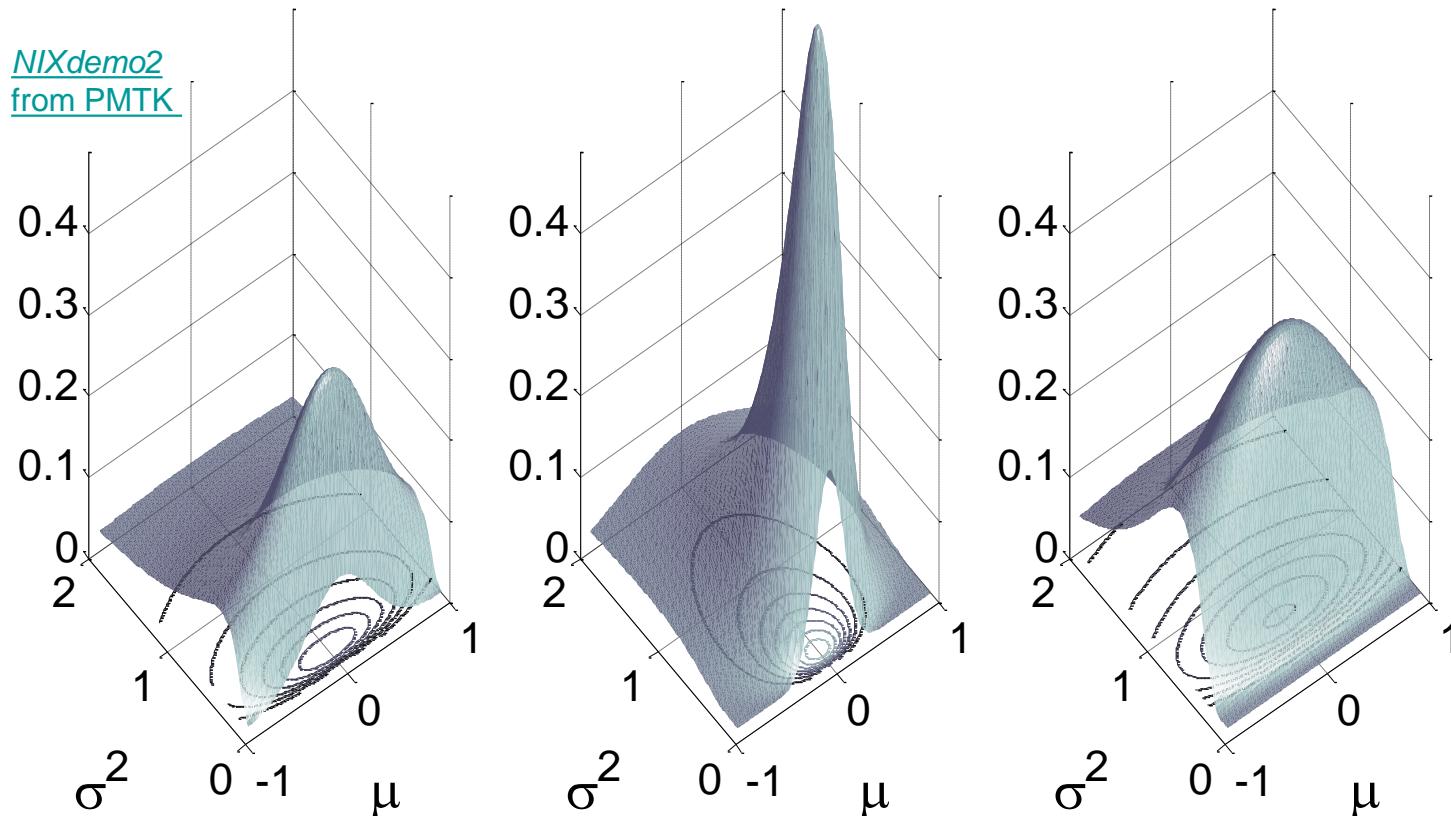
- The posterior marginal for  $\sigma^2$  and posterior expectation are:

$$p(\sigma^2 | \mathcal{D}) = \int p(\mu, \sigma^2 | \mathcal{D}) d\mu = \chi^{-2}(\sigma^2 | v_N, \sigma_N^2), \mathbb{E}[\sigma^2 | \mathcal{D}] = \frac{v_N}{v_N - 2} \sigma_N^2$$

- K. Murphy, Conjugate Bayesian Analysis of the Gaussian Distribution, 2007 (provides additional results for posterior marginals, posterior predictive, and reference results for an uninformative prior, also Section 6 provides the analysis for a normal-inverse-Gamma prior)

# Normal Inverse $\chi^2$ Distribution

$NIX(\mu_0=0, k_0=1, \nu_0=1, \sigma_0^2=1)$     $NIX(\mu_0=0, k_0=5, \nu_0=1, \sigma_0^2=1)$     $NIX(\mu_0=0, k_0=1, \nu_0=5, \sigma_0^2=1)$



- The  $NI\chi^2(\mu_0, \kappa_0, \nu_0, \sigma_0^2)$  distribution.  $\mu_0$  is the prior mean and  $\kappa_0$  is how strongly we believe this;  $\sigma_0^2$  is the prior variance and  $\nu_0$  is how strongly we believe this. (a) The contour plot (underneath the surface) is shaped like a “squashed egg”. (b) We increase the strength of our belief in the mean, so it gets narrower (c) We increase the strength of our belief in the variance, so it gets narrower.

# Posterior for $\mu$ and $\sigma$ for Scalar Data

- The posterior for  $\mu$  is Students'  $\mathcal{T}$ :

$$p(\mu | \mathcal{D}) = \int p(\mu, \sigma^2 | \mathcal{D}) d\sigma^2 = \mathcal{T}(\mu | m_N, \sigma_N^2 / \kappa_N, v_N), \mathbb{E}[\mu | \mathcal{D}] = m_N$$

$$p(x | \mu, \lambda, v) = \frac{\Gamma(\frac{v}{2} + \frac{1}{2})}{\Gamma(\frac{v}{2})} \left( \frac{\lambda}{\pi v} \right)^{1/2} \left[ 1 + \frac{\lambda(x - \mu)^2}{v} \right]^{-v/2-1/2}$$

Mean:  $\mu, v > 1$

Mode:  $\mu$

- Let us revisit these results with *an uninformative prior*:

$$p(\mu, \sigma^2) \propto p(\mu) p(\sigma^2) \propto \sigma^{-2} \propto \mathcal{NI}\chi^2(\mu, \sigma^2 | \mu_0 = 0, \kappa_0 = 0, v_0 = -1, \sigma_0^2 = 0)$$

$$\text{Var: } \frac{v\sigma^2}{v-2} =$$

$$\frac{v}{\lambda(v-2)}, v > 2$$

$$\lambda = \sigma^{-2}$$

- With this prior, the posterior becomes:

$$p(\mu, \sigma^2 | \mathcal{D}) = \mathcal{NI}\chi^2(\mu, \sigma^2 | m_N = \bar{x}, \kappa_N = N, v_N = N-1, \sigma_N^2 = s^2), s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{N}{N-1} \bar{\sigma}_{MLE}^2$$

- $s$  is the sample std. Thus the marginal posterior for  $\mu$  becomes:

$$p(\mu | \mathcal{D}) = \mathcal{T}(\mu | \bar{x}, s^2 / N, N-1), \text{var}[\mu | \mathcal{D}] = \frac{v_N}{v_N - 2} \sigma_N^2 = \frac{N-1}{N-3} \frac{s^2}{N} \rightarrow \frac{s^2}{N}$$

- The standard error of the mean is defined as  $\sqrt{\text{var}[\mu | \mathcal{D}]} \approx \frac{s}{\sqrt{N}}$

- An approximate 95% *posterior credible interval* is thus:  $I_{0.95}[\mu | \mathcal{D}] \approx \bar{x} \pm 2 \frac{s}{\sqrt{N}}$

# Multivariate Gaussian: Posterior of $\mu$

- Consider a known variance  $\Sigma$  and a Gaussian prior  $\mathcal{N}(\mu, \Sigma_0)$  with the posterior for the unknown mean  $\mu$  taking the form:

$$p(\mu | X) \propto p(\mu) \prod_{n=1}^N p(x_n | \mu, \Sigma)$$

- This posterior is the exponential of a quadratic in  $\mu$ :

$$\begin{aligned} -\frac{1}{2}(\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0) - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) = \\ -\frac{1}{2} \mu^T \underbrace{\left( \Sigma_0^{-1} + N \Sigma^{-1} \right)}_{\Sigma_N^{-1}} \mu + \mu^T \underbrace{\left( \Sigma_0^{-1} \mu_0 + \Sigma^{-1} \sum_{n=1}^N x_n \right)}_{\Sigma_N^{-1} \mu_N} + const \end{aligned}$$

- So the variance and mean of the posterior  $p(\mu | X, \Sigma) = \mathcal{N}(\mu | \mu_N, \Sigma_N)$  are:

$$\Sigma_N^{-1} = \Sigma_0^{-1} + N \Sigma^{-1},$$

$$\mu_N = \left( \Sigma_0^{-1} + N \Sigma^{-1} \right)^{-1} \left( \Sigma_0^{-1} \mu_0 + \Sigma^{-1} \sum_{n=1}^N x_n \right) = \left( \Sigma_0^{-1} + N \Sigma^{-1} \right)^{-1} \left( \Sigma_0^{-1} \mu_0 + N \Sigma^{-1} \mu_{ML} \right)$$

- For uninformative prior,  $\Sigma_0 = \infty I \Rightarrow p(\mu | X, \Sigma) \rightarrow \mathcal{N}\left(\mu | \mu_{ML}, \frac{1}{N} \Sigma\right)$

# Posterior Distribution of Precision $\Lambda$

- We now discuss how to compute  $p(\Lambda | \mathcal{D}, \mu)$  for a  $D$ -dimensional Gaussian. The likelihood has the form

$$p(\mathcal{D} | \mu, \Lambda) \propto |\Lambda|^{N/2} \exp\left(-\frac{1}{2} \text{Tr}(\mathbf{S}_\mu \Lambda)\right), \mathbf{S}_\mu = \sum_n (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T$$

- The corresponding conjugate prior is known as [the Wishart distribution](#)

$$\mathcal{W}(\Lambda | \mathbf{W}, v) = B |\Lambda|^{(v-D-1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\mathbf{W}^{-1} \Lambda)\right), v > D-1 \text{ (dof)}, \mathbf{W} \text{ sym pos. def. } (D \times D)$$

$$B(\mathbf{W}, v) = |\mathbf{W}|^{-v/2} \left( 2^{vD/2} \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{v+1-i}{2}\right) \right)^{-1}$$

$$\Gamma_D\left(\frac{v}{2}\right) = \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{v+1-i}{2}\right) \text{(multivariate Gamma Function)}$$

- The following slide shows the similarities of this distribution with the [Gamma prior for  \$\lambda\$  used earlier for univariate Gaussian distributions](#).

$$\text{For } D=1, \mathcal{W}(\lambda | v, s^{-1}) = \text{Gamma}\left(\lambda | \frac{v}{2}, \frac{s}{2}\right)$$

# Wishart Distribution

Wishart	$W \sim \text{Wishart}_\nu(S)$ $p(W) = \text{Wishart}_\nu(W S)$ (implicit dimension $k \times k$ )	degrees of freedom $\nu$ symmetric, pos. definite $k \times k$ scale matrix $S$
---------	--	---

Wishart	$p(W) = \left( 2^{\nu k/2} \pi^{k(k-1)/4} \prod_{i=1}^k \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1}$ $\times  S ^{-\nu/2}  W ^{(\nu-k-1)/2}$ $\times \exp\left(-\frac{1}{2} \text{tr}(S^{-1}W)\right), W \text{ pos. definite}$	$E(W) = \nu S$ mode $= (\nu - k - 1)S$ for $\nu \geq k + 1$
---------	---	---

Gamma	$\theta \sim \text{Gamma}(\alpha, \beta)$ $p(\theta) = \text{Gamma}(\theta \alpha, \beta)$	shape $\alpha > 0$ inverse scale $\beta > 0$
-------	---	---

$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad \theta > 0$	$E(\theta) = \frac{\alpha}{\beta}$ $\text{var}(\theta) = \frac{\alpha}{\beta^2}$ $\text{mode}(\theta) = \frac{\alpha-1}{\beta}, \text{ for } \alpha \geq 1$
--	---

For  $x_i \sim \mathcal{N}(0, \Sigma)$ ,  $S = \sum_{i=1}^N x_i x_i^T$  (scatter matrix)  $\sim \text{Wishart}(S | \Sigma, N) \Rightarrow \mathbb{E}[S] = N\Sigma$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2018, N. Zabaras)



# Posterior Distribution of $\Sigma$

- We similarly discuss computing  $p(\Sigma|\mathcal{D}, \mu)$ . The likelihood has the form

$$p(\mathcal{D}|\mu, \Sigma) \propto |\Sigma|^{-N/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\mathbf{S}_\mu \Sigma^{-1}\right)\right), \quad \mathbf{S}_\mu = \sum_n^N (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T$$

- The corresponding conjugate prior is known as the inverse Wishart

$$\text{InvWi}(\Sigma | S_0, v_0) \propto |\Sigma|^{-(v_0 + D + 1)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\mathbf{S}_0 \Sigma^{-1}\right)\right), \quad \mathbf{S}_0 \text{ sym pos. def.}$$

- $v_0 + D + 1$  controls the strength of the prior, and hence plays a role analogous to the sample size  $N$ .

- The prior scatter matrix is here  $S_0$ .

$$\text{If } \Lambda = \Sigma^{-1} \sim \text{Wi}(\Lambda | S, v) \text{ then } \Sigma \sim \text{InvWi}(S^{-1}, v)$$

- Note: There are many parametrizations of the InvWi. We here follow the notation from Gelman et al. with the same  $v$  for both Wi and InvWi in the Eq. above. In some literature (e.g. K. Murphy's book), the distribution is denoted as  $\text{InvWi}(\Sigma | S_0^{-1}, v_0)$

- Steven W. Nydick, The Wishart and Inverse Wishart Distributions, Report, 2012.
- A. Gelman, J. Carlin, H. Stern and D. Rubin, Bayesian Data Analysis, 2004



# Inverse Wishart Distribution

Inverse-Wishart

$$W \sim \text{Inv-Wishart}_\nu(S)$$

$$p(W) = \text{Inv-Wishart}_\nu(W | S)$$

degrees of freedom  $\nu$   
symmetric, pos. definite  
 $k \times k$  scale matrix  $S$

$$\begin{aligned} p(W) &= \left( 2^{\nu k/2} \pi^{k(k-1)/4} \prod_{i=1}^k \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1} \\ &\times |S|^{\nu/2} |W|^{-(\nu+k+1)/2} \\ &\times \exp\left(-\frac{1}{2} \text{tr}(SW^{-1})\right), W \text{ pos. definite} \end{aligned}$$

$$\begin{aligned} E(W) &= (\nu - k - 1)^{-1} S \\ \text{mode} &= (\nu + k + 1)^{-1} S \end{aligned}$$

Inverse-gamma

$$\theta \sim \text{Inv-gamma}(\alpha, \beta)$$

$$p(\theta) = \text{Inv-gamma}(\theta | \alpha, \beta)$$

shape  $\alpha > 0$   
scale  $\beta > 0$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$$

$$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$$

$$\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$$

$$\text{mode}(\theta) = \frac{\beta}{\alpha+1}$$

$$\text{For } k=1, \text{InvWi}\left(\sigma^2 | \nu, S\right) = \text{InvGamma}\left(\sigma^2 | \frac{\nu}{2}, \frac{S}{2}\right)$$

$$\begin{aligned} \text{If } \lambda \sim \text{Gamma}(a, b) \Rightarrow \frac{1}{\lambda} &\sim \text{InvGamma}(a, b) \\ \text{If } \Sigma^{-1} \sim \mathcal{W}(v, S) \Rightarrow \Sigma &\sim \text{InvWi}\left(v, S^{-1}\right) \end{aligned}$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004

[Statistical Computing, University of Notre Dame, Notre Dame, IN, USA \(Fall 2018, N. Zabaras\)](#)



# Posterior Distribution of $\Sigma$

---

- Multiplying the likelihood and prior, we find that the posterior is also inverse Wishart:

$$\begin{aligned} p(\Sigma | \mathcal{D}, \mu) &\propto |\Sigma|^{-N/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\mathbf{S}_\mu \Sigma^{-1}\right)\right) |\Sigma|^{-(v_0+D+1)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\mathbf{S}_0 \Sigma^{-1}\right)\right), \\ &\propto |\Sigma|^{-[N+(v_0+D+1)]/2} \exp\left(-\frac{1}{2} \text{Tr}\left((\mathbf{S}_\mu + \mathbf{S}_0) \Sigma^{-1}\right)\right) \\ p(\Sigma | \mathcal{D}, \mu) &= \text{InvWi}(\Sigma | N + v_0, \mathbf{S}_N), \quad \mathbf{S}_N = \mathbf{S}_\mu + \mathbf{S}_0 \end{aligned}$$

- The posterior strength  $v_N = v_0 + N$ , is the prior strength  $v_0$  plus the number of observations  $N$ .
- The posterior scatter matrix  $\mathbf{S}_N$  is the prior scatter matrix  $\mathbf{S}_0$  plus the data scatter matrix  $\mathbf{S}_\mu$ .

# MAP Estimation

- From the mode of the inverse Wishart and

$$p(\Sigma | \mathcal{D}, \mu) = \text{InvWi}(\Sigma | v_N, S_N), S_N = S_\mu + S_0, v_N = N + v_0$$

we conclude that the MAP estimate is:

$$\bar{\Sigma}_{MAP} = \frac{S_N}{v_N + D + 1} = \frac{S_\mu + S_0}{\underbrace{N + v_0 + D + 1}_{N_0}} = \frac{S_\mu + S_0}{N + N_0}$$

- For an improper prior,  $S_0 = \mathbf{0}$  and  $N_0 = 0$ ,  $\bar{\Sigma}_{MAP} \rightarrow \frac{S_\mu}{N} = \bar{\Sigma}_{MLE}$
- Consider now the use of a proper informative prior, which is necessary whenever  $D/N$  is large. Let  $\mu = \bar{x} \Rightarrow S_\mu = S_{\bar{x}}$ . Rewrite the *MAP estimate as a convex combination of the prior mode and MLE*

$$\bar{\Sigma}_{MAP} = \frac{S_{\bar{x}} + S_0}{N + N_0} = \underbrace{\frac{N_0}{N + N_0}}_{\lambda} \frac{S_0}{N_0} + \underbrace{\frac{N}{N + N_0}}_{1-\lambda} \frac{S_{\bar{x}}}{N} = \lambda \Sigma_0 + (1 - \lambda) \bar{\Sigma}_{MLE}, \Sigma_0 \equiv \frac{S_0}{N_0} \text{ (prior mode)}$$

where  $\lambda$  controls the amount of shrinkage towards the prior.



# MAP Estimation

$$\bar{\Sigma}_{MAP} = \lambda \Sigma_0 + (1 - \lambda) \bar{\Sigma}_{MLE}, \quad \Sigma_0 \equiv \frac{\mathbf{S}_0}{N_0} \quad (\textit{prior mode})$$

- Can set  $\lambda$  by cross validation.

Alternatively, we can use the formula provided in Ledoit & Wolf and Schaefer & Strimmer which is the optimal frequentist estimate (for squared loss).

This loss function for covariance matrices ignores the positive definite constraint but results in a simple estimator (see PMTK function [shrinkcov](#)).

- For the prior covariance matrix,  $\Sigma_0$ , it is common to use the following (*data dependent*) prior:  $\Sigma_0 = \text{diag}\left(\bar{\Sigma}_{MLE}\right)$

- Ledoit, O. and M. Wolf (2004b). [A well conditioned estimator for large dimensional covariance matrices](#). *J. of Multivariate Analysis* 88(2), 365– 411.
- Ledoit, O. and M. Wolf (2004a). [Honey, I Shrunk the Sample Covariance Matrix](#). *J. of Portfolio Management* 31(1).
- Schaefer, J. and K. Strimmer (2005). [A shrinkage approach to largescale covariance matrix estimation and implications for functional genomics](#). *Statist. Appl. Genet. Mol. Biol* 4(32).



# MAP Shrinkage Estimation

$$\Sigma_0 = \text{diag}(\bar{\Sigma}_{MLE})$$

- In this case, the MAP estimate is:

$$\bar{\Sigma}_{MAP} = \begin{cases} \bar{\Sigma}_{MLE}(i, j) & \text{if } i = j \\ (1 - \lambda)\bar{\Sigma}_{MLE}(i, j) & \text{otherwise} \end{cases}$$

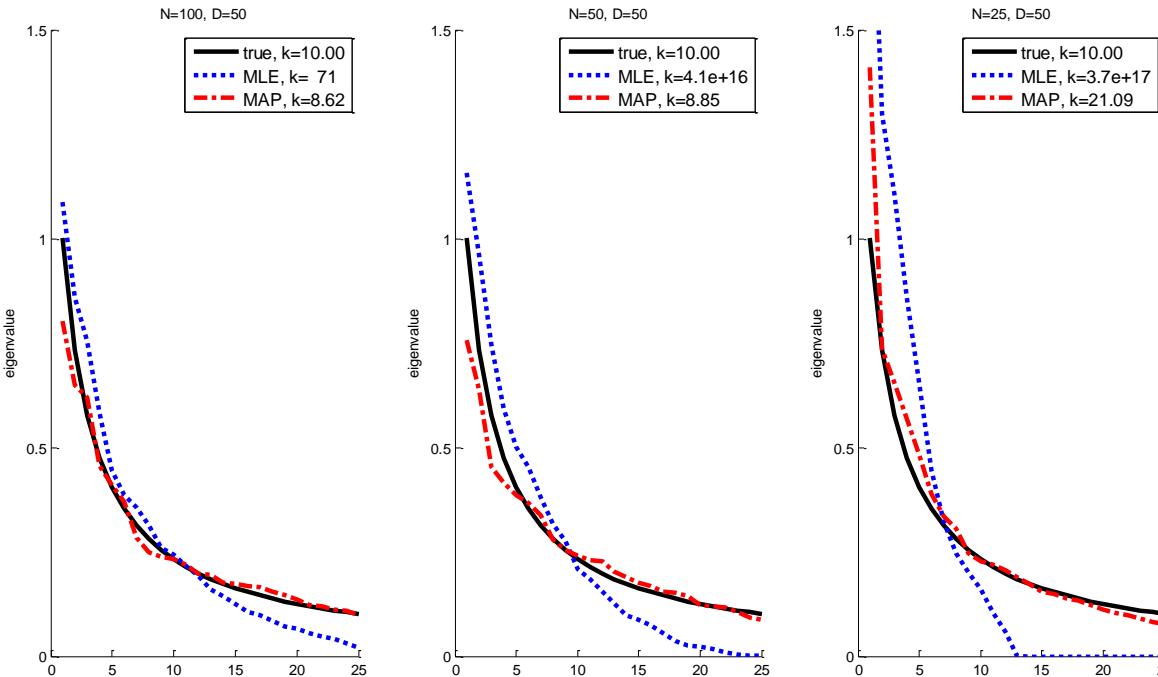
- Thus we see that the diagonal entries are equal to their MLE estimates, and *the off diagonal elements are “shrunk” somewhat towards 0 (shrinkage estimation, or regularized estimation)*.
- The benefits of MAP estimation are illustrated next. We consider fitting a 50 – dim Gaussian to  $N = 100$ ,  $N = 50$  and  $N = 25$  data points.
  - The MAP estimate is always well-conditioned, unlike the MLE.
  - The *eigenvalue spectrum* of the MAP estimate is much closer to that of the true matrix than the MLE's.
  - The eigenvectors, however, are unaffected.



# Posterior Distribution of $\Sigma$

- Estimating a covariance matrix in  $D = 50$  dimensions using  $N \in \{100, 50, 25\}$  samples.

- Eigenvalues in descending order for the true covariance matrix (solid black), MLE (dotted blue) and MAP estimates (dashed red) with  $\lambda = 0.9$ .



- The condition number  $k$  of each matrix is also given in the legend.

[shrinkcovDemo](#)  
from [PMTK](#)

# Inference for Both $\mu$ and $\Lambda$

---

- Suppose  $x_1, x_2, \dots, x_N \sim (\text{i.i.d}) \mathcal{N}(\mu, \Lambda^{-1})$ . We do not know  $\mu$  or  $\Lambda$
- When both sets of parameters are unknown, a conjugate family of priors is one in which

$$\Lambda \sim \mathcal{W}(\Lambda | \nu, T)$$

and

$$\mu | \Lambda \sim \mathcal{N}(\mu_0, (\kappa \Lambda)^{-1})$$

- The Wishart distribution is the multivariate analog of the Gamma distribution (*extension to positive definite matrices*). If matrix  $U$  has the Wishart distribution, then  $U^{-1}$  has the inverse-Wishart distribution. The resulting  $p(\mu, \Lambda | \mu_0, \kappa, T, \nu)$  is the **Gaussian-Wishart distribution**.
- The quantity  $\nu$  is a positive scalar, while  $T$  is a positive definite matrix. They play roles analogous to  $a$  and  $\beta$ , respectively, in the Gamma distribution.
- Other parameters of the prior are the mean vector  $\mu_0$  and  $\kappa$  which represents the ‘a priori number of observations’.



# Inference for Both $\mu$ and $\Lambda$

- The likelihood and prior distributions are given explicitly as:

$$\begin{aligned} p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-ND/2} |\boldsymbol{\Lambda}|^{N/2} \exp\left(-\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x}_n - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-ND/2} |\boldsymbol{\Lambda}|^{N/2} \exp\left(-\frac{N}{2} (\boldsymbol{\mu} - \bar{\mathbf{x}})^T \boldsymbol{\Lambda} (\boldsymbol{\mu} - \bar{\mathbf{x}})\right) \exp\left(-\frac{1}{2} \text{Tr}(\boldsymbol{\Lambda} S_{\bar{\mathbf{x}}})\right) \end{aligned}$$

$$\begin{aligned} p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) &= \mathcal{N}\mathcal{W}\mathcal{I}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{\mu}_0, \kappa_0, v_0, \mathbf{T}_0) = \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_0, (\kappa_0 \boldsymbol{\Lambda})^{-1}) \mathcal{W}\mathcal{I}(\boldsymbol{\Lambda} | \mathbf{T}_0, v_0) = \\ &= \frac{1}{Z} |\boldsymbol{\Lambda}|^{1/2} \exp\left(-\frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right) |\boldsymbol{\Lambda}|^{(v_0 - D - 1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\mathbf{T}_0^{-1} \boldsymbol{\Lambda})\right) \end{aligned}$$

$$Z_{NIW} = \left(\frac{\kappa_0}{2\pi}\right)^{D/2} |\mathbf{T}_0|^{v_0/2} 2^{Dv_0/2} \Gamma_D\left(\frac{v_0}{2}\right), \text{ } \Gamma_D \text{ multivariate Gamma function}$$

- Combining gives:

$$\begin{aligned} p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) &\propto \exp\left(-\frac{N}{2} (\boldsymbol{\mu} - \bar{\mathbf{x}})^T \boldsymbol{\Lambda} (\boldsymbol{\mu} - \bar{\mathbf{x}}) - \frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right) \times \\ &\quad |\boldsymbol{\Lambda}|^{(N + v_0 - D)/2} \exp\left(-\frac{1}{2} \text{Tr}((\mathbf{T}_0^{-1} + S_{\bar{\mathbf{x}}}) \boldsymbol{\Lambda})\right) \end{aligned}$$

- M. DeGroot. *Optimal Statistical Decisions*. McGraw-Hill, 1970.
- K. Murphy, *Conjugate Bayesian Analysis of the Gaussian Distribution*, 2007 ([Section 8](#))

# Inference for Both $\mu$ and $\Lambda$

---

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) \propto \exp\left(-\frac{N}{2}\left(\boldsymbol{\mu} - \bar{\mathbf{x}}\right)^T \boldsymbol{\Lambda} \left(\boldsymbol{\mu} - \bar{\mathbf{x}}\right) - \frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right) \times \\ |\boldsymbol{\Lambda}|^{(N+v_0-D)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\left(\mathbf{T}_0^{-1} + \mathbf{S}_{\bar{\mathbf{x}}}\right) \boldsymbol{\Lambda}\right)\right)$$

□ Can close the square in  $\boldsymbol{\mu}$  as follows:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) \propto \exp\left(-\frac{1}{2}\left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}\right)^T (\kappa_0 + N) \boldsymbol{\Lambda} \left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}\right)\right) \times \\ |\boldsymbol{\Lambda}|^{(v_0+N-D-1)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\left(\mathbf{T}_0^{-1} + \mathbf{S}_{\bar{\mathbf{x}}} + N \bar{\mathbf{x}} \bar{\mathbf{x}}^T + \kappa_0 \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^T - \frac{1}{\kappa_0 + N} (\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}})(\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}})^T\right) \boldsymbol{\Lambda}\right)\right)$$

□ We can simplify as:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) \propto \exp\left(-\frac{1}{2}\left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}\right)^T (\kappa_0 + N) \boldsymbol{\Lambda} \left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}\right)\right) \times \\ |\boldsymbol{\Lambda}|^{(v_0+N-D-1)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\left(\mathbf{T}_0^{-1} + \mathbf{S}_{\bar{\mathbf{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\mathbf{x}})(\boldsymbol{\mu}_0 - \bar{\mathbf{x}})^T\right) \boldsymbol{\Lambda}\right)\right)$$

# Inference for Both $\mu$ and $\Lambda$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) \propto \exp\left(-\frac{1}{2}\left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\boldsymbol{x}}}{\kappa_0 + N}\right)^T (\kappa_0 + N) \boldsymbol{\Lambda} \left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\boldsymbol{x}}}{\kappa_0 + N}\right)\right) \times \\ |\boldsymbol{\Lambda}|^{(v_0 + N - D - 1)/2} \exp\left(-\frac{1}{2} Tr\left(\left(\boldsymbol{T}_0^{-1} + \boldsymbol{S}_{\bar{\boldsymbol{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})(\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})^T\right) \boldsymbol{\Lambda}\right)\right)$$

□ This can be written as:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) = \mathcal{NWi}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{\mu}_N, \kappa_N, v_N, \boldsymbol{T}_N) = \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_N, (\kappa_N \boldsymbol{\Lambda})^{-1}) \mathcal{Wi}(\boldsymbol{\Lambda} | \boldsymbol{T}_N, v_N)$$

$$\boldsymbol{\mu}_N = \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\boldsymbol{x}}}{\kappa_0 + N}$$

$$\boldsymbol{T}_N^{-1} = \boldsymbol{T}_0^{-1} + \boldsymbol{S}_{\bar{\boldsymbol{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})(\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})^T, \text{ where } \boldsymbol{S}_{\bar{\boldsymbol{x}}} = \sum_{i=1}^N (\boldsymbol{x}_i - \bar{\boldsymbol{x}})(\boldsymbol{x}_i - \bar{\boldsymbol{x}})^T$$

$$v_N = v_0 + N, \kappa_N = \kappa_0 + N$$



# Normalization Factors

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) \propto \exp\left(-\frac{1}{2}\left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\boldsymbol{x}}}{\kappa_0 + N}\right)^T (\kappa_0 + N) \boldsymbol{\Lambda} \left(\boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\boldsymbol{x}}}{\kappa_0 + N}\right)\right) \times \\ |\boldsymbol{\Lambda}|^{(v_0 + N - D - 1)/2} \exp\left(-\frac{1}{2} \text{Tr}\left(\left(\boldsymbol{T}_0^{-1} + \boldsymbol{S}_{\bar{\boldsymbol{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})(\boldsymbol{\mu}_0 - \bar{\boldsymbol{x}})^T\right) \boldsymbol{\Lambda}\right)\right)$$

- Using the normalization constants of the multivariate Gaussian and the Wishart:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) = \frac{1}{Z_N} \exp\left(-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T \kappa_N \boldsymbol{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_N)\right) |\boldsymbol{\Lambda}|^{(v_N - D - 1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\boldsymbol{T}_N^{-1} \boldsymbol{\Lambda})\right)$$

$$Z_N = (2\pi)^{D/2} (\kappa_N)^{-D/2} 2^{v_N D/2} \Gamma_D\left(\frac{v_N}{2}\right) |\boldsymbol{T}_N|^{v_N/2}$$

- Similarly for the prior the normalization factor is:

$$Z_0 = (2\pi)^{D/2} (\kappa_0)^{-D/2} 2^{v_0 D/2} \Gamma_D\left(\frac{v_0}{2}\right) |\boldsymbol{T}_0|^{v_0/2}$$

# Inference for Both $\mu$ and $\Lambda$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) = \mathcal{N}\mathcal{W}\mathcal{i}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{\mu}_N, \boldsymbol{\kappa}_N, v_N, \mathbf{T}_N) = \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_N, (\boldsymbol{\kappa}_N \boldsymbol{\Lambda})^{-1}) \mathcal{W}\mathcal{i}(\boldsymbol{\Lambda} | \mathbf{T}_N, v_N)$$

$$\boldsymbol{\mu}_N = \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}, \mathbf{T}_N^{-1} = \mathbf{T}_0^{-1} + \mathbf{S}_{\bar{\mathbf{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\mathbf{x}})(\boldsymbol{\mu}_0 - \bar{\mathbf{x}})^T, \text{ where } \mathbf{S}_{\bar{\mathbf{x}}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

$$v_N = v_0 + N, \kappa_N = \kappa_0 + N$$

□ The posterior marginals can be derived as:

$$p(\boldsymbol{\Lambda} | \mathcal{D}) = \mathcal{W}\mathcal{i}(\boldsymbol{\Lambda} | \mathbf{T}_N, v_N)$$

$$p(\boldsymbol{\mu} | \mathcal{D}) = \mathcal{J}_{v_N-D+1} \left( \boldsymbol{\mu} | \boldsymbol{\mu}_N, \underbrace{\boldsymbol{\kappa}_N v_N \mathbf{T}_N}_{precision} \right)$$

$$\mathcal{J}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}, v) = \frac{\Gamma(\frac{D}{2} + \frac{(v_N - D + 1)}{2})}{\Gamma(\frac{(v_N - D + 1)}{2})} \frac{|\boldsymbol{\kappa}_N v_N \mathbf{T}_N|^{1/2}}{(\pi(v_N - D + 1))^{D/2}} \left[ 1 + \frac{(\mathbf{x} - \boldsymbol{\mu}_N)^T \boldsymbol{\kappa}_N v_N \mathbf{T}_N (\mathbf{x} - \boldsymbol{\mu}_N)}{(v_N - D + 1)} \right]^{-(v_N - D + 1)/2 - D/2}$$

\* Refer to [this report](#) for these results based on Bernardo and Smith.



# Inference for Both $\mu$ and $\Lambda$

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) = \mathcal{N}\mathcal{W}\mathcal{i}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{\mu}_N, \kappa_N, v_N, \mathbf{T}_N) = \mathcal{N}\left(\boldsymbol{\mu} | \boldsymbol{\mu}_N, (\kappa_N \boldsymbol{\Lambda})^{-1}\right) \mathcal{W}\mathcal{i}(\boldsymbol{\Lambda} | \mathbf{T}_N, v_N)$$

$$\boldsymbol{\mu}_N = \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}, \mathbf{T}_N^{-1} = \mathbf{T}_0^{-1} + S_{\bar{\mathbf{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\boldsymbol{\mu}_0 - \bar{\mathbf{x}})(\boldsymbol{\mu}_0 - \bar{\mathbf{x}})^T, \text{ where } S_{\bar{\mathbf{x}}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

$$v_N = v_0 + N, \kappa_N = \kappa_0 + N$$

- Differentiating the Eq. [on the top of this slide](#), we can also derive the MAP estimates as:

$$(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Lambda}}) = \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Lambda}} p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D})$$

$$\bar{\boldsymbol{\mu}} = \frac{\sum_{i=1}^N \mathbf{x}_i + \kappa_0 \boldsymbol{\mu}_0}{N + \kappa_0}$$

$$\bar{\boldsymbol{\Lambda}} = \frac{\sum_{i=1}^N (\mathbf{x}_i - \bar{\boldsymbol{\mu}})(\mathbf{x}_i - \bar{\boldsymbol{\mu}})^T + \kappa_0 (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)(\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)^T + \mathbf{T}_0^{-1}}{N + v_0 - D}$$

- These are reduced to the MLE by setting

$$\kappa_0 = 0, v_0 = D, |\mathbf{T}_0^{-1}| = 0$$



# Inference for both $\mu$ and $\Lambda$

- The posterior predictive is:  $p(\mathbf{x} | \mathcal{D}) = \mathcal{J}_{v_N-D+1} \left( \mu_N, \underbrace{\frac{\kappa_N(v_N - D + 1)T_N}{(\kappa_N + 1)}}_{precision} \right)$
- The marginal likelihood is computed as a ratio of normalization constants using  $p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathcal{D}) = p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) / p(\mathcal{D})$ :  
$$p(\mathcal{D}) = \frac{Z_N}{Z_0} \frac{1}{(2\pi)^{ND/2}} = \frac{1}{\pi^{ND/2}} \frac{\Gamma_D(v_N/2)}{\Gamma_D(v_0/2)} \frac{|T_N|^{v_N/2}}{|T_0|^{v_0/2}} \left( \frac{\kappa_0}{\kappa_N} \right)^{D/2}$$

- A useful reference analysis considers

$$\mu_0 = 0, \kappa_0 = 0, v_0 = -1, |T_0^{-1}| = 0$$

- This results in the following for the prior:

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \propto |\boldsymbol{\Lambda}|^{-(D+1)/2}$$

- The posterior parameters are also simplified as:

$$\mu_N = \bar{\mathbf{x}}, T_N^{-1} = S_{\bar{\mathbf{x}}}^{-1}, \kappa_N = N, v_N = N - 1$$

- The posterior marginals and posterior predictive are given as:

$$p(\boldsymbol{\Lambda} | \mathcal{D}) = \mathcal{W}_{N-D}(\boldsymbol{\Lambda} | S_{\bar{\mathbf{x}}}^{-1}), p(\boldsymbol{\mu} | \mathcal{D}) = \mathcal{J}_{N-D}\left(\boldsymbol{\mu} | \bar{\mathbf{x}}, \frac{S_{\bar{\mathbf{x}}}^{-1}}{N(N-D)}\right)$$

$$p(\mathbf{x} | \mathcal{D}) = \mathcal{J}_{N-D}\left(\bar{\mathbf{x}}, \frac{S_{\bar{\mathbf{x}}}^{-1}(N+1)}{N(N-D)}\right)$$



# Inference for $\mu$ and $\Sigma$

- For the case of the multivariate Gaussian of a  $D$ -dim variable  $\mathbf{x}, \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with both the mean and variance unknowns, the likelihood is

$$\begin{aligned}\prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-ND/2} |\boldsymbol{\Sigma}|^{-N/2} \exp\left(-\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-ND/2} |\boldsymbol{\Sigma}|^{-N/2} \exp\left(-\frac{N}{2} (\boldsymbol{\mu} - \bar{\mathbf{x}})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \bar{\mathbf{x}})\right) \exp\left(-\frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_{\bar{\mathbf{x}}})\right)\end{aligned}$$

- The conjugate prior is given as *the product of a Gaussian and the Inverse Wishart distribution:*

$$\begin{aligned}\mathcal{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\mu}_0, \kappa_0, \mathbf{S}_0, v_0) &= \mathcal{N}\left(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \frac{1}{\kappa_0} \boldsymbol{\Sigma}\right) \mathcal{IW}(\boldsymbol{\Sigma} | \mathbf{S}_0, v_0) = \\ &= \frac{1}{Z_{NIW}} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right) |\boldsymbol{\Sigma}|^{-(v_0 + D + 1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right) \\ &= \frac{1}{Z_{NIW}} \exp\left(-\frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_0)\right) |\boldsymbol{\Sigma}|^{-(v_0 + D + 2)/2} \\ Z_{NIW} &= 2^{v_0 D/2} \Gamma_D\left(\frac{v_0}{2}\right) \left(\frac{2\pi}{\kappa_0}\right)^{D/2} |\mathbf{S}_0|^{-v_0/2}, \Gamma_D \text{ multivariate Gamma function}\end{aligned}$$

# Inference for $\mu$ and $\Sigma$

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) \propto |\boldsymbol{\Sigma}|^{-\frac{N-1}{2} - \frac{v_0+D+1}{2}} \exp\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\left(\mathbf{S}_{\bar{x}} + \mathbf{S}_0 + N(\boldsymbol{\mu} - \bar{\mathbf{x}})(\boldsymbol{\mu} - \bar{\mathbf{x}})^T + \kappa_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right)\right)$$

$$\propto |\boldsymbol{\Sigma}|^{-\frac{v_0+D+2+N}{2}} \exp\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\left(\mathbf{S}_{\bar{x}} + \mathbf{S}_0 + N(\boldsymbol{\mu} - \bar{\mathbf{x}})(\boldsymbol{\mu} - \bar{\mathbf{x}})^T + \kappa_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right)\right)$$

➤ One can show by completing the square in  $\boldsymbol{\mu}$  that:

$$N(\boldsymbol{\mu} - \bar{\mathbf{x}})(\boldsymbol{\mu} - \bar{\mathbf{x}})^T + \kappa_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T =$$

$$\underbrace{(\kappa_0 + N)}_{\kappa_N} \left( \boldsymbol{\mu} - \underbrace{\frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N}}_{\boldsymbol{\mu}_N} \right) \left( \boldsymbol{\mu} - \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_0 + N} \right)^T + \frac{\kappa_0 N}{\kappa_0 + N} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T$$

➤ Thus:

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) \propto |\boldsymbol{\Sigma}|^{-\frac{v_N+D+2}{2}} \exp\left(-\frac{1}{2}\boldsymbol{\Sigma}^{-1}\left(\mathbf{S}_{\bar{x}} + \mathbf{S}_0 + \kappa_N(\boldsymbol{\mu} - \boldsymbol{\mu}_N)(\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T + \frac{\kappa_0 N}{\kappa_N} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T\right)\right),$$

where:  $v_N = v_0 + N$

# The Posterior of $\mu$ and $\Sigma$

➤ The posterior is  $\mathcal{NIW}$  given as:

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}, \boldsymbol{\mu}_0, \kappa_0, \mathbf{S}_0, v_0) = \mathcal{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\mu}_N, \kappa_N, \mathbf{S}_N, v_N)$$

$$\boldsymbol{\mu}_N = \frac{\kappa_0 \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\kappa_N} = \frac{\kappa_0}{\kappa_0 + N} \boldsymbol{\mu}_0 + \frac{N}{\kappa_0 + N} \bar{\mathbf{x}}$$

$$\kappa_N = \kappa_0 + N, v_N = v_0 + N$$

$$\mathbf{S}_N = \mathbf{S}_0 + \mathbf{S}_{\bar{\mathbf{x}}} + \frac{\kappa_0 N}{\kappa_0 + N} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0) (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T = \mathbf{S}_0 + \mathbf{S} + \kappa_0 \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^T - \kappa_N \boldsymbol{\mu}_N \boldsymbol{\mu}_N^T, \mathbf{S} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$$

- The posterior mean is a convex combination of the prior mean and the MLE with strength  $\kappa_0 + N$ .
- The posterior scatter matrix  $\mathbf{S}_N$  is the prior scatter matrix  $\mathbf{S}_0$  plus the empirical scatter matrix  $\mathbf{S}_{\bar{\mathbf{x}}}$  plus an extra term due to the uncertainty in the mean which creates its own scatter matrix.

- Minka, T. (2000). [Inferring a Gaussian distribution](#). Technical report, MIT.
- [Chipman, H., E. George, and R. McCulloch](#) (2001). [The practical implementation of Bayesian Model Selection. Model Selection](#). IMS Lecture Notes.
- Fraley, C. and A. Raftery (2007). [Bayesian Regularization for Normal Mixture Estimation and Model-Based Clustering](#). *J. of Classification* 24, 155–181
- K. Murphy, [Conjugate Bayesian Analysis of the Gaussian Distribution](#), 2007



# MAP Estimate of $\mu$ and $\Sigma$

$$p(\mu, \Sigma | \mathcal{D}, \mu_0, \kappa_0, S_0, v_0) = \mathcal{NIW}(\mu, \Sigma | \mu_N, \kappa_N, v_N, S_N)$$

$$\mu_N = \frac{\kappa_0 \mu_0 + N \bar{x}}{\kappa_N} = \frac{\kappa_0}{\kappa_0 + N} \mu_0 + \frac{N}{\kappa_0 + N} \bar{x}, \kappa_N = \kappa_0 + N, v_N = v_0 + N$$

$$S_N = S_0 + S_{\bar{x}} + \frac{\kappa_0 N}{\kappa_0 + N} (\bar{x} - \mu_0) (\bar{x} - \mu_0)^T = S_0 + S + \kappa_0 \mu_0 \mu_0^T - \kappa_N \mu_N \mu_N^T, S = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$$

- The mode of the joint posterior is:

$$(\bar{\mu}, \bar{\Sigma}) = \arg \max_{\mu, \Sigma} p(\mu, \Sigma | \mathcal{D}, \mu_0, \kappa_0, S_0, v_0) = \left( \mu_N, \frac{S_N}{v_N + D + 2} \right)$$

- For  $k_0 = 0$ , this becomes:

$$(\bar{\mu}, \bar{\Sigma}) = \arg \max_{\mu, \Sigma} p(\mu, \Sigma | \mathcal{D}, \mu_0, \kappa_0 = 0, S_0, v_0) = \left( \bar{x}, \frac{S_0 + S_{\bar{x}}}{v_0 + N + D + 2} \right)$$

- It is interesting to note that this mode is almost the same as [the MAP estimate shown next](#) – it differs by 1 in the denominator as the mode above is the mode of the joint posterior and not of the marginal.

# The Posterior Marginals of $\mu$ and $\Sigma$

- The posterior marginal for  $\Sigma$  and  $\mu$  are:

$$p(\Sigma | \mathcal{D}) = \int p(\mu, \Sigma | \mathcal{D}) d\mu = \mathcal{IW}(\Sigma | S_N, v_N)$$
$$\bar{\Sigma}_{MAP} = \frac{S_N}{v_N + D + 1}, \mathbb{E}[\Sigma] = \frac{S_N}{v_N - D - 1}$$
$$p(\mu | \mathcal{D}) = \int p(\mu, \Sigma | \mathcal{D}) d\Sigma = \mathcal{T}_{v_N - D + 1}\left(\mu | \mu_N, \frac{1}{\kappa_N(v_N - D + 1)} S_N\right)$$

- It is not surprising that the last marginal is Student's  $\mathcal{T}$  that we know can be represented as a mixture of Gaussians.
- To see the connection for the scalar case ( $D = 1$ ), note that  $S_N$  plays the role of the posterior sum of squares  $v_N \sigma_N^2$  (you may want to revisit the results for the scalar case of simultaneously estimating  $\mu$  and  $\sigma^2$ ):

$$\frac{1}{\kappa_N(v_N - D + 1)} S_N = \frac{S_N}{\kappa_N v_N} = \frac{\sigma_N^2}{\kappa_N}$$



# The Posterior Predictive of $\mu$ and $\Sigma$

- The posterior predictive  $p(\mathbf{x}|\mathcal{D}) = p(\mathbf{x}, \mathcal{D})/p(\mathcal{D})$  can be evaluated as:

$$\begin{aligned} p(\mathbf{x} | \mathcal{D}) &= \int \int \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathcal{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\mu}_N, \kappa_N, v_N, S_N) d\boldsymbol{\mu} d\boldsymbol{\Sigma} \\ &= \mathcal{T}_{v_N - D + 1} \left( \mathbf{x} | \boldsymbol{\mu}_N, \frac{\kappa_N + 1}{\kappa_N(v_N - D + 1)} S_N \right) \end{aligned}$$

- Recall that the Student's  $\mathcal{T}$  distribution has heavier tails than the Gaussian but rapidly becomes Gaussian like.
- To see the connection of the above expression with the scalar case, note:

$$\frac{\kappa_N + 1}{\kappa_N(v_N - D + 1)} S_N = \frac{(\kappa_N + 1)v_N \sigma_N^2}{\kappa_N v_N} = \frac{(\kappa_N + 1)\sigma_N^2}{\kappa_N}$$

- K. Murphy, [Conjugate Bayesian Analysis of the Gaussian Distribution](#), 2007 ([Section 9](#))



# Marginal Likelihood

➤ The posterior is given as:

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}, \boldsymbol{\mu}_0, \boldsymbol{\kappa}_0, \mathbf{S}_0, v_0) = \frac{1}{p(\mathcal{D})} \frac{1}{Z_0} \mathcal{NIW}'(\boldsymbol{\mu}, \boldsymbol{\Sigma} | a_0) \frac{1}{(2\pi)^{ND/2}} \mathcal{N}'(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{Z_N} \mathcal{NIW}'(\boldsymbol{\mu}, \boldsymbol{\Sigma} | a_N)$$

$$\mathcal{NIW}'(\boldsymbol{\mu}, \boldsymbol{\Sigma} | a_0) = |\boldsymbol{\Sigma}|^{-(v_0+D)/2+1} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \boldsymbol{\Sigma}^{-1}) - \frac{\kappa_0}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right)$$

$$\mathcal{N}'(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |\boldsymbol{\Sigma}|^{-N/2} \exp\left(-\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathcal{D})\right)$$

➤ In the last two expressions ()' give the *unnormalized likelihood and prior*.

➤ The marginal likelihood  $p(\mathcal{D}) = \frac{Z_N}{Z_0} \frac{1}{(2\pi)^{ND/2}}$  is then:

$$p(\mathcal{D}) = \frac{\frac{2^{v_N D/2} \Gamma_D\left(\frac{v_N}{2}\right) \left(\frac{2\pi}{\kappa_N}\right)^{D/2}}{|\mathbf{S}_N|^{v_N/2}}}{\frac{1}{2^{v_0 D/2} \Gamma_D\left(\frac{v_0}{2}\right) \left(\frac{2\pi}{\kappa_0}\right)^{D/2}}} \frac{1}{(2\pi)^{ND/2}} = \frac{1}{(2\pi)^{ND/2}} \frac{2^{v_N D/2} \left(\frac{2\pi}{\kappa_N}\right)^{D/2}}{\left(\frac{2\pi}{\kappa_0}\right)^{D/2}} \frac{|\mathbf{S}_0|^{v_0/2} \Gamma_D\left(\frac{v_N}{2}\right)}{\Gamma_D\left(\frac{v_0}{2}\right)} = \frac{1}{\pi^{ND/2}} \frac{\Gamma_D\left(\frac{v_N}{2}\right)}{\Gamma_D\left(\frac{v_0}{2}\right)} \frac{|\mathbf{S}_0|^{v_0/2}}{|\mathbf{S}_N|^{v_N/2}} \left(\frac{\kappa_0}{\kappa_N}\right)^{D/2}$$

➤ Note that for  $D = 1$ , this reduces to the familiar Equ.

- K. Murphy, [Conjugate Bayesian Analysis of the Gaussian Distribution](#), 2007 (See calculation of the [marginal likelihood for 1D analysis of the Normal-Inverse-Chi-Squared prior on Section 5](#))



# Non Informative Prior

- The uninformative Jeffrey's prior is  $p(\mu, \Sigma) \propto |\Sigma|^{-(D+1)/2}$ . This is obtained in the limit

$$\kappa_0 \rightarrow 0, v_0 \rightarrow -1, |S_0| \rightarrow 0$$

$$p(\mu, \Sigma | \mathcal{D}) \propto |\Sigma|^{-\frac{v_0+D+2}{2}} \exp\left(-\frac{1}{2}\Sigma^{-1} \left(S_0 + \kappa_0(\mu - \mu_0)(\mu - \mu_0)^T\right)\right) = |\Sigma|^{-\frac{D+1}{2}}$$

- In this case, we have:

$$\mu_N = \bar{x}, \kappa_N = N, v_N = N-1, S_N = S_{\bar{x}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

- The posterior marginals are then given as:

$$p(\Sigma | \mathcal{D}) = \mathcal{IW}(\Sigma | S, N-1), p(\mu | \mathcal{D}) = \mathcal{T}_{N-D}\left(\mu | \bar{x}, \frac{1}{N(N-D)} S\right)$$

- Also the posterior predictive is:

$$p(\mathbf{x} | \mathcal{D}) = \mathcal{T}_{N-D}\left(\mathbf{x} | \bar{\mathbf{x}}, \frac{N+1}{N(N-D)} S\right)$$

- Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004 (pp. 88)
- K. Murphy, [Conjugate Bayesian Analysis of the Gaussian Distribution](#), 2007 ([See Section 9](#))



# Non Informative Prior

- Based on the report of Minka below, the uninformative prior should be instead

$$\lim_{\kappa \rightarrow 0} \mathcal{N}\left(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \frac{1}{\kappa} \boldsymbol{\Sigma}\right) \mathcal{INWIS}(\boldsymbol{\Sigma} | \boldsymbol{S}_0, \kappa)$$

$v_0 = 0$  instead of  $v_0 \rightarrow -1$

$$\propto |2\pi\boldsymbol{\Sigma}|^{-1/2} |\boldsymbol{\Sigma}|^{-(D+1)/2} \propto |\boldsymbol{\Sigma}|^{-(D/2+1)} \propto \mathcal{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{0}, 0, 0, \mathbf{I}, 0)$$

- Often, a *data-dependent weakly informative prior is recommended* (see Chipman et al. and Fraley and Raftery):

Set :  $\boldsymbol{S}_0 = \frac{\text{diag} \bar{\boldsymbol{x}}}{N}$ ,  $v_0 = D + 2$  to ensure  $\mathbb{E}[\boldsymbol{\Sigma}] = \boldsymbol{S}_0$ , and  
 $\boldsymbol{\mu}_0 = \bar{\boldsymbol{x}}$ ,  $\kappa_0 = 0.01$

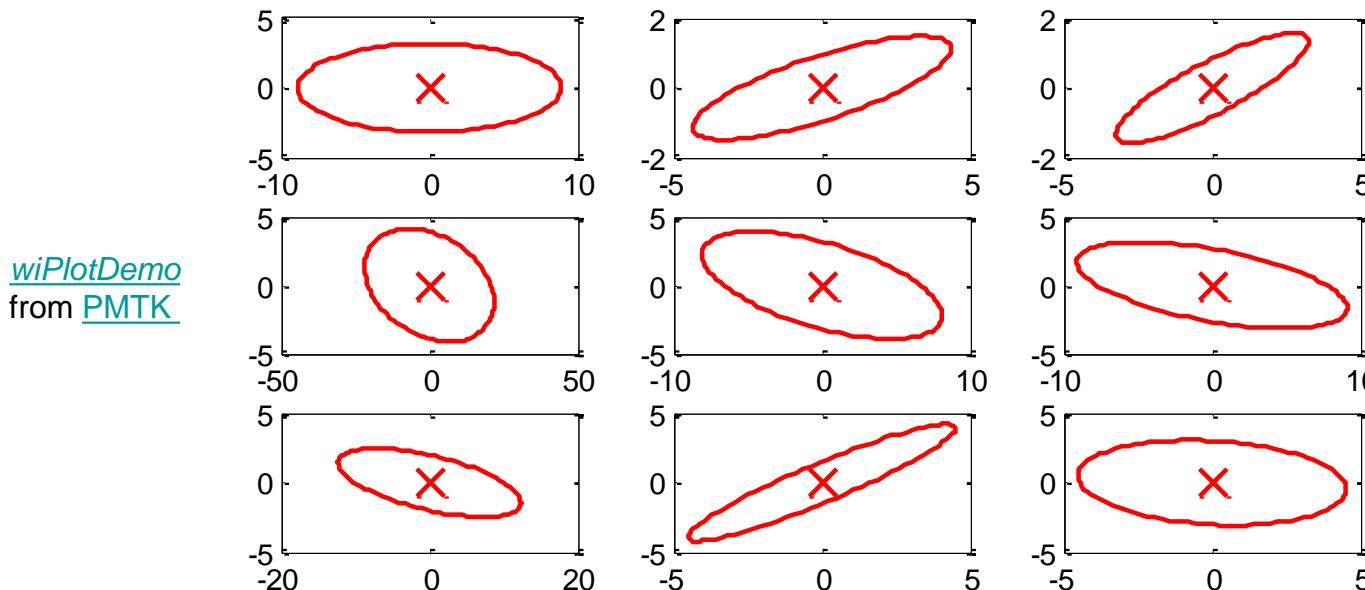
- Minka, T. (2000). [Inferring a Gaussian distribution](#). Technical report, MIT.
- [Chipman, H., E. George, and R. Mc-Culloch](#) (2001). [The practical implementation of Bayesian Model Selection. Model Selection. IMS Lecture Notes](#).
- Fraley, C. and A. Raftery (2007). [Bayesian Regularization for Normal Mixture Estimation and Model-Based Clustering](#). *J. of Classification* 24, 155–181



# Visualization of the Wishart Distribution

- Wi is a distribution over matrices thus difficult to plot the PDF. However, one can sample from it and in 2d use the eigen-vectors of the resulting sample to define an ellipse as we have done for the 2D Gaussian.

Wi(dof=3.0, S), E=[9.5, -0.1; -0.1, 1.9],  $\rho=-0.018$



- Above: Samples from  $\Sigma \sim Wi(S, \nu)$ , where  $S = [3.1653, -0.0262; -0.0262, 0.6477]$  and  $\nu = 3$ .
- The sampled matrices are highly variable, and some are nearly singular. As  $\nu$  increases, the sampled matrices are more concentrated on the prior  $S$ .

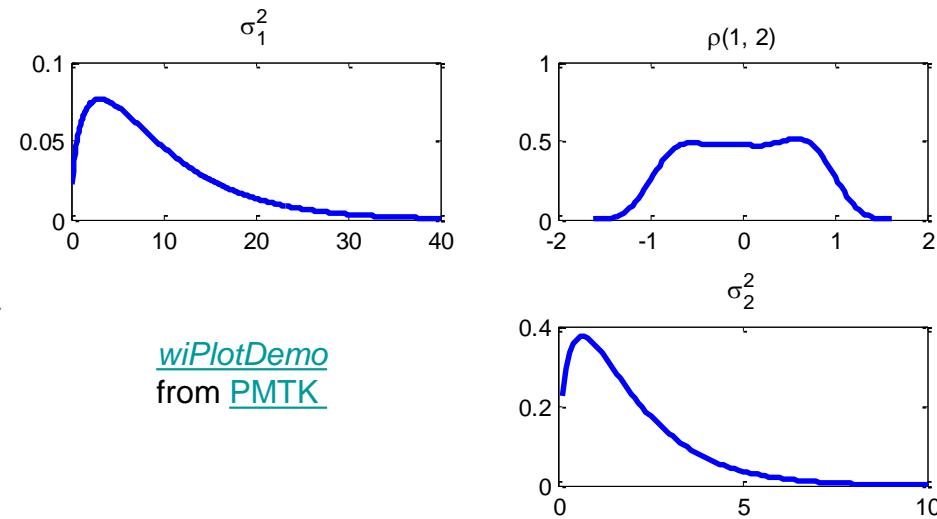
# Visualization of the Wishart Distribution

- For off-diagonal elements, one can sample matrices from the distribution, and then compute their distribution empirically.
- We can *convert each sampled matrix to a correlation matrix*, and thus compute a Monte Carlo approximation

$$\mathbb{E}[R_{ij}] \approx \frac{1}{S} \sum_{s=1}^S R(\Sigma^{(s)})_{ij}, \Sigma^{(s)} \sim \mathcal{W}(\Sigma, v), R(\Sigma)_{ij} = \frac{\Sigma_{ij}}{\sqrt{\Sigma_{ii}\Sigma_{jj}}}$$

- We can then use kernel density estimation to produce for plotting purposes a smooth approximation to the univariate density  $\mathbb{E}[R_{ij}]$ .

- On the right: Plots of the marginals (which are *Gamma*), and the sample-based marginal on the correlation coefficient



# **Exponential Family**

---

- Large family of useful distributions with common properties
  - Bernoulli, beta, binomial, chi-square, Dirichlet, gamma, Gaussian, geometric, multinomial, Poisson, Weibull, . . .
- Not in the family:
  - ✓ Uniform,
  - ✓ Student's T,
  - ✓ Cauchy,
  - ✓ Laplace,
  - ✓ mixture of Gaussians,
  - ✓ . . .
- Variable can be discrete/continuous (or vectors thereof)

# Exponential Family

- The exponential family of distributions over  $x$ , given parameters  $\eta$ , is defined to be the set of distributions of the form

$$p(x | \eta) = h(x)g(\eta)\exp\{\eta^T u(x)\} \text{ or}$$

$$p(x | \eta) = h(x)\exp\{\eta^T u(x) - A(\eta)\}, \text{ where } A(\eta) = -\log g(\eta)$$

$x$  is scalar/vector, discrete/continuous.  **$\eta$  are the natural parameters and  $u(x)$  is referred to as a sufficient statistic.**

- $g(\eta)$  ensures that the distribution is normalized and satisfies

$$g(\eta) \int h(x) \exp\{\eta^T u(x)\} dx = 1$$

- The normalization factor  $Z$  and the log of it  $A$  are defined as:

$$Z(\eta) = \frac{1}{g(\eta)}, A(\eta) = \ln Z(\eta) = -\ln g(\eta) = \ln \int h(x) \exp\{\eta^T u(x)\} dx$$

$$p(x | \eta) = h(x) \exp\{\eta^T u(x)\} / Z(\eta)$$

- The space of  $\eta$  for which  $\int h(x) \exp\{\eta^T u(x)\} dx < \infty$  is the **natural parameter space**.



# **Canonical or Natural Parameters**

---

- When the parameter  $\theta$  enters the exponential family as  $\eta(\theta)$ , we write the probability density of the exponential family as follows:

$$p(x | \theta) = h(x)g(\eta(\theta))\exp\{\eta^T(\theta)u(x)\} \text{ or}$$

$$p(x | \theta) = h(x)\exp\{\eta^T(\theta)u(x) - A(\eta(\theta))\},$$

$$\text{where : } A(\eta(\theta)) = -\log g(\eta(\theta))$$

- $\eta(\theta)$  are the canonical or natural parameters,
- $\theta$  is the parameter vector of some distribution that can be written in the exponential family format

# Exponential Family: The Bernoulli Distribution

- Consider the Bernoulli distribution:

$$\begin{aligned} p(x | \mu) &= \mathcal{B}\text{ern}(x | \mu) = \mu^x (1 - \mu)^{1-x} = \exp \left\{ x \ln \mu + (1-x) \ln(1-\mu) \right\} = \\ &= \underbrace{(1-\mu)}_{g(\eta)} \exp \left\{ \ln \left( \underbrace{\frac{\mu}{1-\mu}}_{\eta} \right) x \right\} \end{aligned}$$
$$\begin{aligned} p(x | \boldsymbol{\eta}) &= h(x) g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T u(x) \right\} \\ &= h(x) \exp \left\{ \boldsymbol{\eta}^T u(x) - A(\boldsymbol{\eta}) \right\} \end{aligned}$$

- From this we see that (note that *the relation  $\mu(\eta)$  is invertible*)

$$\eta = \ln \left( \frac{\mu}{1-\mu} \right) \Rightarrow$$

$$\mu = \sigma(\eta) = \frac{1}{1+e^{-\eta}}$$

Logistic sigmoid function

and

$$g(\eta) = 1 - \mu = 1 - \sigma(\eta) = \sigma(-\eta)$$

- Finally:

$$\begin{aligned} p(x | \boldsymbol{\eta}) &= g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T u(x) \right\}, u(x) = x, h(x) = 1, g(\boldsymbol{\eta}) = \sigma(-\boldsymbol{\eta}), \\ A(\boldsymbol{\eta}) &= -\ln g(\boldsymbol{\eta}) = -\log(1 - \mu) = \log(1 + e^{\boldsymbol{\eta}}) \end{aligned}$$

# Exponential Family: The Beta Distribution

- Consider the Beta distribution

$$\text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \exp[(a-1) \ln \mu + (b-1) \ln(1-\mu)]$$

- Comparing this with our exponential family:

$$\begin{aligned} p(x | \boldsymbol{\eta}) &= h(x)g(\boldsymbol{\eta})\exp\{\boldsymbol{\eta}^T u(x)\} \\ &= h(x)\exp\{\boldsymbol{\eta}^T u(x) - A(\boldsymbol{\eta})\} \end{aligned}$$

we can easily identify:

$$u(\mu) = (\ln \mu, \ln(1-\mu))^T, \boldsymbol{\eta} = (a-1, b-1)^T, h(\mu) = 1, g(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)},$$

$$A(a, b) = -\ln g(\boldsymbol{\eta}) = \ln \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

# Exponential Family: The Gaussian

- Consider the univariate Gaussian

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 + \frac{\mu}{\sigma^2}x\right\}$$

- Comparing this with our exponential family:

$$p(x | \boldsymbol{\eta}) = h(x)g(\boldsymbol{\eta})\exp\left\{\boldsymbol{\eta}^T u(x)\right\} = h(x)\exp\left\{\boldsymbol{\eta}^T u(x) - A(\boldsymbol{\eta})\right\}$$

we can identify (this is a two parameter distribution):

$$u(x) = (x, x^2)^T, \boldsymbol{\eta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)^T, h(x) = \frac{1}{\sqrt{2\pi}}, g(\boldsymbol{\eta}) = (-2\eta_2)^{1/2} \exp \frac{\eta_1^2}{4\eta_2}$$

$$A(\boldsymbol{\eta}) = -\ln g(\boldsymbol{\eta}) = -\frac{1}{2} \ln(-2\eta_2) - \frac{\eta_1^2}{4\eta_2}$$

# Conjugate Priors

- In general, for a given probability distribution  $p(x|\boldsymbol{\eta})$ , we can seek a prior  $p(\boldsymbol{\eta})$  that is conjugate to the likelihood function, so that the posterior distribution has the same functional form as the prior.
  - For the Bernoulli, the conjugate prior is the Beta distribution
  - For the Gaussian, the conjugate prior for the mean is a Gaussian, and the conjugate prior for the precision is the Wishart distribution



# Conjugate Priors

- For any member of the exponential family with likelihood

$$p(x | \theta) = h(x)g(\eta(\theta))\exp\{\eta^T(\theta)u(x)\}$$

there exists a conjugate prior that can be written in the form

$$p(\theta | v_0, \tau_0) \propto g(\eta(\theta))^{v_0} \exp\{\eta^T(\theta)\tau_0\} = \exp\{v_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))v_0\}, \text{ where: } \tau_0 \equiv v_0 \bar{\tau}_0$$

- In normalized form, we write:

$$p(\theta | v_0, \tau_0) = \frac{1}{Z(v_0, \tau_0)} g(\eta(\theta))^{v_0} \exp\{\eta^T(\theta)\tau_0\} = \frac{1}{Z(v_0, \tau_0)} \exp\{v_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))v_0\}$$

$$\text{where: } Z(v_0, \tau_0) = \int \exp\{v_0 \eta^T(\theta)\bar{\tau}_0 - A(\eta(\theta))v_0\} d\theta$$

# Conjugate Priors

$$p(X | \boldsymbol{\theta}) = \prod_{n=1}^N \left( h(\mathbf{x}_n) g(\boldsymbol{\eta}(\boldsymbol{\theta})) \exp\left\{ \boldsymbol{\eta}^T(\boldsymbol{\theta}) u(\mathbf{x}_n) \right\} \right) = \prod_{n=1}^N \left( h(\mathbf{x}_n) \right) g(\boldsymbol{\eta}(\boldsymbol{\theta}))^N \exp\left\{ \boldsymbol{\eta}^T(\boldsymbol{\theta}) \sum_{n=1}^N u(\mathbf{x}_n) \right\}$$

$$p(\boldsymbol{\theta} | \nu_0, \boldsymbol{\tau}_0) = \frac{1}{Z(\nu_0, \boldsymbol{\tau}_0)} g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_0} \exp\{\boldsymbol{\eta}^T(\boldsymbol{\theta}) \boldsymbol{\tau}_0\} = \frac{1}{Z(\nu_0, \boldsymbol{\tau}_0)} \exp\{\nu_0 \boldsymbol{\eta}^T(\boldsymbol{\theta}) \bar{\boldsymbol{\tau}}_0 - A(\boldsymbol{\eta}(\boldsymbol{\theta})) \nu_0\}$$

□ Using  $\bar{\mathbf{u}} = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$ , the posterior becomes (this form justifies  $\bar{\boldsymbol{\tau}}_0$ ):

$$p(\boldsymbol{\theta} | X, \chi, \nu) \propto g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_0 + N} \exp\left\{ \boldsymbol{\eta}^T(\boldsymbol{\theta}) \left( \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + \nu_0 \bar{\boldsymbol{\tau}}_0 \right) \right\} = g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_0 + N} \exp\{\boldsymbol{\eta}^T(\boldsymbol{\theta})(N \bar{\mathbf{u}} + \nu_0 \bar{\boldsymbol{\tau}}_0)\}$$

□ The parameter  $\nu_0$  can be interpreted as *effective number of fictitious observations* in the prior each of which has a value for the sufficient statistic equal to  $\bar{\boldsymbol{\tau}}_0$ .

$$p(\boldsymbol{\theta} | X, \nu_N, \boldsymbol{\tau}_N) = \frac{1}{Z(\nu_N, \boldsymbol{\tau}_N)} g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_N} \exp\left\{ (N + \nu_0) \boldsymbol{\eta}^T(\boldsymbol{\theta}) \frac{N \bar{\mathbf{u}} + \nu_0 \bar{\boldsymbol{\tau}}_0}{N + \nu_0} \right\} = \frac{1}{Z(\nu_N, \boldsymbol{\tau}_N)} g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_N} \exp\{\nu_N \boldsymbol{\eta}^T(\boldsymbol{\theta}) \bar{\boldsymbol{\tau}}_N\},$$

$$\text{where } \nu_N = \nu_0 + N, \bar{\boldsymbol{\tau}}_N = \frac{N \bar{\mathbf{u}} + \nu_0 \bar{\boldsymbol{\tau}}_0}{N + \nu_0}, \boldsymbol{\tau}_N = \nu_N \bar{\boldsymbol{\tau}}_N = N \bar{\mathbf{u}} + \nu_0 \bar{\boldsymbol{\tau}}_0 = \sum_{i=1}^N \mathbf{u}(\mathbf{x}_i) + \boldsymbol{\tau}_0$$



# Posterior Predictive

- Let  $u(X) = \sum_{i=1}^N u(x_i)$ ,  $u(X') = \sum_{i=1}^{N'} u(x'_i)$ , the posterior predictive is then:

$$\begin{aligned} p(X' | X) &= \int p(X' | \theta) p(\theta | X) d\theta \\ &= \prod_{i=1}^{N'} h(x'_i) \int g(\eta)^{N'} \exp\{\eta^T(\theta) u(X')\} \frac{1}{Z(v_0 + N, u(X) + \tau_0)} g(\eta(\theta))^{\nu_N} \exp\{\eta^T(\theta)(u(X) + \tau_0)\} d\theta \end{aligned}$$

- This is simplified as follows:

$$\begin{aligned} p(X' | X) &= \prod_{i=1}^{N'} h(x'_i) \frac{1}{Z(v_0 + N, u(X) + \tau_0)} \int g(\eta(\theta))^{N' + \nu_N} \exp\{\eta^T(\theta)(u(X') + u(X) + \tau_0)\} d\theta \\ &= \prod_{i=1}^{N'} h(x'_i) \frac{Z(v_0 + N + N', u(X') + u(X) + \tau_0)}{Z(v_0 + N, u(X) + \tau_0)} \end{aligned}$$

- If  $N = 0$ , this becomes the marginal likelihood of  $X'$ , which reduces to the normalizer of the posterior divided by the normalizer of the prior multiplied by a constant.

# Beta/Bernoulli: Posterior Predictive

- Consider a Bernoulli likelihood with a Beta prior. The likelihood takes the familiar exponential distribution form:

$$p(\mathcal{D} | \theta) = \theta^{\sum_i x_i} (1-\theta)^{N - \sum_i x_i} = (1-\theta)^N \exp\left(\log \frac{\theta}{1-\theta} \sum_i x_i\right)$$

- The conjugate prior is a Beta:  $p(\theta | \nu_0, \tau_0) \propto (1-\theta)^{\nu_0} \exp\left(\log\left(\frac{\theta}{1-\theta}\right)\tau_0\right) = \theta^{\tau_0} (1-\theta)^{\nu_0 - \tau_0}$   
 $p(\theta | \nu_0, \tau_0) = \text{Beta}(\alpha, \beta), \alpha = \tau_0 + 1, \beta = \nu_0 - \tau_0 + 1,$

- Thus the posterior becomes:  $p(\theta | \mathcal{D}) \propto \theta^{\tau_0 + s} (1-\theta)^{\nu_0 - \tau_0 + N - s} \Rightarrow$

$$p(\theta | \mathcal{D}) = \text{Beta}(\alpha_N, \beta_N), \alpha_N = \alpha + s, \beta_N = \beta + (N - s), s = \sum_i \mathbb{I}(x_i = 1)$$

- Let  $s'$  the number of heads in the past data. The probability of  $s' = \sum_{i=1}^m \mathbb{I}(x'_i = 1)$  future heads in  $m$  trials is then:

$$p(s' | \mathcal{D}, m) = \int \theta^{s'} (1-\theta)^{m-s'} \frac{\Gamma(\alpha_N + \beta_N)}{\Gamma(\alpha_N)\Gamma(\beta_N)} \theta^{\alpha_N-1} (1-\theta)^{\beta_N-1} d\theta = \frac{\Gamma(\alpha_N + \beta_N)}{\Gamma(\alpha_N)\Gamma(\beta_N)} \frac{\Gamma(\alpha_{N+m})\Gamma(\beta_{N+m})}{\Gamma(\alpha_{N+m} + \beta_{N+m})}$$
$$\alpha_{N+m} = \alpha_N + s', \beta_{N+m} = \beta_N + (m - s')$$

