## ECE285/SIO209, Machine learning for physical applications, Spring 2017

Peter Gerstoft, 534-7768, gerstoft@ucsd.edu

We meet Wednesday from 5 to 6:20pm in Spies Hall 330

**Text Bishop** 

Grading A or maybe S

#### Classes

First 4 weeks: Focus on theory/implementationMiddle 3 weeks: 50% Applications plus 50% theoryLast 3 weeks: 30% Final Project, 30% Applications plus 50% theory

#### **Applications**

Graph theory for source localization: Gerstoft Source tracking in ocean acoustics: Grad Student Emma Reeves Aramco Research: Weichang Li, Group leader Seismic network using mobile phones: Berkeley Eric Orenstein: identifying plankton Plus more - Dictionary learning

\_

Homework: Both matlab/python will be used, Just email the code to me (I dont need anything else).
Homework is due 11am on Wednesday. That way we can discuss in class.
Hw 1:

matlab/ python/ ipython/ jupyter ?

Tritoned? https://tritoned.ucsd.edu/

#### May 31 - June 2 Big Data and The Earth Sciences: Grand Challenges Workshop

#### by John Graham — published Feb 07, 2017 06:31 PM, last modified Mar 28, 2017 08:14 PM

Big Data and The Earth Sciences: Grand Challenges Workshop May 31 - June 2

The goal of the The Big Data and Earth Sciences: Grand Challenges Workshop is to bring thought leaders in Big Data and Earth Sciences together for a three day, intensive workshop to discuss what is needed to advance our understanding and predictability of the Earth systems and to highlight key technological advances and methods that are readily available (or will be soon) to assist this advancement. With the ever growing quantity and quality of hyper-dimensional earth science data (satellite and ground based observations and cutting–edge Numerical Weather Prediction (NWP) models), the advancements in machine learning (e.g. supervised, unsupervised and semi-supervised learning techniques), and the progress made in the application of Graphical Processing Units (GPUs) and GPU clusters, we now have an unprecedented opportunity and challenge to engage these computational advances to improve our understanding of the complex nature of the interactions between various earth science events, their variables and their impacts on society (flooding, drought, agriculture, etc.).

#### Grand Challenges Lectures (CONFIRMED):

**Dr. Larry Smarr**, Founding Director of the California Institute for Telecommunications and Information Technology (Calit2), a UC San Diego/UC Irvine partnership, holds the Harry E. Gruber professorship in Computer Science and Engineering (CSE) at UC San Diego's Jacobs School.

**Dr. Vipin Kumar**, Regents Professor at the University of Minnesota, holds the William Norris Endowed Chair in the Department of Computer Science and Engineering, University of Minnesota.

Dr. Padhraic Smyth, Professor, Director, UCI Data Science Initiative and Associate Director, Center for Machine Learning and Intelligent Systems, UC Irvine.

Dr. Michael Wehner, Senior Staff Scientists, Computational Research Division at the Lawrence Berkeley National Laboratory.

Hotel accommodations:

- · La Jolla Sheraton (nice, economical, close by): http://www.sheratonlajolla.com/
- · Estancia (Closest location and most beautiful): http://meritagecollection.com/estancialajolla/
- La Jolla Shores (beach front property farther away): <a href="http://www.ljshoreshotel.com/?src=ppc\_google\_ljshores\_brand\_expanded&NCK=ppc\_google\_ljshores\_brand&gclid=CNTF8JTmqdICFQmlfgodfTMP7A">http://www.ljshoreshotel.com/?src=ppc\_google\_ljshores\_brand\_expanded&NCK=ppc\_google\_ljshores\_brand&gclid=CNTF8JTmqdICFQmlfgodfTMP7A</a>

Please send abstracts to scottsellars@ucsd.edu

Please register here: Workshop Registration Form

#### Download the call for papers HERE

When	May 31, 2017 08:00 AM to Jun 02, 2017 06:00 PM
Where	Qualcomm Institute, Calit2 - University of California, (UCSD) Atkinson Hall Auditorium 9500 Gilman La Jolla, CA 92093
Contact Name	Scott L. Sellars
Add event to calendar	<mark>⊠,</mark> vCal <mark>∐,</mark> iCal

# Entropy

$$\mathbf{H}[x] = -\sum_{x} p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning



## **Parametric Distributions**

- $p(\mathbf{x}|\boldsymbol{\theta})$ Basic building blocks: ullet
- Need to determine  $oldsymbol{ heta}_{ extsf{given}}$   $\{\mathbf{x}_1,\ldots,\mathbf{x}_N\}$ Representation:  $oldsymbol{ heta}^{\star}$  or  $p(oldsymbol{ heta})$ •
- ullet
- **Recall Curve Fitting** ٠

#### The Gaussian Distribution



## **Central Limit Theorem**

•The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N grows.

•Example: N uniform [0,1] random variables.



#### Geometry of the Multivariate Gaussian



### Moments of the Multivariate Gaussian (1)

$$\mathbb{E}[\mathbf{x}] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \mathbf{x} \, \mathrm{d}\mathbf{x}$$
$$= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{z}\right\} (\mathbf{z}+\boldsymbol{\mu}) \, \mathrm{d}\mathbf{z}$$

thanks to anti-symmetry of z

$$\mathbb{E}[\mathbf{x}] = oldsymbol{\mu}$$

### Moments of the Multivariate Gaussian (2)

$$\mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] = \boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}} + \boldsymbol{\Sigma}$$
$$\operatorname{cov}[\mathbf{x}] = \mathbb{E}\left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}}\right] = \boldsymbol{\Sigma}$$

A Gaussian requires D\*(D-1)/2 +D parameters. Often we use D +D or Just D+1 parameters.



## **Partitioned Gaussian Distributions**

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

$$\mathbf{x} = egin{pmatrix} \mathbf{x}_a \ \mathbf{x}_b \end{pmatrix} \qquad \qquad oldsymbol{\mu} = egin{pmatrix} oldsymbol{\mu}_a \ oldsymbol{\mu}_b \end{pmatrix} \qquad \qquad oldsymbol{\Sigma} = egin{pmatrix} oldsymbol{\Sigma}_{aa} & oldsymbol{\Sigma}_{ab} \ oldsymbol{\Sigma}_{ba} & oldsymbol{\Sigma}_{bb} \end{pmatrix}$$

$$oldsymbol{\Lambda} \equiv oldsymbol{\Sigma}^{-1} \qquad oldsymbol{\Lambda} = egin{pmatrix} oldsymbol{\Lambda}_{aa} & oldsymbol{\Lambda}_{ab} \ oldsymbol{\Lambda}_{ba} & oldsymbol{\Lambda}_{bb} \end{pmatrix}$$

## Partitioned Conditionals and Marginals

$$egin{aligned} p(\mathbf{x}_a | \mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a | oldsymbol{\mu}_{a|b}, oldsymbol{\Sigma}_{a|b}) \ \mathbf{\Sigma}_{a|b} &= & \mathbf{\Lambda}_{aa}^{-1} = oldsymbol{\Sigma}_{aa} - oldsymbol{\Sigma}_{ab} oldsymbol{\Sigma}_{bb}^{-1} oldsymbol{\Sigma}_{ba} \ oldsymbol{\mu}_{a|b} &= & oldsymbol{\Sigma}_{a|b} \left\{ \mathbf{\Lambda}_{aa} oldsymbol{\mu}_{a} - oldsymbol{\Lambda}_{ab} (\mathbf{x}_b - oldsymbol{\mu}_b) 
ight\} \ &= & oldsymbol{\mu}_a - oldsymbol{\Lambda}_{aa}^{-1} oldsymbol{\Lambda}_{ab} (\mathbf{x}_b - oldsymbol{\mu}_b) 
ight\} \ &= & oldsymbol{\mu}_a + oldsymbol{\Sigma}_{ab} oldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - oldsymbol{\mu}_b) \end{aligned}$$

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

## Partitioned Conditionals and Marginals



## Bayes' Theorem for Gaussian Variables

• Given

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

• we have

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
  
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

• where

$$\mathbf{\Sigma} = (\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}$$

## Maximum Likelihood for the Gaussian (1)

• Given i.i.d. data

, the log likeli-hood function is given by

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^{\mathrm{T}}$$

• Sufficient statistics

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

$$\sum_{n=1}^{N} \mathbf{x}_n \qquad \qquad \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^{\mathrm{T}}$$

# Maximum Likelihood for the Gaussian (2)

• Set the derivative of the log likelihood function to zero,

$$\frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) = 0$$

• and solve to obtain

• Similarly 
$$oldsymbol{\mu}_{ ext{ML}} = rac{1}{N}\sum_{n=1}^N \mathbf{x}_n.$$

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

Maximum Likelihood for the Gaussian (3)

# Under the true distribution

$$egin{array}{rcl} \mathbb{E}[oldsymbol{\mu}_{ ext{ML}}] &=& oldsymbol{\mu} \ \mathbb{E}[oldsymbol{\Sigma}_{ ext{ML}}] &=& rac{N-1}{N}oldsymbol{\Sigma}. \end{array}$$

Hence define

$$\widetilde{\mathbf{\Sigma}} = rac{1}{N-1} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}_{\mathrm{ML}}) (\mathbf{x}_n - \boldsymbol{\mu}_{\mathrm{ML}})^{\mathrm{T}}.$$

## **Sequential Estimation**

# Contribution of the $N^{th}$ data point, $x_N$



## Bayesian Inference for the Gaussian (1)

• Assume  $\sigma^2$  is known. Given i.i.d. data  $\mathbf{x} = \{x_1, \ldots, x_N\}$  the likelihood function for  $\mu$  is given by

$$p(\mathbf{x}|\mu) = \prod_{n=1}^{N} p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2\right\}.$$

- This has a Gaussian shape as a function of  $\mu$  (but it is *not* a distribution over  $\mu$ ).

## Bayesian Inference for the Gaussian (2)

• Combined with a Gaussian prior over μ,

$$p(\mu) = \mathcal{N}\left(\mu|\mu_0, \sigma_0^2\right).$$

• this gives the posterior

 $p(\mu|\mathbf{x}) \propto p(\mathbf{x}|\mu)p(\mu).$ 

• Completing the square over  $\mu$ , we see that

$$p(\mu|\mathbf{x}) = \mathcal{N}\left(\mu|\mu_N, \sigma_N^2\right)$$

$$\mu_{N} = \frac{\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{0} + \frac{N\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{ML}, \qquad \mu_{ML} = \frac{1}{N}\sum_{n=1}^{N}x_{n}$$
$$\frac{1}{\sigma_{N}^{2}} = \frac{1}{\sigma_{0}^{2}} + \frac{N}{\sigma^{2}}.$$

$$\begin{array}{c|ccc} N = 0 & N \to \infty \\ \hline \mu_N & \mu_0 & \mu_{\rm ML} \\ \sigma_N^2 & \sigma_0^2 & 0 \end{array}$$

## Bayesian Inference for the Gaussian (4)

• Example: for N = 0, 1, 2 and 10.

$$p(\mu|\mathbf{x}) = \mathcal{N}\left(\mu|\mu_N, \sigma_N^2\right)$$



## Bayesian Inference for the Gaussian (5)

• Sequential Estimation

$$p(\mu|\mathbf{x}) \propto p(\mu)p(\mathbf{x}|\mu)$$

$$= \left[p(\mu)\prod_{n=1}^{N-1}p(x_n|\mu)\right]p(x_N|\mu)$$

$$\propto \mathcal{N}\left(\mu|\mu_{N-1},\sigma_{N-1}^2\right)p(x_N|\mu)$$

- The posterior obtained after observing  $N \ \{ \ 1 \ \text{data points becomes the prior when we observe the } N^{th} \ \text{data point.}$ 

## Bayesian Inference for the Gaussian (6)

• Now assume  $\mu$  is known. The likelihood function for  $\lambda = 1/\sigma^2$  is given by

$$p(\mathbf{x}|\lambda) = \prod_{n=1}^{N} \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp\left\{-\frac{\lambda}{2} \sum_{n=1}^{N} (x_n - \mu)^2\right\}.$$

- This has a Gamma shape as a function of  $\lambda$ .
- The Gamma distribution:

$$\operatorname{Gam}(\lambda|a,b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \qquad \mathbb{E}[\lambda] = \frac{a}{b} \qquad \operatorname{var}[\lambda] = \frac{a}{b^2}$$



## Bayesian Inference for the Gaussian (8)

- Now we combine a Gamma prior,  $\operatorname{Gam}(\lambda|a_0,b_0)$  with the likelihood function for  $\lambda$  to obtain

$$p(\lambda|\mathbf{x}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp\left\{-b_0 \lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

• which we recognize as  $\operatorname{Gam}(\lambda|a_N,b_N)$  with

$$a_N = a_0 + \frac{N}{2}$$
  
 $b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{\text{ML}}^2.$ 

## Bayesian Inference for the Gaussian (9)

• If both  $\mu$  and  $\lambda$  are unknown, the joint likelihood function is given by

$$p(\mathbf{x}|\mu,\lambda) = \prod_{n=1}^{N} \left(\frac{\lambda}{2\pi}\right)^{1/2} \exp\left\{-\frac{\lambda}{2}(x_n-\mu)^2\right\}$$
$$\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu\sum_{n=1}^{N} x_n - \frac{\lambda}{2}\sum_{n=1}^{N} x_n^2\right\}.$$

• We need a prior with the same functional dependence on  $\mu$  and  $\lambda$ .

## Bayesian Inference for the Gaussian (10)

• The Gaussian-gamma distribution

$$p(\mu, \lambda) = \mathcal{N}(\mu|\mu_0, (\beta\lambda)^{-1}) \operatorname{Gam}(\lambda|a, b)$$

$$\propto \exp\left\{-\frac{\beta\lambda}{2}(\mu - \mu_0)^2\right\} \lambda^{a-1} \exp\left\{-b\lambda\right\}$$
• Quadratic in  $\mu$ .  
• Linear in  $\lambda$ .  
• Independent of  $\mu$ .  

$$\mu_0=0, \beta=2, a=5, b=6$$

$$\lambda \prod_{\substack{\substack{a = 0 \\ a = 2 \\ a = 2$$

## Bayesian Inference for the Gaussian (12)

- Multivariate conjugate priors
- $\mu$  unknown,  $\Lambda$  known:  $p(\mu)$  Gaussian.
- $\Lambda$  unknown,  $\mu$  known:  $p(\Lambda)$  Wishart,

$$\mathcal{W}(\mathbf{\Lambda}|\mathbf{W},\nu) = B|\mathbf{\Lambda}|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\mathrm{Tr}(\mathbf{W}^{-1}\mathbf{\Lambda})\right)$$

•  $\Lambda$  and  $\mu$  unknown:  $p(\mu, \Lambda)$  Gaussian-Wishart,

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \boldsymbol{\mu}_0, \boldsymbol{\beta}, \mathbf{W}, \nu) = \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_0, (\boldsymbol{\beta} \boldsymbol{\Lambda})^{-1}) \mathcal{W}(\boldsymbol{\Lambda} | \mathbf{W}, \nu)$$

# Mixtures of Gaussians (1)

#### Old Faithful geyser:

The time between eruptions has a <u>bimodal distribution</u>, with the mean interval being either 65 or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of ±10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less than  $2\frac{1}{2}$  minutes, or 91 minutes after an eruption lasting more than  $2\frac{1}{2}$  minutes.



## Mixtures of Gaussians (2)



$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
  
Component  
Mixing coefficient

$$\forall k : \pi_k \ge 0 \qquad \sum_{k=1}^K \pi_k = 1$$

## Mixtures of Gaussians (3)



## Mixtures of Gaussians (4)

• Determining parameters  $\pi$ ,  $\mu$ , and  $\Sigma$  using maximum log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Log of a sum; no closed form maximum.

• Solution: use standard, iterative, numeric optimization methods or the *expectation maximization* algorithm (Chapter 9).