

---

# ***The Gaussian Distribution – MLE Estimators and Introduction to Bayesian Estimation***

*Prof. Nicholas Zabaras  
School of Engineering  
University of Warwick  
Coventry CV4 7AL  
United Kingdom*

*Email: [nzabaras@gmail.com](mailto:nzabaras@gmail.com)  
URL: <http://www.zabaras.com/>*

*August 7, 2014*

# Contents

---

- [The Gaussian Distribution](#), [Standard Normal](#), [Degenerate Gaussian Distribution](#), [Multivariate Gaussian](#), [the Gaussian and Maximum Entropy](#), [the CLT and the Gaussian Distribution](#), [Convolution of Gaussians](#), [MLE for the Gaussian](#), [MLE for the Multivariate Gaussian](#)
  - [Sequential MLE Estimation for the Gaussian](#), [Robbins-Monro Algorithm](#)
  - [Bayesian Inference for the Gaussian with Known Variance](#), [Bayesian Inference for the Gaussian with Known Mean](#), [Bayesian Inference for the Gaussian with unknown Mean and Variance](#)
  - [Normal-Gamma Distribution](#), [Gaussian-Wishart Distribution](#)
- 
- Following closely [Chris Bishops' PRML book](#), Chapter 2
  - Kevin Murphy's, [Machine Learning: A probabilistic perspective](#), Chapter 2

# *The Gaussian Distribution*

- A random variable  $X \in \mathbb{R}$  is Gaussian or normally distributed  $X \sim \mathcal{N}(\mu, \sigma^2)$  if:

$$P\{X \leq t\} = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^t \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx$$

- The following can be shown easily with direct integration:

$$\mathbb{E}[X] = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) x dx = \mu,$$

$$\mathbb{E}[X^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) x^2 dx = \mu^2 + \sigma^2, \quad \text{var}[X] = \mathbb{E}[(X-\mu)^2] = \sigma^2$$

- The following integrals are useful in these derivations :

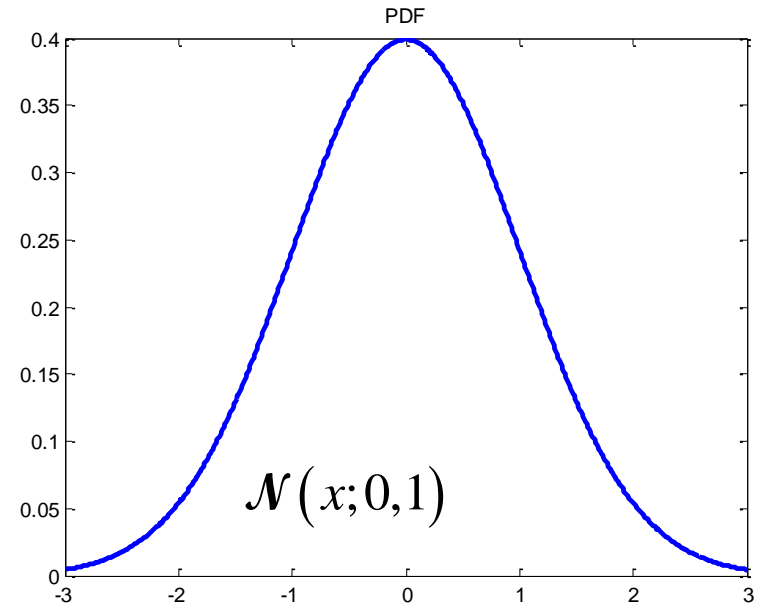
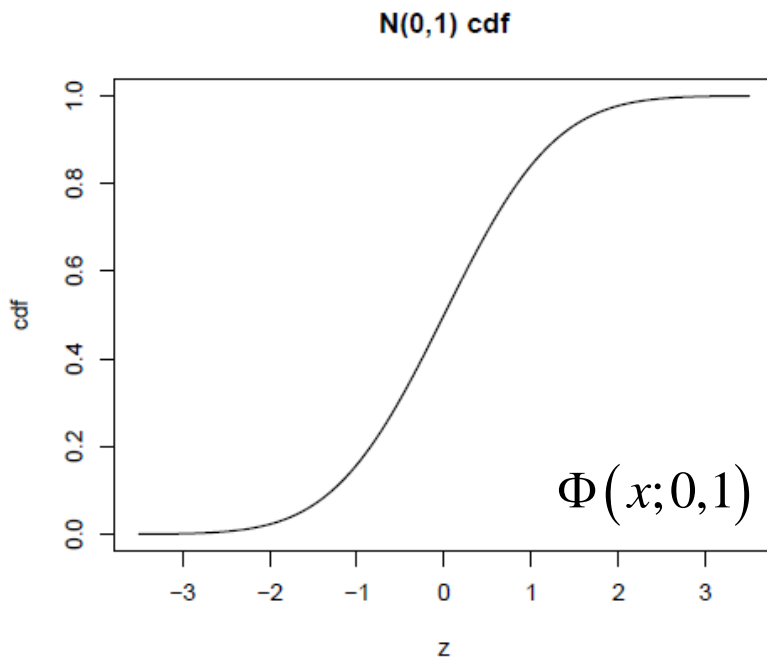
$$\int_{-\infty}^{+\infty} \exp(-u^2) du = \sqrt{\pi}, \quad \int_{-\infty}^{+\infty} u \exp(-u^2) du = 0, \quad \int_{-\infty}^{+\infty} u^2 \exp(-u^2) du = \frac{\sqrt{\pi}}{2}$$

- *We often work with the precision of a Gaussian  $\lambda=1/\sigma^2$ . The higher  $\lambda$  the narrower the distribution is.*

# Standard Normal, CDF, Error Function

- Plot of the Standard Normal  $\mathcal{N}(0,1)$  and CDF. Let  $\Phi(x;0,1)$  the corresponding CDF.

Run [gaussPlotDemo](#)  
from [PMTK](#)



$$\int_{-\infty}^x \mathcal{N}(z | \mu, \sigma^2) dz = \Phi(z; 0, 1), \quad z = (x - \mu) / \sigma$$

$$\Phi(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt = \frac{1}{2} \left[ 1 + \operatorname{erf}\left(z / \sqrt{2}\right) \right]$$

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

# *Degenerate Gaussian Distribution*

---

- Note that as  $\sigma^2 \rightarrow 0$ , the Gaussian becomes a delta function centered at the mean  $\mu$ :

$$\lim_{\sigma^2 \rightarrow 0} \mathcal{N}(x | \mu, \sigma^2) = \delta(x - \mu)$$

# Multivariate Gaussian

---

- A multivariate  $X \in \mathbb{R}^D$  is Gaussian if its probability density is

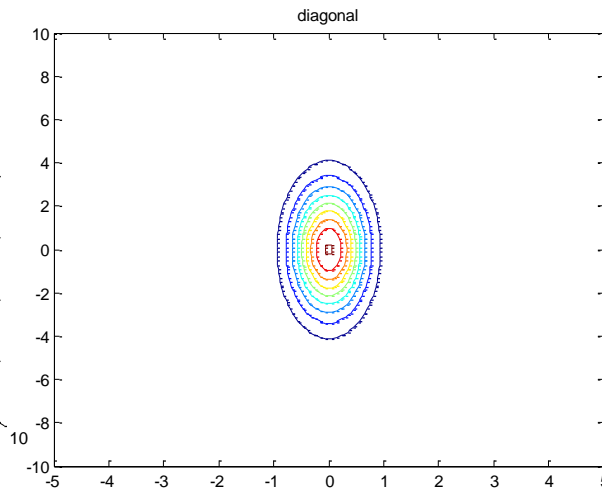
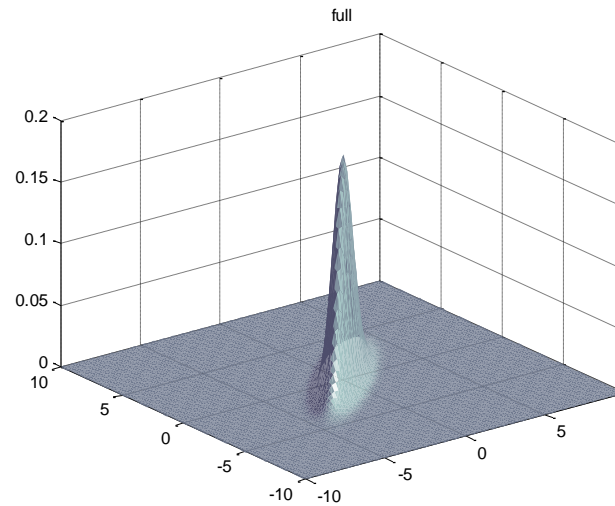
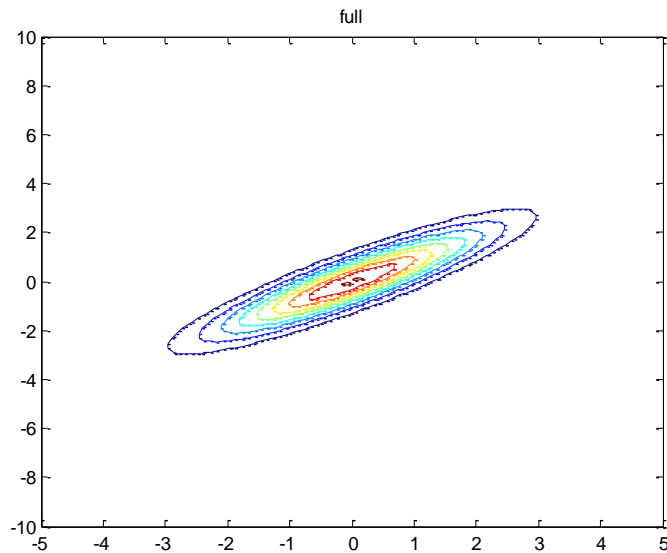
$$\mathcal{N}(x / \mu, \Sigma) = \left( \frac{1}{(2\pi)^D \det \Sigma} \right)^{1/2} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

where  $\mu \in \mathbb{R}^D$ ,  $\Sigma \in \mathbb{R}^{D \times D}$  is symmetric positive definite matrix (*covariance matrix*).

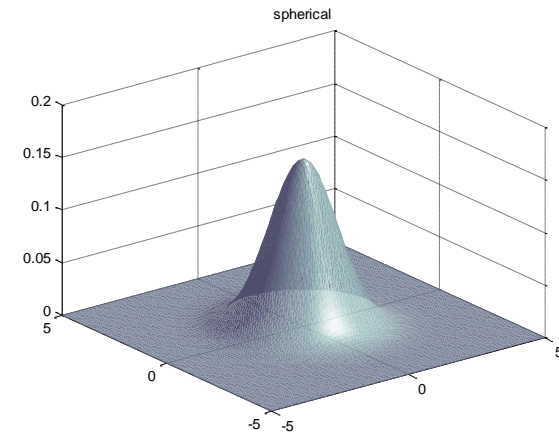
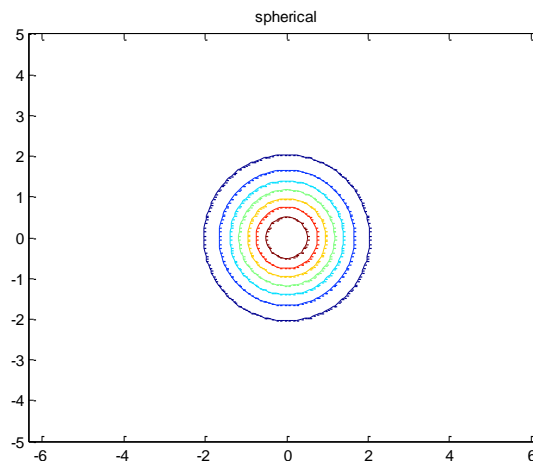
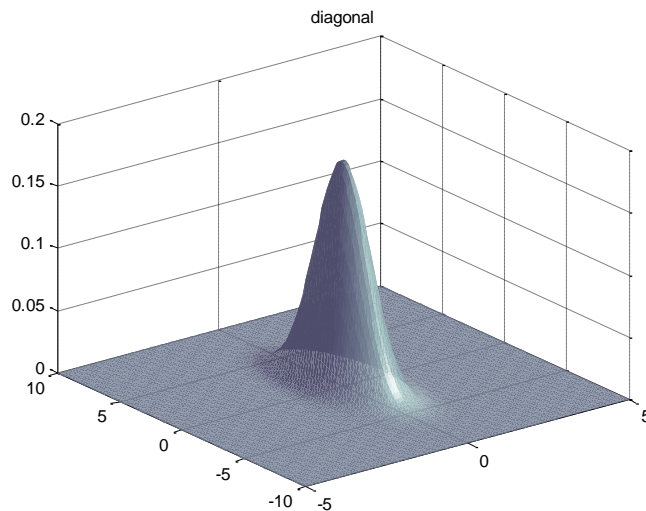
- We often work with *the precision matrix*  $\Lambda = \Sigma^{-1}$

# 2D Gaussian

- Level sets of 2D Gaussians (full, diagonal and spherical covariance matrix)



[gaussPlot2DDemo](#)  
from [PMTK](#)



# Multivariate Gaussian: Maximum Entropy

- We can *show that the multivariate Gaussian maximizes the entropy  $H$  with the constraints of normalization with given mean  $\mu$  and given variance  $\Sigma$ :*

$$\begin{aligned} \max_{p(\mathbf{x}), \lambda, \mathbf{m}, \mathbf{L}} = & -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \lambda \left( \int p(\mathbf{x}) d\mathbf{x} - 1 \right) + \mathbf{m}^T \left( \int \mathbf{x} p(\mathbf{x}) d\mathbf{x} - \mu \right) \\ & + \text{Tr} \left( \mathbf{L} \left( \int p(\mathbf{x}) (\mathbf{x} - \mu)(\mathbf{x} - \mu)^T d\mathbf{x} - \Sigma \right) \right) \end{aligned}$$

- Setting the derivative wrt  $p(\mathbf{x})$  to zero gives:

$$0 = -1 - \ln p(\mathbf{x}) + \lambda + \mathbf{m}^T \mathbf{x} + \text{Tr} \left( \mathbf{L} (\mathbf{x} - \mu)(\mathbf{x} - \mu)^T \right)$$

$$p(\mathbf{x}) = e^{-1 + \lambda + \mathbf{m}^T \mathbf{x} + (\mathbf{x} - \mu)^T \mathbf{L} (\mathbf{x} - \mu)}$$

- The coefficients can be found by satisfying the constraints. We start by completing the square.



# Multivariate Gaussian: Maximum Entropy

$$\begin{aligned} p(\mathbf{x}) &= e^{-1+\lambda+\mathbf{m}^T \mathbf{x}+(\mathbf{x}-\boldsymbol{\mu})^T \mathbf{L}(\mathbf{x}-\boldsymbol{\mu})} = \\ &= e^{-1+\lambda+\boldsymbol{\mu}^T \mathbf{m}-\frac{1}{4}\mathbf{m}^T \mathbf{L}^{-1} \mathbf{m}+(\mathbf{x}-\boldsymbol{\mu}+\frac{1}{2}\mathbf{L}^{-1} \mathbf{m})^T \overbrace{\mathbf{L}(\mathbf{x}-\boldsymbol{\mu}+\frac{1}{2}\mathbf{L}^{-1} \mathbf{m})}^{\mathbf{y}}} \end{aligned}$$

➤ Satisfying the mean constraint:

$$\int e^{-1+\lambda+\boldsymbol{\mu}^T \mathbf{m}-\frac{1}{4}\mathbf{m}^T \mathbf{L}^{-1} \mathbf{m}+\mathbf{y}^T \mathbf{L} \mathbf{y}} \left( \mathbf{y} + \boldsymbol{\mu} - \frac{1}{2} \mathbf{L}^{-1} \mathbf{m} \right) d\mathbf{y} = \boldsymbol{\mu}$$

➤ The 1<sup>st</sup> term drops from symmetry, the 2<sup>nd</sup> gives  $\boldsymbol{\mu}$  from normalization, thus we need to have:

$$-\frac{1}{2} \mathbf{L}^{-1} \mathbf{m} = \mathbf{0} \Rightarrow \mathbf{m} = \mathbf{0}$$

# Multivariate Gaussian: Maximum Entropy

$$p(\mathbf{x}) = e^{-1+\lambda+\mathbf{z}^T \mathbf{L}(\mathbf{x}-\boldsymbol{\mu})}$$

- Satisfying the variance constraint:

$$\int e^{-1+\lambda+\mathbf{z}^T \mathbf{L} \mathbf{z}} \mathbf{z} \mathbf{z}^T d\mathbf{z} = \boldsymbol{\Sigma}$$

- Note that with  $\mathbf{L} = -\boldsymbol{\Sigma} / 2$ , the 3<sup>rd</sup> term from the exponential when integrated gives:

$$\int e^{\mathbf{z}^T \mathbf{L} \mathbf{z}} \mathbf{z} \mathbf{z}^T d\mathbf{z} = \boldsymbol{\Sigma} (2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}$$

- It remains to select  $\lambda$  such that:

$$e^{-1+\lambda} = (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{-1/2} \Rightarrow \lambda - 1 = \ln \left\{ \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \right\}$$

- The optimizing  $p(\mathbf{x})$  is now clearly the Gaussian.

# Multivariate Gaussian: Maximum Entropy

- The entropy of the multivariate Gaussian is now computed as follows:

$$\begin{aligned} H[x] &= - \int \mathcal{N}(x / \mu, \Sigma) \ln \mathcal{N}(x / \mu, \Sigma) dx \\ &= \int \mathcal{N}(x / \mu, \Sigma) \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| + (x - \mu)^T \Sigma^{-1} (x - \mu) \right) dx \\ &= \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| \right) + \int \mathcal{N}(x / \mu, \Sigma) \frac{1}{2} \text{tr} \left( (x - \mu)(x - \mu)^T \Sigma^{-1} \right) dx \\ &= \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| \right) + \frac{1}{2} \text{tr} \left( \left( \int \mathcal{N}(x / \mu, \Sigma) (x - \mu)(x - \mu)^T dx \right) \Sigma^{-1} \right) \\ &= \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| \right) + \frac{1}{2} \text{tr} \left( \Sigma \Sigma^{-1} \right) \\ &= \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| + \text{tr} \left( \Sigma^{-1} \Sigma \right) \right) \\ &= \frac{1}{2} \left( D \ln(2\pi) + \ln |\Sigma| + D \right) \end{aligned}$$

# Multivariate Gaussian: Maximum Entropy

- Using also the KL distance definition, one can show that *the Gaussian has the largest entropy from any other distribution satisfying the mean and 2<sup>nd</sup> moment constraints*. To make the presentation simple, consider

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} / \boldsymbol{\theta}, \Sigma), \int q(\mathbf{x}) \mathbf{x} \mathbf{x}^T d\mathbf{x} = \Sigma$$

- Then:

$$\begin{aligned} 0 \leq KL(q \parallel p) &= -\int q(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} = -\int q(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \int q(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} \\ &= -\int q(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} - H[q] = -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} - H[q] \\ &= H[p] - H[q] \Rightarrow H[p] \geq H[q] \end{aligned}$$

- The intermediate step in the proof above accounts for the moment constraint on  $q$  and the fact that  $\log(p)$  is quadratic in  $\mathbf{x}$ !

# *The CLT and the Gaussian Distribution*

- Let  $(X_1, X_2, \dots, X_n)$  be independent and identically distributed (i.i.d.) continuous random variables each with expectation  $\mu$  and variance  $\sigma^2$ .

- Define: 
$$Z_n = \frac{1}{\sigma\sqrt{N}} (X_1 + X_2 + \dots + X_n - N\mu)$$

- As  $N \rightarrow \infty$ , the distribution of  $Z_n$  converges to the distribution of a standard normal random variable

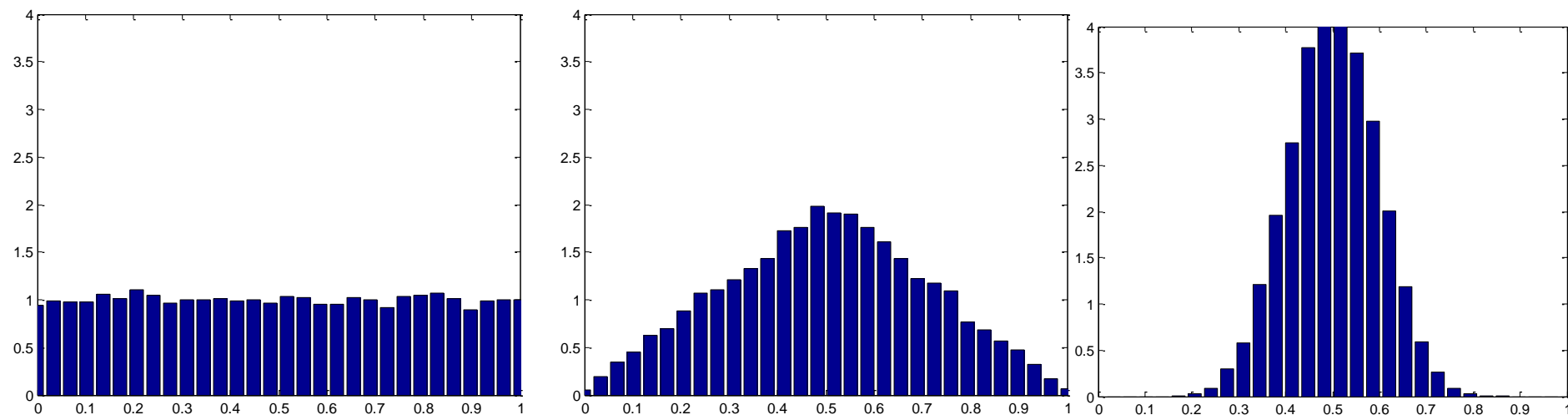
$$\lim_{N \rightarrow \infty} P\{Z_n \leq x\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt$$

- If  $\bar{X}_n = \frac{1}{N} \sum_{j=1}^N X_j$ , for  $N$  large,  $\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right)$  as  $N \rightarrow \infty$

- Somewhat of a justification for assuming that Gaussian noise is common

# The CLT and the Gaussian Distribution

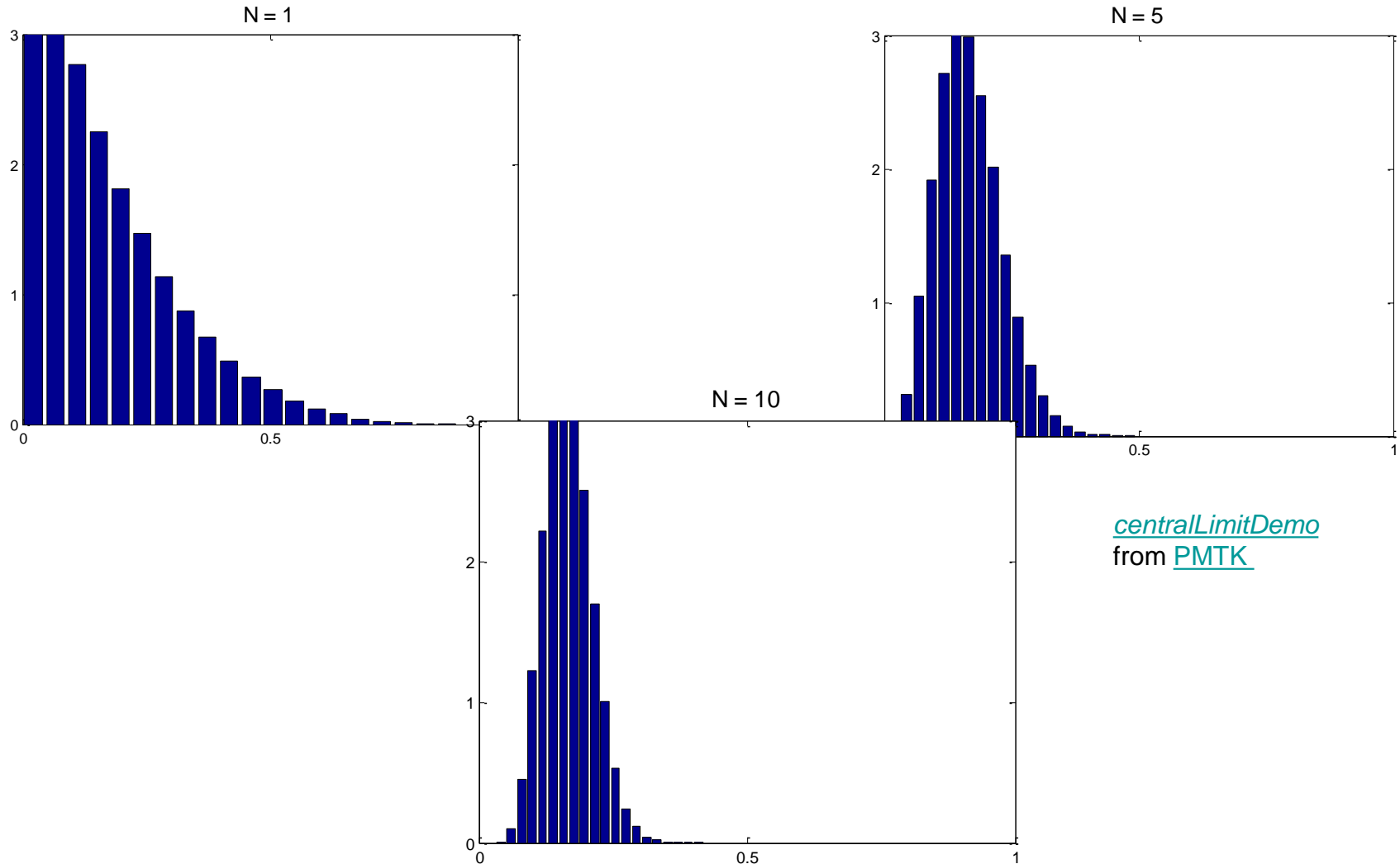
- As an example, assume  $N$  variables  $(X_1, X_2, \dots, X_n)$  each of which has a uniform distribution over  $[0, 1]$  and then consider the distribution of the mean  $(X_1 + X_2 + \dots + X_n)/N$ . For large  $N$ , this distribution tends to a Gaussian. The convergence as  $N$  increases can be rapid.



[MatLab Code](#)

# The CLT and the Gaussian Distribution

□ Histogram of  $\frac{1}{N} \sum_{j=1}^{10000} x_{ij}$  where  $x_{ij} \sim \text{Beta}(1,5)$



[centralLimitDemo](#)  
from [PMTK](#)

# ***The CLT and the Gaussian Distribution***

---

- One consequence of this result is that the binomial distribution which is a distribution over  $m$  *defined by the sum of  $N$  observations of the random binary variable  $x$* , will tend to a Gaussian as  $N \rightarrow \infty$ .



# Example of the Convolution of Gaussians

- Consider 2 Gaussians  $x_1 \sim \mathcal{N}(\mu_1, \tau_1^{-1})$ ,  $x_2 \sim \mathcal{N}(\mu_2, \tau_2^{-1})$ . We want to compute the entropy of the distribution of  $x = x_1 + x_2$ .

- $p(x)$  can be computed from the convolution of two Gaussians

$$p(x) = \int \underbrace{p(x | x_2)}_{\mathcal{N}(\mu_1 + x_2, \tau_1^{-1})} p(x_2) dx_2$$

- We need to *complete the square in the exponential in  $x_2$* :

$$\begin{aligned} & -\frac{1}{2} \tau_1 (x - (\mu_1 + x_2))^2 - \frac{1}{2} \tau_2 (x_2 - \mu_2)^2 = \\ & -\frac{1}{2} (\tau_1 + \tau_2) \left( x_2 - \frac{\tau_1 (x - \mu_1) + \tau_2 \mu_2}{\tau_1 + \tau_2} \right)^2 - \frac{1}{2} \tau_1 (x - \mu_1)^2 + \frac{1}{2} \frac{(\tau_1 (x - \mu_1) + \tau_2 \mu_2)^2}{\tau_1 + \tau_2} \end{aligned}$$

- The 1<sup>st</sup> term is integrated out and the precision of  $x$  is:

$$\tau_1 - \frac{\tau_1^2}{\tau_1 + \tau_2} = \frac{\tau_1 \tau_2}{\tau_1 + \tau_2}$$

- Thus the entropy of  $x$  is:  $H[x] = \frac{1}{2} \ln(2\pi e \sigma^2) = \frac{1}{2} \ln \left( 2\pi e \frac{\tau_1 + \tau_2}{\tau_1 \tau_2} \right)$

# 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.
- *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

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

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

# Max Likelihood for a Gaussian Distribution

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

- 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:

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

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

# Maximum Likelihood for a Gaussian Distribution

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

- ❑ 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 :*

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

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

# Maximum Likelihood for a Gaussian Distribution

$$\sigma_{ML}^2 = \frac{N-1}{N} \sigma^2$$

We use :

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

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

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

# Maximum Likelihood for a Gaussian Distribution

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

- ❑ 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:

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

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 MLE*.

# MLE for the Multivariate Gaussian

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

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

- 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

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

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

$$-\frac{N}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \ln |\boldsymbol{\Sigma}| = \frac{N}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \ln |\boldsymbol{\Sigma}^{-1}| = \frac{N}{2} \boldsymbol{\Sigma}^T = \boxed{\frac{N}{2} \boldsymbol{\Sigma}} \quad \text{A useful trick!}$$

$$\begin{aligned} -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) &= -\frac{1}{2} N \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \text{Tr} \left( \boldsymbol{\Sigma}^{-1} \sum_{n=1}^N \frac{1}{N} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T \right) \\ &= -\frac{1}{2} N \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \text{Tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}) = -\frac{1}{2} N \mathbf{S} = \quad \text{S symmetric} \end{aligned}$$

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

□ So finally  $\boldsymbol{\Sigma}_{ML} = \mathbf{S}$

□ 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*

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

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

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

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

$$\mathbb{E}[\mathbf{x}_n \mathbf{x}_m^T] = \boldsymbol{\mu} \boldsymbol{\mu}^T + \delta_{mn} \boldsymbol{\Sigma}$$

# Sequential MLE Estimation for Gaussians

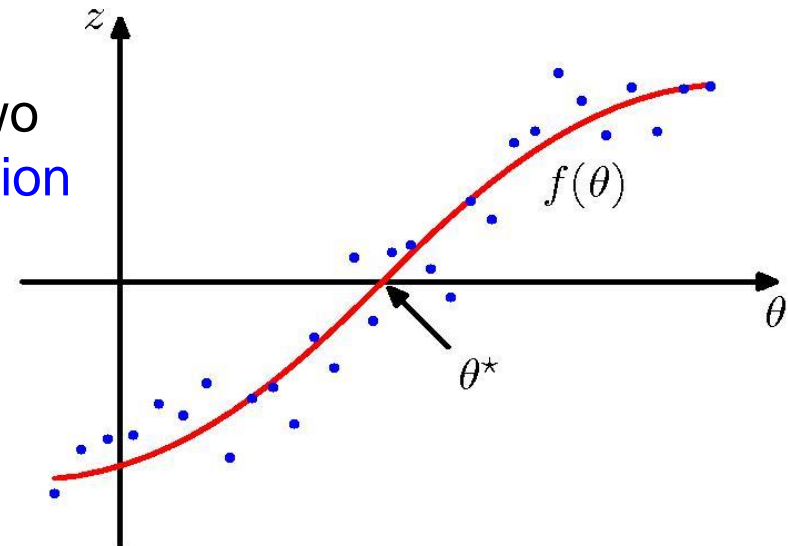
- Often we are interested to compute sequentially an estimate of  $\mu_{ML}$  as more data arrive. This can easily be done:

$$\begin{aligned}\mu_{ML}^{(N)} &= \frac{1}{N} \sum_{n=1}^N x_n = \frac{x_N}{N} + \frac{1}{N} \sum_{n=1}^{N-1} x_n = \\ &= \frac{x_N}{N} + \frac{N-1}{N} \frac{1}{N-1} \sum_{n=1}^{N-1} x_n = \\ &= \frac{x_N}{N} + \frac{N-1}{N} \mu_{ML}^{(N-1)} = \underbrace{\mu_{ML}^{(N-1)}}_{\text{Learning rate}} + \underbrace{\frac{1}{N} (x_N - \mu_{ML}^{(N-1)})}_{\text{Error signal}}\end{aligned}$$

- This sequential approach cannot easily be generalized to other cases (non-Gaussians, etc.)

# Robbins-Monro Algorithm

- A more powerful approach to computing sequentially the MLE estimates is via the [Robbins-Monro algorithm](#).
- We review the algorithm by considering the calculation of the zero of a regression function.\*
- Consider the joint distribution  $p(z, \theta)$  of two random variables and define the regression function as:
$$f(\theta) = \mathbb{E}(z | \theta) = \int z p(z | \theta) dz$$
- Assume we are given samples from  $p(z, \theta)$  one at a time.



\* Effectively, we don't know the regression function  $f(\theta)$  but we have data of a noisy version  $z$  of that. We take the regression function to be the expectation  $\mathbb{E}(z | \theta)$ .

- Robbins, H. and S. Monro (1951). [A stochastic approximation method](#). *Annals of Mathematical Statistics* **22**, 400–407.
- Fukunaga, K. (1990). [Introduction to Statistical Pattern Recognition](#) (Second ed.). Academic Press.

# Robbins-Monro Algorithm

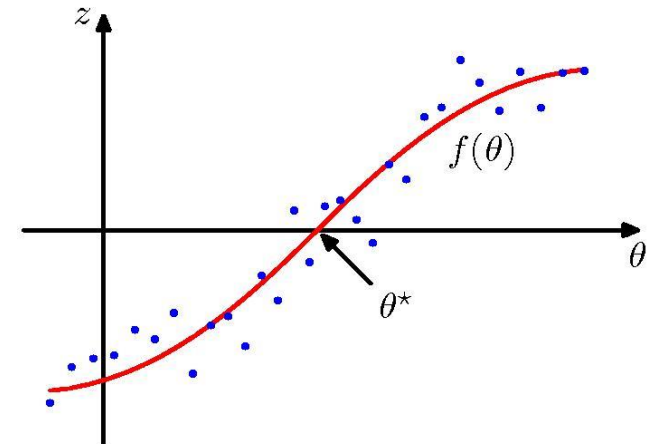
$$f(\theta) = \mathbb{E}(z | \theta) = \int z p(z | \theta) dz$$

□ We want to find the root  $f(\theta^*) = 0$  in a sequential manner: The Robbins-Monro algorithm proceeds as:

$$\theta^{(N)} = \theta^{(N-1)} - a_{N-1} z(\theta^{(N-1)})$$

□ The *learning coefficients*  $\{a_N\}$  should satisfy:

$$\lim_{N \rightarrow \infty} a_N = 0, \sum_{n=1}^{\infty} a_n = \infty, \sum_{n=1}^{\infty} a_n^2 < \infty$$



# Robbins-Monro Algorithm

- We can state the MLE calculation  $\mu_{ML}$  for our Gaussian example as finding the root of a regression function:

$$-\frac{\partial}{\partial \mu} \left\{ \frac{1}{N} \sum_{n=1}^N \ln p(x_n | \mu) \right\} \Big|_{ML} = 0 \Rightarrow -\sum_{n=1}^N \frac{1}{N} \frac{\partial}{\partial \mu} \ln p(x_n | \mu) \Big|_{ML} = 0 \xrightarrow[N \rightarrow \infty]{CLT} \mathbb{E} \left[ \underbrace{-\frac{\partial}{\partial \mu} \ln p(x | \mu)}_{z|\mu} \right]_{ML} = 0$$

- In the context of the Robbins-Monro algorithm,

$$z = \frac{\partial}{\partial \mu} [-\ln p(x | \mu)] \Big|_{\mu=\mu_{ML}} = -\frac{x - \mu_{ML}}{\sigma^2}, \text{ } z \text{ is a Gaussian, } f(\mu) = \mathbb{E}(z | \mu) = -\frac{\mu - \mu_{ML}}{\sigma^2}$$

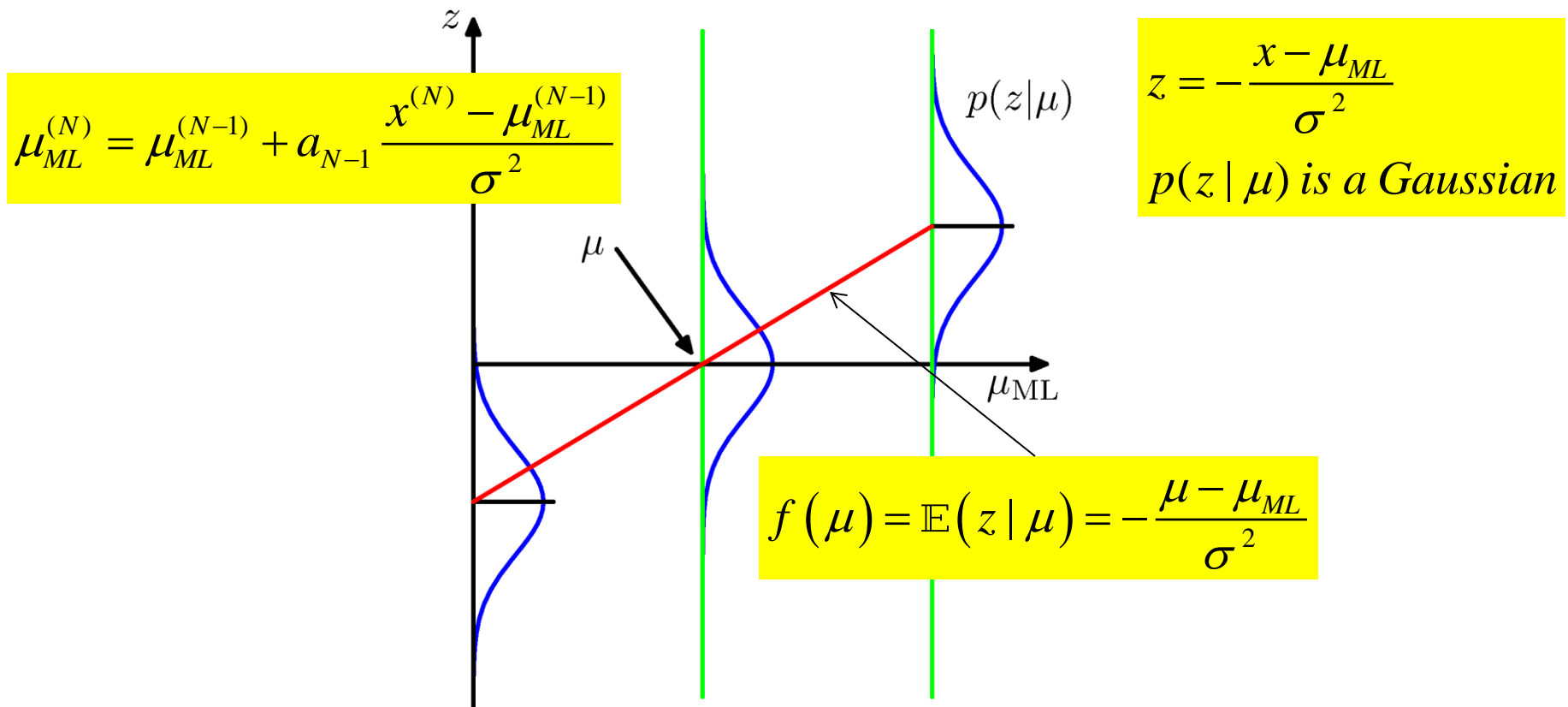
- The algorithm takes the form:

$$\mu_{ML}^{(N)} = \mu_{ML}^{(N-1)} + a_{N-1} \frac{x^{(N)} - \mu_{ML}^{(N-1)}}{\sigma^2}$$

- Substituting  $a_{N-1} = \frac{\sigma^2}{N}$  gives the estimate discussed earlier.

# Robbins-Monro Algorithm

- A graphical interpretation of the algorithm is shown here.



- The Robbin-Monro algorithm computes the zero of the regression function.

- Blum, J. A. (1965). [Multidimensional stochastic approximation methods](#). *Annals of Mathematical Statistics* **25**, 737–744.



# Sequential MLE Estimation for Gaussians

- Let us now repeat the same calculations but for the MLE estimate of  $\sigma^2$ :

$$\begin{aligned}\sigma_{(N)}^2 &= \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 = \frac{1}{N} \sum_{n=1}^{N-1} (x_n - \mu)^2 + \frac{(x_N - \mu)^2}{N} \\ &= \frac{N-1}{N} \sigma_{(N-1)}^2 + \frac{(x_N - \mu)^2}{N} = \\ &= \sigma_{(N-1)}^2 + \frac{1}{N} \left\{ (x_N - \mu)^2 - \sigma_{(N-1)}^2 \right\}\end{aligned}$$

- If we substitute the expression for the Gaussian likelihood into the Robbins-Monro procedure for maximizing likelihood:

$$\sigma_{(N)}^2 = \sigma_{(N-1)}^2 + a_{N-1} \frac{\partial}{\partial \sigma_{(N-1)}^2} \left\{ -\frac{1}{2} \ln \sigma_{(N-1)}^2 - \frac{(x_N - \mu)^2}{2\sigma_{(N-1)}^2} \right\} = \sigma_{(N-1)}^2 + a_{N-1} \frac{1}{2\sigma_{(N-1)}^4} \left\{ (x_N - \mu)^2 - \sigma_{(N-1)}^2 \right\}$$

- The 2 formulas are identical for:  $a_{N-1} = 2\sigma_{(N-1)}^4 / N$ .

# Sequential MLE: Multivariate Gaussian

- To simplify things, *assume that*  $\mu_{ML} = \mu$  and thus:

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

- From this equation we can derive:

$$\Sigma_{ML}^{(N)} = \Sigma_{ML}^{(N-1)} + \frac{1}{N} \left( (\mathbf{x}_N - \mu)(\mathbf{x}_N - \mu)^T - \Sigma_{ML}^{(N-1)} \right)$$

- To apply the Robbins-Monro algorithm, *assume that*  $\Sigma$  is *diagonal* and as before compute the derivative

$$\frac{\partial}{\partial \Sigma_{ML}^{(N-1)}} \left( -\ln p(\mathbf{x}_N | \mu, \Sigma_{ML}^{(N-1)}) \right) = -\frac{1}{2} \left( \Sigma_{ML}^{(N-1)} \right)^{-2} \left( (\mathbf{x}_N - \mu)(\mathbf{x}_N - \mu)^T - \Sigma_{ML}^{(N-1)} \right)$$

- Substituting into the RM algorithm:

$$\Sigma_{ML}^{(N)} = \Sigma_{ML}^{(N-1)} + A_{N-1} \frac{1}{2} \left( \Sigma_{ML}^{(N-1)} \right)^{-2} \left( (\mathbf{x}_N - \mu)(\mathbf{x}_N - \mu)^T - \Sigma_{ML}^{(N-1)} \right)$$

- Thus from the RM algorithm, we can obtain the exact update by selecting

$$A_{N-1} = \frac{2}{N} \left( \Sigma_{ML}^{(N-1)} \right)^2$$

# Bayesian Inference for the Gaussian: Known Variance

➤ Consider  $X_1 | \mu \sim \mathcal{N}(\mu, \sigma^2)$ , with prior  $\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$ . We want to infer  $\mu$  with the variance  $\sigma^2$  taken as known. The case with multiple data points will be considered later on.

➤ Then we can derive the following:

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

$\mu | x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  with

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

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

# Bayesian Inference: Predictive distribution

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

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

- We can complete the square by treating the integrand above as a bivariate Gaussian in  $(x, \mu)$ . One can verify that:

$$\frac{1}{2}\left(\frac{(x-\mu)^2}{\sigma^2} + \frac{(\mu-\mu_1)^2}{\sigma_1^2}\right) = \frac{1}{2} \begin{pmatrix} x - \mu_1 & \mu - \mu_1 \end{pmatrix} \underbrace{\begin{pmatrix} \frac{1}{\sigma^2} & -\frac{1}{\sigma^2} \\ -\frac{1}{\sigma^2} & \frac{1}{\sigma^2} + \frac{1}{\sigma_1^2} \end{pmatrix}}_{\Sigma^{-1}} \begin{pmatrix} x - \mu_1 \\ \mu - \mu_1 \end{pmatrix} + \text{const.}$$

- From the above expression note that:  $\Sigma = \begin{pmatrix} \sigma^2 + \sigma_1^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma_1^2 \end{pmatrix}$

# Bayesian Inference: Predictive distribution

- We will see [at a follow up lecture](#) that if we partition the mean and variance of a multivariate Gaussian as:

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

then, the marginal

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

- In our predictive distribution we need to integrate out  $\boldsymbol{\mu}$ .  
Thus based on the above result and  $\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_1 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2 + \sigma_1^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma_1^2 \end{pmatrix}$ , we have:

$$f(x | x_1) = \int \underbrace{f(x | \mu)}_{\text{Likelihood}} \underbrace{\pi(\mu | x_1)}_{\text{Posterior}} d\mu = \mathcal{N}(x | \mu_1, \sigma^2 + \sigma_1^2)$$

- Note *the variance is the sum of model variance + variance of posterior uncertainty in  $\mu$ .*

# 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:

$$\begin{aligned} \pi(\mu | \mathbf{X}) &= \prod_{n=1}^N f(x_n | \mu) \pi(\mu) \propto \exp\left(-\frac{\sum_{n=1}^N (x_n - \mu)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \Rightarrow \\ \pi(\mu | \mathbf{X}) &\propto \exp\left(-\frac{\mu^2}{2} \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right) + \mu \left(\frac{\sum_{n=1}^N x_n}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right)\right) \propto \exp\left(-\frac{1}{2\sigma_1^2} (\mu - \mu_N)^2\right) \end{aligned}$$

# Bayesian Inference for the Gaussian

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

➤ So the posterior is a Gaussian as before with

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

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

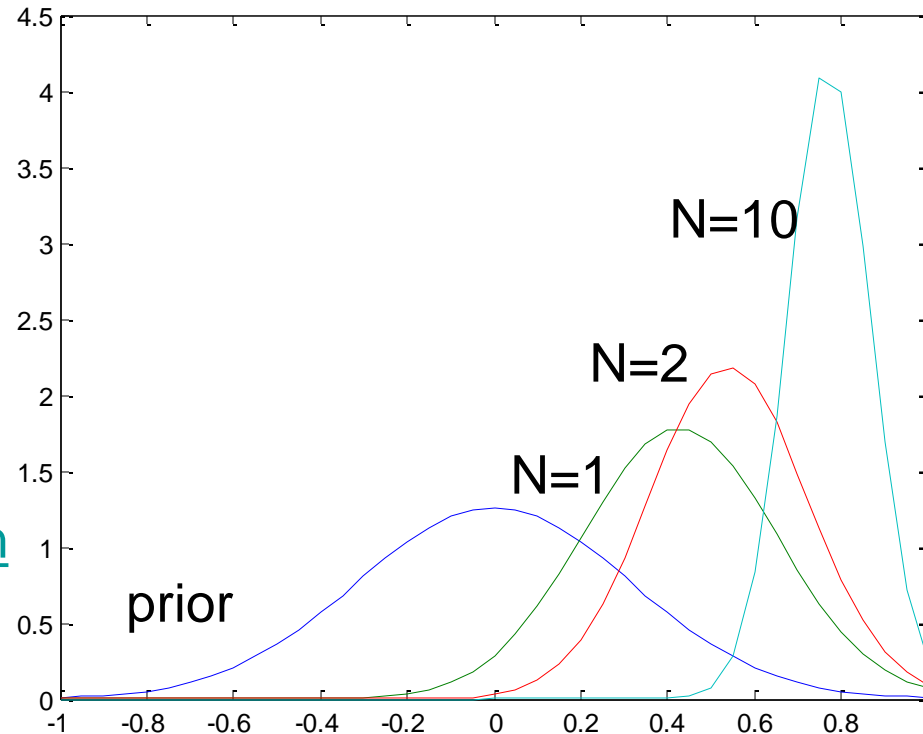
- Observe the posterior mean for  $N \rightarrow \infty$  and  $N \rightarrow 0$ .
- *The posterior precision is the sum of the precision of the prior plus one contribution of the data precision for each observed data point.* As we have seen before for  $N \rightarrow \infty$  the posterior peaks around 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  $\mu_N \rightarrow \mu_{ML}$



# 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](#)



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

# 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 easily derive sequential estimates of the MLE.  
They are as follows:

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