# Bayesian Linear Regression Pt. 2, Gaussian Properties

Lecturer: Drew Bagnell  Scribe: Laura Lindzey (2009) Hans Pirnay(2008)

# 1 Parameterizations for Gaussians

There are two common parameterizations for Gaussians, the moment parameterization and the natural parameterization. It is often most practical to switch back and forth between representations, depending on which calculations are needed.

## 1.1 Moment Parameterization

$$\mathcal{N}(\mu, \Sigma) = p(\theta) \quad = \quad \frac{1}{z} \exp\left(-\frac{1}{2}\left(\theta - \mu\right)\Sigma^{-1}\left(\theta - \mu\right)\right) \tag{1}$$

Given:

$$\mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right) \tag{2}$$

**Marginal:** computing $p(x_2)$

$$\mu_2^{\text{marg}} \quad = \quad \mu_2 \tag{3}$$
$$\Sigma_2^{\text{marg}} \quad = \quad \Sigma_{22} \tag{4}$$

**Conditional:** computing $p(x_1|x_2)$

$$\mu_{1|2} \quad = \quad \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}\left(x_2 - \mu_2\right) \tag{5}$$
$$\Sigma_{1|2} \quad = \quad \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \tag{6}$$

## 1.2 Natural Parameterization

$$\tilde{\mathcal{N}}(J, P) = \tilde{p}(\theta) \quad = \quad \frac{1}{z} \exp\left(J^T\theta - \frac{1}{2}\theta^T P\theta\right) \tag{7}$$

Given:

$$\mathcal{N}\left(\begin{bmatrix} J_1 \\ J_2 \end{bmatrix}, \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix}\right) \tag{8}$$

**Marginal:** computing $p(x_2)$

$$J_2^{\text{marg}} \quad = \quad J_2 - P_{21}P_{11}^{-1}J_1 \tag{9}$$
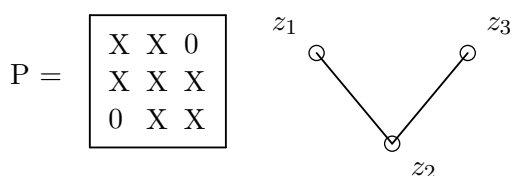$$P_1^{\text{marg}} \quad = \quad P_{12} - P_{21}P_{11}^{-1}P_{12} \tag{10}$$

**Conditional:** computing $p(x_1|x_2)$

$$J_{1|2} = J_1 - P_{12}x_2 \tag{11}$$
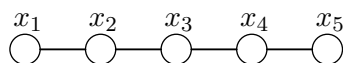$$P_{1|2} = P_{11} \tag{12}$$

# 2 Graphical Models

The matrix $P$ of the natural parameterization has a graphical model interpretation. If and only if there is a non-zero entry for $(z_i, z_j)$, then there is a lik between $z_i$ and $z_j$ in the MRF that corresponds to $\tilde{N}(J, P)$.

$$P = \begin{array}{|ccc|} \hline X & X & 0 \\ X & X & X \\ 0 & X & X \\ \hline \end{array}$$

$z_1$  $z_3$

$z_2$

Following the graphical model interpretation, $P$ is in many cases highly structured. Consider for example the graphical model of a markov chain:

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

The corresponding matrix $P$ will be non-zero only on the diagonal and off-diagonal:

$$P = \begin{pmatrix} X & X & 0 & 0 & 0 \\ X & X & X & 0 & 0 \\ 0 & X & X & X & 0 \\ 0 & 0 & X & X & X \\ 0 & 0 & 0 & X & X \end{pmatrix} \tag{13}$$

**Note:** $P$ describes which variables **directly** affect each other.

**Note:** $P^{-1}$ is, in general, not sparse! (this makes intuitve sense since $P^{-1} = \Sigma$ the covariance matrix, and the covariance of two states along the markov chain are not independent.)

# 3 Update Rules for Bayes Linear Regression

We want to derive how to update our representation given a series of new data $D = \{x_1, y_1, ..., x_n, y_n\}$. Let the prior be $\theta$, where $\theta \sim N(\mu_0, \Sigma_0)$. Then, we want to calculate

$$p(y_t|D) = \int p(y_t|\theta, x_t)p(\theta|D_t) = E_{p(\theta|D)}[p(y_t|\theta, x] \tag{14}$$

2

So, we need to calculate $p(\theta|D)$:

$$p(y|x,\theta) \quad = \quad \mathcal{N}(\theta^T x_t, \sigma_t^2) = \frac{1}{z}\exp\left(\frac{-(\theta^T x - y)(\theta^T x - y)}{2\sigma^2}\right) \tag{15}$$

(Don't worry about the weird notation of $\mathcal{N}$ as a function of $\sigma^2$. This is an arbitrary definition)

## 3.1 Deriving the update rules

Apply Bayes' Rule to the probability of a weight vector $\theta$ given a datapoint $D$.

$$p(\theta|D) \quad = \quad \frac{p(D|\theta)p(\theta)}{z} \tag{16}$$

This results in the multiplication of two exponential functions. Adding the exponent of the prior to that of the likelihood yields

$$-\frac{1}{2\sigma^2}\left(\theta^T x - y\right)^2 + J^T \theta - \frac{1}{2}\theta^T P\theta \tag{17}$$

collecting terms to find updates $J'_\theta$ and $P'_\theta$:

$$= \quad -\frac{1}{2\sigma^2}\left(\theta^T xx^T \theta - 2\theta^T xy + y^2\right) + J^T\theta - \frac{1}{2}\theta^T P\theta \tag{18}$$

$$= \quad \left(\frac{x^T y}{\sigma^2} + J^T\right)\theta - \frac{1}{2}\theta^T\left(\frac{xx^T}{\sigma^2} + P\right)\theta - \frac{y^2}{2\sigma^2} \tag{19}$$

Since this all happens in the exponent of an exponential function, the constat $y^2$-term can be shifted into the regularizing $z$. Thus, the update rules for $J'_\theta$ and $P'_\theta$ are

$$J'_\theta \quad = \quad \frac{xy}{\sigma^2} + J \tag{20}$$

$$P'_\theta \quad = \quad \frac{xx^T}{\sigma^2} + P \tag{21}$$

1. in a gaussian model, a new datapoint <u>always</u> lowers the variance - this downgrading of the variance does not always make sense

2. if you believe there are outliers, this model won't work for you

3. the variance is not a function of $y$. The precision if only affected by <u>input not output</u>. This is a consequence of having the same $\sigma$ (observation error) everywhere in space.

## 3.2 Transfer to moment parameterization

The update rules for the natural parameterization at timestep $t$ are

$$J \quad += \quad \frac{y_t x_t}{\sigma^2} \tag{22}$$

$$P \quad += \quad \frac{1}{\sigma^2}xx^T. \tag{23}$$

Given the transfer rules to the moment parameterization

$$\Sigma = P^{-1} \tag{24}$$
$$\mu = P^{-1}J \tag{25}$$

the moment parameterization after $N$ timesteps is then

$$\Sigma_{\theta|D} = \left[ \Sigma_0^{-1} + \sum_{i=1}^{N} \frac{x_i x_i^T}{\sigma^2} \right]^{-1} \tag{26}$$

$$\mu_\theta = \left[ \Sigma_0^{-1} + \sum_{i=1}^{N} \frac{x_i x_i^T}{\sigma^2} \right] \sum_{t=1}^{N} \frac{y_t x_t}{\sigma^2} \tag{27}$$

1. If $\mu_0 \neq 0$, the update for $\mu_\theta$ will be more complicated

2. $\sum_{t=1}^{N} \frac{y_t x_t}{\sigma^2}$ is the gradient of Bayes online linear regression

3. this looks just like Newton's method

4. Computation time: $o(d^2)$ for update, $o(d^3)$ for mean (that can be reduced to $o(d^2)$ with tricks)

## 3.3   Making predictions

Given all data $D$ up to timestep $t$ and $x_{t+1}$, the probability of an observation $\tilde{y}_{t+1}$ is

$$p(\tilde{y}_{t+1}|x_{t+1}, D) = \int p(\tilde{y}_{t+1}|x_{t+1}, D, \theta) \cdot p(\theta|D)d\theta \tag{28}$$

$$= \int p(\tilde{y}_{t+1}|x_{t+1}, \theta) \cdot p(\theta, D)d\theta \tag{29}$$

To know $p(\tilde{y}_{t+1}|x_{t+1}, D)$, we only need $y$ and $\sigma^2$, because these parameters determine the gaussian.