

# Gaussians and Continuous-Density HMMs

Mark Hasegawa-Johnson  
These slides are in the public domain

ECE 417: Multimedia Signal Processing

- 1 Gaussians, Brownian motion, and white noise
- 2 Gaussian Random Vector
- 3 HMM with Gaussian Observation Probabilities
- 4 Summary

# Outline

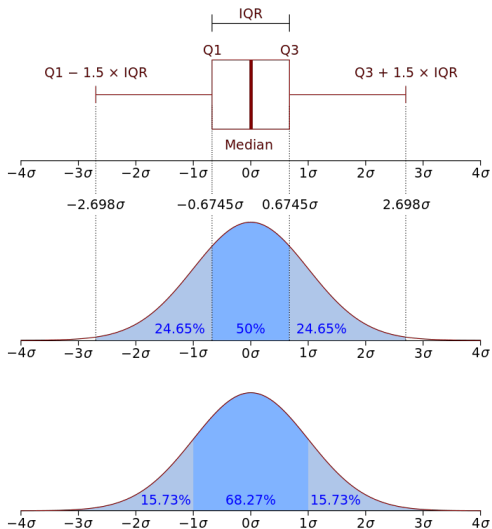
- 1 Gaussians, Brownian motion, and white noise
- 2 Gaussian Random Vector
- 3 HMM with Gaussian Observation Probabilities
- 4 Summary

# Gaussian (Normal) pdf

- Gauss considered this problem: under what circumstances does it make sense to estimate the mean of a distribution,  $\mu$ , by taking the average of the experimental values,  
$$m = \frac{1}{n} \sum_{i=1}^n x_i?$$
- He demonstrated that  $m$  is the maximum likelihood estimate of  $\mu$  if (not only if!)  $X$  is distributed with the following probability density:

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

# Gaussian pdf



[https://commons.wikimedia.org/wiki/File:](https://commons.wikimedia.org/wiki/File:Boxplot_vs_PDF.svg)

Boxplot vs PDF.svg

# Unit Normal pdf

Suppose that  $X$  is normal with mean  $\mu$  and standard deviation  $\sigma$  (variance  $\sigma^2$ ):

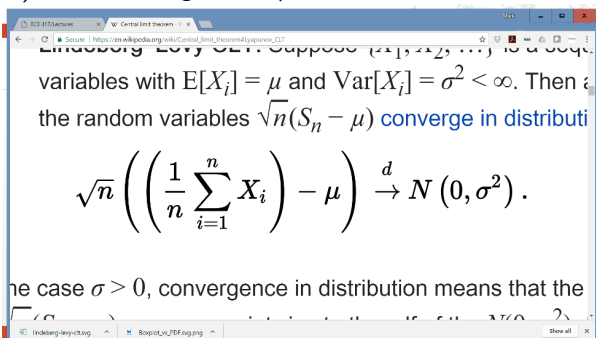
$$p_X(x) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Then  $U = \left(\frac{X-\mu}{\sigma}\right)$  is normal with mean 0 and standard deviation 1:

$$p_U(u) = \mathcal{N}(u; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2}$$

# Central Limit Theorem

The Gaussian pdf is important because of the Central Limit Theorem. Suppose  $X_i$  are i.i.d. (independent and identically distributed), each having mean  $\mu$  and variance  $\sigma^2$ . Then



The screenshot shows a web browser window with the URL [https://en.wikipedia.org/wiki/Central\\_limit\\_theorem#/media/File:Indenberg-Levy-CLT](https://en.wikipedia.org/wiki/Central_limit_theorem#/media/File:Indenberg-Levy-CLT). The text on the page reads: "Suppose  $(X_1, X_2, \dots)$  is a sequence of independent and identically distributed random variables with  $E[X_i] = \mu$  and  $\text{Var}[X_i] = \sigma^2 < \infty$ . Then as  $n \rightarrow \infty$ , the random variables  $\sqrt{n}(S_n - \mu)$  converge in distribution to a normal distribution with mean 0 and variance  $\sigma^2$ . In the case  $\sigma > 0$ , convergence in distribution means that the

$$\sqrt{n} \left( \left( \frac{1}{n} \sum_{i=1}^n X_i \right) - \mu \right) \xrightarrow{d} N(0, \sigma^2).$$

The browser's taskbar at the bottom shows two open files: "Indenberg-levy-clt.svg" and "Boiprot\_vn\_PDF.svg.png".

# Brownian motion

The Central Limit Theorem matters because Einstein showed that the movement of molecules, in a liquid or gas, is the sum of  $n$  i.i.d. molecular collisions.

In other words, the position after  $t$  seconds is Gaussian, with mean 0, and with a variance of  $Dt$ , where  $D$  is some constant.

`https://commons.wikimedia.org/wiki/File:Brownianmotion5particles150frames.gif`

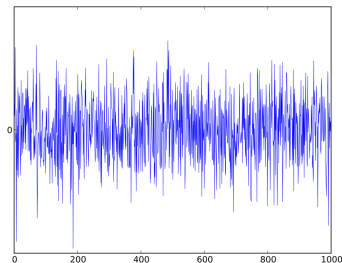


# White Noise

- Sound = air pressure fluctuations caused by velocity of air molecules
- Velocity of warm air molecules without any external sound source = Gaussian

Therefore:

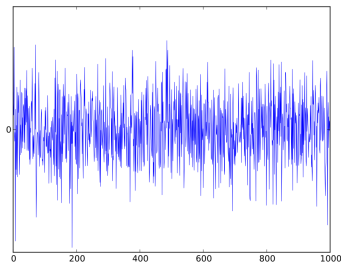
- Sound produced by warm air molecules without any external sound source = Gaussian noise
- Electrical signals: same.



[https://commons.wikimedia.org/wiki/File:White\\_noise.svg](https://commons.wikimedia.org/wiki/File:White_noise.svg)

# White Noise

- White Noise = noise in which each sample of the signal,  $x_n$ , is i.i.d.
- Why “white”? Because the Fourier transform,  $X(\omega)$ , is a zero-mean random variable whose variance is independent of frequency (“white”)
- Gaussian White Noise:  $x[n]$  are i.i.d. and Gaussian



[https://commons.wikimedia.org/wiki/File:White\\_noise.svg](https://commons.wikimedia.org/wiki/File:White_noise.svg)



# Vector of Independent Gaussian Variables

Suppose we have a frame containing  $D$  samples from a Gaussian white noise process,  $x_1, \dots, x_D$ . Let's stack them up to make a vector:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}$$

This whole frame is random. In fact, we could say that  $\mathbf{x}$  is a sample value for a Gaussian random vector called  $X$ , whose elements are  $X_1, \dots, X_D$ :

$$X = \begin{bmatrix} X_1 \\ \vdots \\ X_D \end{bmatrix}$$

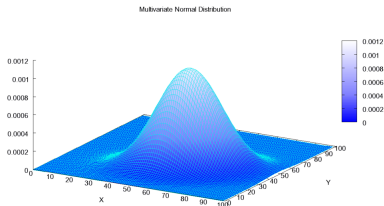
# Vector of Independent Gaussian Variables

Suppose that the  $N$  samples are i.i.d., each one has the same mean,  $\mu$ , and the same variance,  $\sigma^2$ . Then the pdf of this random vector is

$$p_{\mathbf{X}}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 \mathbf{I}) = \prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left( \frac{x_i - \mu}{\sigma} \right)^2}$$

# Vector of Independent Gaussian Variables

Here's an example from Wikipedia with a mean of about 50 and a standard deviation of about 12.



[https://commons.wikimedia.org/wiki/File:Multivariate\\_Gaussian.png](https://commons.wikimedia.org/wiki/File:Multivariate_Gaussian.png)

# Independent Gaussians that aren't identically distributed

Suppose that the  $N$  samples are independent Gaussians that aren't identically distributed, i.e.,  $X_i$  has mean  $\mu_i$  and variance  $\sigma_i^2$ . Then the pdf of this random vector is

$$p_{\mathbf{X}}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2}\left(\frac{x_i - \mu_i}{\sigma_i}\right)^2}$$

where  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are the mean vector and covariance matrix:

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_D \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots \\ 0 & \sigma_2^2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

# Independent Gaussians that aren't identically distributed

Another useful form is:

$$\begin{aligned}
 p_{\mathbf{X}}(\mathbf{x}) &= \prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-\frac{1}{2}\left(\frac{x_i - \mu_i}{\sigma_i}\right)^2} \\
 &= \frac{1}{(2\pi)^{D/2} \prod_{i=1}^D \sigma_d} e^{-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_d - \mu_d}{\sigma_d}\right)^2}
 \end{aligned}$$



## Example

Suppose that  $\mu_1 = 1$ ,  $\mu_2 = -1$ ,  $\sigma_1^2 = 1$ , and  $\sigma_2^2 = 4$ . Then

$$p_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^2 \frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-\frac{1}{2}\left(\frac{x_i - \mu_i}{\sigma_i}\right)^2} = \frac{1}{4\pi} e^{-\frac{1}{2}\left((x_1 - 1)^2 + \left(\frac{x_2 + 1}{2}\right)^2\right)}$$

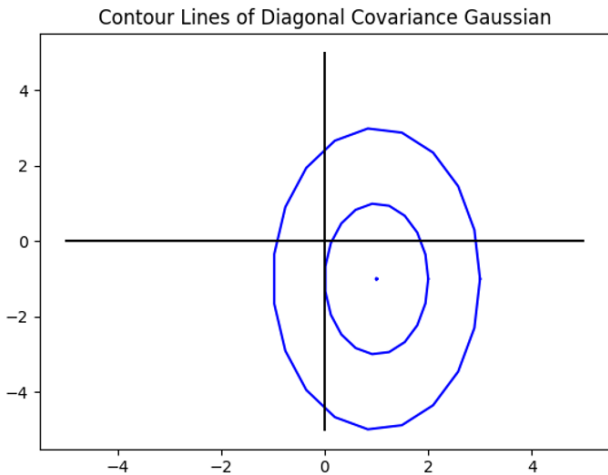
The pdf has its maximum value,  $p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{4\pi}$ , at  $\mathbf{x} = \boldsymbol{\mu} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ .

It drops to  $p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{4\pi\sqrt{e}}$  at  $\mathbf{x} = \begin{bmatrix} \mu_1 \pm \sigma_1 \\ \mu_2 \end{bmatrix}$  and at

$\mathbf{x} = \begin{bmatrix} \mu_1 \\ \mu_2 \pm \sigma_2 \end{bmatrix}$ . It drops to  $p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{4\pi e^2}$  at  $\mathbf{x} = \begin{bmatrix} \mu_1 \pm 2\sigma_1 \\ \mu_2 \end{bmatrix}$

and at  $\mathbf{x} = \begin{bmatrix} \mu_1 \\ \mu_2 \pm 2\sigma_2 \end{bmatrix}$ .

# Example



# Facts about linear algebra #1: determinant of a diagonal matrix

Suppose that  $\Sigma$  is a diagonal matrix, with variances on the diagonal:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots \\ 0 & \sigma_2^2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Then its determinant is

$$|\Sigma| = \prod_{i=1}^D \sigma_d^2$$

So we can write the Gaussian pdf as

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{|2\pi\Sigma|^{1/2}} e^{-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_d - \mu_d}{\sigma_d}\right)^2}$$

## Facts about linear algebra #2: inverse of a diagonal matrix

Suppose that  $\Sigma$  is a diagonal matrix, with variances on the diagonal:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots \\ 0 & \sigma_2^2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Then its inverse is:

$$\Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \cdots \\ 0 & \frac{1}{\sigma_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

## Facts about linear algebra #3: weighted distance

Suppose that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_D \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots \\ 0 & \sigma_2^2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Then

$$\begin{aligned} \sum_{i=1}^D \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2 &= [x_1 - \mu_1, x_2 - \mu_2, \dots] \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \cdots \\ 0 & \frac{1}{\sigma_2^2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \end{bmatrix} \\ &= (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \end{aligned}$$



# Independent Gaussians that aren't identically distributed

So if we have independent Gaussians that aren't identically distributed, we can write the pdf as

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \prod_{i=1}^D \sigma_i} e^{-\frac{1}{2} \sum_{i=1}^D \left(\frac{x_i - \mu_i}{\sigma_i}\right)^2}$$

or as

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{|2\pi \mathbf{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})}$$

or as

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{|2\pi \mathbf{\Sigma}|^{1/2}} e^{-\frac{1}{2} d_{\mathbf{\Sigma}}^2(\mathbf{x}, \boldsymbol{\mu})}$$





# Review: HMM with Discrete Observations

## 1 Initial State Probabilities:

$$\pi'_i = \frac{E[\# \text{ state sequences that start with } q_1 = i]}{\# \text{ state sequences in training data}}$$

## 2 Transition Probabilities:

$$\pi'_i = \frac{E[\# \text{ frames in which } q_{t-1} = i, q_t = j]}{E[\# \text{ frames in which } q_{t-1} = i]}$$

## 3 Observation Probabilities:

$$b'_j(k) = \frac{E[\# \text{ frames in which } q_t = j, k_t = k]}{E[\# \text{ frames in which } q_t = j]}$$

# Baum-Welch with Gaussian Probabilities

The requirement that we vector-quantize the observations is a problem. It means that we can't model the observations very precisely.

It would be better if we could model the observation likelihood,  $b_j(\mathbf{x})$ , as a probability density in the space  $\mathbf{x} \in \mathbb{R}^D$ . One way is to use a parameterized function that is guaranteed to be a properly normalized pdf. For example, a Gaussian:

$$b_i(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

# Diagonal-Covariance Gaussian pdf

Let's assume the feature vector has  $D$  dimensions,

$\mathbf{x}_t = [x_{t,1}, \dots, x_{t,D}]$ . The Gaussian pdf is

$$b_i(\mathbf{x}_t) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_i|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_i)}$$

The logarithm of a Gaussian is

$$\ln b_i(\mathbf{x}_t) = -\frac{1}{2} \left( (\mathbf{x}_t - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_i) + \ln |\boldsymbol{\Sigma}_i| + C \right)$$

where the constant is  $C = D \ln(2\pi)$ .

# Baum-Welch

Baum-Welch maximizes the expected log probability, i.e.,

$$E_{\mathbf{q}|\mathbf{X}} [\ln b_i(\mathbf{x}_t)] = -\frac{1}{2} \sum_{i=1}^N \gamma_t(i) \left( (\mathbf{x}_t - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_i) + \ln |\boldsymbol{\Sigma}_i| + C \right)$$

If we include all of the frames, then we get

$$\begin{aligned} E_{\mathbf{q}|\mathbf{X}} [\ln p(\mathbf{X}, \mathbf{q}|\Lambda)] &= \text{other terms} \\ &- \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^N \gamma_t(i) \left( (\mathbf{x}_t - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_i) + \ln |\boldsymbol{\Sigma}_i| + C \right) \end{aligned}$$

where the “other terms” are about  $a_{i,j}$  and  $\pi_i$ , and have nothing to do with  $\boldsymbol{\mu}_i$  or  $\boldsymbol{\Sigma}_i$ .

# M-Step: optimum $\mu$

First, let's optimize  $\mu$ . We want

$$0 = \frac{\partial}{\partial \mu_q} \sum_{t=1}^T \sum_{i=1}^N \gamma_t(i) (\mathbf{x}_t - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_t - \mu_i)$$

Re-arranging terms, we get

$$\mu'_q = \frac{\sum_{t=1}^T \gamma_t(q) \mathbf{x}_t}{\sum_{t=1}^T \gamma_t(q)}$$

M-Step: optimum  $\Sigma$ 

Second, let's optimize  $\Sigma_i$ . For this, it's easier to express the log likelihood as

$$E_{\mathbf{q}|\mathbf{X}} [\ln p_{\mathbf{X}}(\mathbf{X}, \mathbf{q})] = \text{other stuff} - \frac{1}{2} \sum_{t=1}^T \gamma_t(i) \sum_{d=1}^D \left( \ln \sigma_{i,d}^2 + \frac{(x_{t,d} - \mu_{i,d})^2}{\sigma_{i,d}^2} \right)$$

Its scalar derivative is

$$\frac{\partial E_{\mathbf{q}|\mathbf{X}} [\ln p_{\mathbf{X}}(\mathbf{X}, \mathbf{q})]}{\partial \sigma_{i,d}^2} = -\frac{1}{2} \sum_{t=1}^T \gamma_t(i) \left( \frac{1}{\sigma_{i,d}^2} - \frac{(x_{t,d} - \mu_{i,d})^2}{\sigma_{i,d}^4} \right)$$

Which we can solve to find

$$\sigma_{i,d}^2 = \frac{\sum_{t=1}^T \gamma_t(i) (x_{t,d} - \mu_{t,d})^2}{\sum_{t=1}^T \gamma_t(i)}$$

# Minimizing the cross-entropy: optimum $\sigma$

Arranging all the scalar derivatives into a matrix, we can write

$$\Sigma'_i = \frac{\sum_{t=1}^T \gamma_t(i) (\mathbf{x}_t - \boldsymbol{\mu}_i) (\mathbf{x}_t - \boldsymbol{\mu}_i)^T}{\sum_{t=1}^T \gamma_t(i)}$$

- Actually, the above formula holds even if the Gaussian has a non-diagonal covariance matrix, but Gaussians with non-diagonal covariance matrices work surprisingly badly in HMMs.
- For a diagonal-covariance Gaussian, we evaluate only the diagonal elements of the vector outer product  $(\mathbf{x}_t - \boldsymbol{\mu}_i) (\mathbf{x}_t - \boldsymbol{\mu}_i)^T$

# Summary: Gaussian Observation PDFs

So we can use Gaussians for  $b_j(\mathbf{x})$ :

- **E-Step:**

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{\sum_{i'} \alpha_t(i')\beta_t(i')}$$

- **M-Step:**

$$\mu'_i = \frac{\sum_{t=1}^T \gamma_t(i)\mathbf{x}_t}{\sum_{t=1}^T \gamma_t(i)}$$
$$\Sigma'_i = \frac{\sum_{t=1}^T \gamma_t(i)(\mathbf{x}_t - \mu_i)(\mathbf{x}_t - \mu_i)^T}{\sum_{t=1}^T \gamma_t(i)}$$





# Summary: Independent Gaussians that aren't identically distributed

$$\begin{aligned} p_X(\mathbf{x}) &= \frac{1}{(2\pi)^{D/2} \prod_{i=1}^D \sigma_i} e^{-\frac{1}{2} \sum_{i=1}^D \left(\frac{x_i - \mu_i}{\sigma_i}\right)^2} \\ &= \frac{1}{|2\pi \boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \\ &= \frac{1}{|2\pi \boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} d_{\boldsymbol{\Sigma}}^2(\mathbf{x}, \boldsymbol{\mu})} \end{aligned}$$

# Summary: Gaussian Observation PDFs

So we can use Gaussians for  $b_j(\mathbf{x})$ :

- **E-Step:**

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{\sum_{i'} \alpha_t(i')\beta_t(i')}$$

- **M-Step:**

$$\mu'_i = \frac{\sum_{t=1}^T \gamma_t(i)\mathbf{x}_t}{\sum_{t=1}^T \gamma_t(i)}$$
$$\Sigma'_i = \frac{\sum_{t=1}^T \gamma_t(i)(\mathbf{x}_t - \mu_i)(\mathbf{x}_t - \mu_i)^T}{\sum_{t=1}^T \gamma_t(i)}$$