

# Machine Learning 4771

Instructor: Tony Jebara

# Topic 9

- Continuous Probability Models
- Gaussian Distribution
- Maximum Likelihood Gaussian
- Sampling from a Gaussian

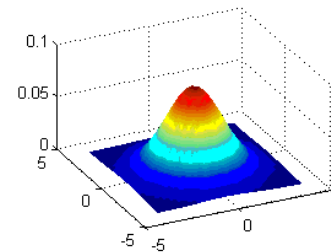
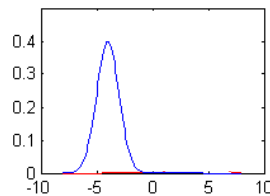
# Continuous Probability Models

- Probabilities can have both discrete & continuous variables
- We will discuss:
  - 1) discrete probability tables

x=T	x=H
0.4	0.6

x=1	x=2	x=3	x=4	x=5	x=6
0.1	0.1	0.1	0.1	0.1	0.5

## 2) continuous probability distributions

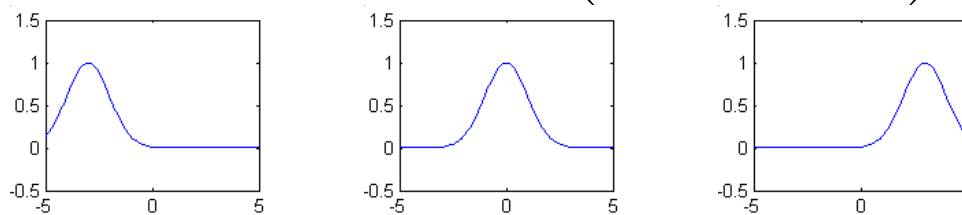


- Most popular continuous distribution = Gaussian

# Gaussian Distribution

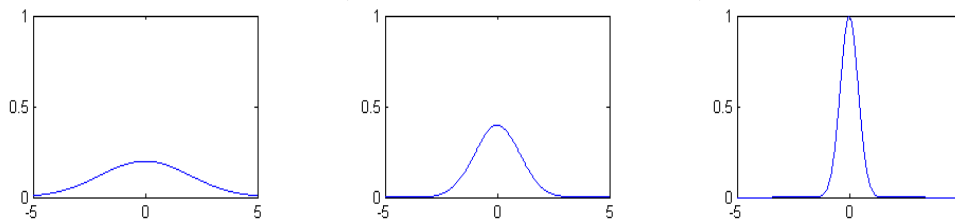
- Recall 1-dimensional Gaussian with mean parameter  $\mu$  translates Gaussian left & right

$$p(x | \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right)$$



- Can also have variance parameter  $\sigma^2$  widens or narrows the Gaussian

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$



Note:  $\int_{x=-\infty}^{\infty} p(x) dx = 1$

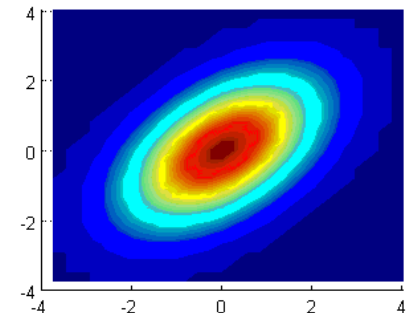
# Multivariate Gaussian

- Gaussian can extend to D-dimensions
- Gaussian mean parameter  $\mu$  vector, it translates the bump
- Covariance matrix  $\Sigma$  stretches and rotates bump

$$p(\vec{x} \mid \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)$$

- Mean is any real vector
- Max and expectation =  $\mu$
- Variance parameter is now  $\Sigma$  matrix
- Covariance matrix is positive definite
- Covariance matrix is symmetric
- Need matrix **inverse** (inv)
- Need matrix **determinant** (det)
- Need matrix **trace** operator (trace)

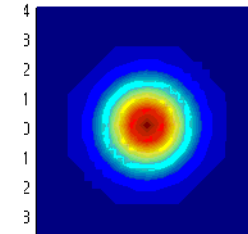
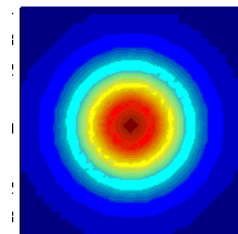
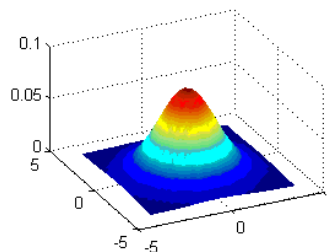
$$\vec{x} \in \mathbb{R}^D, \vec{\mu} \in \mathbb{R}^D, \Sigma \in \mathbb{R}^{D \times D}$$



# Multivariate Gaussian

- Spherical:

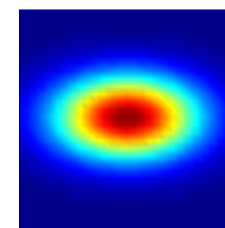
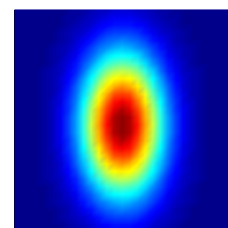
$$\Sigma = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$$



- Diagonal Covariance:

dimensions of  $\mathbf{x}$  are independent  
product of multiple 1d Gaussians

$$p(\vec{x} \mid \vec{\mu}, \Sigma) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi}\vec{\sigma}(d)} \exp\left(-\frac{(\vec{x}(d) - \vec{\mu}(d))^2}{2\vec{\sigma}(d)^2}\right)$$



$$\Sigma = \begin{bmatrix} \vec{\sigma}(1)^2 & 0 & 0 & 0 \\ 0 & \vec{\sigma}(2)^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \vec{\sigma}(D)^2 \end{bmatrix}$$

# Max Likelihood Gaussian

- Have IID samples as vectors  $i=1..N$ :  $\mathcal{D} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$
- How do we recover the mean and covariance parameters?
- Standard approach: Maximum Likelihood (IID)
- Maximize probability of data given model (likelihood)

$$\begin{aligned} p(\mathcal{D} | \theta) &= p(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N | \theta) \\ &= \prod_{i=1}^N p(\vec{x}_i | \vec{\mu}_i, \Sigma_i) \quad \text{independent Gaussian samples} \\ &= \prod_{i=1}^N p(\vec{x}_i | \vec{\mu}, \Sigma) \quad \text{identically distributed} \end{aligned}$$

- Instead, work with maximum of log-likelihood

$$\sum_{i=1}^N \log p(\vec{x}_i | \vec{\mu}, \Sigma) = \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu})\right)$$

# Max Likelihood Gaussian

•Max over  $\mu$   $\frac{\partial}{\partial \mu} \left( \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp \left( -\frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right) \right) = 0$

$$\frac{\partial}{\partial \mu} \left( \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right) = 0$$

$$\frac{\partial \vec{x}^T \vec{x}}{\partial \vec{x}} = 2\vec{x}^T$$

$$\sum_{i=1}^N (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} = \vec{0}$$

see Jordan Ch. 12, get sample mean...

$$\vec{\mu} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$$

•For  $\Sigma$  need Trace operator:  $tr(A) = tr(A^T) = \sum_{d=1}^D A(d, d)$

and several properties:

$$tr(AB) = tr(BA)$$

$$tr(BAB^{-1}) = tr(A)$$

$$tr(\vec{x}\vec{x}^T A) = tr(\vec{x}^T A \vec{x}) = \vec{x}^T A \vec{x}$$



# Max Likelihood Gaussian

- Likelihood rewritten in trace notation:

$$\begin{aligned}
 l &= \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right] \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} \right] \\
 &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |A| - \frac{1}{2} \sum_{i=1}^N \text{tr} \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T A \right]
 \end{aligned}$$

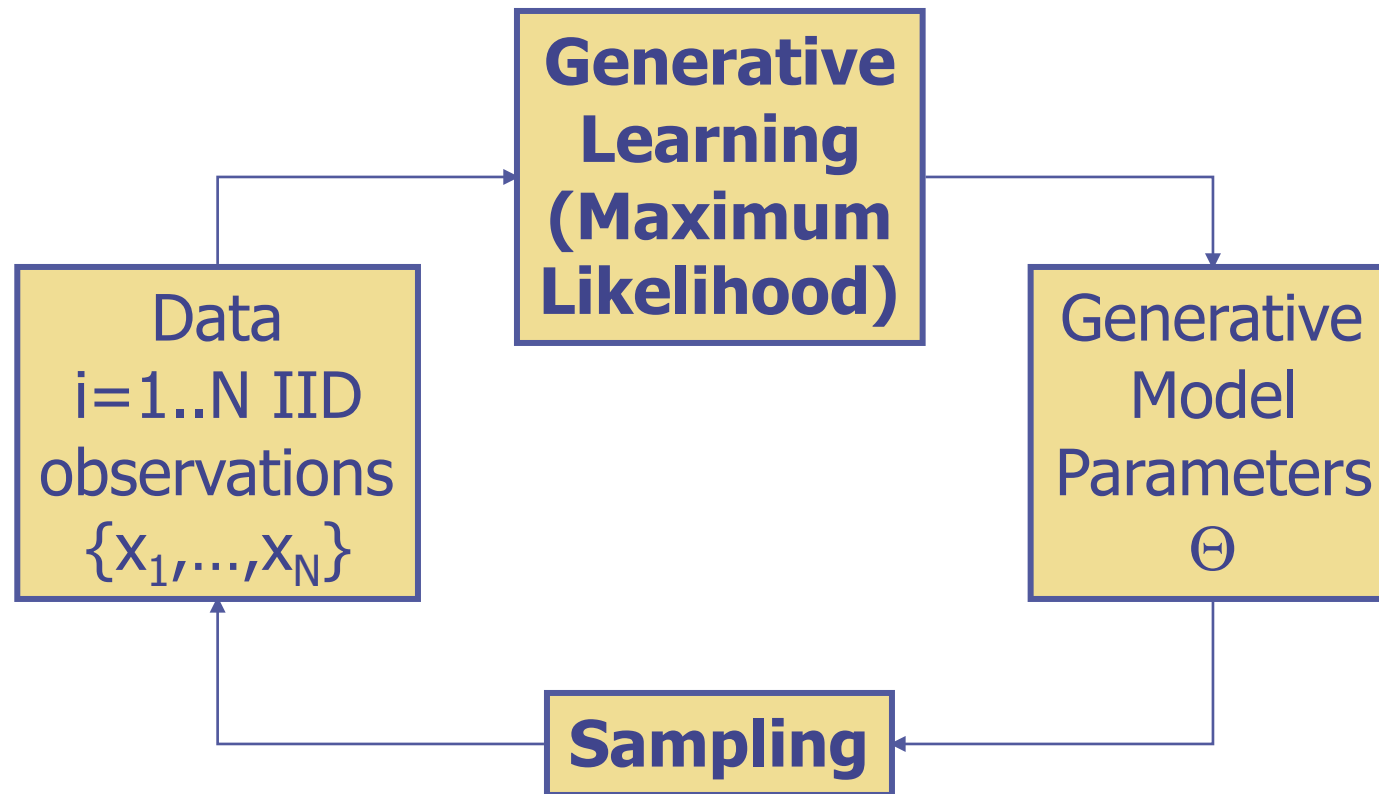
- Max over  $A = \Sigma^{-1}$   
use properties:

$$\begin{aligned}
 \frac{\partial l}{\partial A} &= -0 + \frac{N}{2} \left( A^{-1} \right)^T - \frac{1}{2} \sum_{i=1}^N \left[ (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T \right]^T \\
 &= \frac{N}{2} \Sigma - \frac{1}{2} \sum_{i=1}^N (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T
 \end{aligned}$$

$\frac{\partial \log |A|}{\partial A} = (A^{-1})^T$        $\frac{\partial \text{tr}[BA]}{\partial A} = B^T$

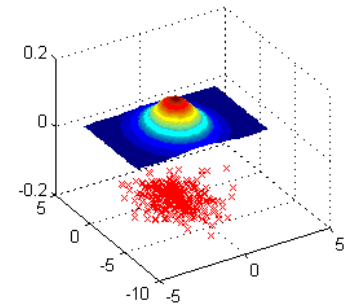
- Get sample covariance:  $\frac{\partial l}{\partial A} = 0 \rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{\mu}) (\vec{x}_i - \vec{\mu})^T$

# Sampling & Max Likelihood

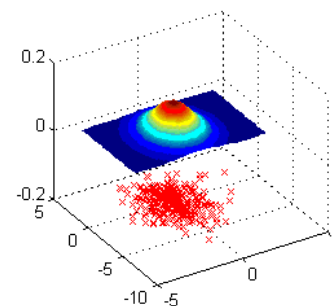


# Sampling from a Gaussian

- Fit Gaussian to data, how is this Generative?



# Sampling from a Gaussian



- Fit Gaussian to data, how is this Generative?

- Sampling! Generating discrete data easy:

0.73	0.1	0.17
------	-----	------

- Assume we can do uniform sampling:

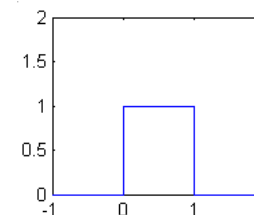
- i.e. rand between (0,1)

- if  $0.00 \leq \text{rand} < 0.73$  get A

- if  $0.73 \leq \text{rand} < 0.83$  get B

- if  $0.83 \leq \text{rand} < 1.00$  get C

- What are we doing?



# Sampling from a Gaussian

- Fit Gaussian to data, how is this Generative?

- Sampling! Generating discrete data easy:

0.73	0.1	0.17
------	-----	------

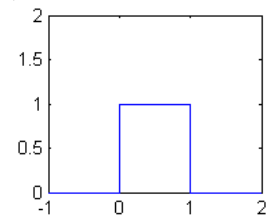
- Assume we can do uniform sampling:

i.e. rand between (0,1)

if  $0.00 \leq \text{rand} < 0.73$  get A

if  $0.73 \leq \text{rand} < 0.83$  get B

if  $0.83 \leq \text{rand} < 1.00$  get C

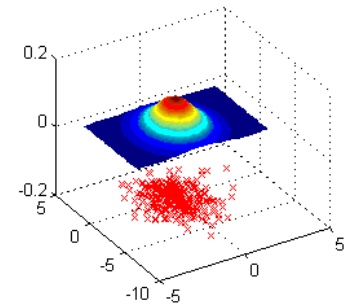
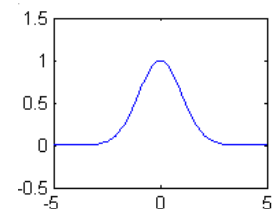


- What are we doing?

Sum up the Probability Density Function (PDF)  
to get Cumulative Density Function (CDF)

0.73	0.83	1.00
------	------	------

- For 1d Gaussian, Integrate Probability Density Function get Cumulative Density Function  
Integral is like summing many discrete bars



# Sampling from a Gaussian

- Integrate 1d Gaussian to get CDF:

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

$$F(x) = \int_{-\infty}^x p(t) dt = \frac{1}{2} \operatorname{erf}\left(\frac{1}{\sqrt{2}}x\right) + \frac{1}{2}$$

- If sample from uniform, get:  $u \sim \operatorname{uniform}(0,1)$

- Compute mapping:  $x = F^{-1}(u) = \sqrt{2} \operatorname{erfinv}(2u - 1)$

- This is a Gaussian sample:  $x \sim N(x | 0, 1)$

- For D-dimensional Gaussian  $N(\mathbf{z} | 0, I)$  concatenate samples:

$$\vec{x} = [\vec{x}(1) \dots \vec{x}(D)]^T \sim p(\vec{x} | 0, I) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \vec{x}(d)^2\right)$$

- For  $N(\mathbf{z} | \vec{\mu}, \Sigma)$ , add mean & multiply by root cov

$$\vec{z} = \Sigma^{1/2} \vec{x} + \vec{\mu} \sim p(\vec{z} | \vec{\mu}, \Sigma)$$

- Example code: `gendata.m`

