Gaussian Distribution

Sargur N. Srihari

The Gaussian Distribution

- For single real-valued variable x $N(x \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$
- Parameters:
 - Mean μ , variance σ^2 ,
- Standard deviation σ

Cai 177

Carl Friedrich Gauss 1777-1855

68% of data lies within σ of mean 95% within 2σ

- Precision $\beta = 1/\sigma^2$, $E[x] = \mu$, $Var[x] = \sigma^2$
- For *D*-dimensional vector x, multivariate Gaussian

 $\mathcal{N}(x|\mu, \sigma^2)$

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

 μ is a mean vector, Σ is a *D* x *D* covariance matrix, $|\Sigma|$ is the determinant of Σ Σ⁻¹ is also referred to as the precision matrix

Covariance Matrix

- Gives a measure of the dispersion of the data
- It is a *D* x *D* matrix
 - Element in position i,j is the covariance between the i^{th} and j^{th} variables.
 - Covariance between two variables x_i and x_j is defined as $E[(x_i \mu_i)(y_i \mu_j)]$
 - Can be positive or negative
 - If the variables are independent then the covariance is zero.
 - Then all matrix elements are zero except diagonal elements which represent the variances

Importance of Gaussian

- Gaussian arises in many different contexts, e.g.,
 - For a single variable,
 Gaussian maximizes entropy (for given mean and variance)
 - Sum of set of random variables becomes increasingly Gaussian



Mean of two variables

The two values could be 0.8 and 0.2 whose average is 0.5 More ways of getting 0.5 than say 0.1



Mean of ten variables



Geometry of Gaussian

 Functional dependence of Gaussian on x is through

 $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$

- Called Mahanalobis Distance
- reduces to Euclidean distance when Σ is an identity matrix
- Matrix Σ is symmetric
 - Has an Eigenvector equation

$$\Sigma u_i = \lambda_i u_i$$

 u_i are Eigen vectors λ_i are Eigen values

Two dimensional Gaussian $x = (x_1, x_2)$



Red: Elliptical contour of constant density Major axes: eigenvectors u_i

Contours of Constant Density

 Determined by Covariance Matrix

 Covariances represent
 how features vary together

(b) Diagonal matrix(aligned withcoordinate axes)





(c) Proportional toIdentity matrix(concentric circles)



Machine Learning

Joint Gaussian implies that Marginal and Conditional are Gaussian

- If two sets of variables x_a,x_b are jointly Gaussian then the two conditional densities and the two marginals are also Gaussian
- Given joint Gaussian $N(x|\mu,\Sigma)$ with $\Lambda = \Sigma^{-1}$ and $x = [x_a, x_b]^T$ where x_a are first *m* components of x and x_b are next *D*-*m* components
- Conditionals

 $p(\mathbf{x}_a | \mathbf{x}_b) = N(\mathbf{x} | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \text{ where } \boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$

• Marginals

$$p(x_a) = N(x_a \mid \mu_a, \Sigma_{aa}) \text{ where } \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

Joint $p(x_a, x_b)$



Marginal $p(x_a)$ and Conditional $p(x_a|x_b)$



If Marginals are Gaussian, Joint need not be Gaussian

- Constructing such a joint pdf:
 - Consider 2-D Gaussian, zero-mean uncorrelated



х

v

Due to symmetry about x- and x- axes, we can write marginals:

i.e., we only need to integrate over hatched regions

$$p_{X}(x) = \begin{cases} 2\int_{0}^{\infty} p_{XY}(x, y)dy, & x > 0\\ 2\int_{-\infty}^{0} p_{XY}(x, y)dy, & x \le 0 \end{cases}$$
$$p_{Y}(y) = \begin{cases} 2\int_{0}^{\infty} p_{XY}(x, y)dx, & y > 0\\ 2\int_{-\infty}^{0} p_{XY}(x, y)dx, & y \le 0 \end{cases}$$

 Take original 2-D Gaussian and set it to zero over non-hatched quadrants and multiply remaining by 2 we get a 2-D pdf that is definitely NOT Gaussian

srihari

Maximum Likelihood for the Gaussian

- Given a data set $X = (x_1, ..., x_N)^T$ where the observations $\{x_n\}$ are drawn independently
- Log-likelihood function is given by $\ln p(X \mid \mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$
- Derivative wrt μ is $\frac{\partial}{\partial \mu} \ln p(X | \mu, \Sigma) = \sum_{n=1}^{N} \Sigma^{-1}(\mathbf{x}_n - \mu)$
- Whose solution is $\left| \mu_{ML} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_{n} \right|$
- Maximization w.r.t. Σ is more involved. Yields

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

Bias of M. L. Estimate of Covariance Matrix

• For $N(\mu, \Sigma)$, m.l.e. of Σ for samples $\mathbf{x}_1, ..., \mathbf{x}_N$ is $\sum_{ML} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \mu_{ML}) (\mathbf{x}_n - \mu_{ML})^T$

• arithmetic average of N matrices: $[(\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T]$

• Since
$$E[\Sigma_{ML}] = \frac{1}{N-1} \sum_{n=1}^{N} (\mathbf{x}_n - \mu_{ML}) (\mathbf{x}_n - \mu_{ML})^T$$
 we have $E[\Sigma_{ML}] = \frac{N-1}{N} \Sigma$

– m.l.e. is smaller than the true value of Σ

- Thus m.l.e. is biased
 - irrespective of no of samples does not give exact value.
- For large N inconsequential.
 - Rule of thumb: use 1/N for known mean and 1/(N-1) for estimated mean.
- Bias does not exist in Bayesian solution.

Sequential Estimation

- In *on-line* applications and large data sets batch processing of all data points in infeasible
 - Real-time learning scenario where steady stream of data is arriving and predictions must be made before all data is seen
- Sequential methods allow data points to be processed one-at-a-time and then discarded
 - Sequential learning arises naturally with Bayesian viewpoint
- M.L.E. for parameters of Gaussian gives a convenient opportunity to discuss more general discussion of sequential estimation for maximum likelihood

Sequential Estimation of Gaussian Mean

By dissecting contribution of final data point

- Same as earlier batch result
- Nice interpretation:
 - After observing $\mathit{N-1}$ data points we have estimated $\,\mu\,by\,\,\mu_{ML}^{\mathit{N-1}}$
 - We now observe data point \boldsymbol{x}_{N} and we obtain revised estimate by moving old estimate by small amount
 - As N increases contribution from successive points smaller

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_{n}$$
$$= \mu_{ML}^{N-1} + \frac{1}{N} (\mathbf{x}_{N} - \mu_{ML}^{N-1})$$

General Sequential Estimation

- Sequential algorithms cannot always be factored out
- Robbins and Monro (1951) gave a general solution
- Consider pair of random variables θ and z with joint distribution $p(z, \theta)$
- Conditional expectation of *z* given θ is $f(\theta) = E[z | \theta] = \int zp(z | \theta) dz$



- Which is called a regression function
 - Same as one that minimizes expected squared loss seen earlier
- It can be shown that maximum likelihood solution is equivalent to finding the root of the regression function
 - Goal is to find θ^* at which $f(\theta^*)=0$

Robbins-Monro Algorithm

- Defines sequence of successive estimates of root θ^{\star} as follows

$$\boldsymbol{\theta}^{(N)} = \boldsymbol{\theta}^{(N-1)} + \boldsymbol{a}_{N-1} \boldsymbol{z}(\boldsymbol{\theta}^{(N-1)})$$

- Where $z(\theta^{(N)})$ is observed value of z when θ takes the value $\theta^{(N)}$
- Coefficients $\{a_N\}$ satisfy reasonable conditions

$$\lim_{N \to \infty} a_N = 0, \quad \sum_{N=1}^{\infty} a_N = \infty, \quad \sum_{N=1}^{\infty} {a_N}^2 < \infty$$

- Solution has a form where z involves a derivative of $p(x|\theta)$ wrt θ
- Special case of Robbons-Monro is solution for Gaussian mean

Bayesian Inference for the Gaussian

- MLE framework gives point estimates for parameters μ and Σ
- Bayesian treatment introduces prior distributions
 over parameters
- Case of known variance
- Likelihood of *N* observations $X = \{x_1, ..., x_N\}$ is

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

- Likelihood function is not a probability distribution over $\mu \, \text{and}$ is not normalized
- Note that likelihood function is quadratic in μ

Bayesian formulation for Gaussian mean

Likelihood function

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

- Note that likelihood function is quadratic in μ
- Thus if we choose a prior p(θ) which is Gaussian it will be a conjugate distribution for the likelihood because product of two exponentials will also be a Gaussian

$$p(\boldsymbol{\mu}) = N(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\sigma}_0^2)$$

16

Bayesian inference: Mean of Gaussian

Given Gaussian prior

 $p(\boldsymbol{\mu}) = N(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\sigma}_0^2) -$

- Posterior is given by $p(\mu|X) \ \alpha \ p(X|\mu)p(\mu)$
- Simplifies to

 $\frac{1}{\sigma_{\rm w}^2} = \frac{1}{\sigma_{\rm o}^2} + \frac{1}{\sigma^2}$

$$P(\mu|X) = N(\mu|\mu_N, \sigma_N^2)$$
 where

Prior and posterior have same form: conjugacy

N = 10

N=2

N = 1

Data points from mean=0.8 and known variance=0.1

 $\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML}$

If N=0 reduces to prior mean If $N \rightarrow \infty$ posterior mean is ML solution

N = 0

-1

Precision is additive: sum of precision of prior plus one contribution of data precision from each observed data point 17

Bayesian Inference of the Variance

- Known Mean
- Wish to infer variance
- Analysis simplified if we choose conjugate form for prior distribution
- Likelihood function with precision $\lambda = 1/\sigma^2$

$$p(X \mid \lambda) = \prod_{n=1}^{N} N(x_n \mid \mu, \lambda - 1)$$
$$\alpha \quad \lambda^{N/2} \exp\left\{-\frac{\lambda}{2} \sum_{n=1}^{N} (x - \mu)^2\right\}$$

Conjugate prior is given by Gamma distribution

Gamma Distribution

- Gaussian with known mean but unknown variance
- Conjugate prior for the precision of a Gaussian is given by a Gamma distribution

– Precision $I = 1/\sigma^2$

$$Gam(\lambda \mid a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \quad \Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$$

– Mean and Variance

$$E[\lambda] = \frac{a}{b}, \quad \operatorname{var}[\lambda] = \frac{a}{b^2}$$

Gamma distribution $Gam(\lambda|a,b)$ for various values of *a* and *b*





19

Gamma Distribution Inference

- Given prior distribution $Gam(\lambda | a_0, b_0)$
- Multiplying by likelihood function
- The posterior distribution has the form $Gam(\lambda | a_N, b_N)$ where

$$a_N = a_0 + \frac{N}{2}$$

Effect of *N* observations is to increase *a* by N/2Interpret a_0 as $2a_0$ effective prior observations

$$b_{N} = b_{0} + \frac{1}{2} \sum_{n=1}^{N} (x_{n} - \mu)^{2}$$
$$= b_{0} + \frac{N}{2} \sigma_{ML}^{2}$$

Both Mean and Variance are Unknown

- Consider dependence of likelihood function on μ and λ

$$p(X \mid \mu, \lambda) = \prod_{n=1}^{N} \left(\frac{\lambda}{2\pi}\right)^{1/2} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\}$$

- Identify a prior distribution $p(\mu, \lambda)$ that has same functional dependence on μ and λ as likelihood function
- Normalized prior takes the form $p(\mu, \lambda) = N(\mu \mid \mu_0, (\beta \lambda)^{-1})Gam(\lambda \mid a, b)$
 - $p(\mu, \lambda) N(\mu | \mu_0, (p\lambda))$ $fram(\lambda | u, b)$
 - Called *normal-gamma* or *Gaussian-gamma* distribution

Normal Gamma

- Both mean and precision unknown
- Contour plot with

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



Estimation for Multivariate Case

- For a *multivariate Gaussian* distribution $N(\mathbf{x} | \mu, \Lambda^{-1})$ for a *D*-dimensional variable \mathbf{x}
 - Conjugate prior for mean µ assuming known precision is Gaussian
 - For known mean and *unknown precision* matrix Λ , conjugate prior is *Wishart* distribution
 - If both mean and precision are unknown conjugate prior is Gaussian-Wishart

Machine Learning

Student's t-distribution

- Conjugate prior for precision of Gaussian is given by Gamma
- If we have a univariate Gaussian N(x|μ, τ⁻¹) together with Gamma prior Gam(τ|a,b) and we integrate out the precision we obtain marginal distribution of x

$$p(x \mid \mu, a, b) = \int_{0}^{0} N(x \mid \mu, \tau^{-1}) Gam(\tau \mid a, b) d\tau$$

Has the form

$$St(x \mid \mu, \lambda, v) = \frac{\Gamma(v \mid 2 + 1 \mid 2)}{\Gamma(v \mid 2)} \left(\frac{\lambda}{\pi v}\right)^{1/2} \left[1 + \frac{\lambda(x - \mu)^2}{v}\right]^{-v/2 - 1/2}$$



• Parameter v=2a is called *degrees of freedom and* $\lambda=a/b$ as the *precision* of the t distribution

 $v \rightarrow \infty$ becomes Gaussian

- Infinite mixture of Gaussians with same mean but different precisions
 - Obtained by adding many Gaussian distributions
 - Result has longer tails than Gaussian

Robustness of Student's t

- Has longer tails than Gaussian
- ML solution can be found by expectation maximization algorithm
- Effect of small no of outliers much less on t distributions
- Can also obtain multivariate
 form of t-distribution



Periodic Variables

- Gaussian inappropriate for continuous variables that are periodic or angular
 - Wind direction on several days
 - Calendar time
 - Fingerprint minutiae direction
- If we choose standard Gaussian, results depend on choice of origin
 - With 0^* as origin two observations $\theta_1 = 1^*$ and $\theta_2 = 359^*$ will have mean at 180^* and std dev 179^*
 - With 180^* as origin mean= 0^* , std dev= 1^*
- Quantity represented by polar coordinates $0 \le \theta \le 2\rho$

Conversion to Polar Coords

- Observations $D = \{\theta_1, \dots, \theta_N\}, \theta$ measured in radians
- Viewed as points on unit circle
 - Represent data as 2-d vectors $x_1, ..., x_N$ where $||x_n|| = 1$

- Mean
$$\overline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n$$

- Cartesian coords: $x_n = (\cos \theta_n, \sin \theta_n)$
- Coords of sample mean $\overline{\mathbf{x}} = (\overline{r}\cos\overline{\theta}, \overline{r}\sin\overline{\theta})$
- Simplifying and solving gives $\overline{\theta} = \tan^{-1} \left\{ \frac{\sum_{n} \sin \theta_{n}}{\nabla \cos \theta_{n}} \right\}$





Von Mises Distribution

- Periodic generalization of Gaussian

 p(θ) satisfies
 p(θ) ≥ 0, ∫₀^{2π} *p*(θ)*d*θ = 1, *p*(θ + 2π) = *p*(θ)

 Consider 2-d Gaussian x = (x₁,x₂), μ = (μ₁,μ₂), Σ = σ²I
- **Consider 2-d Gaussian** $\stackrel{0}{x} = (x_1, x_2), \ \mu = (\mu_1, \mu_2), \ \Sigma = \sigma^2 I$ $p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2}\right\}$
 - Contours of constant density are circles
- Transforming to polar coordinates (r, θ)
 - $x_1 = r \cos \theta$, $x_2 = r \sin \theta$ and $\mu_1 = r_0 \cos \theta_0$, $\mu_2 = r_0 \sin \theta_0$



• Defining $\mu = r_0 / s^2$ distribution of $p(\theta)$ along unit circle $p(\theta | \theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\{m \cos(\theta - \theta_0)\}$ Von Mises distribution

Where $I_0(m)$ is the Bessel Function $I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp\{m\cos\theta\} d\theta$

28

 $I_0(m)$







Two different parameter values

For large μ distribution Is approximately Gaussian

Zeroth-order Bessel Function of the first kind



ML estimates of von Mises parameters

- Parameters are θ_0 and m
- Log-likelihood function $\ln p(D \mid \theta_0, m) = -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^{N} \cos(\theta_n - \theta_0)$
- Setting derivative wrt q_0 equal to zero gives

$$\theta_0^{ML} = \tan^{-1} \left\{ \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right\}$$

• Maximizing wrt μ gives solution for $A(\mu_{ML})$ which can be inverted to get μ



Mixtures of Gaussians

- Gaussian has limitations in modeling real data sets
- Old Faithful (Hydrothermal Geyser in Yellowstone)
 - 272 observations
 - Duration (mins, horiz axis) vs Time to next eruption (vertical axis)
 - Simple Gaussian unable to capture structure
 - Linear superposition of two Gaussians is 100 better
- Linear combinations of Gaussians can give very complex densities

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k N(x \mid \mu_k, \Sigma_k)$$

 π_k are mixing coefficients that sum to one

- One –dimension
 - Three Gaussians in blue
 - Sum in red









srihari

Mixture of 3 Gaussians

- Contours of constant density for components
- Contours of mixture density
 p(x)
- Surface plot of distribution
 p(x)



Estimation for Gaussian Mixtures

Log likelihood function is

$$\ln p(X \mid \pi, \mu, \Sigma) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k N(\mathbf{x}_n \mid \mu_k, \Sigma_k) \right\}$$

- Situation is more complex
- No closed form solution
- Use either iterative numerical optimization techniques or *Expectation Maximization*