

classical extensions to higher dimensions

- multi-variate standard normal $p(x) = N(\mu=0, \Sigma=I) = \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{x x^T}{2}\right)$
 $D = \dim(x)$ x : row vector

can be expressed as a product of 1-D normal distributions

$$x x^T = \sum_{j=1}^D x_j^2 \Rightarrow \exp\left(-\frac{x x^T}{2}\right) = \exp\left(-\frac{\sum_j x_j^2}{2}\right) = \prod_{j=1}^D \exp\left(-\frac{x_j^2}{2}\right)$$

\Rightarrow sampling in D dimensions: $x \sim N(0, I) \Leftrightarrow x_j \sim N(0, 1) \leftarrow$ 1-D gaussian
 \uparrow "random number generation"

- multivariate Gaussian distribution with mean μ , covariance Σ

$$p(x) = N(\mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{(x-\mu)\Sigma^{-1}(x-\mu)^T}{2}\right) \quad \det(\gamma A) = \gamma^D \det(A)$$

eigen decomposition of $\Sigma = U \Lambda U^T$ $\Sigma^{-1} = (U^T)^{-1} \Lambda^{-1} U^{-1} = U \Lambda^{-1} U^T$
 eigenvalues (diagonal $D \times D$)
 eigenvectors (orthonormal $D \times D$)

$$\Lambda^{-1/2} = \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} & & \\ & \ddots & \\ & & \frac{1}{\sqrt{\lambda_D}} \end{pmatrix}$$

define new coordinates. $\tilde{x} = (x-\mu) \cdot U \Lambda^{-1/2}$

$$\tilde{x} \tilde{x}^T = (x-\mu) \Sigma^{-1} (x-\mu)^T$$

sampling from $x \sim p(x) = N(\mu, \Sigma)$
 - sample $\tilde{x} \sim N(0, I)$

calculate $x = \tilde{x} \cdot \Lambda^{1/2} \cdot U^T + \mu$

"reparameterization trick"
 \Rightarrow reduce sampling from arbitrary Gaussian to sampling from standard normal, followed by linear transf.

in our language:

$$z = f(x) = (x-\mu) U \Lambda^{-1/2} \sim N(0, I) \Rightarrow x = f^{-1}(z) = z \Lambda^{1/2} U^T + \mu \sim N(\mu, \Sigma)$$

- similar to other analytic multi-variate distributions (scipy.stats offers 16)
if $p(x)$ must be learned. mixture model idea naturally generalize to D dimensions

- histograms: if bins are defined by a regular grid with k levels per dimension
 $\Rightarrow L = k^D$ exponential growth, does not scale

(requires $N = O(L^D)$ as well to "fill" the bins)

\Rightarrow define bins by recursive subdivision, i.e. density tree
quality of models is only medium

the # of correlations to be considered bounded by tree depth.

tree depth = $O(\log N)$

- GMM. work as before, but need to learn co-variance instead of variance ^{outer product}
 \Rightarrow change EM alg: $\Sigma_e^{(+)} = \frac{\sum_{i=1}^N \gamma_{ie} (X_i - \mu_e^{(+)})^T (X_i - \mu_e^{(+)})}{\sum_{i=1}^N \gamma_{ie}}$

- kernel density estimation: unchanged
but finding a bandwidth σ^2 that works equally well for all X_i is very difficult

\rightarrow become more expensive for growing D GMM = $O(L \cdot D^2)$, KDE = $O(N D)$

compute $p(x)$ efficiently: use (approximate) nearest neighbor search to find the non-negligible component $p_i(x)$

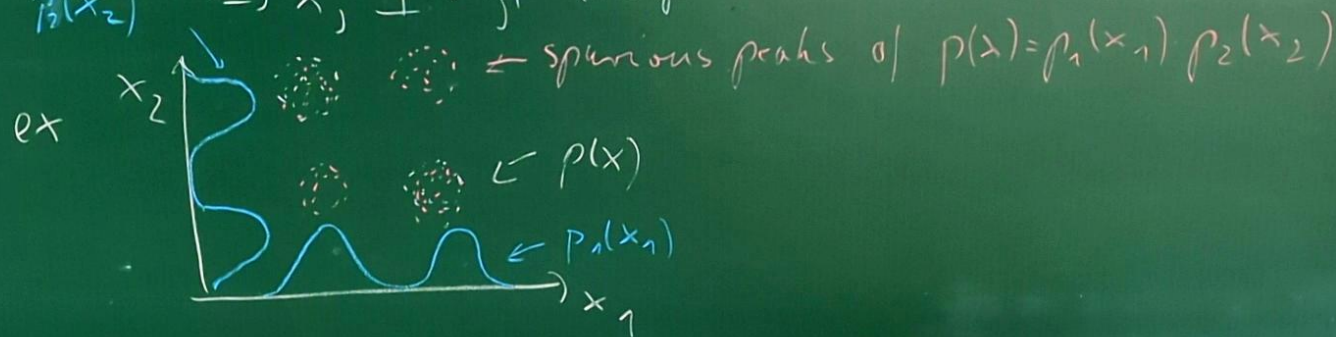
- sampling from a mixture model

- sample $l \sim \text{discrete}(\pi_1, \dots, \pi_L)$, sample $x \sim p_l(x)$
 $= N_1/N$ $= N_L/N$

methods to reduce $p(x)$ to a collection of 1-D distributions

- naive. $p(x) = \prod_{j=1}^D p_j(x_j)$, learn a 1-D model for each coordinate x_j

$p_1(x_1) \Rightarrow x_1 \perp x_2$ (independent coordinates \Rightarrow all correlations lost \Rightarrow highly unrealistic)



- exact: auto-regressive models : decompose $p(x)$ by Bayesian chain rule

$$p(x) = P_1(x_1) P_2(x_2 | x_1) P_3(x_3 | x_1, x_2) \dots P_D(x_D | x_{1:(D-1)})$$

$$= P_2(x_2) P_D(x_D | x_2) P_3(x_3 | x_2, x_D) \dots P_j(\dots)$$

any ordering of chain rule is valid for $p(x)$

each $P_j(x_j | X_{<j})$ is a collection of 1-D distributions (one distribution per value of $X_{<j}$)

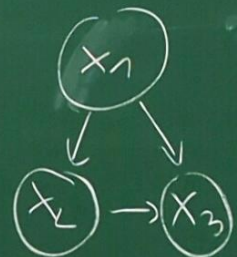
\Rightarrow use conditional inverse transform method

$$x_j \sim P(x_j | X_{<j}) \Leftrightarrow z_j \sim P(z_j), \quad x_j = f_j^{-1}(z_j | X_{<j})$$

$$f^{-1}(x) = \begin{pmatrix} x_1 = f_1^{-1}(z_1) \\ x_2 = f_2^{-1}(z_2 | x_1) \\ \vdots \\ x_D = f_D^{-1}(z_D | x_1, \dots, x_{D-1}) \end{pmatrix}$$

auto-regressive model

\rightarrow visualize as directed acyclic graph (DAG)



- problem: $f_j^{-1}(z_j, x_{<j})$ is a different 1-D function for each value of $x_{<j}$

\Rightarrow if $x_{<j}$ is defined on a regular grid with h levels per dimension

$\Rightarrow x_{<j}$ can take $h^{(j-1)}$ different values \Rightarrow does not scale well

• general solution: learn f_j^{-1} with a neural network, which generalizes from a few seen (TS) values of $x_{<j}$ to all possible values

• if $x_j \perp x_{j'}$ for $j' < j$, then $x_{j'}$ can be dropped from $f_j^{-1}(\cdot)$ \Rightarrow drop arrow in DAG

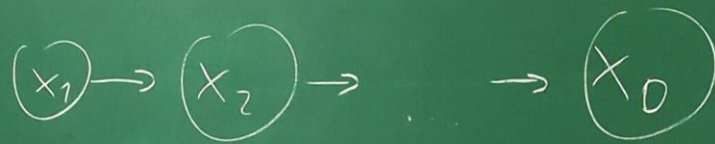
$x_j \perp x_{j''} \mid x_{j'}$ for $j'' < j$, $x_{j''}$ can be dropped \Rightarrow —||—

$x_j \perp x_{j'''} \mid x_{j'}, x_{j''}$

etc.

\Rightarrow drop as many arrows from DAG as possible $\Rightarrow f_j^{-1}(\cdot)$ become easier to learn

① Markov chain: $X_j \perp X_{j'} \mid X_{j-1}$ for all $j' < j-1$

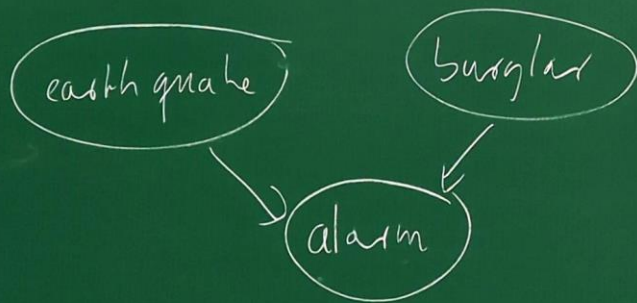


typical: $j = \text{time step}$
"future only depends on present, not the past"

$f_j^{-1}(z_j, X_{j-1})$ for all $j \Rightarrow$ easy to learn

② causal model: $p(x) = \prod_{j=1}^D p_j(x_j \mid \text{PA}(x_j))$

causal parents of x_j (set of direct causes)



property: causal DAG has fewest edges
 \Rightarrow easier to learn

but: finding the causal DAG is very hard
hot research topic "causal discovery"