

Static Facts:
 $x, y \in \mathbb{R}^m, A \in \mathbb{R}^{m \times m}$
 $\frac{dx^T y}{dx} = y \quad \frac{dx^T x}{dx} = 2x$
 $\frac{dx^T A}{dx} = A \quad \frac{d(A^T x)}{dx} = (A + A^T)x$

Covariance:
 $\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])^T]$
 For random vectors $X \in \mathbb{R}^m, Y \in \mathbb{R}^n$
 then the $m \times n$ covariance matrix
 $\text{Cov}(X, Y) = E[XY^T] - E[X]E[Y]^T$

correlation:
 $\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$

Variance $= E[X^2] - E[X]^2$
 $E[f(x)] = \int f(x)p(x)dx$
 $E[f(x)] = \int f(x)q(x|z)dx$

Exponential families of x given param η :
 $p(x|\eta) = h(x)g(\eta) \exp\{\eta^T u(x)\}$
 $u(x)$: natural statistic (sufficient)
 $h(x)$: base measure (often constant)
 $g(\eta)$: normalizer.

$g(\eta) \int h(x) \exp\{\eta^T u(x)\} dx = 1$
 E.g. Bernoulli: $p(x|\mu) = \mu^x (1-\mu)^{1-x}$
 $h(x) = 1, u(x) = x, \eta = \ln \frac{\mu}{1-\mu}$
 $g(\eta) = \sigma(-\eta) = \frac{1}{1 + e^\eta}$

Multinomial: $p(x_1, \dots, x_M | \mu) = \frac{N!}{x_1! \dots x_M!} \prod_{k=1}^M \mu_k^{x_k}$
 $h(x) = \frac{N!}{\prod x_k!}, u(x) = \sum x_k \eta_k, \eta = [\ln \mu_1, \dots, \ln \mu_M]$
 $g(\eta) = \sum e^{\eta_i}$, where $\sum e^{\eta_i} = 1$.

Gaussian: $p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\{-\frac{1}{2\sigma^2}(x-\mu)^2\}$
 $h(x) = \frac{1}{\sqrt{2\pi}}, g(\eta) = \frac{1}{\sqrt{2\pi}} \exp\{\frac{\eta^2}{2}\}$
 $\eta^T = [\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2}] \quad u(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$

Multivariate normal μ, Σ (cov)
 $p(x_1, \dots, x_K | \mu, \Sigma) = \frac{1}{(2\pi)^{K/2} |\Sigma|^{1/2}} \exp\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\}$
 $h(x) = \frac{1}{(2\pi)^{K/2}}, g(\eta) = \frac{1}{(2\pi)^{K/2}} \exp\{-\frac{1}{2}\eta^T \Sigma^{-1} \eta\}$
 $\eta = \begin{bmatrix} \Sigma^{-1} \mu \\ -\frac{1}{2} \Sigma^{-1} \end{bmatrix} \quad u(x) = \begin{bmatrix} x \\ x x^T \end{bmatrix}$

Poisson: $h(x) = \frac{1}{x!}, u(x) = x, \eta = \ln \lambda$
 $g(\eta) = e^\eta$

Beta: $h(x) = \frac{1}{x! (1-x)!}, u(x) = [\ln x, \ln(1-x)]$
 $\eta = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}, g(\eta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$

MLE for Exponential family:
 $-\nabla \ln g(\eta|\mu) = \frac{1}{N} \sum_{n=1}^N u(x_n)$
 $P(\theta|d) = \frac{P(d|\theta)P(\theta)}{P(d)}$ Likelihood.
 $P(\theta)$ prior. $P(\theta|d)$ posterior.

$\text{MLE} = \underset{\theta}{\text{argmax}} P(d|\theta)$
 $= \underset{\theta}{\text{argmax}} \sum_{i=1}^N \ln P(x_i|\theta)$
 $\frac{d \sum \log P(x_i|\theta)}{d\theta} = 0$, find θ
 (e.g. μ and σ for normal).
MAP $\theta_{\text{MAP}} = \underset{\theta}{\text{argmax}} P(\theta|d)$
 $= \underset{\theta}{\text{argmax}} \sum \log P(x_i|\theta)P(\theta)$

Conjugate prior: Conjugate priors induce a posterior of the same form: Gaussian for Gaussian, Gamma for poisson likelihood, Beta for Bernoulli/binomial.
 & Exp family \exists conjugate prior of $p(\eta|x, y) = f(x, y)g(\eta) \exp\{\eta^T u(x)\}$
 multiply by likelihood, posterior $p(\eta|D, x, y) \propto g(\eta)^{N+1} \exp\{\eta^T (u(x) + \sum_{n=1}^N u(x_n))\}$

Discriminative:
 $P(y|x) = P(k|x) = \frac{f(x, k, \theta)}{\sum_k f(x, k, \theta)}$
Generative:
 $P(y|x) = \frac{P(x, y)}{P(x)} = \frac{P(x|k)P(k)}{P(x)}$

D-separation: $(X) \rightsquigarrow (Z) \rightarrow (Y)$ blocked
 $(X) \rightsquigarrow (Z) \rightarrow (Y)$
 $(X) \rightsquigarrow (Z) \rightarrow (Y)$ if z is not and no descent of z is observed, then $X \perp Y$.

Regression: $p(y|x) = \frac{p(y, x)}{p(x)}$
classification: $p(c|x) = \frac{p(c, x)}{p(x)}$
clustering $p(c|x) = \frac{p(c, x)}{p(x)}$, c unknown.
Density estimation: $p(y|x) = \frac{p(y, x)}{p(x)}$
 y unobserved.
 If something is always unobserved, they are called hidden/latent.

Complexity: DAG: m variables, each variable takes k different values.
Naive inference: Cost $O(k)$ operations to update each of $O(k^{m-1})$ table entries.
Markov blanket (\Leftarrow) Independence $P(Y|E=\emptyset) = \frac{P(Y, e)}{P(e)}$

Markov blanket: parents, children, spouses.
 $mb(8) = \{3, 7, 9, 13\} \cup \{12, 4\}$
 Undirected graph also called Markov random field (MRF)
 In UGM: all neighbours.
 $X_A \perp X_B | X_C$ iff C separates A and B

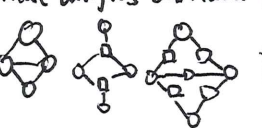
Global Markov Property:
 $\sum \log P(x_i|\theta) = 0$, find θ
 marry unmarried parents together pairwise.
 DAG \rightarrow ALM: moralization: marry unmarried parents together pairwise.
 ALM \rightarrow DAG: triangulation.
 $\psi_c(y_c|\theta_c)$ is a non negative potential function or factor.

Ham: $P(y|\theta) = \frac{1}{Z(\theta)} \prod_{c \in C} \psi_c(y_c|\theta_c)$
 C is maximal cliques, y is all variables, $Z(\theta) = \sum_y \prod_{c \in C} \psi_c(y_c|\theta_c)$ called partition function.

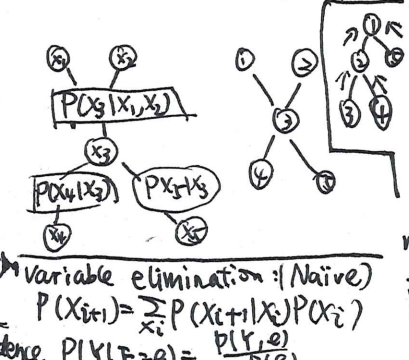
$p(y) \propto \psi_{123}(y_1, y_2, y_3) \psi_{235}(y_2, y_3, y_5) \psi_{345}(y_3, y_4, y_5)$
 $p(A, B, C, D) = \frac{1}{Z} \psi_{AB}(A, B) \psi_{BC}(B, C) \psi_{CD}(C, D) \psi_{AD}(A, D)$

$p(x, y, z) = \psi_{x,y}(x, y) \psi_{y,z}(y, z) \psi_{x,z}(x, z)$
 then $\psi_{x,y}(x, y, z) = \psi_{x,y}(x, y) \psi_{y,z}(y, z)$

Markov net \mathcal{H} over $X, U=U$ be the context. Residue reduced net $\mathcal{H}(U)$ is a Markov net over nodes $W=X-U$ where exists an edge between i, j there is an edge between in \mathcal{H} .

a factor graph is a graphical model that unifies DAG and UGM.


for directed model: one factor per (conditional distribution) CPD. And connect the factor to all the variables that use the CPD.



Complexity: m variables, each has k different values.
 $P(x_i) = \sum_{x_{-i}} P(x_i, x_{-i})$: k^2 multiplication and $k(k-1)$ add cost of total chain: $O(nk^2)$
 full joint: $O(k^n)$.

Sum-product inference:
 $\pi(y) = \sum_{\phi \in \Phi} \pi(\phi) \psi(\text{scope}(\phi) \cap Y, \text{scope}(\phi) \cap \bar{Y})$
 ψ is the set of potentials/factors. For DAG, Φ is CPD. N_i here: $\phi_c(A, D, I), Z = X - Y$

For UGM, Φ is the set of potentials. $Z = \sum_y J(y)$.

Variable Elim:
 let N_i be the # of entries in the factor ϕ_i , let $N_{\max} = \max N_i$
 Total cost: $O(mk N_{\max})$. (here n : # of variables, m : # of initial factors. N_{\max} can be as large as n .)

In UGM, N_{\max} = Size of largest clique.

$p(x)$	$x, z \sim p(x, z)$	$x \sim p(x)$	$p(x z)$	$p(z x)$
DAG	✓	✓	✓	✗
UGM	✗	✗	✗	✓

$p(x) = \int p(x, z) dz$
 $p(x) = \sum_z p(x, z)$

Message Belief Propagation:
 Message from j to $i \in N(i)$ is:
 $m_{j \rightarrow i}(x_i) = \sum_{x_j} \phi_j(x_j) \phi_{ij}(x_i, x_j)$

$\pi_{k \in N(i)}(x_i) = \sum_{x_j} \phi_j(x_j) \phi_{ij}(x_i, x_j)$. Each message is a vector with one value for each state of x_i . In order to have $m_{j \rightarrow i}$, we must have $m_{k \rightarrow j}$ for $k \in N(j) \setminus i$. Thus we need a specific ordering of messages.

We want $P(x)$. τ is the tree.

$p(x) \propto \sum_{x_1, \dots, x_n} \prod_{i,j \in \tau} \phi_{ij}(x_i, x_j)$
 $m_{5 \rightarrow 1}(x_1) = \sum_{x_5} \phi_5(x_5) \phi_{15}(x_1, x_5)$
 $m_{3 \rightarrow 2}(x_2) = \sum_{x_3} \phi_3(x_3) \phi_{13}(x_1, x_3)$
 $m_{2 \rightarrow 1}(x_1) = \sum_{x_2} \phi_2(x_2) \phi_{12}(x_1, x_2)$
 $m_{3 \rightarrow 2}(x_2) m_{4 \rightarrow 2}(x_2)$
 $P(x) \propto \phi_1(x_1) m_{5 \rightarrow 1}(x_1) m_{3 \rightarrow 1}(x_1) \sum_{x_1} \phi_1(x_1)$
 Choose root r . Leaf i to r to leaf

Running time is 2 times the cost of $VE = O(nk^2)$, $n = \#$ of nodes, $k = \#$ states per node.

KL-Divergence:
 $D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$, measures the expected number of bits required to describe samples from $p(x)$ using a code based on q instead of p .

$D(p||q) \geq 0 \forall p, q. = 0 \iff p = q$. Asymmetric.

Punish when high, q low. So, see the unnormalized distn.

$J(q) = \sum_x q(x) \log \frac{q(x)}{p(x)} = \dots = \frac{q(x)}{p(x)}$
 $= KL(q||p) - \log Z$. Since Z is constant, we force q to become p by minimize $J(q)$.

$J(q) \geq -\log(Z) = -\log p(D)$. negative log likelihood. ∇
 $J(q) = E_q[\log q(x)] + E_q[\log \tilde{p}(x)]$
 $= -H(q) + E_q[E_{\tilde{p}}]$
 $= E_q[-\log p(D)] + KL(q||p)$

Monte Carlo: Problem: Sampling from $p(x)$. (we can only compute $p(x)$ for a given x !) calculate $\hat{\Phi} = E_{x \sim p(x)} [\phi(x)] = \int \phi(x) p(x) dx$.

Simple Monte Carlo: if $x^{(r)} \sim p(x)$, then to find $\hat{\Phi} = \frac{1}{R} \sum \phi(x^{(r)})$. Then $E[\hat{\Phi}] = E[\Phi]$, $var[\hat{\Phi}] \propto \frac{\sigma^2}{R}$.

where $\sigma^2 = \int (\phi(x) - \Phi)^2 p(x) dx$. Note: accuracy of MC depend on var of $\phi(x)$, not dim of x . But, obtaining independent samples in high dim. 在哪里都很稀薄!

Why is sampling from $\phi(x)$? assume $\tilde{p}(x) \propto \phi(x)$, but $Z = \int \tilde{p}(x) dx$. Even if know Z , correct samples from p will tend to come from places in x -space where $p(x)$ is big, but how can we find those places without eval $p(x)$ everywhere? (In high dim?) $Z = \int \tilde{p}(x) dx$. Compute Z cost \propto (dimension of x).

Rejection Sampling: has a $Q(x)$. assume \exists constant c , s.t. $cQ(x) \geq p(x)$ for all x . (ZQ is the multiplicative factor, $Q(x)$ is the proposal density). generate 2 random numbers. 1st: $x \sim Q$

\Rightarrow eval $cQ(x)$. 2nd: $u \sim [0, cQ(x)]$ uniform. eval $p(x)$.

If $\tilde{p}(x) < u$, reject
 If $\tilde{p}(x) > u$, accept, add x

to our samples $\{x^{(n)}\}$, then $p(x)$ work well if Q is a good approx of p . Else, rejection rate will be large. In high dim, c will be forced too large. rejection rate so high, find c is hard, too.

Importance sampling. Not soln to problem 1, it is soln to problem 2. ($\phi(x)$). Still, $Q(x) = \tilde{Q}(x)$. Q is called sampler density. $x \sim Q$, sample $\{x^{(n)}\}_{n=1}^R$. $w_r = \frac{p(x^{(n)})}{Q(x^{(n)})}$ (weights). $\hat{\Phi} = \sum_r w_r \phi(x^{(n)})$

problem: it is hard to estimate how reliable the estimator $\hat{\Phi}$ is (the variation of $\hat{\Phi}$ is hard to estimate) w_r and $w_r \phi(x^{(n)})$ may not be a good guide).

dramatically wrong! Also, in high dim, likely $\hat{\Phi}$ will be dominated by a few large w_r . The hyper.

Sphere case: all on surface. Metropolis-Hastings: Q proposed depending on current state $x^{(t)}$. $Q(x'; x^{(t)})$ may be simple as gaussian centered on $x^{(t)}$.

pending new state: $x' \sim Q(x'; x^{(t)})$. $a = \frac{p(x') Q(x^{(t)}; x')}{p(x^{(t)}) Q(x'; x^{(t)})}$. if $a \geq 1$, accept. Else, accept with probability a .

If accepted, set $x^{(t+1)} = x'$. If rejected, set $x^{(t+1)} = x^{(t)}$. Add x' to $\{x^{(n)}\}$ no matter what. Samples are dependent. As an e.g. of MCMC, x^t depend on x^{t-1} .

need long time to be "independent" difficult to tell whether converged. Given enough time, get all! Use for high dim. (Where the scale ϵ is small, so that it won't jump to low-prob area eazly). But long-time random walk.

HMM: stationary: generative process does not depend evolve through time. Non-stationary: ... does... In general, $p(x_1, \dots, x_T) = \prod_{t=1}^T P(x_t | x_{t-1}, \dots, x_1)$.

FOMC: $P(x_t | x_{t-1})$. Each pair of output is a training case. $P(x_t = B | x_{t-1} = A) = \#(t, st. x_t = B, x_{t-1} = A) / \#(x_{t-1} = A)$. Higher order: Many counts may be zero in training data set. # of parameters:

1st order: $K(K-1)$. N th order: $K^M(K-1)$, $x_t \in \{1, \dots, K\}$. $0 \rightarrow 0 \rightarrow 0 \dots$ initial state: $\pi_i = P(z_t = i)$. $0 \rightarrow 0 \rightarrow 0$ (Hidden is 1st order).

Transition Prob: $T(i, j) = P(z_t = j | z_{t-1} = i)$. Emission Prob: $E_i(x_t) = P(x_t | z_t = i)$. initial: 打电话前天气. Transition: 雨 \rightarrow 晴 or 雨 \rightarrow 雨. Emission: 雨天出去的概率. output is not any order. Joint: $P(\vec{x}, \vec{z}) = P(\vec{z}) P(x_1 | z_1) \prod_{t=2}^T P(z_t | z_{t-1}) P(x_t | z_t)$ change.

SVI: Inference: $P(\vec{z} | \vec{x}) = \frac{P(\vec{x}, \vec{z})}{\int P(\vec{x}, \vec{z}) d\vec{z}}$ \downarrow true posterior, \downarrow eval density hard.

Introduce a variational distn: $q_\phi(\vec{z} | \vec{x}) = q(\vec{z} | \vec{x}; \phi)$. 2. optimize ϕ to minimize distance between q_ϕ and p .

$KL(q(y)||p(y)) = E_{y \sim q(y)} [\log \frac{q(y)}{p(y)}]$. ELBO: ① $P(\vec{z} | \vec{x}) \approx \frac{P(\vec{x}, \vec{z})}{P(\vec{x})}$. $KL(q_\phi(\vec{z} | \vec{x}) || P(\vec{z} | \vec{x})) = E_{\vec{z} \sim q_\phi(\vec{z} | \vec{x})} [\log \frac{q_\phi(\vec{z} | \vec{x})}{P(\vec{z} | \vec{x})}]$

$= E_q [\log (q_\phi(\vec{z} | \vec{x}) \frac{P(\vec{x})}{P(\vec{x}, \vec{z})})]$ ① $= E_q \log (\frac{q_\phi(\vec{z} | \vec{x})}{P(\vec{x}, \vec{z})}) + E_q [P(\vec{x})]$ \downarrow not depend on q !

$= KL(q_\phi(\vec{z} | \vec{x}) || P(\vec{x}, \vec{z})) + \log(P(\vec{x}))$. Evidence Lower BO and.

KLBO = $-E_q [\log P(\vec{z} | \vec{x}) - \log q_\phi(\vec{z} | \vec{x})]$ so, maximizing ELBO is minimizing KL posterior ①. = minimizing distance. Also, ELBO = $E_q [\log p(x|\vec{z})] - KL(q(\vec{z}||p(\vec{z}))$

Reparametrization trick: In order to do backprop (BP) cannot go through random node original: $z \sim q_\phi(z|x)$ Reparam: $z = \mu + \sigma \epsilon$, $\epsilon \sim p(\epsilon)$

we can learn $q(z|x)$ in a very flexible way using NN to output param ϕ for q . NN takes x as input, output ϕ . combine deep learning and graphical model. Remember we estimate ELBO by sampling from q .

ELBO: $\mathcal{L}(\phi) = (1/MC) H(B)$ KL. $E_q [\log p(x|\vec{z}) + \log p(\vec{z}) - \log q_\phi(\vec{z} | \vec{x})]$ $E_q [\log p(x|\vec{z})] + H[q_\phi]$ $E_q [\log p(\vec{z})] - KL(q_\phi(\vec{z} | \vec{x}) || p(\vec{z}))$

when $q_\phi(\vec{z}) = p(\vec{z} | \vec{x})$, MC has zero variance. When optimizing, $\nabla \mathcal{L}(\phi)$ use reparam trick (Back prop) $\nabla_{\phi} \mathcal{L}(\phi) = \text{reinforce with score func.}$

high variance score function $E [\log p(x|\vec{z}) - \log q_\phi(\vec{z} | \vec{x})] \nabla_{\phi} \log q_\phi(\vec{z} | \vec{x})$ $\approx \nabla_{\phi} \log q_\phi(\vec{z} | \vec{x})$

if \vec{z} is deterministic fn of ϕ and unbiased, push p_ϕ inside. $E [p_\phi \log p(q(\vec{z}, \phi), x) - \log p(q(\vec{z}, \phi) | x)]$ $\approx \nabla_{\phi} \log p(q(\vec{z}, \phi) | x)$

lower variance ($q(\vec{z}, \phi) = \vec{z}$). MC sampling unbiased estimator of \int . iid sample $\vec{z}^{(n)} \sim q$, take E . Alg: sample $\vec{z}^{(n)} \sim q_\phi(\vec{z} | \vec{x})$ compute $\log P(\vec{z}^{(n)}, x) - \log q_\phi(\vec{z}^{(n)} | \vec{x})$

R times. $\nabla_{\phi} \mathcal{L} = \frac{1}{R} \sum \nabla_{\phi} \log p(q(\vec{z}, \phi) | x)$ P ϕ by Auto diff. Sample from MC \Rightarrow stochastic.

VAE Encoder decoder. encoder: $q(x) \Rightarrow z$ decoder: $z \Rightarrow x$ latent. $x \in \mathbb{R}^{x_1 \times x_2}$ $z \in \mathbb{R}^{z_1 \times z_2}$ $\dim(F) < \dim(L)$. MC is asymptotically exact. SVI simple form underestimates posterior variance, but it can measure inference progress, use fancy tools (Adam).