

Probability Distributions

(30)

Prob. distributions form the building blocks for more complex models.

Density Estimation: model $p(x)$ given a finite set $\{x_1, \dots, x_N\}$ of observations. This is an ill-posed problem, since any $p(x)$ which is non-zero at each data point is a potential candidate.

Parametric: the distr. is governed by a small number of adaptive parameters. Ex.: Gaussian depends on mean μ , and covariance Σ . One limitation is that assumes a specific functional form.

Nonparametric: the form of the distr. depends on the size of the data set. The parameters control the model complexity rather than the form of the prob. distr.

Binary Variables

$$x \in \{0, 1\}. \quad p(x=1|\mu) = \mu \quad 0 \leq \mu \leq 1$$

$$p(x=0|\mu) = 1 - \mu$$

$$\boxed{\text{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}}$$

Bernoulli Distr.

$$\begin{aligned} E[X] &= p(x=1|\mu) \cdot 1 + p(x=0|\mu) \cdot 0 \\ &= \mu \end{aligned}$$

$$\text{Var}[X] = E[X^2] - E[X]^2 = \mu - \mu^2 = \mu(1-\mu)$$

Suppose we have data $D = \{x_1, \dots, x_N\}$. The likelihood function is

$$P(D|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n}$$

In the frequentist approach we maximize the likelihood, or equivalently, its logarithm:

$$\log P(D|\mu) = \sum_{n=1}^N x_n \log \mu + (1-x_n) \log(1-\mu)$$

(31)

$$\begin{aligned} \frac{\partial \log P}{\partial \mu} = 0 &= \sum_{n=1}^N \left\{ \frac{x_n}{\mu} + \frac{(1-x_n)(-1)}{1-\mu} \right\} \\ &= \sum_{n=1}^N \frac{x_n - \cancel{\mu x_n} - \cancel{\mu} + \cancel{x_n \mu}}{\mu(1-\mu)} \end{aligned}$$

$$\boxed{\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n} \quad \text{or} \quad \boxed{\mu_{ML} = \frac{m}{N}} \quad \begin{array}{l} \text{if } m \\ \text{is the number} \\ \text{of } x=1 \text{ occurrences.} \end{array}$$

What is the distribution for the number of $x=1$ observations, denoted m ? It is proportional to $\text{Bern}(x|\mu)$, and of N possibilities, we can get m times $x=1$ in $\binom{N}{m}$ distinct ways. Thus the answer is

$$\boxed{\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}} \quad \text{Binomial Distr.}$$

$$m = x_1 + x_2 + \dots + x_N$$

$$E[m] = N E[x] = N\mu \quad \text{since they are Bernoulli iid.}$$

Another way:

$$E[m] = \sum_{m=0}^N m \frac{N!}{(N-m)! m!} \mu^m (1-\mu)^{N-m}$$

$$= \sum_{m=1}^N \frac{N!}{(N-m)! (m-1)!} \mu^m (1-\mu)^{N-m}$$

$$m \rightarrow m+1 = \sum_{m=0}^{N-1} \frac{N!}{(N-1-m)! m!} \mu \mu^m (1-\mu)^{N-1-m}$$

$$N \rightarrow N+1 = \mu N \underbrace{\sum_{m=0}^N \frac{N!}{(N-m)! m!} \mu^m (1-\mu)^{N-m}}_{=1} = \mu N //$$

$$\begin{aligned}
\text{Var}[m] &= E[(m - E[m])^2] \\
&= E\left[\left(\sum_{i=1}^N x_i - N\mu\right)^2\right] = E\left[\left(\sum_{i=1}^N (x_i - \mu)\right)^2\right] \\
&= E\left[\sum_{i,j=1}^N (x_i - \mu)(x_j - \mu)\right] \\
&= E\left[\sum_{i=1}^N (x_i - \mu)^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^N (x_i - \mu)(x_j - \mu)\right] \xrightarrow{\text{iid}} \\
&= \sum_{i=1}^N E[(x_i - \mu)^2] + \sum_{\substack{i,j=1 \\ i \neq j}}^N E(x_i - \mu) E(x_j - \mu) \\
&= N \text{Var}[x] + 0 \\
&= N\mu(1-\mu)
\end{aligned}$$

Another way:

$$\begin{aligned}
E[m^2] &= \sum_{m=0}^N m^2 \frac{N!}{(N-m)! m!} \mu^m (1-\mu)^{N-m} \\
&= \sum_{m=1}^N \frac{N! m}{(N-m)! (m-1)!} \mu^m (1-\mu)^{N-m} \\
m \rightarrow m+1 &= \sum_{m=0}^{N-1} \frac{N! (m+1)}{(N-1-m)! m!} \mu \mu^m (1-\mu)^{N-1-m} \\
N-1 \rightarrow N' &= N\mu \sum_{m=0}^{N-1} \frac{(N-1)! (m+1)}{(N-1-m)! m!} \mu^m (1-\mu)^{N-1-m} \\
&= N\mu \left\{ \sum_{m=0}^{N'} \frac{N'! (m+1)}{(N'-m)! m!} \mu^m (1-\mu)^{N'-m} \right\} \\
&= N\mu \left\{ \sum_{m=0}^{N'} \frac{N'!}{(N'-m)! m!} \mu^m (1-\mu)^{N'-m} \right\} = 1 \\
&\quad + \sum_{m=0}^{N'} \frac{m N'!}{(N'-m)! m!} \mu^m (1-\mu)^{N'-m} \left\{ \right. \\
&\quad \left. \right\} \mu N' = \mu(N-1)
\end{aligned}$$

$$E[m^2] = N\mu(1 + \mu(N-1))$$

$$\begin{aligned} \text{Var}[m] &= E[m^2] - (E[m])^2 = N\mu + N\mu^2(N-1) \\ &= N\mu + \cancel{N^2\mu^2} - \cancel{N\mu^2} - \cancel{N^2\mu^2} \\ &= \underline{\underline{N\mu(1-\mu)}} \end{aligned}$$

The Beta Distribution

We saw that ML overfits for small N . In a Bayesian treatment we need a prior $p(\mu)$ such that

$$p(\mu|D) = \frac{p(D|\mu) p(\mu)}{p(D)}$$

In the Bernoulli case, if $p(\mu) \sim \mu^a (1-\mu)^b$ the prior and the posterior will have the same form. This is known as conjugacy. So we introduce

$$\boxed{\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}}$$

One can check $\int_0^1 \text{Beta}(\mu|a, b) d\mu = 1$.

Now

$$\begin{aligned} E[\mu] &= \left(\int_0^1 d\mu \mu^a (1-\mu)^{b-1} \right) \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \\ &= \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+1+b)} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rightarrow \Gamma(n+1) = n\Gamma(n) \\ &= \frac{a \cancel{\Gamma(a)} \cancel{\Gamma(b)} \Gamma(a+b)}{(a+b) \cancel{\Gamma(a+b)} \Gamma(a) \Gamma(b)} \\ &= \underline{\underline{\frac{a}{a+b}}} \end{aligned}$$

$$E[\mu^2] = \int_0^1 d\mu \mu^{a+1} (1-\mu)^{b-1} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$= \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$= \frac{(a+1)a \cancel{\Gamma(a)} \cancel{\Gamma(b)} \Gamma(a+b)}{(a+b+1)(a+b) \cancel{\Gamma(a+b)} \cancel{\Gamma(a)} \cancel{\Gamma(b)}}$$

$$= \frac{a(a+1)}{(a+b)(a+b+1)}$$

$$\therefore \text{Var}[\mu] = \frac{a(a+1)}{(a+b)(a+b+1)} - \frac{a^2}{(a+b)^2}$$

$$= \frac{a(a+1)(a+b) - a^2(a+b+1)}{(a+b)^2(a+b+1)}$$

$$= \frac{(a^2+a)(a+b) - a^3 - a^2b - a^2}{(a+b)^2(a+b+1)}$$

$$= \frac{\cancel{a^3} + \cancel{a^2b} + \cancel{a^2} + \cancel{ab} - \cancel{a^3} - \cancel{a^2b} - \cancel{a^2}}{(a+b)^2(a+b+1)}$$

$$= \frac{ab}{(a+b)^2(a+b+1)}$$

(a, b) are called Hyperparameters because they control the distribution of the parameter μ .

The likelihood function for Bernoulli is

$$P(D|\mu) = \prod_{i=1}^N \mu^{x_i} (1-\mu)^{1-x_i} = \mu^{\sum_{i=1}^N x_i} (1-\mu)^{N - \sum_{i=1}^N x_i}$$

if $m = \sum_{i=1}^N x_i$ is the # times we get $x=1$,
and $l = N - m$ the # times we get $x=0$, then

$$P(D|\mu) = \mu^m (1-\mu)^l = p(m, l | \mu)$$

The posterior is the

$$P(\mu | m, l, a, b) \propto p(m, l | \mu) p(\mu | a, b) \\ = \mu^{m+a-1} (1-\mu)^{l+b-1}$$

which is simply another Beta distr. Introducing the normalization:

$$P(\mu | m, l, a, b) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a) \Gamma(l+b)} \mu^{m+a-1} (1-\mu)^{l+b-1}$$

Note that this exactly as the original Beta with $a \rightarrow a+m$, $b \rightarrow b+l$. So the prior corresponds to an effective number of observations of $x=1$, and $x=0$, and when more data is accessible we get the posterior, which is just an update of our current information. This is a sequential approach to "learning".

Only assumption: iid data.

This is useful where we have a real-time stream of data, for instance.

If we want to predict the outcome of the next trial we move from the sum rule:

$$P(x=1 | D) = \int d\mu P(x=1 | \mu) P(\mu | D) \\ = \int d\mu \mu P(\mu | D) = E[\mu | D] \\ = \frac{m+a}{m+a+l+b}$$

When $m, l \rightarrow \infty$, this reproduces μ_{ML} from (36) maximum likelihood. It's a general property that ML and Bayes agree in the limit of infinitely large data.

When $a, b \rightarrow \infty$, $\text{Var}[\mu] \rightarrow 0$. Is it a general property of Bayesian learning that, as we observe more data, the uncertainty contained in the posterior will decrease? Let us consider the following:

$$\begin{aligned} \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]] &= \mathbb{E}_D \left\{ \int d\theta p(\theta|D) \theta \right\} \\ &= \int dD p(D) \left\{ \int d\theta p(\theta|D) \theta \right\} \\ &= \int d\theta \theta \left\{ \int dD p(\theta|D) p(D) \right\} \\ &= \int d\theta \theta p(\theta) = \mathbb{E}_\theta[\theta] \end{aligned}$$

The posterior mean, averaged over the data, is equal to the prior. In other words, the mean of the posterior is the same as the prior, considering the distr. generating the data.

Now let's consider the variance. First the mean, with respect to the data, of the posterior variance:

$$\begin{aligned} \mathbb{E}_D[\text{Var}_\theta[\theta|D]] &= \mathbb{E}_D \left\{ \mathbb{E}_\theta[\theta^2|D] - (\mathbb{E}_\theta[\theta|D])^2 \right\} \\ &= \mathbb{E}_D[\mathbb{E}_\theta[\theta^2|D]] - \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]^2] \\ &= \mathbb{E}_\theta[\theta^2] - \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]^2] \end{aligned}$$

Now consider the variance of the posterior mean:

$$\begin{aligned} \text{Var}_D[\mathbb{E}_\theta[\theta|D]] &= \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]^2] - \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]]^2 \\ &= \mathbb{E}_D[\mathbb{E}_\theta[\theta|D]^2] - \mathbb{E}_\theta[\theta]^2 \end{aligned}$$

Thus summing these results we get

(37)

$$E_D[\text{Var}_\theta[\theta|D]] + \text{Var}_D[E_\theta[\theta|D]] = \text{Var}_\theta[\theta]$$

Since $\text{Var}_\theta[\theta] > 0$, it means that, on average, the posterior variance decreases.

Multinomial Variables

Consider a variable that can take K distinct states. One way to represent a state is through the vector $x \in \mathbb{R}^K$ such that $x_k = 1$ for only one entry, and $x_j = 0$ for $j \neq k$.

Ex.: $x = (0, 0, 1, 0, 0, \dots)^T$

Note that $\sum_{k=1}^K x_k = 1$. Let $p(x_k = 1) = \mu_k$, then

the distribution of x is

$$P(x|\mu) = \prod_{k=1}^K \mu_k^{x_k} \quad \mu = (\mu_1, \dots, \mu_K)^T.$$

$\mu_k \geq 0, \sum_k \mu_k = 1.$

Analogous to Bernoulli, but for K possibilities.

$\sum_x p(x|\mu) = \sum_k \mu_k = 1$, since for each vector x we have $p(x|\mu) = \mu_k$ for some k . Summing over all possibilities gives the result. By analogous reasoning

$$E[x|\mu] = \sum_x p(x|\mu) x = (\mu_1, \dots, \mu_K)^T = \mu$$

Now consider data $D = \{x_1, \dots, x_N\}$. The likelihood function is

$$P(D|\mu) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{\left(\sum_{n=1}^N x_{nk}\right)} = \prod_{k=1}^K \mu_k^{m_k}$$

where $m_k = \sum_{n=1}^N x_{nk}$ is the number of observations (38)
 where $x_k = 1$. Maximizing this we have

$$l = \sum_{k=1}^K m_k \log \mu_k + d \left(\sum_{k=1}^K \mu_k - 1 \right)$$

$$\frac{\partial l}{\partial \mu_j} = 0 = \frac{m_j}{\mu_j} + d \therefore \mu_j = -\frac{m_j}{d}$$

$$\sum_{k=1}^K \mu_k = 1 = \frac{1}{d} \sum_{k=1}^K m_k = -\frac{N}{d} \therefore d = -\frac{1}{N} \therefore \boxed{\mu_k = \frac{m_k}{N}}$$

which is the fraction where $x_k = 1$ occurs.

For the joint distribution of m_k 's we have:

$$\text{Mult}(m_1, \dots, m_K | \mu, N) = \binom{N}{m_1, m_2, \dots, m_K} \prod_{k=1}^K \mu_k^{m_k}$$

which is the multinomial distr.

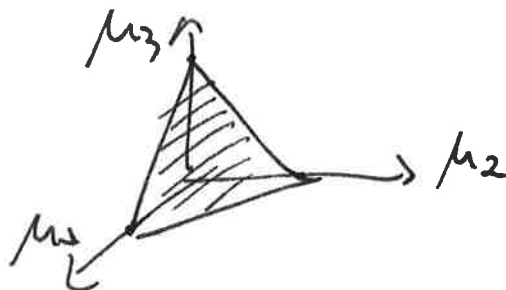
The Dirichlet Distribution

Prior for $\{\mu_k\}$. The conjugate prior is

$$p(\mu | \alpha) \propto \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

$\alpha = (\alpha_1, \dots, \alpha_K)^T$ are the hyperparameters. Because of $\sum_{k=1}^K \mu_k = 1$, the distr. is confined to a

Simplex (bounded linear manifold) of dimension $K-1$.



The normalized form is

(39)

$$\text{Dir}(\mu|x) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

Dirichlet
Distr.

$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

Now the posterior has the form

$$\begin{aligned} P(\mu|D, \alpha) &\propto P(D|\mu) P(\mu|\alpha) \\ &\propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} \end{aligned}$$

Normalizing:

$$\begin{aligned} P(\mu|D, \alpha) &= \text{Dir}(\mu|\alpha+m) \\ &= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} \end{aligned}$$

The Gaussian Distribution

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

$x \in \mathbb{R}$. For D-dimensional case

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

It's the continuous distr. that maximizes the entropy. The sum of random variables, under mild assumptions, also follows a Gaussian distribution (Central Limit Theorem).

The binomial distr. when $N \rightarrow \infty$ tends to a Gaussian.

$$\Delta^2 = (x - \mu)^T \Sigma^{-1} (x - \mu) \quad \text{Mahalanobis distance} \quad (40)$$

$N = cte$ on surfaces where $\Delta^2 = cte$.

Σ can be assumed symmetric, without loss of generality. In effect, we have

$$\Sigma = \Sigma_S + \Sigma_A = \frac{\Sigma + \Sigma^T}{2} + \frac{\Sigma - \Sigma^T}{2}$$

$$\text{Now } z^T \Sigma z = z_i \Sigma_{ij} z_j = z_i \Sigma_{Sij} z_j + z_i \Sigma_{Aij} z_j.$$

$$z_i \Sigma_{Aij} z_j = z_i \Sigma_{Aji} z_j = -z_i \Sigma_{Aij} z_j \therefore z_i \Sigma_{Aij} z_j = 0.$$

Consider the eigenvalue eq:

$$\Sigma u_i = \lambda_i u_i$$

Since $\Sigma = \Sigma^T \Rightarrow \lambda_i$ are real, and its eigenvectors form an orthonormal set, $u_i^T u_j = \delta_{ij}$.

Let's show these things:

$$\begin{aligned} \Sigma u_i = \lambda_i u_i \quad u_j^T \Sigma^T = \lambda_j^* u_j^T &\Rightarrow u_j^T \Sigma^T \Sigma u_i = \lambda_i \lambda_j^* u_j^T u_i \\ u_j^T \Sigma^2 u_i = \lambda_i \lambda_j^* u_j^T u_i \\ \lambda_i^2 u_j^T u_i = \lambda_i \lambda_j^* u_j^T u_i \\ (\lambda_i^2 - \lambda_i \lambda_j^*) u_j^T u_i = 0 \end{aligned}$$

$$\text{if } i=j \Rightarrow \lambda_i^2 = |\lambda_i|^2 \therefore \lambda_i \in \mathbb{R}$$

$$\text{if } i \neq j \Rightarrow u_j^T u_i = 0 \quad \text{orthogonal.}$$

So by appropriate normalization, $\boxed{u_i^T u_j = \delta_{ij}}$

Now multiply this by u_i : $u_i u_i^T u_j = \delta_{ij} u_i$

$$\text{or } \left(\sum_i u_i u_i^T \right) u_j = \left(\sum_i \delta_{ij} u_i \right) = u_j \therefore \boxed{I = \sum_i u_i u_i^T}$$

24

$$\underline{\Sigma} = \underline{Z} \underline{I} = \underline{Z} \sum_i u_i u_i^T = \underline{\sum_i \lambda_i u_i u_i^T}$$

Thus

$$\Delta^2 = (x - \mu)^T \left(\sum_i \frac{1}{\lambda_i} u_i u_i^T \right) (x - \mu)$$

$$= \sum_i \frac{1}{d_i} (x - \mu)^T u_i \quad u_i^T (x - \mu)$$

Define $y_i = \frac{1}{d_i} u_i^T (X - \mu) \Rightarrow \sigma^2 = \frac{1}{n} \sum_i y_i^2$

$$s^2 = \frac{\sum y_i^2}{n}$$

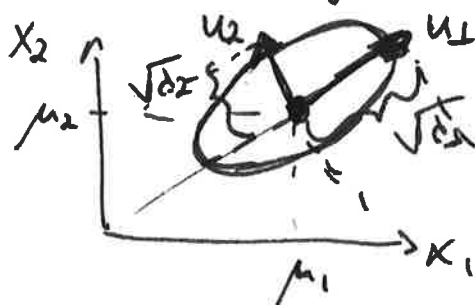
Let $y = (y_1, \dots, y_D)^T \in \mathbb{R}^D$, and U the matrix $U = \begin{pmatrix} u_1^T \\ \vdots \\ u_D^T \end{pmatrix}$. Then $y = U(x - \mu)$. It

follows that ~~$V^T V = (u_1 \dots u_n) \begin{pmatrix} u_1^T \\ \vdots \\ u_n^T \end{pmatrix} = V^T V$~~

$$U^T U = \begin{pmatrix} \underset{1}{u_1} & \underset{1}{u_2} & \dots & \underset{1}{u_D} \end{pmatrix} \begin{pmatrix} -u_1^T - \\ -u_2^T - \\ \vdots \\ -u_D^T - \end{pmatrix} = I = U U^T$$

So U is orthogonal.

N will be constant when y_i 's are constant.
If all $\lambda_i > 0$, these surfaces are ellipsoids with center at μ , and axes oriented along u_i , with scaling factor $\sqrt{\lambda_i}$.



$d_i > 0$ otherwise it's not possible to normalize. Thus Σ must be positive definite.
 If some $d_i = 0$ (singular distr.) then Σ is positive semi-definite.

(42)

Change of basis: $y = U(x - \mu)$. $U^T y = x - \mu$

$$x = U^T y + \mu. \quad x_i = U_{ij}^T (y_j) + \mu_i \\ = U_{ji} y_j + \mu_i$$

$$J = \frac{\partial x}{\partial y}, \quad J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji} \therefore \boxed{J = U^T}$$

$$|J|^2 = |J| |J^T| = |J J^T| = |I| = 1.$$

$$|J| = 1.$$

$$\text{Also, } |\Sigma|^{1/2} = \prod_{j=1}^D d_j^{1/2}. \text{ Therefore}$$

$$p(y) = p(x) |J| = \prod_{j=1}^D \frac{1}{(2\pi d_j)^{1/2}} e^{-\frac{y_j^2}{2d_j}}$$

which is a product of independent Gaussians.
 So the eigenvectors of the covariance Σ define a system of coordinates such that the joint $p(y)$ is factorized.

Now we compute the moments.

$$\begin{aligned} E[x] &= \int dx \, x \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \\ &= \int dz \, (z + \mu) \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} z^T \Sigma^{-1} z} \\ &= \mu + \int dz \, z \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} z^T \Sigma^{-1} z} \quad \text{K odd} = 0 \end{aligned}$$

$$E[x] = \mu$$

(43)

There are D^2 terms $E[x_i x_j]$ which can be grouped in the matrix $E[xx^T]$, thus

$$E[xx^T] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int dz (z+\mu)(z+\mu)^T e^{-\frac{1}{2} z^T \Sigma^{-1} z}$$

The terms $z\mu^T$ and μz^T vanish by symmetry. The term $\mu\mu^T$ is constant. It remains zz^T .

$$z = x - \mu = \left(\sum_{i=1}^D u_i u_i^T \right) (x - \mu) = \sum_{i=1}^D u_i y_i \quad \text{scalar}$$

$$zz^T = \sum_{i,j} y_i y_j u_i u_j^T$$

$$z^T \Sigma^{-1} z = \sum_{k=1}^D \frac{y_k^2}{\lambda_k}$$

$$z = U^T y \quad \text{and} \quad |U| = 1, \text{ thus}$$

$$\begin{aligned} E[xx^T] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int dy \sum_{i,j} y_i y_j u_i u_j^T e^{-\frac{1}{2} \sum_k \frac{y_k^2}{\lambda_k}} \\ &\quad + \mu\mu^T \quad \text{the integral vanishes unless } i=j \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \sum_{i=1}^D \int dy u_i u_i^T y_i^2 e^{-\frac{1}{2} \sum_k \frac{y_k^2}{\lambda_k}} \\ &\quad + \mu\mu^T \end{aligned}$$

The integral over the components that are not y_i will just give the normalization factor $(2\pi)^{D-1} \prod_{j \neq i} \lambda_j^{1/2}$. And the integral $\int dy_i y_i^2 e^{-\frac{1}{2} \frac{y_i^2}{\lambda_i}} = \sqrt{2\pi \lambda_i} \lambda_i$

Thus

$$\begin{aligned}
 \mathbb{E}[x x^T] &= \mu \mu^T + \prod_{j=1}^D \frac{1}{d_j^{1/2}} \sum_{i=1}^D \left(\prod_{k \neq i} d_k^{1/2} \right) u_i u_i^T d_i^{1/2} d_i \quad (44) \\
 &= \mu \mu^T + \sum_{i=1}^D d_i u_i u_i^T \\
 &= \mu \mu^T + \Sigma
 \end{aligned}$$

$$\begin{aligned}
 \text{Cov}[x] &= \mathbb{E}[(x - \mu)(x - \mu)^T] \\
 &= \mathbb{E}[x x^T - x \mu^T - \mu x^T + \mu \mu^T] \\
 &= \mathbb{E}[x x^T] - \mu \mu^T \\
 &= \Sigma
 \end{aligned}$$

Number of parameters:

$$\Sigma \rightarrow \frac{D(D+1)}{2}$$

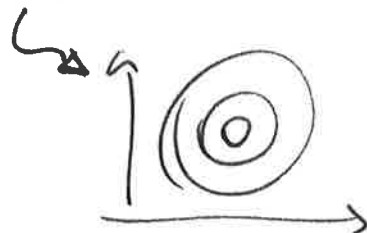
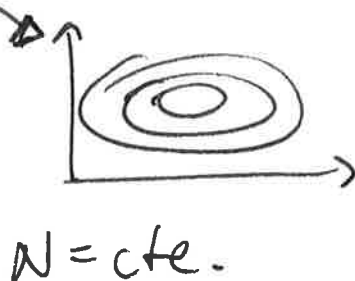
$$\mu \rightarrow D$$

$$\frac{D(D+3)}{2}$$

For large D , grows quadratically with D .

Inverting large matrices are expensive!

One way to deal with this is to assume some restriction on the form of Σ . For instance, $\Sigma = \text{diag}(\sigma_i^2) \Rightarrow$ 2D params, $\Sigma = \sigma^2 I$



conditional gaussian

(45)

If $p(x, y) \sim \mathcal{N}(\mu, \Sigma) \Rightarrow \begin{cases} \text{conditional distr. is } \mathcal{N} \\ \text{marginal distr. is } \mathcal{N} \end{cases}$

$x \in \mathbb{R}^D \sim \mathcal{N}(x | \mu, \Sigma)$ and we partition x into

$$\{x_a, x_b\} \quad x_a \cap x_b = \emptyset$$

$$x_a = (x_1, \dots, x_M)^T$$

$$x_b = (x_{M+1}, \dots, x_D)^T$$

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}; \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Sigma_{ba} = \Sigma_{ab}^T$$

$\Lambda = \Sigma^{-1}$ is the precision matrix

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \quad \Lambda_{ba} = \Lambda_{ab}^T$$

We then have

$$-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu) = -\frac{1}{2} \begin{pmatrix} (x_a - \mu_a)^T & (x_b - \mu_b)^T \end{pmatrix} \cdot \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \begin{pmatrix} x_a - \mu_a \\ x_b - \mu_b \end{pmatrix}$$

$$= -\frac{1}{2} \begin{pmatrix} (x_a - \mu_a)^T & (x_b - \mu_b)^T \end{pmatrix} \cdot$$

$$\begin{pmatrix} \Lambda_{aa}(x_a - \mu_a) + \Lambda_{ab}(x_b - \mu_b) \\ \Lambda_{ba}(x_a - \mu_a) + \Lambda_{bb}(x_b - \mu_b) \end{pmatrix}$$

$$= -\frac{1}{2} (x_a - \mu_a)^T \Lambda_{aa} (x_a - \mu_a)$$

$$- \frac{1}{2} (x_a - \mu_a)^T \Lambda_{ab} (x_b - \mu_b)$$

$$- \frac{1}{2} (x_b - \mu_b)^T \Lambda_{ba} (x_a - \mu_a)$$

$$- \frac{1}{2} (x_b - \mu_b)^T \Lambda_{bb} (x_b - \mu_b)$$

For fixed x_b , this is a quadratic form on x_a .
Thus $p(x_a | x_b) \propto p(x_a, x_b)_{\substack{\text{fixed} \\ x_b}}$ is a Gaussian.

Notice that for a general Gaussian

(46)

$$\begin{aligned} -\frac{1}{2}(X-\mu)^T \Sigma^{-1}(X-\mu) &= -\frac{1}{2}X^T \Sigma^{-1}X + \frac{1}{2}X^T \Sigma^{-1}\mu \\ &\quad + \frac{1}{2}\mu^T \Sigma^{-1}X - \frac{1}{2}\mu^T \Sigma^{-1}\mu \\ &= -\frac{1}{2}X^T \Sigma^{-1}X + X^T \Sigma^{-1}\mu - \underbrace{\frac{1}{2}\mu^T \Sigma^{-1}\mu}_{\text{cte}} \end{aligned}$$

We want the mean and covariance of $p(x_a|x_b)$: $\mu_{a|b}, \Sigma_{a|b}$.

Since x_b is fixed, the only term quadratic in x_a is

$$-\frac{1}{2}x_a^T \Lambda_{aa} x_a$$

Thus $\boxed{\Sigma_{a|b} = \Lambda_{aa}^{-1}}$. Now consider the linear terms on x_a :

$$\begin{aligned} &+ \frac{1}{2}x_a^T \Lambda_{aa} \mu_a + \frac{1}{2}\mu_a^T \Lambda_{aa} x_a \\ &- \frac{1}{2}x_a^T \Lambda_{ab} x_b - \frac{1}{2}x_a^T \Lambda_{ab} \mu_b \\ &- \frac{1}{2}x_b^T \Lambda_{ba} x_a - \frac{1}{2}\mu_b^T \Lambda_{ba} x_a \end{aligned}$$

$$\begin{aligned} &= x_a^T \Lambda_{aa} \mu_a - \frac{1}{2}x_a^T \Lambda_{ab} (x_b - \mu_b) \\ &\quad - \frac{1}{2}(x_b + \mu_b)^T \underbrace{\Lambda_{ba}}_{\Lambda_{ab}^T} x_a \end{aligned}$$

$$= x_a^T \left\{ \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b) \right\}$$

Thus $\Sigma_{ab}^{-1} \mu_{ab} = \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b)$

$$\mu_{ab} = \Sigma_{ab} (\Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b))$$

We already concluded that $\Sigma_{ab} = \Lambda_{aa}^{-1}$, thus

$$\boxed{\mu_{ab} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b)}$$

To express these results in terms of Σ_{ab} we can use (inverse of a partitioned matrix)

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{pmatrix}$$

$M = (A - BD^{-1}C)^{-1}$, M^{-1} is the Schur complement.

Proof.

Just multiply both sides.

$$I = \begin{pmatrix} A & B \\ C & D \end{pmatrix} \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{pmatrix}$$

$$= \begin{pmatrix} AM - BD^{-1}CM & -AMB D^{-1} + BD^{-1} + BD^{-1}CMB D^{-1} \\ CM - CM & -CMB D^{-1} + CMB D^{-1} + I \end{pmatrix}$$

$$-AMB D^{-1} + BD^{-1} + BD^{-1}CMB D^{-1} =$$

$$= (-A + BD^{-1}C)MB D^{-1} + BD^{-1}$$

$$= -M^{-1}MB D^{-1} + BD^{-1} = 0$$

$$\Delta (A - BD^{-1}C)M = M^{-1}M = I$$

Thus $I = I$

□

Using this for

(48)

$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$\Lambda_{aa} = M = (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1}$$

$$\Lambda_{ab} = -M \Sigma_{ab} \Sigma_{bb}^{-1}$$

~~$$\Sigma_{aa} \Sigma_{ab} \Sigma_{bb}^{-1} + \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \Sigma_{ab} \Sigma_{bb}^{-1}$$~~

$$= -(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_{bb}^{-1}$$

Thus $\boxed{\Sigma_{alb} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}}$

~~$$\mu_{alb} = \mu_{ab} (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1}$$~~

$$= \mu_a + (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_{bb}^{-1}$$

$$\cdot (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_{bb}^{-1} \cdot (x_b - \mu_b)$$

$$\boxed{\mu_{alb} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)}$$

Marginal Gaussian

$$p(x_a) = \int p(x_a, x_b) dx_b$$

From the partitioned quadratic form, pick the terms involving x_b :

$$-\frac{1}{2} (x_a - \mu_a)^T \Lambda_{ab} x_b - \frac{1}{2} x_b^T \Lambda_{ba} (x_a - \mu_a)$$

$$-\frac{1}{2} x_b^T \Lambda_{bb} x_b + \frac{1}{2} \mu_b^T \Lambda_{bb} x_b + \frac{1}{2} x_b^T \Lambda_{bb} \mu_b$$

$$-\frac{1}{2} x_b^T \Lambda_{bb} x_b - x_b^T \Lambda_{ba} (x_a - \mu_a) + x_b^T \Lambda_{bb} \mu_b$$

(49)

Thus we have

$$-\frac{1}{2} x_b^T \Lambda_{bb} x_b + x_b^T m //$$

where $m = \Lambda_{bb} \mu_b - \Lambda_{ba} (x_a - \mu_a)$. This is also

$$-\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m) + \frac{1}{2} m^T \Lambda_{bb}^{-1} m$$

So the integral over x_b is

$$\int e^{-\frac{1}{2} (x_b - \Lambda_{bb}^{-1} m)^T \Lambda_{bb} (x_b - \Lambda_{bb}^{-1} m)} dx_b = (2\pi)^{\frac{D-M}{2}} |\Lambda_{bb}^{-1}|$$

Now picking the terms that do not depend on x_b we have

$$= \frac{1}{2} m^T \Lambda_{bb}^{-1} m - \frac{1}{2} x_a^T \Lambda_{aa} x_a + \frac{1}{2} x_a^T \Lambda_{ab} \mu_b + \frac{1}{2} x_a^T \Lambda_{aa} \mu_a + \frac{1}{2} \mu_a^T \Lambda_{aa} x_a + \frac{1}{2} \mu_b^T \Lambda_{ba} x_a + cte$$

$$= \frac{1}{2} m^T \Lambda_{bb}^{-1} m - \frac{1}{2} x_a^T \Lambda_{aa} x_a + \frac{x_a^T \Lambda_{ab} \mu_b}{+ x_a^T \Lambda_{aa} \mu_a} + cte$$

$$= \frac{1}{2} (\Lambda_{bb} \mu_b - \Lambda_{ba} (x_a - \mu_a))^T \Lambda_{bb}^{-1} (\Lambda_{bb} \mu_b - \Lambda_{ba} (x_a - \mu_a)) - \frac{1}{2} x_a^T \Lambda_{aa} x_a + x_a^T (\Lambda_{aa} \mu_a + \Lambda_{ab} \mu_b) + cte$$

$$= -\frac{1}{2} (\Lambda_{bb} \mu_b)^T \Lambda_{bb}^{-1} \Lambda_{ba} x_a - \frac{1}{2} (\Lambda_{ba} x_a)^T \Lambda_{bb}^{-1} \Lambda_{bb} \mu_b + \frac{1}{2} (\Lambda_{ba} (x_a - \mu_a))^T \Lambda_{bb}^{-1} \Lambda_{ba} (x_a - \mu_a) + \dots$$

$$\begin{aligned}
&= -\frac{1}{2} \mu_b^T \cancel{\Lambda_{bb} \Lambda_{bb}^{-1}} \Lambda_{ba} x_a - \frac{1}{2} x_a^T \Lambda_{ab} \cancel{\Lambda_{bb}^{-1} \Lambda_{bb}} \mu_b \\
&\quad + \frac{1}{2} x_a^T \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} x_a - \frac{1}{2} x_a^T \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \mu_a \\
&\quad - \frac{1}{2} \mu_a^T \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} x_a - \frac{1}{2} x_a^T \Lambda_{ac} x_a \\
&\quad + x_a^T (\Lambda_{aa} \mu_a + \Lambda_{cb} \mu_b) + cte
\end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{2} x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) x_a \\
&\quad + x_a^T (\Lambda_{aa} \mu_a + \cancel{\Lambda_{ab} \mu_b} - \cancel{\Lambda_{ab} \mu_b} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \mu_a) + cte \\
&= -\frac{1}{2} x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) x_a \\
&\quad + x_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mu_a + cte.
\end{aligned}$$

Therefore

$$\boxed{\Sigma_a^{-1} = \Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}}$$

$$\Sigma_a^{-1} \mu_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mu_a \quad (\text{just an identity})$$

$$\begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \bar{\Sigma}_{aa} & \bar{\Sigma}_{ab} \\ \bar{\Sigma}_{ba} & \bar{\Sigma}_{bb} \end{pmatrix}$$

$$\boxed{\bar{\Sigma}_{aa} = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1} = \bar{\Sigma}_a}$$

Thus $\mathbb{E}[x_a] = \mu_a$

$$\mathbb{Cov}[x_a] = \bar{\Sigma}_{aa} //$$

Summing up: Give $N(x|\mu, \Sigma)$, $\Lambda = \Sigma^{-1}$ with the following partition

(51)

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

The conditional distribution is

$$p(x_a|x_b) = N(x_a|\mu_{a|b}, \Lambda_{aa}^{-1})$$

$$\mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (x_b - \mu_b)$$

The marginal distribution is \nwarrow linear in x_b

$$p(x_a) = N(x_a|\mu_a, \Sigma_{aa})$$

Bayes' Theorem for Gaussians

Suppose we are given a Gaussian marginal $p(x)$ and a Gaussian conditional $p(y|x)$ - the mean is linear in x (linear Gaussian Model). We want to find $p(y)$ and $p(x|y)$. Thus

$$p(x) = N(x|\mu, \Lambda^{-1})$$

$$p(y|x) = N(y|Ax+b, L^{-1})$$

$$* \quad p(y|x) = \frac{p(x,y)}{p(x)}$$

if $x \in \mathbb{R}^M$, $y \in \mathbb{R}^D$, then $A \in \mathbb{R}^{D \times M}$. First we find the joint dist. in terms of $z = \begin{pmatrix} x \\ y \end{pmatrix}$.

$$\begin{aligned} \log p(z) &= \log p(x) + \log p(y|x) \\ &= -\frac{1}{2}(x-\mu)^T \Lambda (x-\mu) - \frac{1}{2}(y-Ax-b)^T L (y-Ax-b) \\ &\quad + \text{cte.} \end{aligned}$$

This is quadratic in the components of z so it is again a Gaussian.

$$\log p(z) = -\frac{1}{2} x^T \Lambda x - \frac{1}{2} y^T L y - \frac{1}{2} x^T A^T L A x$$

$$+ \frac{1}{2} x^T \Lambda \mu + \frac{1}{2} \mu^T \Lambda x + \frac{1}{2} y^T L A x + \frac{1}{2} y^T L b$$

$$+ \frac{1}{2} x^T A^T L y - \frac{1}{2} x^T A^T L b + \frac{1}{2} b^T L y$$

$$- \frac{1}{2} b^T L A x + \text{cte}$$

$$= -\frac{1}{2} x^T (\Lambda + A^T L A) x - \frac{1}{2} y^T L y + \frac{1}{2} y^T L A x + \frac{1}{2} x^T A^T L y$$

$$+ x^T \Lambda \mu + y^T L b - x^T A^T L b \quad \text{linear}$$

$$+ \text{cte}$$

The quadratic terms can be written as

$$-\frac{1}{2} \begin{pmatrix} x \\ y \end{pmatrix}^T \underbrace{\begin{pmatrix} \Lambda + A^T L A & -A^T L \\ -L A & L \end{pmatrix}}_R \begin{pmatrix} x \\ y \end{pmatrix} = -\frac{1}{2} z^T R z$$

So the joint $p(z)$ has precision R , and

$$\text{cov}[z] = R^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^T \end{pmatrix}$$

use the
inversion
formula

$$\cancel{M = (\Lambda + A^T L A)^{-1} (-A^T L b)}$$

$$\cancel{= (\Lambda + A^T L A A)^{-1}}$$

$$M = (\Lambda + A^T L A + A^T L L^{-1} L A)^{-1}$$

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

$$= \Lambda^{-1}$$

$$+ \Lambda^{-1} A^T L L^{-1}$$

$$- (L^{-1} L - L A) \Lambda^{-1}$$

$$L^{-1} + L^{-1} (-L A) \Lambda^{-1} (-A^T L) L^{-1}$$

$$L^{-1} + A \Lambda^{-1} A^T$$

Doing the same for the linear terms:

(53)

$$x^T(\Lambda\mu - A^T L b) + y^T L b = \begin{pmatrix} x \\ y \end{pmatrix}^T \begin{pmatrix} \Lambda\mu - A^T L b \\ L b \end{pmatrix}$$

$$\bar{z}^{-1} \mu = \text{---} \rightarrow$$

$$\text{So } E[\bar{z}] = R^{-1} \begin{pmatrix} \Lambda\mu - A^T L b \\ L b \end{pmatrix}$$

$$= \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^T \end{pmatrix} \begin{pmatrix} \Lambda\mu - A^T L b \\ L b \end{pmatrix}$$

$$= \begin{pmatrix} \mu - \Lambda^{-1} A^T L b + \Lambda^{-1} A^T L b \\ A \mu - A A^T A^T L b + b + A A^{-1} A^T L b \end{pmatrix}$$

$$E[\bar{z}] = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}$$

To obtain the marginal $p(y)$ we just integrate this last result over x , and according to previous results gives a Gaussian with

$$E[y] = A\mu + b$$

$$\text{cov}[y] = R_{yy} = L^{-1} + A \Lambda^{-1} A^T$$

Particular case: $A = I \Rightarrow E[y] = \mu + b$

$$\text{cov}[y] = L^{-1} + \Lambda^{-1}$$

Convolution

For the conditional $p(x|y)$ we have:

$$p(x|y) = N(x | \mu_{x|y}, \Lambda_{xx}^{-1})$$

$$\mu_{x|y} = \mu_x - \Lambda_{xx}^{-1} \Lambda_{xy} (y - \mu_y)$$

$$= \mu - (\Lambda + A^T L A)^{-1} (-A^T L) (y - A\mu - b)$$

$$= \mu + (\Lambda + A^T L A)^{-1} (A^T L (y - b) + A^T L A \mu)$$

moving this inside

$$(\Lambda + A^T L A) \mu + A^T L A \mu = \Lambda \mu$$

$$\begin{cases} E[X|y] = (\Lambda + A^T L A)^{-1} \{ A^T L (y - b) + \Lambda \mu \} \\ \text{cov}[X|y] = (\Lambda + A^T L A)^{-1} \end{cases}$$

Summing up:

$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$

$$p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$$

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

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

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

Maximum Likelihood for Gaussian

$$\mathcal{N}(x | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

Given iid data drawn from this distribution $X = (x_1, \dots, x_N)^T$ we form the log likelihood function

$$\begin{aligned} \log p(X | \mu, \Sigma) &= -\frac{ND}{2} \log 2\pi - \frac{N}{2} \log |\Sigma| \\ &\quad - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) \end{aligned}$$

using that $\frac{\partial (X^T A X)}{\partial X} = A X + A^T X$
 $= 2 A X$ (symmetric)

we have $\frac{\partial \log p}{\partial \mu} = 0 = \sum_{n=1}^N \Sigma^{-1} (x_n - \mu)$

$$\boxed{\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n}$$

Now we use the following identities

$$\frac{\partial \log |X|}{\partial X} = (X^{-1})^T$$

$$\frac{\partial a^T X^{-1} b}{\partial X} = -(X^{-1})^T a b^T (X^{-1})^T = -(X^{-1} b a^T X^{-1})^T$$

Thus

$$\frac{\partial \log p}{\partial \Sigma} = 0 = -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} \sum_{n=1}^N \Sigma^{-1} (x_n - \mu) (x_n - \mu)^T \Sigma^{-1}$$

$$\boxed{\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML}) (x_n - \mu_{ML})^T}$$

Now taking the expectation over the true distribution we have

$$\boxed{E[\mu_{ML}] = \frac{1}{N} \sum_{n=1}^N E[x_n] = \frac{1}{N} N \mu = \mu}$$

To compute $E[\Sigma_{ML}]$ use μ_{ML}

$$\Sigma_{ML} = \frac{1}{N} \sum_n x_n x_n^T - \left(\frac{1}{N} \sum_n x_n \right) \mu_{ML}^T - \mu_{ML} \left(\frac{1}{N} \sum_n x_n^T \right) + \mu_{ML} \mu_{ML}^T$$

$$\begin{aligned}\Sigma_{ML} &= \frac{1}{N} \sum_n X_n X_n^T - \mu_{ML} \mu_{ML}^T \\ &= \frac{1}{N} \sum_n X_n X_n^T - \frac{1}{N^2} \sum_n \sum_m X_n X_m^T\end{aligned}$$

We computed before that $E[xx^T] = \mu\mu^T + \Sigma$.
Now if $n \neq m$, since the data is iid, we have
 $E[X_n X_m^T] = E[X_n] E[X_m^T] = \mu\mu^T$. Therefore,

$$E[X_n X_m^T] = \mu\mu^T + \delta_{nm} \Sigma$$

From this we have

$$\begin{aligned}E[\Sigma_{ML}] &= \frac{1}{N} \sum_n (\mu\mu^T + \Sigma) - \frac{1}{N^2} \sum_{n,m} (\mu\mu^T + \delta_{nm} \Sigma) \\ &= \cancel{\mu\mu^T} + \Sigma - \cancel{\mu\mu^T} - \frac{1}{N} \Sigma \\ &= \frac{N-1}{N} \Sigma \quad \text{biased}\end{aligned}$$

We can correct this by the unbiased estimator

$$\tilde{\Sigma} = \frac{N}{N-1} \Sigma_{ML} = \frac{1}{N-1} \sum_{n=1}^N (X_n - \mu_{ML})(X_n - \mu_{ML})^T$$

thus $E[\tilde{\Sigma}] = \Sigma$.

Sequential Estimation

Allow data points to be processed one at a time. This is important for on-line applications, and for very large data sets.

This can be done through Robbins-Monro algorithm. Consider the joint $p(z, \theta)$.

(57)

$$f(\theta) = \mathbb{E}[z|\theta] = \int z p(z|\theta) dz$$

this is a deterministic function of θ . This is called a regression function. The goal is to find the roots

$$f(\theta) = 0$$

Under the assumption that we receive one z at a time, and

$$\mathbb{E}[(z-f)^2|\theta] < \infty$$

without loss of generality we assume $f(\theta) > 0$ for $\theta > \theta^*$ and $f(\theta) < 0$ for $\theta < \theta^*$. Then Robbins-Monro procedure is

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

where

$$\lim_{N \rightarrow \infty} a_N = 0$$

$$\sum_{N=1}^{\infty} a_N = \infty$$

$$\sum_{N=1}^{\infty} a_N^2 < \infty$$

They can show that the sequence $\{\theta^{(N)}\}$ converges to θ^* with probability one.

In a ML estimation we have

$$\frac{\partial}{\partial \theta} \left\{ \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta) \right\} = 0$$

$$\frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \theta} \log p(x_n|\theta) = \mathbb{E}_x \left[\frac{\partial}{\partial \theta} \log p(x|\theta) \right]$$

$$\text{Now } \theta^{(N)} = \theta^{(N-1)} + a_{N-1} \frac{\partial \log p(x_N | \theta^{(N-1)})}{\partial \theta^{(N-1)}} \quad (58)$$

Consider a 1D gaussian $p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$

Then $\tilde{z} = \frac{\partial \log p(x | \mu, \sigma^2)}{\partial \mu} = \frac{(x-\mu)}{\sigma^2}$ So \tilde{z} is normally distributed with mean $\mu - \mu_{ML}$.

$$\mu^{(N)} = \mu^{(N-1)} + \frac{a_{N-1}}{\sigma^2} (x_N - \mu^{(N-1)})$$

choosing $a_{N-1} = \frac{\sigma^2}{N}$, $\mu^{(N)} = \mu^{(N-1)} + \frac{1}{N} (x_N - \mu^{(N-1)})$.

It's not hard to see that the solution to this recursive equation is

$$\mu^{(N)} = \frac{1}{N} \sum_{n=1}^N x_n$$

which is the ML estimate.

Bayesian Inference for the Gaussian

We need to introduce priors to the parameters. Suppose we have $N(x | \mu, \sigma^2)$ (1D), and σ is known, and we want to estimate μ .

The data likelihood func. is

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

Choose the conjugate prior for μ as

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

The posterior is the

(59)

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

both functions are quadratic in μ^2 so the posterior is a Gaussian. We can compute its parameters by considering

$$\begin{aligned} & -\frac{1}{2\sigma^2} \sum_{n=1}^N \{x_n^2 - 2x_n\mu + \mu^2\} - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 = \\ & = -\frac{1}{2\sigma^2} \sum_{n=1}^N x_n^2 + \frac{\mu}{\sigma^2} N\mu_{ML} - \frac{1}{2\sigma^2} N\mu^2 - \frac{1}{2\sigma_0^2} \mu^2 + \frac{\mu\mu_0}{\sigma_0^2} + \frac{1}{2\sigma_0^2} \mu_0^2 \\ & = \mu^2 \left(-\frac{N}{2\sigma^2} - \frac{1}{2\sigma_0^2} \right) + \mu \left(\frac{N\mu_{ML}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) + \text{cte} \end{aligned}$$

General form $-\frac{1}{2\sigma_N^2} (\mu - \mu_N)^2 = -\frac{1}{2\sigma_N^2} \mu^2 + \frac{1}{\sigma_N^2} \mu \mu_N + \text{cte}$

$$-\frac{1}{2\sigma_N^2} = -\frac{1}{2} \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \therefore \boxed{\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}$$

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

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

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

so $p(\mu|X) = \mathcal{N}(\mu | \mu_N, \sigma_N^2)$

Notice that $\begin{cases} N=0 \Rightarrow \mu_N = \mu_0, \text{ thus} \\ N \rightarrow \infty \Rightarrow \mu_N \rightarrow \mu_{ML} \end{cases}$ (60)
 $\mu_N \in [\mu_0, \mu_{ML}]$

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

$N=0 \Rightarrow \sigma_N = \sigma_0$
 $N \rightarrow \infty \Rightarrow \sigma_N \rightarrow 0$ so the posterior becomes infinitely peaked around the ML solution

So ML is recovered when $N \rightarrow \infty$.

Notice also, if N is fixed, and $\sigma_0 \rightarrow \infty$ the $\mu_N \rightarrow \mu_{ML}$, $\sigma_N^2 \rightarrow \frac{\sigma^2}{N}$.

Notice that we can view sequential view

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

↑ likelihood
 \propto prior with $N-1$ points.

Now assume the mean is known and we estimate the variance.

$$p(X|d) \propto d^{N/2} e^{-\frac{d}{2} \sum_{n=1}^N (x_n - \mu)^2}$$

precision
 $d = \frac{1}{\sigma^2}$

The conjugate prior must be $e^{ad} d^b$ so we can use

$$\text{Gam}(d|a, b) = \frac{1}{\Gamma(a)} b^a d^{a-1} e^{-bd}$$

$$\int dd \text{Gam}(d|a, b) = 1$$

(6)

$$\begin{aligned} \overline{E[d]} &= \int dd d \text{Gam}(d|a, b) \\ &= \int dd \frac{b^a d^a e^{-bd}}{\Gamma(a)} = \frac{a}{b} \end{aligned}$$

$$\int dd \frac{b^{a+1} d^a e^{-bd}}{\Gamma(a+1)} = 1 = \frac{b}{a} \int dd \frac{b^a d^a e^{-bd}}{\Gamma(a)}$$

$$E[d^2] = \int dd \frac{b^a d^{a+1} e^{-bd}}{\Gamma(a)} = \frac{a+1}{b} \cdot \frac{a}{b}$$

$$\frac{a}{b} = \int dd \frac{b^a d^a e^{-bd}}{\Gamma(a)} \rightarrow \frac{a+1}{b} \int dd \frac{b b^a d^{a+1} e^{-bd}}{a \Gamma(a)}$$

$$\text{Thus } \overline{\text{Var}[d]} = \frac{a+1}{b} \cdot \frac{a}{b} - \frac{a^2}{b^2} = \frac{a}{b^2}$$

The posterior is

$$\begin{aligned} p(d|x) &\propto p(x|d) p(d) \\ &= d^{N/2} e^{-\frac{1}{2} \sum_n (x_n - \mu)^2} d^{a_0-1} e^{-b_0 d} \\ &= d^{N/2 + a_0 - 1} e^{-b_0 d - \frac{1}{2} \sum_n (x_n - \mu)^2} \end{aligned}$$

which is also a $\text{Gam}(d|a_N, b_N)$ distribution with

$$a_N = \frac{N}{2} + a_0 - 1, \quad \boxed{a_N = a_0 + \frac{N}{2}}$$

$$-b_N = -b_0 - \frac{1}{2} \sum_n (x_n - \mu)^2, \quad \boxed{b_N = b_0 + \frac{N}{2} \sigma_{ML}^2}$$

Now suppose that both, the precision and the mean are unknown.

(62)

$$\begin{aligned} P(X|\mu, d) &= \prod_{n=1}^N \left(\frac{d}{2\pi}\right)^{1/2} e^{-\frac{d}{2}(x_n - \mu)^2} \\ &= \prod_{n=1}^N \frac{d^{1/2}}{(2\pi)^{1/2}} e^{-\frac{d}{2}(x_n^2 - 2x_n\mu + \mu^2)} \\ &= \left(\frac{d^{1/2}}{(2\pi)^{1/2}} e^{-\frac{d\mu^2}{2}}\right)^N e^{d\mu \sum_n x_n - \frac{d}{2} \sum_n x_n^2} \end{aligned}$$

So the conjugate prior should be of the form

$$\begin{aligned} P(\mu, d) &\propto \left(d^{1/2} e^{-\frac{d\mu^2}{2}}\right)^\beta e^{c\mu - dd} \\ &= d^{\beta/2} e^{-\frac{d\beta}{2}\left(\mu - \frac{c}{\beta}\right)^2} e^{+\frac{dc^2}{2\beta} - dd} \\ &= e^{-\frac{d\beta}{2}\left(\mu - \frac{c}{\beta}\right)^2} d^{\beta/2} e^{-d\left(d - \frac{c^2}{2\beta}\right)} \end{aligned}$$

$$P(\mu, d) = \underbrace{P(\mu|d)}_{\text{gaussian}} \underbrace{P(d)}_{\text{gamma}}$$

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

$$\mu_0 = \frac{c}{\beta}$$

$$a = 1 + \frac{\beta}{2}$$

$$b = d - \frac{c^2}{2\beta}$$

Normal-gamma

For $N(x|\mu, \Lambda) = \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} e^{-\frac{1}{2}(x-\mu)^T \Lambda (x-\mu)}$ (63)

if $\Sigma^{-1} = \Lambda$ is known and μ is unknown, the conjugate prior $p(\mu)$ is also a Gaussian. If μ is known, and Λ is unknown then

$$p(X|\Lambda) = \frac{1}{(2\pi)^{ND/2}} |\Lambda|^{N/2} e^{-\frac{1}{2} \sum_n (x_n - \mu)^T \Lambda (x_n - \mu)}$$

Notice that $\sum_n (x_n - \mu)^T \Lambda (x_n - \mu) = \sum_n \text{Tr}(z_n^T \Lambda z_n)$

where $z_n \equiv x_n - \mu$. $= \sum_n \text{Tr}(z_n z_n^T \Lambda)$

Since $\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$ $\Rightarrow = \text{Tr}((\sum_n z_n z_n^T) \Lambda)$
 $= \text{Tr}(Z \Lambda)$

Thus the conjugate prior is the Wishart distr.

$$W(\Lambda|W, \nu) = B |\Lambda|^{(\nu-D-1)/2} e^{-\frac{1}{2} \text{Tr}(W^{-1} \Lambda)}$$

ν is the number of degrees of freedom of the distr.

$W \in \mathbb{R}^{D \times D}$ the normalization is given by

$$B(W, \nu) = |\Lambda|^{-1/2} \left(2^{\nu D/2} \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1}$$

If both the mean and precision are unknown we have the conjugate prior

$$p(\mu, \Lambda | \mu_0, \beta, W, \nu) = N(\mu | \mu_0, (\beta \Lambda)^{-1}) W(\Lambda | W, \nu)$$

which is the normal-Wishart distr.

Student's t-Distribution

(64)

If μ is known, but $\sigma^2 = \frac{1}{d}$ is not, the conjugate prior to the gaussian is the gamma distr.

Here we consider:

$$\begin{aligned} p(x|\mu, a, b) &= \int_0^\infty dd N(x|\mu, d^{-1}) \text{Gam}(d|a, b) \\ &= \int_0^\infty dd \left(\frac{d}{2\pi}\right)^{1/2} e^{-\frac{d}{2}(x-\mu)^2} \frac{b^a d^{a-1} e^{-bd}}{\Gamma(a)} \\ &= \frac{b^a}{(2\pi)^{1/2} \Gamma(a)} \int_0^\infty dd d^{a-1/2} e^{-d(b + \frac{(x-\mu)^2}{2})} \end{aligned}$$

solving the integral.

$$\int_0^\infty dd d^{a-1/2} e^{-dc} = c^{-a-1/2} \int_0^\infty dt t^{a-1/2} e^{-t} = c^{-a-1/2} \Gamma(a+1/2)$$

$$\begin{aligned} dc = t & \quad \frac{1}{c} dt \cdot c^{-a+1/2} = dt c^{-a-1/2} \quad \Gamma(z) \equiv \int_0^\infty t^{z-1} e^{-t} dt \\ \frac{dt}{c} = \frac{1}{c} dt & \quad \left| \frac{1}{c} dt \cdot c^{-a+1/2} = dt c^{-a-1/2} \right| \end{aligned}$$

Thus

$$p(x|\mu, a, b) = \frac{b^a}{(2\pi)^{1/2} \Gamma(a)} \left(b + \frac{(x-\mu)^2}{2}\right)^{-a-1/2} \Gamma(a+1/2)$$

Now introduce parameters $\nu = 2a$, $d = a/b$:

$$\begin{aligned} p(x|\mu, \nu, d) &= \text{St}(x|\mu, d, \nu) \\ &= \frac{1}{(2\pi)^{1/2}} \frac{\Gamma(\frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{d^{\nu/2}}{b^{\nu/2}} b^{-\nu/2} \left(1 + \frac{(x-\mu)^2}{2b}\right)^{-\nu/2} \end{aligned}$$

$$b^{-1/2} = \left(\frac{a}{\lambda}\right)^{-1/2} = \left(\frac{\nu}{2\lambda}\right)^{-1/2} = 2^{1/2} \left(\frac{d}{\nu}\right)^{1/2}$$

Thus

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

 d is called the precision ν is called the ~~number~~ degrees of freedom $\nu=1 \Rightarrow St \rightarrow$ Cauchy distr. $\nu \rightarrow \infty \Rightarrow St \rightarrow N(x|\mu, d^{-1})$. Let us show this:

$$= \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2}) \nu^{1/2}} \left(\frac{d}{\pi}\right)^{1/2} \left(1 + \frac{d(x-\mu)^2}{\nu}\right)^{-\frac{\nu}{2}-\frac{1}{2}}$$

$$\stackrel{\nu \rightarrow 2z}{=} \frac{\Gamma(z + \frac{1}{2})}{\Gamma(z) z^{1/2}} \left(\frac{d}{2\pi}\right)^{1/2} \left(1 + \frac{d(x-\mu)^2}{2z}\right)^{-z-\frac{1}{2}}$$

Now using $\lim_{z \rightarrow \infty} \frac{\Gamma(z+\alpha)}{\Gamma(z) z^\alpha} = 1$ and $\lim_{z \rightarrow \infty} \left(1 + \frac{x}{z}\right)^{-z} =$

$$= \lim_{z \rightarrow \infty} \left(\left(1 + \frac{x}{z}\right)^z\right)^{-1}$$

$$= (e^x)^{-1} = e^{-x}$$

we have

$$\lim_{\nu \rightarrow \infty} St(x|\mu, d, \nu) = \left(\frac{d}{2\pi}\right)^{1/2} e^{-\frac{d(x-\mu)^2}{2}} = N(x|\mu, d^{-1})$$

The t-distribution is obtained by adding an infinite number of Gaussians with the same mean but different precisions. This is an infinite mixture. This gives a distribution which has longer tails than a Gaussian. \Rightarrow Robustness, less sensitive to outliers.

(66)

In higher dim. this generalizes to

$$St(X|\mu, \Lambda, \nu) = \int_0^\infty N(X|\mu, (\eta\Lambda)^{-1}) \text{Gam}(\eta|\nu/2, \nu/2) d\eta$$

Computing this integral gives

$$St(X|\mu, \Lambda, \nu) = \frac{\Gamma(D/2 + \nu/2)}{\Gamma(\nu/2)} \frac{|\Lambda|^{D/2}}{(\pi\nu)^{D/2}} \left(1 + \frac{\Delta^2}{\nu}\right)^{-D/2 - \nu/2}$$

$$\Delta^2 = (X - \mu)^T \Lambda (X - \mu)$$

$$E[X] = \mu$$

$$\text{COV}[X] = \frac{\nu}{\nu - 2} \Lambda^{-1}$$

$$\text{mode}[X] = \mu$$

Periodic Variables

$D = \{\theta_1, \dots, \theta_N\}$ periodic. Mean?

$\frac{1}{N} \sum \theta_i$ will be strongly coordinate dependent.

These points are on the unit circle in 2D
So we can average vectors $\{x_i\}$ and then
invert to find $\bar{\theta}$.

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \rightarrow \begin{array}{c} \bar{x} \\ \nearrow \\ \bar{r} \cos \bar{\theta} \\ \searrow \\ \bar{r} \sin \bar{\theta} \end{array} \quad \bar{\theta} = \arctan \frac{\bar{x}_2}{\bar{x}_1}$$

$$\text{or } \bar{\theta} = \arctan \left(\frac{\sum_i \sin \theta_i}{\sum_i \cos \theta_i} \right)$$

Periodic Generalization of a Gaussian is
the von Mises distn.

(67)

$$\left\{ \begin{array}{l} p(\theta) \geq 0 \text{ conditions} \\ \int_0^{2\pi} p(\theta) d\theta = 1 \\ p(\theta + 2\pi) = p(\theta) \end{array} \right.$$

Gaussian 2D: $P(x_1, x_2) = \frac{1}{2\pi\sigma^2} e^{-\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2}}$
 $\Sigma = \sigma^2 I$

$$\left\{ \begin{array}{l} x_1 = r \cos \theta \\ x_2 = r \sin \theta \\ \mu_1 = r_0 \cos \theta_0 \\ \mu_2 = r_0 \sin \theta_0 \end{array} \right. \Rightarrow \frac{r_0 \cos(\theta - \theta_0)}{\sigma^2} + \text{cte}$$

↑ indep. of θ

$$p(\theta) = \frac{1}{2\pi I_0(m)} e^{m \cos(\theta - \theta_0)}$$

Vo-Mises
circular
Normal

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} e^{m \cos \theta} d\theta$$

For large m this becomes approximately Gaussian
 since $\cos \theta = 1 - \frac{\theta^2}{2} + \dots$

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

$$\frac{\partial}{\partial \theta_0} (\dots) = m \sum_{n=1}^N \sin(\theta_n - \theta_0) = 0$$

$$\sum_{n=1}^N \sin(\theta_n - \theta_0) = 0$$

$$\sum_{n=1}^N (\sin \theta_n \cos \theta_0 - \sin \theta_0 \cos \theta_n) = 0$$

$$\boxed{\frac{\sin \theta_0}{\cos \theta_0} = \frac{\sum \sin \theta_n}{\sum \cos \theta_n}} \quad \bar{\theta} = \theta_{0, \text{MLE}}$$

$$\frac{\partial}{\partial m}(\dots) = 0 = -\frac{N}{I_0(m)} \frac{\partial I_0(m)}{\partial m} + \sum_{n=1}^N \cos(\theta_n - \theta_0)$$

$$\downarrow$$

$$\equiv I_1(m)$$

$$\frac{I_1(m)}{I_0(m)} = \frac{1}{N} \sum_{n=1}^N \cos(\theta_n - \theta_0^{ML}) \equiv A(m)$$

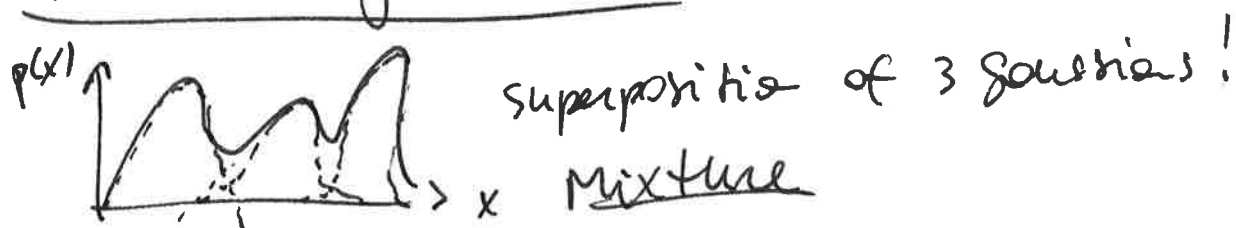
$$A(m) = \cos \theta_0^{ML} \left(\frac{1}{N} \sum \cos \theta_n \right) + \sin \theta_0^{ML} \left(\frac{1}{N} \sum \sin \theta_n \right)$$

$$A^{-1}(m) = m$$

K can be inverted numerically.

We can form mixtures of von Mises distributions to account for multimodality.

Mixture of Gaussians



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

↑
mixing coefficients

K component k of the mixture

$$\int p(x) dx = 1 \Rightarrow \sum_{k=1}^K \pi_k = 1$$

Since both sides are positive $\pi_k \geq 0$,
thus $0 \leq \pi_k \leq 1$. probabilities

$$p(x) = \sum_{k=1}^K p(k) p(x|k)$$

\downarrow π prior \downarrow N cond. likelihood

$p(k|x)$ posterior

$$h_k(x) = p(k|x) = \frac{p(x|k) p(k)}{\sum_n p(n) p(x|n)}$$

responsibilities

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

$$\text{MLE: } P(D|\pi, \mu, \Sigma) = \prod_{n=1}^N \sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k)$$

$$\pi = \{\pi_1, \dots, \pi_K\}$$

$$\mu = \{\mu_1, \dots, \mu_K\}$$

$$\Sigma = \{\Sigma_1, \dots, \Sigma_K\}$$

$$D = \{x_1, \dots, x_N\}$$

$$\log P = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k N_k(x_n) \right\}$$

↑
More complicated
Does not have a closed
form solution.

Exponential Family

Except for Gaussian, all other previous prob. distrs. are members of the exponential family.

$$p(x|\eta) = h(x) g(\eta) e^{\eta^T \phi(x)}$$

η is called natural parameters.
 $g(\eta)$ ensures normalization.

Ex.: Bernoulli

(70)

$$P(x) = \mu^x (1-\mu)^{1-x}$$

$$= e^{x \log \mu + (1-x) \log(1-\mu)}$$

$$= (1-\mu) e^{x \log \mu - x \log(1-\mu)}$$

$$= (1-\mu) e^{x \log \frac{\mu}{1-\mu}}$$

$$\uparrow g(\mu)$$

$$\uparrow \eta = \log \frac{\mu}{1-\mu}$$

$$e^\eta = \frac{\mu}{1-\mu} \therefore \frac{1}{\mu} - 1 = e^{-\eta}$$

$$\frac{1}{\mu} = e^{-\eta} + 1$$

logistic
sigmoid
func.

$$\boxed{\mu = \frac{1}{1+e^{-\eta}} = \sigma(\eta)}$$

$$g(\mu) = 1-\mu$$

$$= 1 - \frac{1}{1+e^{-\eta}}$$

$$= \frac{1+e^{-\eta} - 1}{1+e^{-\eta}}$$

$$= \frac{e^{-\eta}}{1+e^{-\eta}} = \sigma(-\eta)$$

$$\therefore \boxed{\text{Bern}(x|\mu) = \sigma(-\eta) e^{\eta x}}$$

where
 $\mu = \sigma(\eta)$

Ex.: Multinomial

$$P(x|\mu) = \prod_{k=1}^M \mu_k^{x_k} = \prod_{k=1}^M e^{x_k \log \mu_k}$$

$$\sum_k \mu_k = 1$$

$$= e^{\sum_k x_k \log \mu_k} = e^{\eta^T x}$$

$$\eta = (\log \mu_1, \dots, \log \mu_M)^T$$

$$x = (x_1, \dots, x_M)^T$$

Incorporating $\sum_{k=1}^M \mu_k = 1$ and $\sum_{k=1}^M x_k = 1$,

(71)

Since $x = (x_1, \dots, x_M)^T$ is a binary vector, we have

$$\begin{aligned} \sum_{k=1}^M x_k \log \mu_k &= \sum_{k=1}^{M-1} x_k \log \mu_k + x_M \log \mu_M \\ &= \sum_{k=1}^{M-1} x_k \log \mu_k + \left(1 - \sum_{k=1}^{M-1} x_k\right) \log \left(1 - \sum_{k=1}^{M-1} \mu_k\right) \\ &= \sum_{k=1}^{M-1} x_k \log \left(\frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right) + \log \left(1 - \sum_{k=1}^{M-1} \mu_k\right) \end{aligned}$$

Let $\eta_k = \log \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \therefore e^{\eta_k} = \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j}$

Ansatz $\boxed{c e^{\eta_k} = \mu_k} \rightarrow e^{\eta_k} = \frac{c e^{\eta_k}}{1 - c \sum_j e^{\eta_j}}, c = \frac{1}{1 + \sum_{j=1}^{M-1} e^{\eta_j}}$

Thus $\mu_k = \frac{e^{\eta_k}}{1 + \sum_{j=1}^{M-1} e^{\eta_j}}$ Softmax function or normalized exponential

So $p(x) = \left(1 - \sum_{k=1}^{M-1} \mu_k\right) e^{\sum_{k=1}^{M-1} x_k \eta_k}$

$\hookrightarrow \mu_k = \frac{e^{\eta_k}}{1 + \sum_{j=1}^{M-1} e^{\eta_j}}$

$p(x|\eta) = \underbrace{\left(1 + \sum_{j=1}^{M-1} e^{\eta_j}\right)^{-1}}_{g(\eta)} e^{\eta^T x}$

$\eta = (\eta_1, \dots, \eta_M)^T$

Simple Gaussian:

(72)

$$\begin{aligned}
 P(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} \\
 &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} e^{-\frac{1}{2} \frac{x^2 - 2x\mu + \mu^2}{\sigma^2}} \\
 e^{-\frac{1}{2} \frac{\mu^2}{\sigma^2}} &= e^{\frac{\eta_1^2}{4\eta_2}} \\
 \frac{1}{\sigma} &= (-2\eta_2)^{1/2}
 \end{aligned}$$

$$\begin{aligned}
 &\Delta \left(-\frac{1}{2} \frac{x^2}{\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{1}{2} \frac{\mu^2}{\sigma^2} \right) = \eta^T u(x) \\
 &\eta = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix}, u(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}
 \end{aligned}$$

$$\text{Thus } \mathcal{N}(x|\eta) = \underbrace{\frac{1}{\sqrt{2\pi}}}_{h(x)} \underbrace{(-2\eta_2)^{1/2} e^{\frac{\eta_1^2}{4\eta_2}}}_{g(\eta)} e^{\eta^T u(x)}$$

Maximum Likelihood

$$P(x|\eta) = h(x) g(\eta) e^{\eta^T u(x)}$$

$$\int P(x|\eta) dx = 1 = g(\eta) \int h(x) e^{\eta^T u(x)} dx$$

$$0 = \nabla g(\eta) \int h(x) e^{\eta^T u(x)} dx + g(\eta) \int h(x) e^{\eta^T u(x)} u(x) dx$$

$$0 = \frac{\nabla g(\eta)}{g(\eta)} + g(\eta) \int h(x) e^{\eta^T u(x)} u(x) dx$$

$$\boxed{-\nabla \log g(\eta) = \mathbb{E}[u(x)]}$$

The covariance comes from 2nd derivatives.
 As long as we can normalize an exponential distr.
 we can compute its moments in this way.

Likelihood function

(73)

$$P(D|\eta) = \prod_{n=1}^N h(x_n) g(\eta)^N e^{\eta^T \sum_{n=1}^N u(x_n)}$$

$$\log P(D|\eta) = \sum_{n=1}^N \log h(x_n) + N \log g(\eta) + \eta^T \sum_{n=1}^N u(x_n)$$

$$\nabla_{\eta} \log P(D|\eta) = 0 = \frac{1}{g(\eta)} \nabla g(\eta) + \sum_{n=1}^N u(x_n)$$

$$\therefore \boxed{-\nabla \log g(\eta)_{ML} = \frac{1}{N} \sum_{n=1}^N u(x_n)}$$

Sufficient statistic

Conjugate Priors

Given $P(x|\eta)$ we want a prior $p(\eta)$ such that the posterior $p(\eta|x)$ has the same form as the prior. For the exponential family, η prior should be in the form

$$P(\eta|x, \nu) = f(x, \nu) g(\eta)^{\nu} e^{\nu \eta^T x}$$

Thus

$$P(x|\eta) P(\eta|\nu) \Rightarrow P(\eta|x, \nu) \propto g(\eta)^{\nu+N} e^{\nu \eta^T x + \eta^T \sum_{n=1}^N u(x_n)}$$

$$P(\eta|x, \nu) \propto g(\eta)^{\nu+N} e^{\eta^T \left\{ \sum_{n=1}^N u(x_n) + \nu x \right\}}$$

which has the same form as the prior.

When we don't have much prior information, (74)
 we don't want to influence the posterior, so
 we can use uniform distr. as a prior. However,
 in the continuous case this can lead to an
 improper prior, which cannot be normalized.
 This is ok, as long as the posterior is proper.
 A second problem may appear in a change
 of variables: suppose $h(d) = \text{const.}$ $d = \gamma^2 \Rightarrow$
 $h(\gamma^2) = \text{const.}$, but $p_\gamma(\gamma) = p_d(d) \left| \frac{dd}{d\gamma} \right| = p_d(d) \underset{\substack{\uparrow \\ \text{not const.}}}{2\gamma}$
 $\underset{\substack{\uparrow \\ \text{const}}}{p_d(d)}$

Consider the noninformative prior

$$p(x|\mu) = f(x-\mu)$$

$$\hat{x} = x + c \Rightarrow p(\hat{x}|\hat{\mu}) = f(\hat{x} - \hat{\mu}) \quad \text{indep. choice of origin.}$$

prior must assign equal probs. to the
 interval $A \leq \mu \leq B$ as $A-c \leq \mu \leq B-c$

$$\int_A^B p(\mu) d\mu = \int_{A-c}^{B-c} p(\mu) d\mu = \int_A^B p(\mu-c) d\mu$$

$$\therefore \underline{p(\mu-c) = p(\mu)} \Rightarrow p(\mu) = \text{cte.}$$

Ex.: mean μ of a gaussian

$$\text{Consider } p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$$

$$\hat{x} = Cx \Rightarrow p(\hat{x}|\hat{\sigma}) = \frac{1}{\hat{\sigma}} f\left(\frac{\hat{x}}{\hat{\sigma}}\right)$$

$$\int_A^B p(\sigma) d\sigma = \int_{A/c}^{B/c} p(\sigma) d\sigma = \int_A^B p\left(\frac{\sigma}{c}\right) \frac{1}{c} d\sigma$$

(75)

scale
invariant.

$$\therefore p(\sigma) = \frac{1}{c} p\left(\frac{\sigma}{c}\right) \quad p(\sigma) \propto \frac{1}{\sigma}$$

Nonparametric Methods

Parametric Approach: prob. distr. has a specific form depending on few parameters whose values are determined from the data.

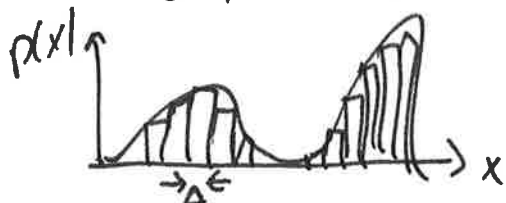
Nonparametric Approach: few assumptions about the form of the distribution.

Histogram method: partition x into bins Δ_i and count the number of observations of x , n_i , falling into Δ_i :

$$p_i = \frac{n_i}{N \Delta_i} = \text{prob}(x \in \Delta_i)$$

$$\sum_i p_i \Delta_i = 1 = \int p(x) dx.$$

Density $p(x)$ is constant inside each Δ_i .



problems:

- discontinuities due to the edges of the bins

- For high dimensions $\#(\Delta_\epsilon) = MD$. Need lots of data. Curse of Dimensionality.

locality



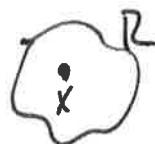
distance
measure around x .

Δ cannot be too small nor too large.
* choice of model complexity.

Kernel Density Estimators

(76)

$x \in \mathbb{R}^D$. Prob mass around x is

 $P = \int_R p(x) dx$

$x_1, \dots, x_N \sim p(x)$. $P(x_i \in R) = P$. The total number K of points that lie in R has prop distribution $\text{Bin}(K|N, P) = \binom{N}{K} P^K (1-P)^{N-K}$

$$E[K/N] = P \quad \xrightarrow{N \rightarrow \infty} K \approx NP$$

$$\text{Var}[K/N] = \frac{P(1-P)}{N}$$

If R is small, $p(x) \approx \text{const.}$ inside R , thus

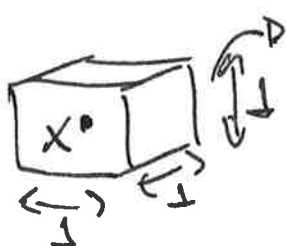
$$P \approx p(x)V \quad V = \text{Vol}(R)$$

Thus $p(x) = \frac{P}{V} = \frac{K}{NV}$

1. Fix K and determine V from data.
K-nearest-neighbour

2. Fix V and determine K from data.
Kernel Approach

Both converge to $p(x)$ when $N \rightarrow \infty$ provided V shrinks with N , and K grows with N .



$$K(u) = \begin{cases} 1, & |u_i| \leq \frac{1}{2} \quad i=1, \dots, D \\ 0, & \text{otherwise} \end{cases}$$

kernel function
Parzen window

$$K\left(\frac{x-x_n}{h}\right) = \begin{cases} 1 & \text{if } x_n \in \text{cube size } h \text{ around } x \\ 0 & \text{otherwise} \end{cases} \quad (77)$$

Total # of points in side this h -cube = $K = \sum_{n=1}^N K\left(\frac{x-x_n}{h}\right)$

Thus $p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} K\left(\frac{x-x_n}{h}\right)$
 $\hookrightarrow V$

We can also see this as a sum of N cubes centered at each data point x_n .

We can "smooth" this model to remove discontinuities in

$$p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{1/2}} e^{-\frac{\|x-x_n\|^2}{2h^2}}$$

$$\Sigma = h^2 I$$

we place a gaussian at each point and add the whole contributions over the data.

h plays the role of smoothing parameter.

choice of h = model complexity.



any $K(u)$ such that $\left\{ \begin{array}{l} \bullet K(u) \geq 0 \\ \bullet \int K(u) du = 1 \end{array} \right\}$ (can be used).


the computation cost for the density grows linearly with N .

Nearest-neighbour Methods

(78)

kernel, h is the same. Bad if data is too concentrated in a region (too smoothing). If h is too small then we get a lot of noise in regions where data is sparse.

$$\text{Recall } p(x) = \frac{k}{NV}$$

 sphere around x
Fix k

We allow the radius to grow until it captures k points

- k controls the smoothing now.
- Not a true density because the integral over all space diverges.

Example (classification): N_k points in class C_k .
 $\sum_k n_k = N$. given a new point x , classify it!

sphere around x containing k points independent of their classes. Suppose K_k points of C_k in this sphere of vol. V . The

$$p(x|C_k) = \frac{K_k}{N_k V}$$

$$p(x) = \frac{k}{NV}$$

$$p(C_k) = \frac{N_k}{N}$$

$$\text{Bayes} \Rightarrow p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)} = \frac{K_k}{k}$$

Assign x to the largest posterior!