Announcements

Plan for today is to continue with basic Bayesian statistics (nearly done).

Example (from Bishop, PRML)

Unknown variance or mean and variance

Similar stories can be told if the mean is known and the variance is not, or both are unknown.

We will simplify things by using the inverse of the covariance matrix which is called the precision matrix.

In the univariate case this is simply: $\lambda = \frac{1}{\sigma^2}$

Review from last lecture

Estimating the mean of a univariate Gaussian

(Case where variance is known)

$$\mu_{ML} = \frac{1}{N} \sum_{i} x_{i}$$

The appropriate conjugate prior for the mean is also Gaussian.

$$\mu_{MAP} = \frac{\sigma^2}{\sigma^2 + \sigma_0^2 N} \mu_0 + \frac{N \sigma_0^2}{\sigma^2 + N \sigma_0^2} \mu_{ML}$$

Example (from Bishop, PRML)

Known mean, unknown variance

$$p(\lbrace x_{i} \rbrace | \lambda) = \prod_{i=1}^{N} \mathbb{N}(x_{i} | \mu, \frac{1}{\lambda})$$

$$= \prod_{i=1}^{N} \left\{ \left(\frac{\lambda}{2\pi} \right)^{\frac{1}{2}} \exp\left(-\frac{\lambda}{2} (x_{i} - \mu)^{2} \right) \right\}$$
(u is constant)
$$\approx \lambda^{\frac{N}{2}} \exp\left\{ -\frac{\lambda}{2} \sum_{i} (x_{i} - \mu)^{2} \right\}$$

Inspection reveals that multiplying this by a gamma distribution

$$\operatorname{Gam}(\lambda \mid a,b) = \left\{ \frac{1}{\Gamma(a)} b^{a} \right\} \lambda^{a-1} \exp(-b\lambda)$$

yields a posterior of the same form.

Example (from Bishop, PRML)

Unknown mean and variance



Indicates optional

$$p(\lbrace x_i \rbrace | \lambda) = \prod_{i=1}^{N} \mathbb{N}(x_i | \mu, \frac{1}{\lambda})$$

(u is variable)

material

Example (from Bishop, PRML)

Quick review of exponentiation

$$\frac{e^a}{e^b} = e^a \cdot e^{-b} = e^{(a-b)}$$

$$\frac{e^a}{e^b} = e^a \cdot e^{-b} = e^{(a-b)}$$
$$e^{\left(\frac{a}{b}\right)} = e^{a \cdot \left(\frac{1}{b}\right)} = \left(e^a\right)^{\left(\frac{1}{b}\right)}$$

Example (from Bishop, PRML)

Quick review of exponentiation

$$e^{(a+b)} = e^a \cdot e^b$$

$$e^{a \cdot b} = \left(e^a\right)^b$$

$$e^{-a} = \frac{1}{e^a}$$

Example (from Bishop, PRML)



Unknown mean and variance

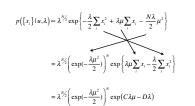
Indicates optional material

$$p(\lbrace x_{i}\rbrace | \lambda) = \prod_{i=1}^{N} \mathbb{N}(x_{i} \mid \mu_{i}) / \lambda$$

$$= \prod_{i=1}^{N} \left\{ \left(\frac{\lambda}{2\pi} \right)^{\frac{1}{2}} \exp\left(-\frac{\lambda}{2} (x_{i} - \mu)^{2} \right) \right\}$$

$$\approx \lambda^{N/2} \exp\left\{ -\frac{\lambda}{2} \sum_{i} (x_{i} - \mu)^{2} \right\}$$

$$= \lambda^{N/2} \exp\left\{ -\frac{\lambda}{2} \sum_{i} x_{i}^{2} + \lambda \mu \sum_{i} x_{i} - \frac{N\lambda}{2} \mu^{2} \right\}$$
(u is variable)



We now manipulate the formula to a more standard form.

$$\begin{split} p(u,\lambda) &\approx \lambda^{\frac{\beta}{2}} \left(\exp(-\frac{\lambda \mu^2}{2}) \right)^{\beta} \exp(c\lambda \mu - d\lambda) \\ &= \lambda^{\frac{\beta}{2}} \left(\exp(-\frac{\lambda \beta}{2} \mu^2) \right) \exp(c\lambda \mu - d\lambda) \\ &= \lambda^{\frac{\beta}{2}} \left(\exp(-\frac{\lambda \beta}{2} \mu^2 + c\lambda \mu - d\lambda) \right) \\ &= \lambda^{\frac{\beta}{2}} \left(\exp\left(-\frac{\lambda \beta}{2} \left(\mu^2 - \frac{2c\mu}{\beta} + \frac{2d}{\beta} \right) \right) \right) \end{split}$$

From the previous slide
$$p(\{x_i\} | u, \lambda) \approx \lambda^{\frac{N_2}{2}} \left(\exp(-\frac{\lambda \mu^2}{2}) \right)^N \exp(C\lambda \mu - D\lambda)$$

So a conjugate prior of the form

$$p(u,\lambda) \propto \lambda^{\beta/2} \left(\exp(-\frac{\lambda \mu^2}{2}) \right)^{\beta} \exp(c\lambda \mu - d\lambda)$$
will do (recall that $\exp(x) \cdot \exp(y) = \exp(x+y)$).

From the previous slide
$$p(u,\lambda) \approx \lambda^{\frac{\beta}{2}} \left(\exp\left(-\frac{\lambda\beta}{2} \left(\mu^2 - \frac{2c\mu}{\beta} + \frac{2d}{\beta} \right) \right) \right)$$

$$\left(\mu^2 - \left(\frac{2c}{\beta} \right) \mu + \frac{2d}{\beta} \right) = \left(\mu - \frac{c}{\beta} \right)^2 + \frac{2d}{\beta} - \frac{c^2}{\beta^2}$$

$$p(u,\lambda) \approx \lambda^{\frac{\beta}{2}} \exp\left(-\frac{\lambda\beta}{2} \left(\mu - \frac{c}{\beta} \right)^2 \right) \exp\left(-\lambda \left(d - \frac{c^2}{2\beta} \right) \right)$$

$$= \exp\left(-\frac{\lambda\beta}{2} \left(\mu - \frac{c}{\beta} \right)^2 \right) \lambda^{\frac{\beta}{2}} \left(\exp\left(-\lambda \left(d - \frac{c^2}{2\beta} \right) \right) \right)$$

From the previous slide

$$\begin{split} p(u,\lambda) &\approx \exp\left(-\frac{\lambda\beta}{2}\left(\mu - \frac{c}{\beta}\right)^2\right) \lambda^{\beta/2} \left(\exp\left(-\lambda\left(d - \frac{c^2}{2\beta}\right)\right)\right) \\ &\sim \mathbb{N}\left(\mu \mid \mu_0, (\lambda\beta)^{-1}\right) Gam(\lambda \mid a,b) \\ \text{where } \mu_0 &= \frac{c}{\beta} \text{ and } a = 1 + \frac{\beta}{2} \text{ (*) and } b = d - \frac{c^2}{2\beta} \\ \text{Recall that } Gam(\lambda \mid a,b) &\approx \lambda^{a-1} \exp(-b\lambda) \\ p(\mu,\lambda) &= p(\mu \mid \lambda) p(\lambda) &= \mathbb{N}\left(\mu \mid \mu_0, (\lambda\beta)^{-1}\right) Gam(\lambda \mid a,b) \\ \text{This is called the Gaussian-Gamma function} \end{split}$$

*According to Bishop, this is how the gamma parameter, a, relates to the Gaussian variance scale beta according to Bishop, but the powers of lambda from the normal do not seem to be accounted for — regardless, the conjugate formula is still correct.

More on priors

If we leave off the prior, then we are completely ignorant.

Note that the prior might be the uniform distribution over all numbers which is not a PDF! (Why?)

Such priors are called improper.

A more interesting example is p(k)=1/k, integral k>0.

Everything can work out fine if the posterior is a PDF.

Unknown mean and variance

To summarize, our conjugate prior is given by

$$p(u,\lambda) = p(u \mid \lambda)p(\lambda)$$
$$= N(u|u_o,(\beta\lambda)^{-1})Gam(\lambda \mid a,b)$$

Here a,b,β are constants. This is the normal-gamma (Gaussian-gamma) distribution.

Bayesian Sequential Update

For independent (conditioned on the model) sequential events

Suppose after N-1 observations we have the posterior $p(\theta|D_{1:N-1})$

This is a natural prior to estimate the posterior $p(\theta|D_{\text{EN}})$ from one more observation, D_{N} .

$$p(\theta \mid D_{1:N}) \propto p(\theta \mid D_{1:N-1}) p(D_N \mid \theta)$$
 (but is it true?)

Bayesian Sequential Update

More formally, for independent sequential events

$$\begin{split} p(\theta \mid D_{::N}) & \approx p(\theta) \, p(D_{::N} \mid \theta) \\ &= p(\theta) \prod_{i=1}^{N} \left(p(D_i \mid \theta) \right) \\ &= \left\{ p(\theta) \prod_{i=1}^{N-1} \left(p(D_i \mid \theta) \right) \right\} p(D_N \mid \theta) \\ &\approx p(\theta \mid D_{::N-1}) \, p(D_N \mid \theta) \end{split}$$
Prior from first

Predictive Distribution

Given training observations, X.

What is the density for another observation, x?

If we know the model parameters then we have:

$$p(x | \theta)$$

Predictive Distribution

Example --- you are tracking points following a curve in time

Question --- where do you think the next point will be?

Your observations up to time N-1 are "training data" (learn the curve).

What is the density for the location of the next point?

Predictive Distribution

If we know the model parameters then we have:

$$p(x \mid \theta)$$

However, we have been developing methods to compute posterior distributions:

$$p(\theta \mid X)$$

Predictive Distribution

 In the most general case, the predictive distribution marginalizes over uncertain model parameters

Model Selection

- There is no real difference between choosing the parameters of a model, which model class instance, and which model class.
- Difficulties
 - Prior densities are of different dimensionality
 - Constructing priors and likelihoods that properly penalize model complexity
 - Standard penalties include AIC, BIC
 - Cross-validation over novel data is more reliable.

Model Selection

- Model selection refers to choosing among different instances within a model class (1) or different model classes (2).
- Examples:
 - The number of clusters (1)
 - The degree of a polynomial to fit a curve to data (1)
 - Polynomials versus other basis functions such as Fourier (2)