

# Gaussian and Linear Gaussian Observation Models

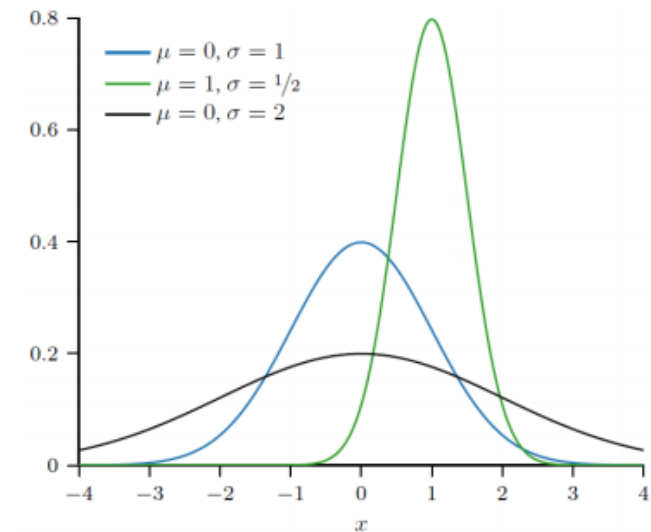
CS772A: Probabilistic Machine Learning

Piyush Rai

# Gaussian Distribution (Univariate)

- Distribution over real-valued scalar random variables  $Y \in \mathbb{R}$ , e.g., height of students in a class
- Defined by a scalar mean  $\mu$  and a scalar variance  $\sigma^2$

$$\mathcal{N}(Y = y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y - \mu)^2}{2\sigma^2} \right]$$



- Mean:  $\mathbb{E}[Y] = \mu$
- Variance:  $\text{var}[Y] = \sigma^2$
- Inverse of variance is called **precision**:  $\beta = \frac{1}{\sigma^2}$ .

$$\mathcal{N}(Y = y | \mu, \beta) = \sqrt{\frac{\beta}{2\pi}} \exp \left[ -\frac{\beta}{2} (y - \mu)^2 \right]$$

Gaussian PDF in terms of precision

# Gaussian Distribution (Multivariate)

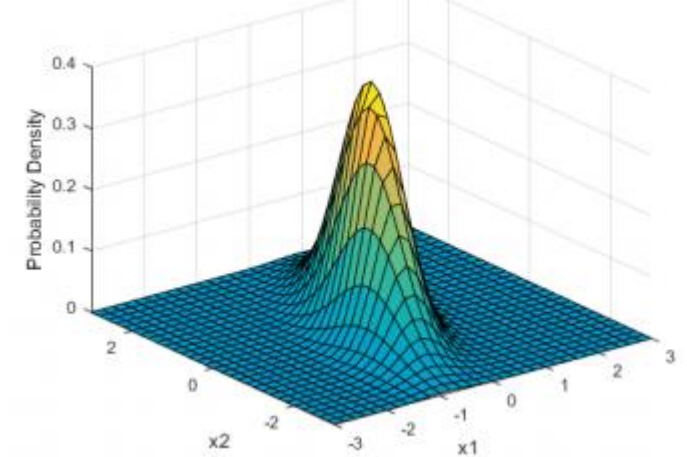
- Distribution over real-valued vector random variables  $\mathbf{Y} \in \mathbb{R}^D$
- Defined by a mean vector  $\boldsymbol{\mu} \in \mathbb{R}^D$  and a covariance matrix  $\boldsymbol{\Sigma}$

$$\mathcal{N}(\mathbf{Y} = \mathbf{y} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp[-(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})]$$

- Note: The cov. matrix  $\boldsymbol{\Sigma}$  must be symmetric and PSD
  - All eigenvalues are positive
  - $\mathbf{z}^\top \boldsymbol{\Sigma} \mathbf{z} \geq 0$  for any real vector  $\mathbf{z}$

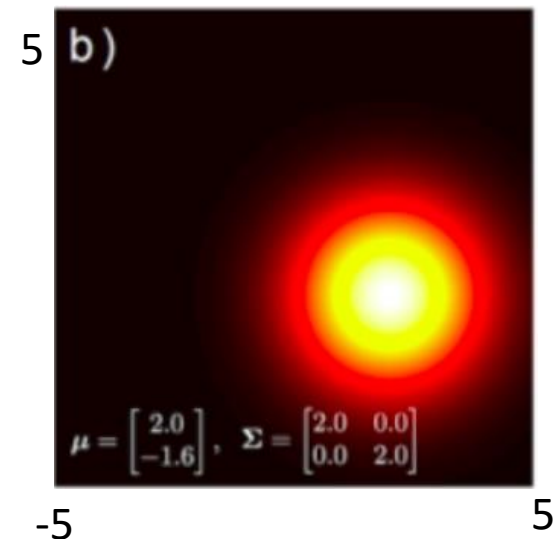
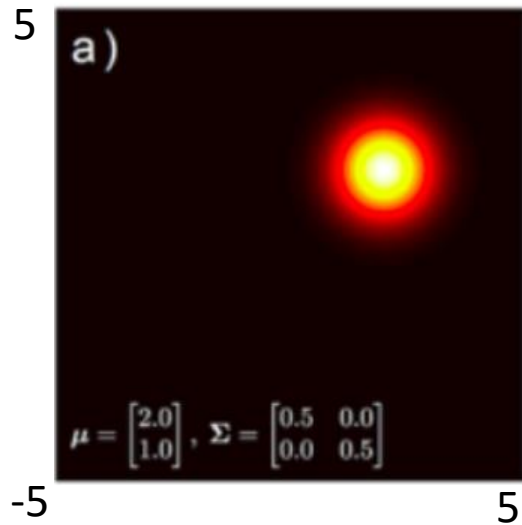
- The covariance matrix also controls the shape of the Gaussian
- Sometimes we work with precision matrix (inverse of covariance matrix)  $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$

A two-dimensional Gaussian

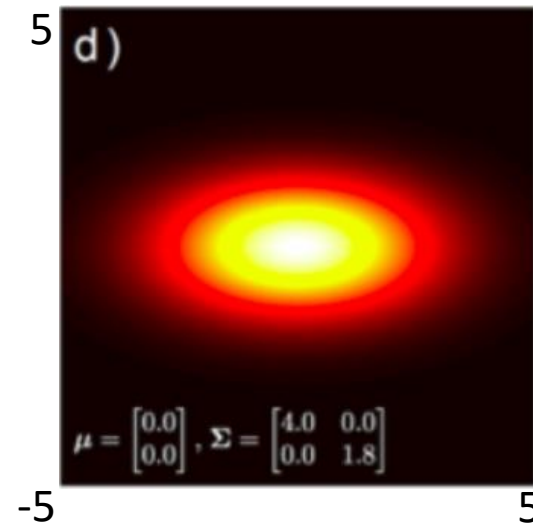
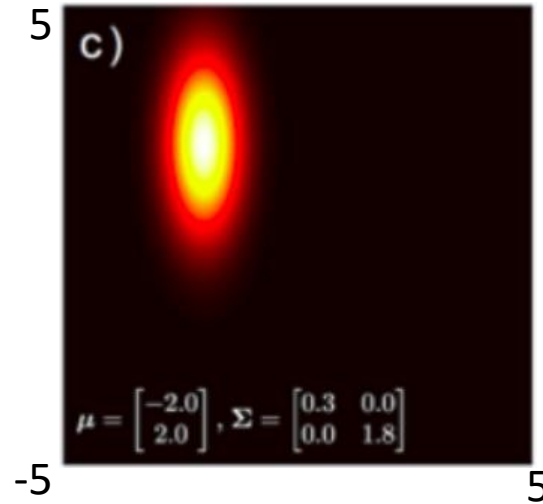


# Covariance Matrix for Multivariate Gaussian

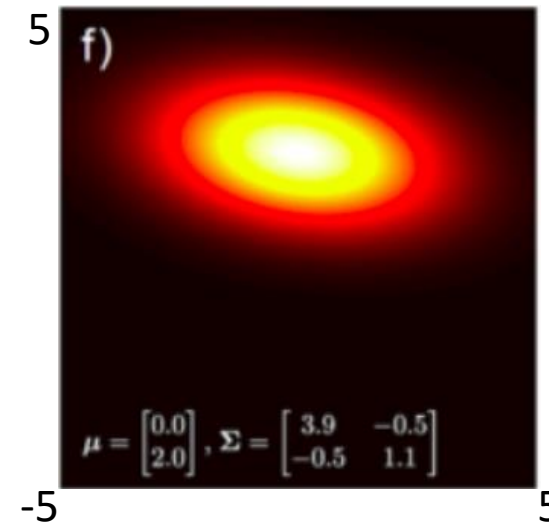
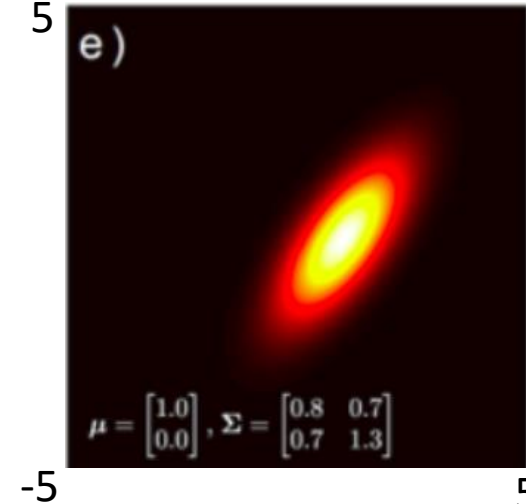
Spherical Covariance



Diagonal Covariance



Full Covariance



Spherical: Equal spreads (variances) along all dimensions

Diagonal: Unequal spreads (variances) along all directions but still axis-parallel

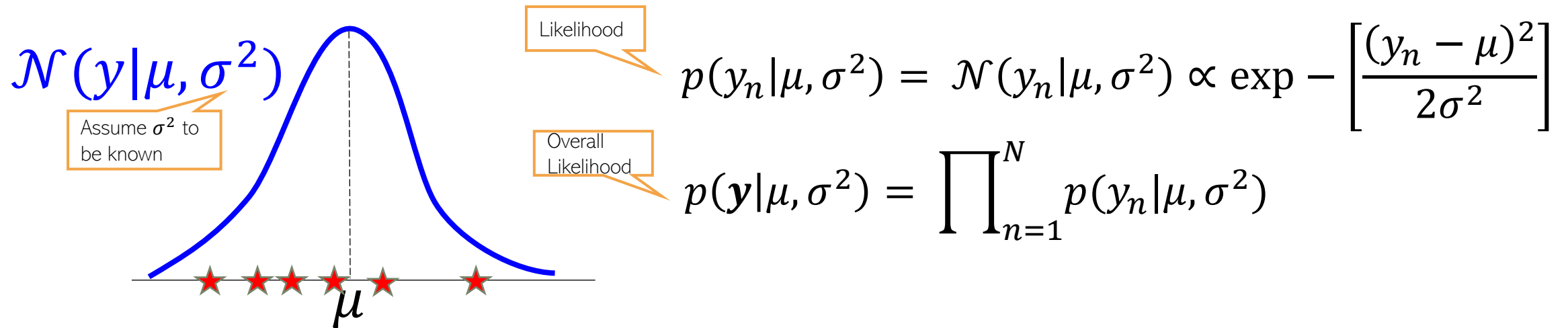
Full: Unequal spreads (variances) along all directions and also spreads along oblique directions



# Posterior Distribution for Gaussian's Mean

Its MLE/MAP estimation left as an exercise

- Given:  $N$  i.i.d. scalar observations  $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$  assumed drawn from  $\mathcal{N}(y|\mu, \sigma^2)$



- Note: Easy to see that each  $y_n$  drawn from  $\mathcal{N}(y|\mu, \sigma^2)$  is equivalent to the following

Thus  $y_n$  is like a noisy version of  $\mu$  with zero mean Gaussian noise added to it

$$y_n = \mu + \epsilon_n \quad \text{where } \epsilon_n \sim \mathcal{N}(0, \sigma^2)$$

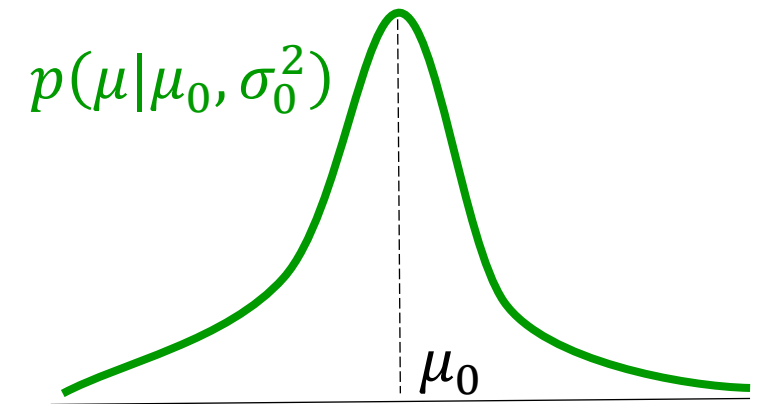
- Let's estimate mean  $\mu$  given  $\mathbf{y}$  using fully Bayesian inference (not point estimation)



# A prior distribution for the mean

- To compute posterior, need a prior over  $\mu$
- Let's choose a Gaussian prior

$$p(\mu|\mu_0, \sigma_0^2) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \\ \propto \exp \left[ -\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$



- The prior basically says that a priori we believe  $\mu$  is close to  $\mu_0$
- The prior's variance  $\sigma_0^2$  denotes how certain we are about our belief
- We will assume that the prior's hyperparameters  $(\mu_0, \sigma_0^2)$  are known
- Since  $\sigma^2$  in the likelihood  $\mathcal{N}(y|\mu, \sigma^2)$  is known, Gaussian prior  $\mathcal{N}(\mu|\mu_0, \sigma_0^2)$  on  $\mu$  is also conjugate to the likelihood (thus posterior of  $\mu$  will also be Gaussian)

# The posterior distribution for the mean

- The posterior distribution for the unknown mean parameter  $\mu$

On conditioning side, skipping all fixed params and hyperparams from the notation

$$p(\mu|\mathbf{y}) = \frac{p(\mathbf{y}|\mu)p(\mu)}{p(\mathbf{y})} \propto \prod_{n=1}^N \exp \left[ -\frac{(y_n - \mu)^2}{2\sigma^2} \right] \exp \left[ -\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$

- Easy to see that the above will be prop. to **exp of a quadratic function** of  $\mu$ . Simplifying:

$$p(\mu|\mathbf{y}) \propto \exp \left[ -\frac{(\mu - \mu_N)^2}{2\sigma_N^2} \right]$$

Gaussian posterior's precision is the sum of the prior's precision and sum of the noise precisions of all the observations

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

Contribution from the prior

Contribution from the data

Also the MLE solution for  $\mu$

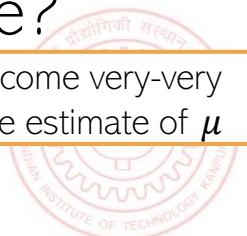
Gaussian posterior's mean is a convex combination of prior's mean and the MLE solution

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \bar{y} \quad (\text{where } \bar{y} = \frac{\sum_{n=1}^N y_n}{N})$$

- What happens to the posterior as  $N$  (number of observations) grows very large?

- Data (likelihood part) overwhelms the prior
- Posterior's variance  $\sigma_N^2$  will approximately be  $\sigma^2/N$  (and goes to 0 as  $N \rightarrow \infty$ )
- The posterior's mean  $\mu_N$  approaches  $\bar{y}$  (which is also the MLE solution)

Meaning, we become very-very certain about the estimate of  $\mu$



# The Predictive Distribution

- If given a point estimate  $\hat{\mu}$ , the plug-in predictive distribution for a test  $\mathbf{y}_*$  would be

This is an approximation of the true PPD  $p(\mathbf{y}_*|\mathbf{y})$

The best point estimate

$$p(\mathbf{y}_*|\hat{\mu}, \sigma^2) = \mathcal{N}(\mathbf{y}_*|\hat{\mu}, \sigma^2)$$

- On the other hand, the posterior predictive distribution of  $\mathbf{x}_*$  would be

$$\begin{aligned} p(\mathbf{y}_*|\mathbf{y}) &= \int p(\mathbf{y}_*|\mu, \sigma^2)p(\mu|\mathbf{y})d\mu \\ &= \int \mathcal{N}(\mathbf{y}_*|\mu, \sigma^2)\mathcal{N}(\mu|\mu_N, \sigma_N^2)d\mu \\ &= \mathcal{N}(\mathbf{y}_*|\mu_N, \sigma^2 + \sigma_N^2) \end{aligned}$$

This “extra” variance  $\sigma_N^2$  in PPD is due to the averaging over the posterior’s uncertainty

If **conditional** is Gaussian then **marginal** is also Gaussian

**A useful fact:** When we have conjugacy, the posterior predictive distribution also has a closed form (will see this result more formally when talking about exponential family distributions)



- For an alternative way to get the above result, note that, for test data

$$\mathbf{y}_* = \mu + \epsilon \quad \mu \sim \mathcal{N}(\mu_N, \sigma_N^2) \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Using the **posterior** of  $\mu$  since we are at test stage now

$$\Rightarrow p(\mathbf{y}_*|\mathbf{y}) = \mathcal{N}(\mathbf{y}_*|\mu_N, \sigma^2 + \sigma_N^2)$$

Since both  $\mu$  and  $\epsilon$  are Gaussian r.v., and are independent,  $\mathbf{y}_*$  also has a Gaussian posterior predictive, and the respective means and variances of  $\mu$  and  $\epsilon$  get added up

PRML [Bis 06], 2.115, and also mentioned in prob-stats refresher slides





# Gaussian Observation Model: Some Other Facts

- MLE/MAP for  $\mu, \sigma^2$  (or both) is straightforward in Gaussian observation models.
- Posterior also straightforward in most situations for such models
  - (As we saw) computing posterior of  $\mu$  is easy (using Gaussian prior) if variance  $\sigma^2$  is known
  - Likewise, computing posterior of  $\sigma^2$  is easy (using **gamma prior** on  $\sigma^2$ ) if mean  $\mu$  is known
- If  $\mu, \sigma^2$  both are unknown, posterior computation requires computing  $p(\mu, \sigma^2 | \mathbf{y})$ 
  - Computing joint posterior  $p(\mu, \sigma^2 | \mathbf{y})$  exactly requires a jointly conjugate prior  $p(\mu, \sigma^2)$
  - “**Gaussian-gamma**” (“Normal-gamma”) is such a conjugate prior – a product of normal and gamma
  - Note: Computing joint posteriors exactly is possible only in rare cases such this one
- If each observation  $\mathbf{y}_n \in \mathbb{R}^D$ , can assume a likelihood/observation model  $\mathcal{N}(\mathbf{y} | \mu, \Sigma)$ 
  - Need to estimate a **vector-valued** mean  $\mu \in \mathbb{R}^D$ . Can use a **multivariate Gaussian prior**
  - Need to estimate a  $D \times D$  positive definite covariance **matrix**  $\Sigma$ . Can use a **Wishart prior**
  - If  $\mu, \Sigma$  both are unknown, can use **Normal-Wishart** as a conjugate prior



# Linear Gaussian Model (LGM)

- LGM defines a noisy **lin. transform** of a Gaussian r.v.  $\boldsymbol{\theta}$  with  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Both  $\boldsymbol{\theta}$  and  $\mathbf{y}$  are vectors (can be of different sizes)

Also assume  $\mathbf{A}, \mathbf{b}, \boldsymbol{\Lambda}, \mathbf{L}$  to be known; only  $\boldsymbol{\theta}$  is unknown

$$\mathbf{y} = \mathbf{A}\boldsymbol{\theta} + \mathbf{b} + \boldsymbol{\epsilon}$$

Noise vector - independently and drawn from  $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{L}^{-1})$

- Easy to see that, conditioned on  $\boldsymbol{\theta}$ ,  $\mathbf{y}$  too has a Gaussian distribution

$$p(\mathbf{y}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\theta} + \mathbf{b}, \mathbf{L}^{-1})$$

- Assume  $p(\boldsymbol{\theta})$  as prior and  $p(\mathbf{y}|\boldsymbol{\theta})$  as the likelihood, and defining  $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$

Posterior of  $\boldsymbol{\theta}$

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})} = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\Sigma}(\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}), \boldsymbol{\Sigma})$$

Marginal likelihood

$$p(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top + \mathbf{L}^{-1})$$

- Many probabilistic ML models are LGMs
- These results are very widely used (PRML Chap. 2 contains a proof)

Closed form expressions for posterior and marginal likelihood (and both Gaussian)

