

PROPERTIES OF MULTIVARIATE GAUSSIANS

2/23/11
JASON L. PACHECO
CS195 F

CONDITIONAL GAUSSIAN DISTRIBUTIONS:

If two RV's are jointly Gaussian then the distribution of one conditional on the other is also Gaussian.

Consider X_a and X_b jointly Gaussian

$$\begin{pmatrix} X \\ \begin{matrix} X_a \\ X_b \end{matrix} \end{pmatrix} \sim N\left(\begin{pmatrix} \mu \\ \begin{matrix} \mu_a \\ \mu_b \end{matrix} \end{pmatrix}, \begin{pmatrix} \Sigma \\ \begin{matrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ab}^T & \Sigma_{bb} \end{matrix} \end{pmatrix}\right)$$

Sometimes it's useful to parameterize by $\Lambda = \Sigma^{-1}$ the precision matrix.

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ab}^T & \Lambda_{bb} \end{pmatrix}$$

NOTE: Λ is also symmetric, but $\Lambda_{aa} \neq \Sigma_{aa}^{-1}$.

Now we solve for the conditional,

$$p(X_a | X_b) = N(X_a | \mu_{a|b}, \Sigma_{a|b})$$

COMPLETING THE SQUARE

To do this we use a common technique. The quadratic form of a Gaussian expands as,

$$\begin{aligned} -\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu) &= \dots \\ &= -\frac{1}{2} x^T \Sigma^{-1} x + x^T \underbrace{\Sigma^{-1} \mu}_{\text{precision } \times \text{mean}} - \underbrace{\frac{1}{2} \mu^T \Sigma^{-1} \mu}_{\text{const.}} \end{aligned}$$

So if we can alter into this form we can "see" the precision directly, then solve for mean μ .

In our example we expand the quadratic form as,

$$\begin{aligned} -\frac{1}{2}(x-\mu)^T \Delta(x-\mu) &= \dots \\ &= -\frac{1}{2}(x_a - \mu_a)^T \Delta_{aa} (x_a - \mu_a) - \frac{1}{2}(x_a - \mu_a)^T \Delta_{ab} (x_b - \mu_b) \\ &\quad - \frac{1}{2}(x_b - \mu_b)^T \Delta_{ba} (x_a - \mu_a) - \frac{1}{2}(x_b - \mu_b)^T \Delta_{bb} (x_b - \mu_b) \\ &= -\frac{1}{2} x_a^T \Delta_{aa} x_a + x_a^T \Delta_{aa} \mu_a - \frac{1}{2} x_a^T \Delta_{ab} x_b + \frac{1}{2} x_a^T \Delta_{ab} \mu_b \\ &\quad - \frac{1}{2} x_b^T \Delta_{ba} x_a + \frac{1}{2} \mu_b^T \Delta_{ba} x_a + \text{const} \end{aligned}$$

Since $\Delta_{ab} = \Delta_{ba}^T$ we have,

$$\begin{aligned} &= -\frac{1}{2} x_a^T \Delta_{aa} x_a + x_a^T \Delta_{aa} \mu_a - \frac{1}{2} x_a^T \Delta_{ab} x_b + x_a^T \Delta_{ab} \mu_b \\ &\quad + \text{const} \end{aligned}$$

Collecting terms we have,

$$-\frac{1}{2}(x - \mu)^T \Delta(x - \mu) = -\frac{1}{2} x_a^T \Delta_{aa} x_a + x_a^T (\Delta_{aa} \mu_a - \Delta_{ab} (x_b - \mu_b)) + \text{const}$$

So,

$$\Sigma_{a|b}^{-1} = \Delta_{aa}^{-1}$$

$$\Sigma_{a|b}^{-1} \mu_{a|b} = (\Delta_{aa} \mu_a - \Delta_{ab} (x_b - \mu_b))$$

And multiplying both sides by $\Sigma_{a|b}$ the mean is,

$$\begin{aligned} \mu_{a|b} &= \Sigma_{a|b} \left\{ \Delta_{ab} \mu_a - \Delta_{ab} (x_b - \mu_b) \right\} \\ &= \mu_a - \Delta_{aa}^{-1} \Delta_{ab} (x_b - \mu_b) \quad \leftarrow \text{Linear in } x_b \end{aligned}$$

So, if $p(x_a, x_b)$ is Gaussian $N(\mu, \Sigma)$ then the conditional $p(x_a | x_b)$ is Gaussian $N(\mu_{a|b}, \Sigma_{a|b})$

MARGINAL GAUSSIANS

What about the marginal?

$$p(x_a) = \int p(x_a, x_b) dx_b$$

This is also Gaussian!

$$p(x_a) = N(x_a; \mu_a, \Sigma_{aa})$$

(3)

BAYE'S THEOREM For GAUSSIAN RV's:

So far we know if x_a and x_b are jointly Gaussian, then the marginal $p(x_a)$ is Gaussian and conditional $p(x_a|x_b)$ is Gaussian.

Furthermore, the conditional mean is a linear function of x_b . Generalizing assume we have

$$p(x) = N(x|\mu, \Lambda^{-1}) \quad (\text{think prior})$$

$$p(y|x) = N(y|Ax + b, L^{-1}) \quad (\text{think likelihood})$$

Multiplying and completing the square we look at the quadratic form,

$$\begin{aligned} & -\frac{1}{2}(x-\mu)^T \Lambda (x-\mu) - \frac{1}{2}(y-Ax-b)^T L (y-Ax-b) \\ &= -\frac{1}{2} \underline{x^T \Delta x} + x^T \Delta \mu - \frac{1}{2} y^T L y + \frac{1}{2} y^T \Delta A x \\ & \quad + \frac{1}{2} y^T L b + \frac{1}{2} x^T A^T L y - \frac{1}{2} x^T A^T L A x - \frac{1}{2} x^T A^T L b \\ & \quad + \frac{1}{2} b^T L y - \frac{1}{2} b^T L A x + \text{const} \\ &= -\frac{1}{2} x^T (\Delta + A^T L A) x - \frac{1}{2} y^T L y + \frac{1}{2} y^T L A x \\ & \quad + \frac{1}{2} x^T A^T L y + x^T \Delta \mu - x^T A^T L b + y^T L b \end{aligned}$$

Collecting terms we have,

$$-\frac{1}{2} \begin{pmatrix} x \\ y \end{pmatrix}^T \underbrace{\begin{pmatrix} I + A^T L A & -A^T L \\ -L A & L \end{pmatrix}}_{\text{the inverse precision matrix}} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} x \\ y \end{pmatrix}^T \begin{pmatrix} I \mu - A^T L b \\ L b \end{pmatrix}$$

Which is the form of the full joint. We "see" the inverse precision matrix

Using our previous results we can compute,

$$p(y) = N(y | A\mu + b, L^{-1} + A\Lambda^{-1} A^T)$$

$$p(x|y) = N(x | \Sigma \{ A^T L(y - b) + \Lambda \mu \}, \Sigma)$$

where,

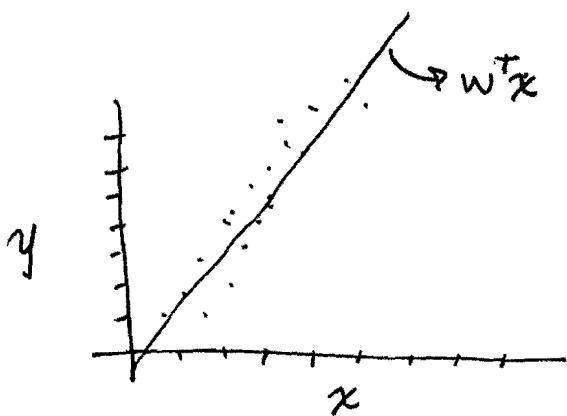
$$\Sigma = (I + A^T L A)^{-1}$$

These are general formulas which can be used!

LINEAR REGRESSION

Assume we have data

$$D \in \{(y_i, x_i), \dots, (y_N, x_N)\}$$



where,

$$y_i \in \mathbb{R}, x_i \in \mathbb{R}^M$$

We want to fit a line w/ coefficients $w \in \mathbb{R}^M$

$$y_i = w^T x_i + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

This is equivalent to the likelihood,

$$p(y_i | x_i, w, \sigma^2) = N(y_i | w^T x_i, \sigma^2)$$

MAXIMUM LIKELIHOOD \hat{w}_{ML}

The negative log-likelihood is,

$$\ell(w) = - \sum_{i=1}^N \log p(y_i | x_i, w, \sigma^2)$$

$$= \sum_{i=1}^N \left\{ \frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right\}$$

$$= \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^T x_i)^2 + \text{const}$$

$$\hat{w}_{ML} = \arg \min_w \ell(w) = \arg \min_w \frac{1}{N} \sum_{i=1}^N (y_i - w^T x_i)^2 \quad (*)$$

↳ Doesn't affect optimum

~~So it is minimum like that~~

LEAST SQUARES:

So ML is equivalent to MMSE for linear regression. Furthermore, objective is convex.

Let,

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{pmatrix} \in \mathbb{R}^{N \times M}$$

We rewrite (*) as,

$$\begin{aligned}
 \hat{\mathbf{w}}_{\text{MMSE}} &= \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{N} (\vec{\mathbf{y}} - \mathbf{X}\vec{\mathbf{w}})^T (\vec{\mathbf{y}} - \mathbf{X}\vec{\mathbf{w}}) \\
 &= \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{N} \mathbf{y}^T \mathbf{y} - \frac{2}{N} \mathbf{w}^T (\mathbf{X}^T \mathbf{y}) + \frac{1}{N} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} \\
 &= \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2N} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \frac{1}{N} \mathbf{w}^T (\mathbf{X}^T \mathbf{y}) \quad (**)
 \end{aligned}$$

ASIDE

Derivative of scalar product: $\nabla_{\mathbf{x}} \mathbf{a}^T \mathbf{x} = \mathbf{a}$

Derivative of quadratic form: $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$

RETURN

The objective is convex which means unique optimum.

We take gradient of (*) w.r.t. \vec{w} and solve,

$$\nabla_{\vec{w}} \hat{W}_{MMSE} = 0 = \frac{1}{2N} (\vec{X}^T \vec{X} + \vec{X}^T \vec{X}) \vec{w} - \frac{1}{N} \vec{X}^T \vec{y}$$

$$0 = \frac{1}{N} (\vec{X}^T \vec{X} \vec{w} - \vec{X}^T \vec{y})$$

$$\vec{X}^T \vec{y} = \vec{X}^T \vec{X} \vec{w}$$

$$\hat{W}_{MMSE} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$$

(NORMAL EQUATIONS)