

Gaussian Distribution

Sargur N. Srihari

The Gaussian Distribution

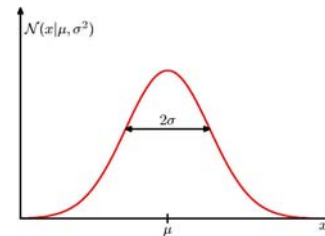
- For single real-valued variable x

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

- Parameters:

– Mean μ , variance σ^2 ,

- Standard deviation σ



Carl Friedrich Gauss
1777-1855

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

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

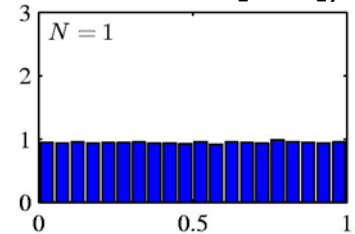
$\boldsymbol{\mu}$ is a mean vector, $\boldsymbol{\Sigma}$ is a $D \times D$ covariance matrix, $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$

$\boldsymbol{\Sigma}^{-1}$ is also referred to as the precision matrix

Importance of Gaussian

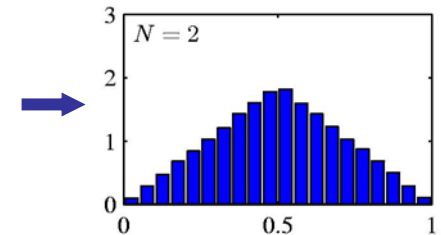
- Gaussian arises in many different contexts, e.g.,
 - For a single variable, Gaussian maximizes entropy (for given mean and variance)
 - Sum of set of random variables becomes increasingly Gaussian

One variable histogram
(uniform over $[0,1]$)

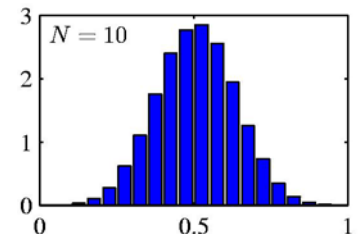


Mean of two variables

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 \mathbf{x} is through

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

- Called Mahalanobis Distance
- reduces to Euclidean distance when $\boldsymbol{\Sigma}$ is an identity matrix

- Matrix $\boldsymbol{\Sigma}$ is symmetric

- Has an Eigenvector equation

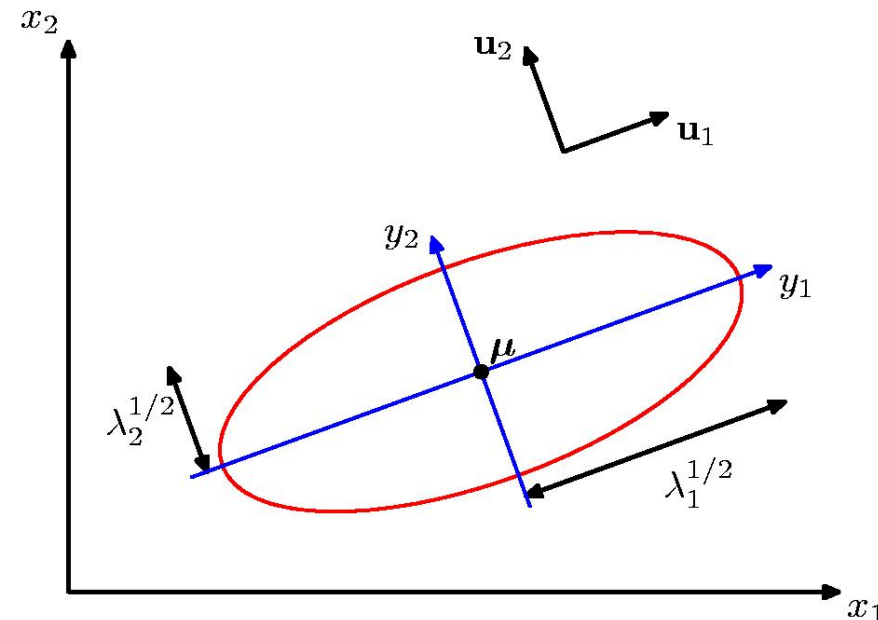
$$\boldsymbol{\Sigma} \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

\mathbf{u}_i are Eigen vectors

λ_i are Eigen values

Two dimensional Gaussian

$$\mathbf{x} = (x_1, x_2)$$

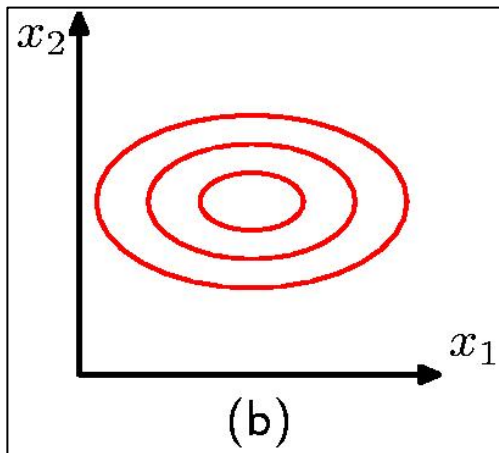


Red: Elliptical contour of constant density
Major axes: eigenvectors \mathbf{u}_i

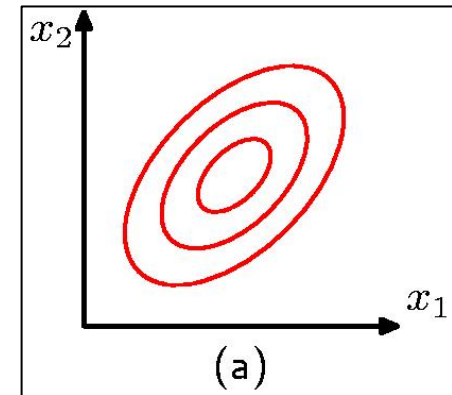
Contours of Constant Density

- Determined by Covariance Matrix
 - Covariances represent how features vary together

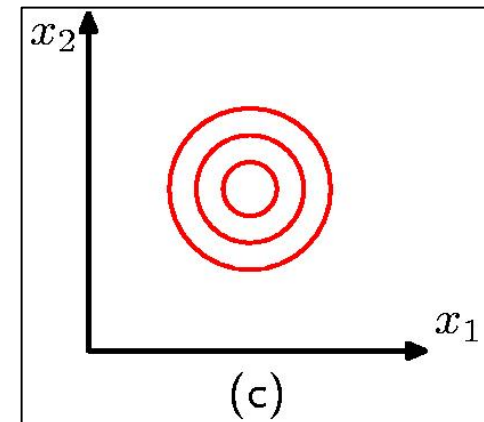
(b) Diagonal matrix
(aligned with coordinate axes)



(a) General form



(c) Proportional to Identity matrix
(concentric circles)



Joint, Marginal and Conditional with Gaussian

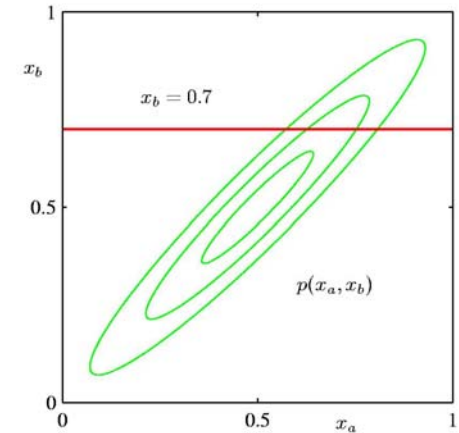
- If two sets of variables $\mathbf{x}_a, \mathbf{x}_b$ are jointly Gaussian then the two conditional densities and the two marginals are also Gaussian
- Given joint Gaussian $N(\mathbf{x}|\mu, \Sigma)$ with $\Lambda = \Sigma^{-1}$ and $\mathbf{x} = [\mathbf{x}_a, \mathbf{x}_b]^T$ where \mathbf{x}_a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

$$p(\mathbf{x}_a | \mathbf{x}_b) = N(\mathbf{x}_a | \mu_{a|b}, \Lambda_{aa}^{-1}) \text{ where } \mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b)$$

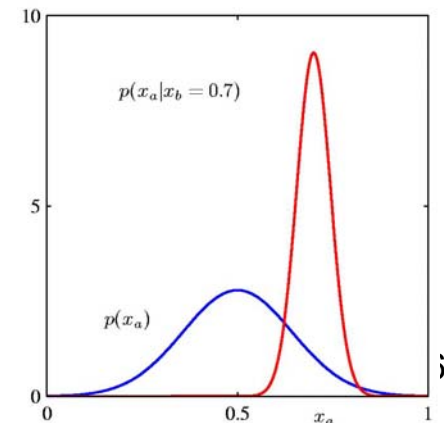
- Marginals

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

Joint $p(x_a, x_b)$



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



Maximum Likelihood for the Gaussian

- Given a data set $X=(x_1, \dots, x_N)^T$ where the observations $\{x_n\}$ are drawn independently
- Log-likelihood function is given by

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

- 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

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Maximization w.r.t. Σ is more involved. Yields

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

Sequential Estimation

- In *on-line* applications and large data sets batch processing of all data points is 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

$$\begin{aligned}\mu_{ML} &= \frac{1}{N} \sum_{n=1}^N x_n \\ &= \mu_{ML}^{N-1} + \frac{1}{N} (x_N - \mu_{ML}^{N-1})\end{aligned}$$

- Same as earlier batch result

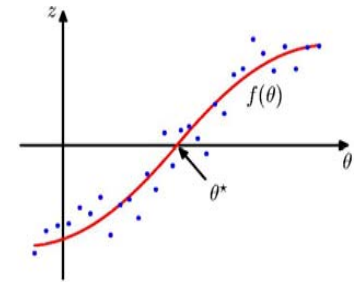
- Nice interpretation:

- After observing $N-1$ data points we have estimated μ by μ_{ML}^{N-1}
- We now observe data point x_N and we obtain revised estimate by moving old estimate by small amount
- As N increases contribution from successive points smaller

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^* as follows

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

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

- Coefficients $\{a_N\}$ satisfy reasonable conditions

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

- Solution has a form where z involves a derivative of $p(x|\theta)$ wrt θ

- Special case of Robbins-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 $X = \{x_1, \dots, x_N\}$ is

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

- Likelihood function is not a probability distribution over μ and is not normalized
- Note that likelihood function is quadratic in μ

Bayesian formulation for Gaussian mean

- Likelihood function

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

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

$$p(\mu) = N(\mu | \mu_0, \sigma_0^2)$$

Bayesian inference: Mean of Gaussian

- Given Gaussian prior

$$p(\mu) = N(\mu | \mu_0, \sigma_0^2)$$

- Posterior is given by

$$p(\mu | X) \propto p(X | \mu) p(\mu)$$

- Simplifies to

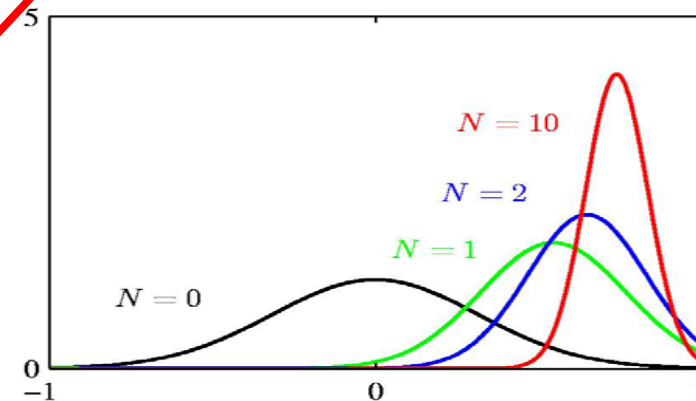
$$P(\mu | X) = N(\mu | \mu_N, \sigma_N^2) \text{ where}$$

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

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}$$

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

Prior and posterior have same form: conjugacy



Data points from mean=0.8 and known variance=0.1

If $N=0$ reduces to prior mean

If $N \rightarrow \infty$ posterior mean is ML solution

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 $\lambda = 1/\sigma^2$

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

- 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

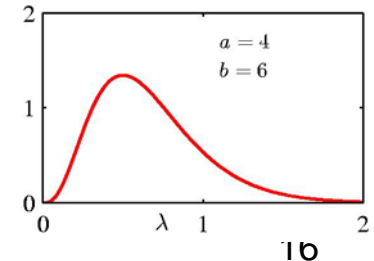
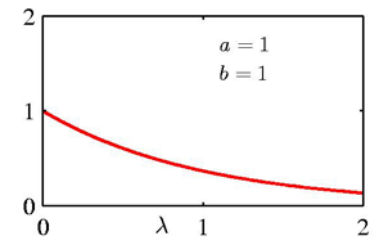
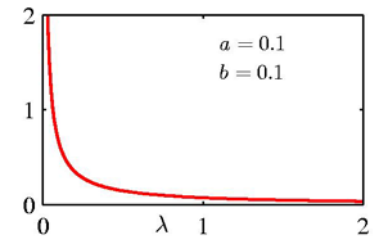
– Precision $\lambda = 1/\sigma^2$

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

– Mean and Variance

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

Gamma distribution
 $\text{Gam}(\lambda/a, b)$ for various values of a and b



Gamma Distribution Inference

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

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

Effect of N observations is to increase a by $N/2$
Interpret 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

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

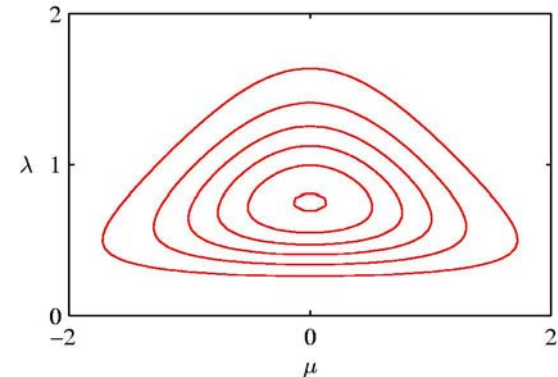
- Identify a prior distribution $p(\mu, \lambda)$ that has same functional dependence on μ and λ as likelihood function
- Normalized prior takes the form

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

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

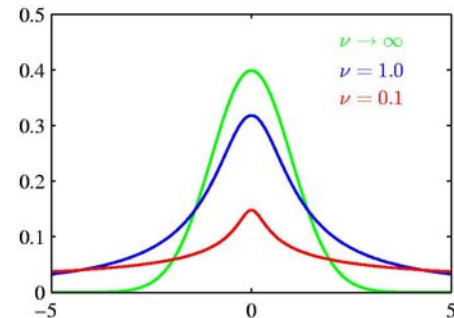
Student's t-distribution

- Conjugate prior for precision of Gaussian is given by Gamma
- If we have a univariate Gaussian $N(x|\mu, \tau^{-1})$ together with Gamma prior $Gam(\tau|a, b)$ and we integrate out the precision we obtain marginal distribution of x

$$p(x | \mu, a, b) = \int_0^{\infty} N(x | \mu, \tau^{-1}) Gam(\tau | a, b) d\tau$$

- Has the form

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



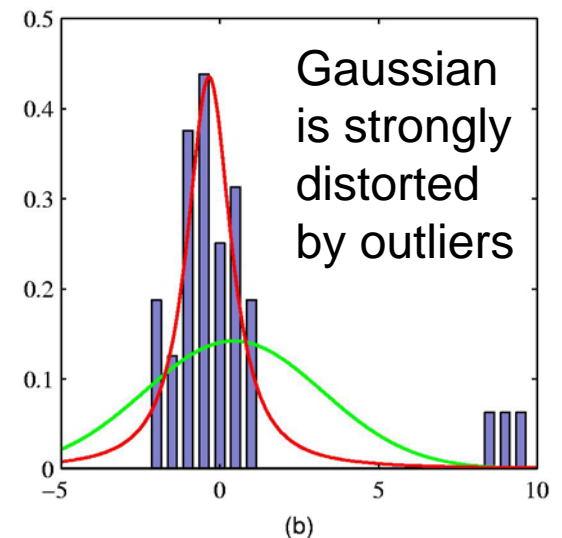
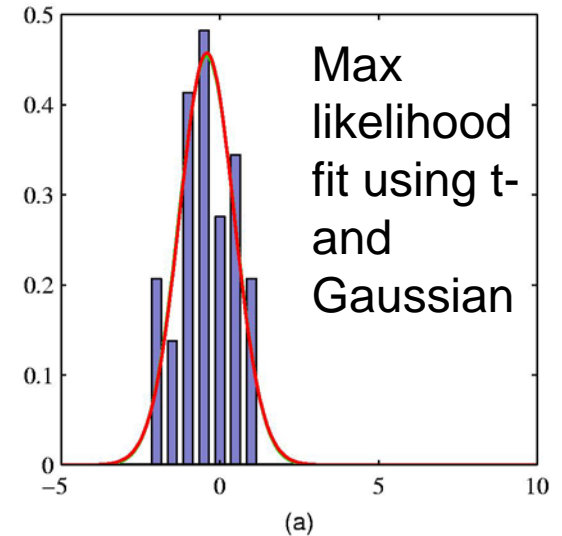
- Parameter $\nu=2a$ is called *degrees of freedom* and $\lambda=a/b$ as the *precision* of the t distribution

$\nu \rightarrow \infty$ 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 $\theta_1=1^\circ$ and $\theta_2=359^\circ$ will have mean at 180° and std dev 179°
 - With 180° as origin mean= 0° , std dev= 1°
- Quantity represented by polar coordinates $0 \leq \theta \leq 2\pi$

Conversion to Polar Coords

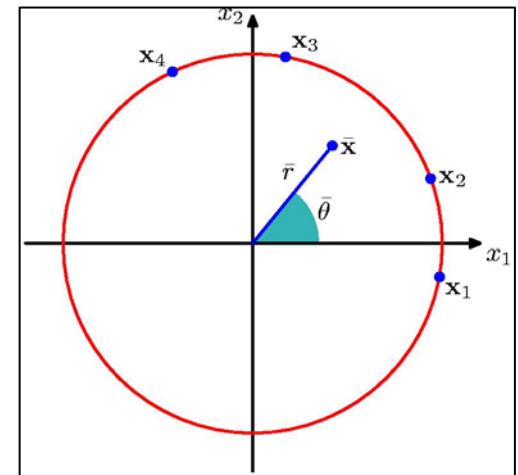
- Observations $D = \{\theta_1, \dots, \theta_N\}$, θ measured in radians
- Viewed as points on unit circle
 - Represent data as 2-d vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ where $\|\mathbf{x}_n\| = 1$

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

- Cartesian coords: $\mathbf{x}_n = (\cos \theta_n, \sin \theta_n)$

- Coords of sample mean

$$\bar{\mathbf{x}} = (\bar{r} \cos \bar{\theta}, \bar{r} \sin \bar{\theta})$$



- Simplifying and solving gives $\bar{\theta} = \tan^{-1} \left\{ \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right\}$

Von Mises Distribution

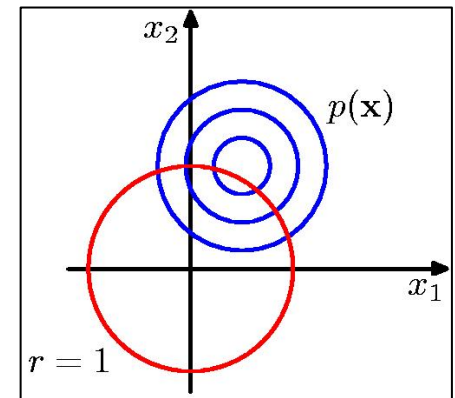
- Periodic generalization of Gaussian

– $p(\theta)$ satisfies $p(\theta) \geq 0$, $\int_0^{2\pi} p(\theta) d\theta = 1$, $p(\theta + 2\pi) = p(\theta)$

- Consider 2-d Gaussian $\mathbf{x} = (x_1, x_2)$, $\boldsymbol{\mu} = (\mu_1, \mu_2)$, $\boldsymbol{\Sigma} = \sigma^2 I$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2}\right\}$$

– Contours of constant density are circles



- Transforming to polar coordinates (r, θ)

$$x_1 = r \cos \theta, \quad x_2 = r \sin \theta \quad \text{and} \quad \mu_1 = r_0 \cos \theta_0, \quad \mu_2 = r_0 \sin \theta_0$$

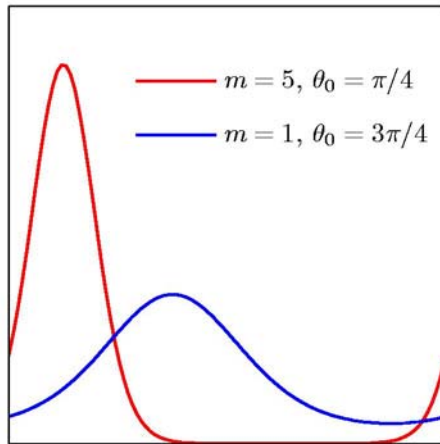
- Defining $m = r_0 / \sigma^2$ distribution of $p(\theta)$ along unit circle

$$p(\theta | \theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\{m \cos(\theta - \theta_0)\} \quad \text{Von Mises distribution}$$

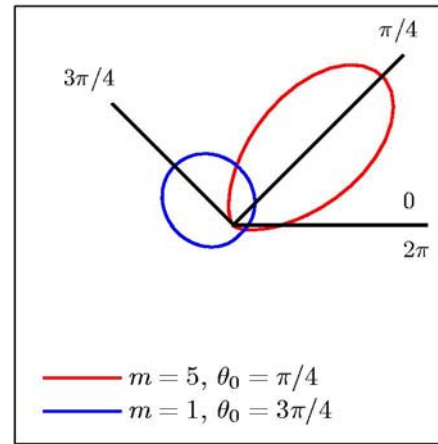
Where $I_0(m)$ is the Bessel Function $I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp\{m \cos \theta\} d\theta$

Von Mises Plots

Cartesian Plot



Polar Plot

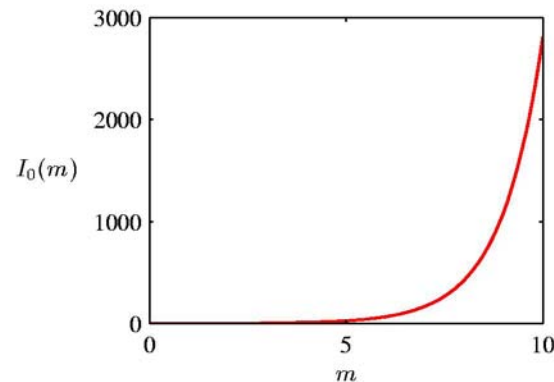


Two different
parameter values

For large m distribution
Is approximately Gaussian

- Zeroth-order Bessel Function of the first kind

$$I_0(m)$$



ML estimates of von Mises parameters

- Parameters are θ_0 and m

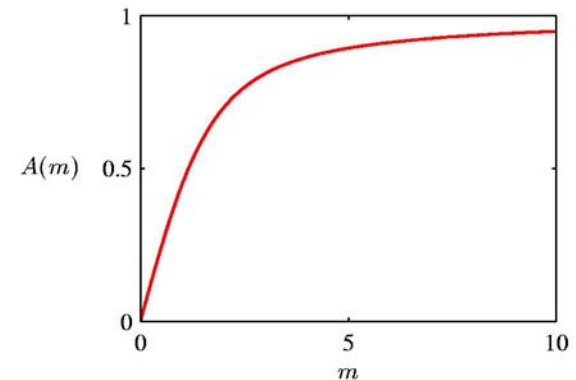
- Log-likelihood function

$$\ln p(D | \theta_0, m) = -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^N \cos(\theta_n - \theta_0)$$

- Setting derivative wrt θ_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 m gives solution for $A(m_{ML})$ which can be inverted to get m



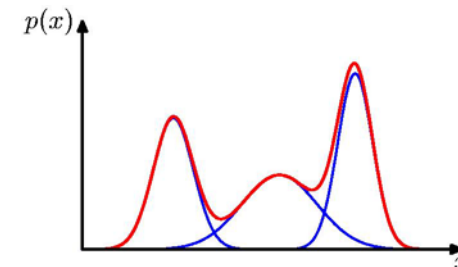
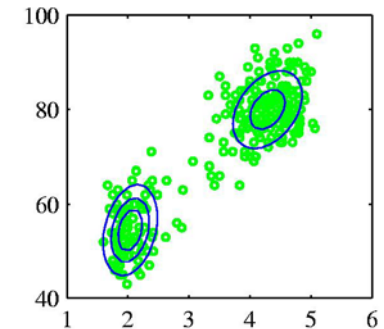
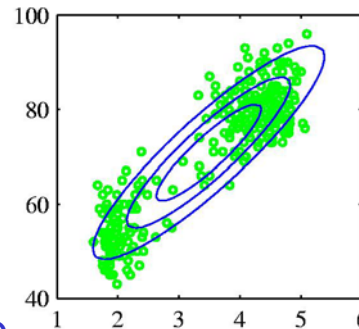
Mixtures of Gaussians

- 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 better
- Linear combinations of Gaussians can give very complex densities

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

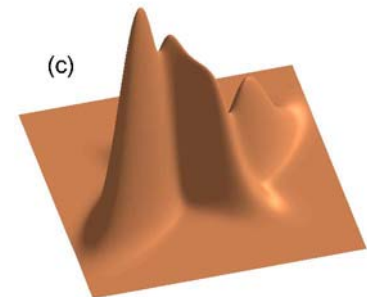
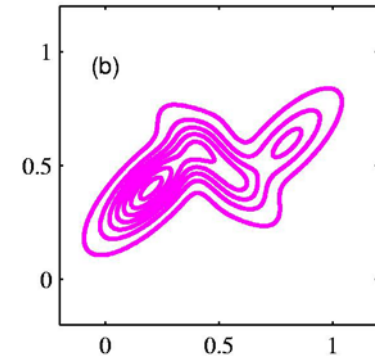
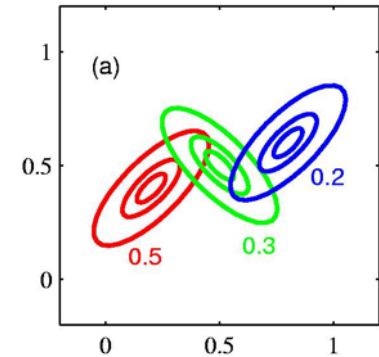
π_k are mixing coefficients that sum to one

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



Mixture of 3 Gaussians

- Contours of constant density for components
- Contours of mixture density $p(\mathbf{x})$
- Surface plot of distribution $p(\mathbf{x})$



Estimation for Gaussian Mixtures

- Log likelihood function is

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

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