

The big picture: - probabilistic model

- [RECAP]
- Specify a likelihood func.: $p(x|\vec{\theta}) \xrightarrow{x: \text{observed data (e.g. samples, D features)}} \vec{\theta}: \text{vector of model params. } \vec{\theta} \in \mathbb{R}^p \quad (p: \# \text{params})$
 - Max. Likelihood estimation (frequentist approach) \rightarrow if likelihood p = Gauss., ML estim. = least squares

$$\hat{\vec{\theta}}_{ML} = \underset{\vec{\theta}}{\operatorname{argmax}} p(x|\vec{\theta}) \quad \text{where, } p(x|\vec{\theta}) = \prod_i p(x_i|\vec{\theta}) \quad \begin{array}{l} [\text{i.i.d. assumption}] \\ \text{indep. \& ident. distrib.} \end{array}$$

0

b) Add regularizer for smoother estimates: MAP (\approx penalized ML)

c) Specify a prior on model params. $p(\vec{\theta})$

$$\hat{\vec{\theta}}_{MAP} = \underset{\vec{\theta}}{\operatorname{argmax}} p(x|\vec{\theta}) = \underset{\vec{\theta}}{\operatorname{argmax}} p(x|\vec{\theta}) p(\vec{\theta}) = \underset{\vec{\theta}}{\operatorname{argmax}} \frac{p(x|\vec{\theta}) p(\vec{\theta})}{p(x)} = \underset{\vec{\theta}}{\operatorname{argmax}} p(x|\vec{\theta}) p(\vec{\theta})$$

\downarrow
normalize joint probab.
a.k.a. prior

= $\operatorname{argmax}_{\vec{\theta}} p(\vec{\theta}|x)$
 $p(x)$
dividing by this does
not affect the 'argmax'

- * Bayesian approach:
- specify prior $p(\vec{\theta})$
 - Calcul. posterior $p(\vec{\theta}|x)$

Supervised learning:

if : $x \in \mathbb{R}^{n \times d}$
o/p : $y \in \mathbb{R}^n$ (regression), $y \in \{c_1, \dots, c_k\}^n$ (classif.)

Conditional likel.: $p(y|x, \vec{\theta}) \xrightarrow{\downarrow} = \prod_{i=1}^N p(y_i|x_i, \vec{\theta})$

\hookrightarrow Learning same as un-conditional case:

- 1) ML or,
- 2) introduce prior and MAP.

3. After learning \rightarrow make preds :

frequentist
 $\hat{y}_{ML} = \underset{\vec{\theta}}{\operatorname{argmax}} p(y|x, \vec{\theta})$

predictive distrib. is: $p(y^*|x^*, \theta_{ML})$
 \downarrow new y \downarrow new x

\rightarrow Best / chosen prediction = $\hat{y} = \underset{y^*}{\operatorname{argmax}} p(\hat{y}^*|x^*, \hat{\theta}_{ML})$
 \downarrow (has max. probab.)

\rightarrow But sometimes we don't want the one that has max. probab. BUT minimizes expected loss.

(Bishop 15) \rightarrow we have a loss func. $L(y, y_{pred})$ (this quantifies loss)
 \rightarrow Min. expected loss func. :

$$\hat{y} = \underset{y}{\operatorname{argmin}} \sum p(y|x^*, \vec{\theta}) L(y, y^*)$$

\downarrow prob. of true y given new x & $\vec{\theta}_{ML}$?

in medicine
loss func. is very asym.
 \hookleftarrow e.g.: TP, FP,
etc.

Bayesian

predic. w/ complete posterior: (over params θ)
 $p(\vec{\theta} | \vec{x}, \vec{y})$

$$\text{predictive dist. } p(y^* | x^*, \vec{x}, \vec{y}) = \int p(y^* | x^*, \vec{\theta}) p(\vec{\theta} | \vec{x}, \vec{y}) d\vec{\theta}$$

Then proceed as before. . . .

* Sometimes we have latent var. in our prob. model:

$$\vec{z} = (z_1, \dots, z_n)$$

so in case ①:

Sum/integrate latent var. out:

$\overset{k}{\underset{z}{\sum}}$

$$p(x | \vec{\theta}) = \sum_z p(x, z | \vec{\theta}) = \sum_{z_1} \dots \sum_{z_n} p(x, z | \vec{\theta}) \quad [\text{discrete}]$$

$$p(x | \vec{\theta}) = \int p(x, z | \vec{\theta}) dz \quad [\text{contin.}]$$

In practice, how do we learn in this case ??

1) EM algo.

2) variational bayes / sampling methods for approx. Bayesian learning & pred.

Exponential Family Distrub. : (Bishop 2.4)

$$\text{Def: } p(x | \vec{\eta}) = h(\vec{x}) g(\vec{\eta}) \exp(\vec{\eta}^T \cdot \vec{u}(\vec{x}))$$

choose 3 func.

$h(\vec{x})$ = natural param. | $h(\vec{x})$ = backgr. measure

$g(\cdot)$ = normal. term | $\vec{u}(\vec{x})$ = sufficient/natural stats.

$$z(\vec{\eta}) = \frac{1}{g(\vec{\eta})} = \int \exp(\vec{\eta}^T \cdot \vec{u}(\vec{x})) h(\vec{x}) d\vec{x}$$

$$\frac{\partial \log(z(\vec{\eta}))}{\partial \vec{\eta}} = -\frac{\partial (\log(g(\vec{\eta})))}{\partial \vec{\eta}} = \frac{1}{z(\vec{\eta})} \int \exp(\vec{\eta}^T \cdot \vec{u}(\vec{x})) h(\vec{x}) d\vec{x} = \mathbb{E}(\vec{u}(\vec{x}) | \vec{\eta})$$

Eg: Gaussian:

$$(\text{multivariate Gauss.}) \quad N(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)$$

$$= \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \vec{\mu}^T \Sigma^{-1} \vec{\mu}\right) \exp\left(-\frac{1}{2} \text{Tr}(\Sigma^{-1} \vec{x} \vec{x}^T) + (\vec{\mu}^T \Sigma^{-1} \vec{x})\right)$$

wrong
actually $\text{Tr}(\vec{x} \vec{x}^T)$

$$\Rightarrow \begin{pmatrix} \vec{\eta}_1 \\ \vec{\eta}_2 \end{pmatrix} = \begin{pmatrix} \vec{\Sigma}^{-1} & \vec{\mu} \\ \vec{\Sigma}^{-1} & -1 \end{pmatrix}$$

$$\parallel \begin{pmatrix} \vec{\mu}_1(\vec{x}) \\ \vec{\mu}_2(\vec{x}) \end{pmatrix} = \begin{pmatrix} \vec{x} \\ \vec{x}^T \end{pmatrix} \parallel \Lambda(\vec{x}) = 1$$

$$\text{Tr}(AB) = \sum_{ij} (AB)_{ij} = \sum_{ij} A_{ij} B_{ij}$$

[here, Σ is symm., so ...]

($D+D^2$ elements)

$$= \text{vec}(A^T) \text{vec}(B)$$

Q.

if $\vec{x} \sim N(\vec{\mu}, \Sigma)$ then $E(x) = \mu$

$$\Rightarrow E(\vec{x} | \vec{\eta}) = E(\mu(\vec{x}) | \vec{\eta}) = \frac{\partial}{\partial \vec{\eta}} \log Z(\vec{\eta}) = \frac{\partial}{\partial \vec{\eta}} \left(\frac{1}{2} \vec{\eta}^T \Sigma \vec{\eta} \right) = \Sigma \vec{\eta} = \Sigma \vec{\eta}^{-1} \vec{\eta} = \underline{\underline{\mu}}$$

Passed

Derive second moment on your own (var.).

2) ML of m.v. Gaussian (2.3 Bishop)

Given Dataset $D = \{\vec{x}_1, \dots, \vec{x}_n\}$, log-likel.

$$L(\vec{\eta}, D) = \sum_{i=1}^n \ln p(\vec{x}_i | \vec{\eta})$$

$$\frac{\partial L}{\partial \vec{\eta}} = \sum_i \ln \left[g(\vec{x}_i) g(\vec{\eta}) \exp(\vec{\eta}^T \vec{m}(\vec{x}_i)) \right] = n \frac{\partial}{\partial \vec{\eta}} \ln(g(\vec{\eta})) + \sum_i \vec{m}(\vec{x}_i)$$

$$\Leftrightarrow \frac{\partial}{\partial \vec{\eta}} \ln(g(\vec{\eta})) = \frac{1}{N} \sum_{i=1}^n \vec{m}(\vec{x}_i) = \underline{\underline{u(\vec{x})}}$$

$\boxed{\vec{\eta}_{ML} = \vec{F}(\vec{x})}$

empirical mean

* For Gaussian:

$$\begin{aligned} \text{(first moment)} \quad \vec{\mu}_{ML} &= E(\vec{x}) = \overline{\vec{x}} = \vec{\bar{x}} = \frac{1}{N} \sum_{i=1}^n \vec{x}_i \\ \text{(2nd moment)} \quad \vec{\Sigma}_{ML} &= E(\vec{x} \vec{x}^T) + \vec{\mu}_{ML} \vec{\mu}_{ML}^T \\ &= \frac{1}{N} \sum_{i=1}^n \vec{x}_i \vec{x}_i^T + \vec{\mu}_{ML} \vec{\mu}_{ML}^T \end{aligned}$$

3) Marginal Distrib. & condit. distrib.:

$$\vec{x} \sim N(\vec{\mu}, \Sigma)$$

$$\text{Suppose } \vec{x} \text{ split into } (\vec{x}_a, \vec{x}_b), \text{ def: } \vec{\mu} = (\vec{\mu}_a, \vec{\mu}_b), \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

Σ is symm. \rightarrow so, $\Sigma_{ab} = \Sigma_{ba}$
 $\Sigma_{aa} = \Sigma_{aa}$
:

$$\text{Marg. Dist. } p(\vec{x}_a) = N(\vec{x}_a | \vec{\mu}_a, \Sigma_{aa})$$

$$\text{conditional. } p(\vec{x}_a | \vec{x}_b) = N(\vec{x}_a | \vec{\mu}_{a|b}, \Sigma_{a|b}) \quad \text{where, } \Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

$$\vec{\mu}_{a|b} = \vec{\mu}_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\vec{x}_b - \vec{\mu}_b)$$

4) Another useful identity

$$\text{product of 2 Gaussians: } N(\vec{x} | \vec{z}, \vec{A}) N(\vec{x} | \vec{y}, \vec{B}) = \underbrace{N(\vec{x} | \vec{B}^{-1} A + \vec{B}^{-1} \vec{y}, \vec{C})}_{\text{normalizing}}$$

$$\text{with } C = (A^{-1} + B^{-1})^{-1}$$

$$\vec{z} = C(A^{-1}\vec{z} + B^{-1}\vec{y})$$

"splitting of square"

Student t-distrib.

→ with heavy tails, density falls rapidly ($\propto x^2$) for $x \rightarrow \pm \infty \Rightarrow N(\mu) \propto e^{-\frac{1}{2}\sigma^2 x^2}$
 → sometimes we need heavier tails., st $\propto |x|^{-\alpha}$ (power law)

Mixt. of Gauss.: $st(x|\mu, \sigma^2, v=2\alpha)$

1. Draw precision $\tau \sim \text{gamma}(a, b)$
2. Draw $x \sim N(\mu, \tau^{-1})$

$$\text{gamma distrib. : } \text{gamma}(\tau|a, b) = \frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau}$$

$$\text{gamma func. : } \Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du, \quad \Gamma(x+1) = \Gamma(x) \cdot x \quad (\text{factorial})$$

$$\star p(x|\mu, a, b) = st(x|\mu, \sigma^2, v) = \frac{b^a}{\Gamma(a) 2\pi} \int_0^\infty r^{a-\frac{1}{2}} \exp\left(-\left(b + \frac{1}{2}(x-\mu)^2\right)r\right) dr$$

Lab 1:

$$(1) \quad \tanh a = \frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{1 - e^{-2a}}{1 + e^{-2a}} \Rightarrow \int \phi(a) da = - \int \frac{1 - e^{-2a}}{1 + e^{-2a}} da$$

$$\text{Let, } y = e^{-2a} \Rightarrow \frac{-1}{2} \ln y = a \quad \frac{da}{dy} = \frac{-1}{2y}$$

$$= \frac{1}{2} \int \frac{1-y}{(1+y)y} dy = \frac{1}{2} \int \frac{1+y-2y}{(1+y)y} dy$$

$$= \frac{1}{2} \int \frac{1}{y} dy - \int \frac{1}{1+y} dy$$

$$= \frac{1}{2} \ln y - \ln(1+y) + C$$

$$= \frac{1}{2} \cdot -2a - \ln(1 + e^{-2a}) + C$$

$$p(a) = e^{-a - \ln(1 + e^{-2a})} + C$$

$$\int \frac{e^a - e^{-a}}{e^a + e^{-a}} da$$

$$x = e^a - e^{-a} \Rightarrow \ln x = 2a \quad \frac{\partial \ln x}{\partial a} = \frac{1}{2x} = \frac{dx}{2x}$$

$$\int \frac{1}{u} du = \ln(u) + C$$

9th Sep

2.4.2 Conjugate priors (read)

2.3.6 Bayesian inference for Gaussian:

1-dim data $X = \{x_1, \dots, x_n\}$

① Var. known, mean estimated
 σ^2 given, $\mu = ?$

$$\text{conj. prior} \Rightarrow p(\mu) = N(\mu | \mu_0, \sigma_0^2)$$

$$\text{posterior} \Rightarrow p(\mu | D) = N(\mu | \mu_n, \sigma_n^2)$$

$$\mu_n = \left(\quad \right), \frac{1}{\sigma_n^2} = \left(\quad \right)$$

② $\sigma^2 = ?$, $\mu = \text{known}$ ($\lambda = \text{precision}$)

$$\text{likelihood } p(x|\lambda) = \prod_{i=1}^n N(x_i | \mu, \lambda^{-1})$$

$$\propto \lambda^{n/2} \exp\left(-\frac{\lambda}{2} \sum_i (x_i - \mu)^2\right)$$

$$\propto \text{Gamma}(\lambda | \frac{n+1}{2}, \frac{1}{2} \sum_i (x_i - \mu)^2)$$

$$\text{prior: } p(\lambda) = \text{Gamma}(\lambda | a_0, b_0)$$

$$\text{posterior: } p(\lambda|x) = \text{Gamma}(\lambda | a_n, b_n)$$

$$a_n = a_0 + \frac{n}{2}, \quad b_n = b_0 + \frac{1}{2} \sum_i (x_i - \mu)^2$$

③ $\sigma^2 = ?$, $\mu = ?$

"normal gamma distn."

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

\downarrow prior on μ $\overbrace{\qquad \qquad}^{\text{not indep.}}$ \uparrow prior on λ
 sample λ then
 use it to sample μ .

Inform. Theory

$$\text{inform} = -\log_2(\text{prob.}) \quad (\text{in bits})$$

$$\therefore \text{inform. of } A \quad h(A) = -\log_2 p(A)$$

\rightarrow
 $= -\ln p(A)$

Shannon entropy:

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \cdot \ln_2 p(x) = -E(\ln_2 p(x))$$

$(x \text{ is in } \mathcal{X}, \mathcal{X} \text{ is discrete domain})$

[Eq]: rain flip : $X = \{0, 1\} = \mathcal{D}_X$

$$H_X = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$= 1 \text{ bit}$

if K indep. fair coins : K bits

Differential Entropy:

$$H(X) = - \int p(x) \log_2 p(x) dx \quad (x \text{ is r.v. in cont. domain, i.e., } x \in \mathbb{R}^d)$$

KL divergence: (relative entropy)

$$KL(p(x) \| q(x)) = - \int p(x) \log \frac{q(x)}{p(x)} dx$$

Imp. properties -

- $\hookrightarrow KL(\cdot) \geq 0$
- $\hookrightarrow KL(p \| q) = 0 \Rightarrow p = q$ (if p, q sufficiently regular)
- $\hookrightarrow KL(p \| q) \neq KL(q \| p)$
- $\hookrightarrow KL(p \| q) + KL(q \| r) \neq KL(p \| r)$

Eg: 2 rain tasses
(indep.)

mutual info = 0
cond. info = same as
info. of
1 coin.
 $H(X) = H(X|Y)$

Conditional entropy:

$$H(Y|X) = - \underbrace{\int p(x) \cdot \int p(y|x) \ln p(y|x) dx}_{\text{average over all values of } X}$$

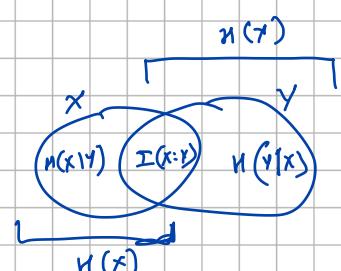
$$\text{(joint) } H(X, Y) = H(X) + H(Y|X)$$

$$p(x, y) = p(x)p(y) \leftarrow \text{Mutual info. : (shared info. betw. 2 dist.)}$$

if they are indep.

If not indep., then we want to know how much info. they share.

$$\begin{aligned} I(X:Y) &= KL(p(x,y) \| p(x) \cdot p(y)) \\ &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \end{aligned}$$



Inf. theory interpretation of MLE

$$\min_{\theta} KL(p(x) \parallel q(x|\theta))$$

↑ ↑
empirical dist. model dist.

(Ch - 8)

Probabilistic Graphical Models

- graphical way of expressing prob. dist.
 - useful when many vars / big data.
 - useful for causal reasoning / modelling
 - design / communicate stat. models.
 - cond. indep. relations are encoded in graph.
 - directed / undirected [DAG / MRF]
 - inference efficiently, calcul. cond. dist. (Eg. Bayes thm.)
- * more useful if have less edges.

Eg: directed graph. model

$$p(A, B, C) = p(C) \cdot p(B|C) \cdot p(A|B, C) \quad \text{or,} \quad p(A) \cdot p(B|A) \cdot p(C|A, B)$$



(nodes are r.v.s.)

- A depends on B & C
- B " " " C

* more useful if have less edges:

$$p(A, B, C) = p(C) \cdot p(B|C) \cdot p(A|B)$$



generalizing for N-variables :

1. determine ordering of var. (topological ordering)
2. In this ordering, call parents of x_i $p_{a,i}$: $(p_a(x_i))$

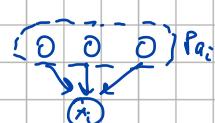
$$P(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | p_{a,i})$$

does not
have to be all,
can be a subset.

if no parent, take marginal:
 $p(x_i | \emptyset) = p(x_i)$

$$\sum p(x_i | p_{a,i}) = 1$$

$$p_{a,i} \subseteq \{1, \dots, i-1\}$$

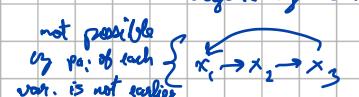


Directed graph
model
or
Bayesian netw.

$\forall j \in p_{a,i}, j \rightarrow i$ in graph
(DAG)

acyclic czg can't have a cycle.

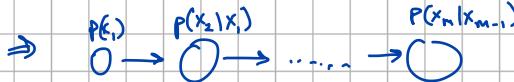
not possible
as par. of each
var. is not earlier.



Eg:

Markov chain $\Theta \rightarrow O \rightarrow \dots \Theta$

$$p(x_1 \dots x_m) = p(x_1) \prod_{i=2}^m p(x_i | x_{i-1})$$



Fully connected:



$$p(x_1 \dots x_m) = \prod_{i=1}^m p(x_i | x_1 \dots x_{i-1})$$

Assume $D_x = \{1, \dots, K\}$ (but $\leq_k = \perp$)

#params

$$p(x_1 \dots x_m) = \prod_{i=1}^m p(x_i | x_1 \dots x_{i-1})$$

 $O(K^m)$

\leq
↑
can be redundancy here

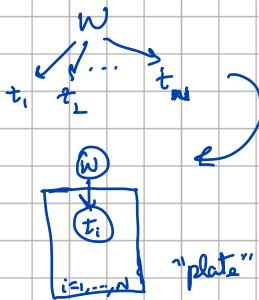
$$\left. \begin{array}{l} p(x_1) = K-1 \\ p(x_1, x_2) = K^2-1 \\ p(x_1 \dots x_m) = K^m-1 \end{array} \right\} \text{#params}$$

1) Draw param w 2) Draw N pts t_1, \dots, t_N

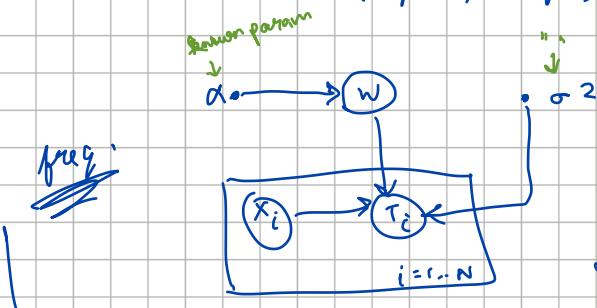
$$p(w) = \prod_{i=1}^N p(t_i | w)$$

prior

likelih.

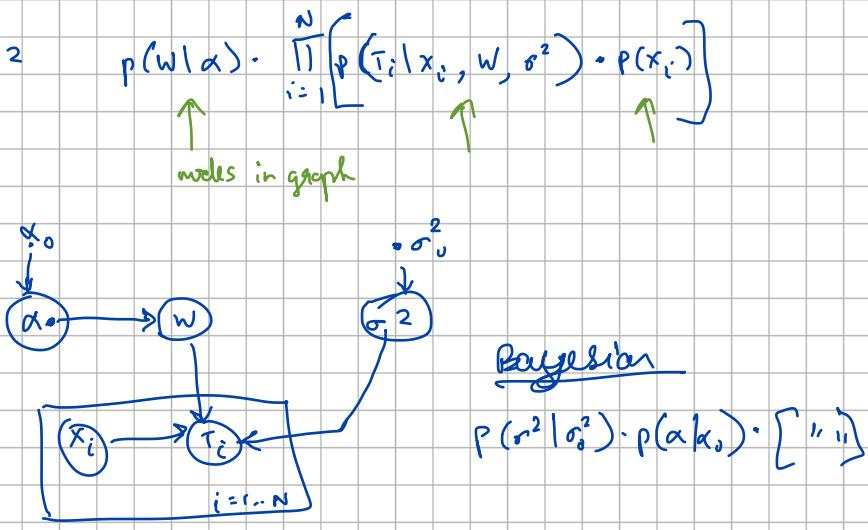


"plate" repetition of part inside plate indexed by i.

Add params, say if, $x = (x_1, \dots, x_n)$, α, σ^2 Learn w from T , conditioning on x , for known α, σ^2
(regress./classif.)

Predictive distz. :

$$p(t^* | x^*, x, \alpha, \sigma^2)$$

* add nodes for t^* & x^* 

For ass./lab :

x, y indep. if $p(x, y) = p(x) \cdot p(y)$

$\boxed{x \perp\!\!\!\perp y}$

Als valid for sets.
 $x, \vec{x}, \vec{y}, \vec{z} \dots$

x conditionally indep. of y given z if $p(x, y | z) = p(x|z) \cdot p(y|z)$

$\boxed{x \perp\!\!\!\perp y | z}$

D-separation:

if x_A is d-separated from x_B & x_C then,

$$p(x_A, x_B | x_C) = p(x_A | x_C) p(x_B | x_C)$$

i.e., $\underline{\underline{x_A \perp\!\!\!\perp x_B | x_C}}$

How to check d-separ. holds (from graph):

1) Consider all paths (seq. of nodes, connected by edges) betw. node A & a node in B

↳ - does not have to be all

↳ such that no node repeats

2) A node is blocked by x_C if:

a) it contains a collider: $\dots \rightarrow u \leftarrow \dots$

such that u is not ancestor of a node in C .

↳ transitive closure of parents
 $x \rightarrow \dots \rightarrow v$
 H is ancs. of V
 u is ult. ancs. of H

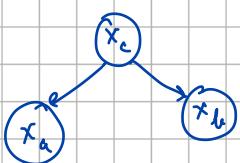
b) it contains a non-collider

$$\begin{array}{c} \dots \rightarrow u \\ \dots \rightarrow u \rightarrow \dots \\ \dots \leftarrow u \leftarrow \dots \\ \dots \leftarrow u \rightarrow \dots \\ u \leftarrow \dots \end{array}$$

such that u is in C
(when C is a set)

else $u = C$

Eg:



A is not d-sep. from $\{B\}$
 $A \leftarrow C \rightarrow B$
↳ non-collider

$\{A\}$ is d-sep. $\{B\}$ given $\{C\}$

$x_A \perp\!\!\!\perp x_B | x_C$

12th Sep

Bayesian Netw.: $\langle G, \{p(x_i | x_{P_G^{(i)}})\} \rangle$

$$\rightarrow \text{Prob. distr.}: p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_{P_G^{(i)}})$$

d-separ.

$$A \perp B | C$$

conditional / marginal independencies

$$x_A \perp\!\!\!\perp x_B | x_C$$

$$\xrightarrow{\text{markov prop.}} x_A \perp\!\!\!\perp x_B | x_C$$

$$\left. \begin{array}{l} * \text{ dsep} \Rightarrow I \\ * \text{ no dsep} \Rightarrow \text{no } I \\ = \text{ dsep} \Leftarrow I \end{array} \right\} \text{ deterministic relations}$$

$$\Rightarrow x \rightarrow y \rightarrow z \perp\!\!\!\perp z | x$$

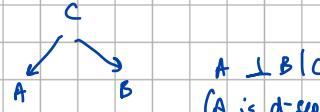
faithfulness
(holds in most cases)

$$\textcircled{1} \quad x \rightarrow y \rightarrow z \quad \text{where, } p(y|x) = \delta_{y,x} \quad (\text{if } y \text{ is a copy of } x)$$

$$y \not\perp\!\!\!\perp z | x$$

$$y \perp\!\!\!\perp z | x$$

Eg.:



$$A \perp\!\!\!\perp B | C$$

(A is d-sep. from B → C ∵ C is a non-collider → C ∨ 2-way arrows & no other path)

$$\therefore x_A \perp\!\!\!\perp x_B | x_C$$

$$p(x_A, x_B | x_C) = \frac{p(x_A, x_B, x_C)}{p(x_C)}$$

$$= \frac{p(x_C) \cdot p(A|C) \cdot p(B|C)}{p(x_C)} \quad \begin{matrix} \nearrow = A \& B \text{ indep. } | C \\ \searrow \end{matrix}$$

Q:

$$\star \quad A \not\perp\!\!\!\perp B | \emptyset \quad \text{in general} \quad x_A \not\perp\!\!\!\perp x_B | x_C$$

Eg.:

$$A \rightarrow B \rightarrow C$$

Here, $A \perp\!\!\!\perp C | B$
 $\therefore A \perp\!\!\!\perp C | \emptyset$

$$p(A, C | B) = \frac{p(A, B, C)}{p(B)}$$

$$= \frac{p(A) p(B|A) p(C|B)}{p(B)}$$

$$= p(A|B) p(C|B)$$

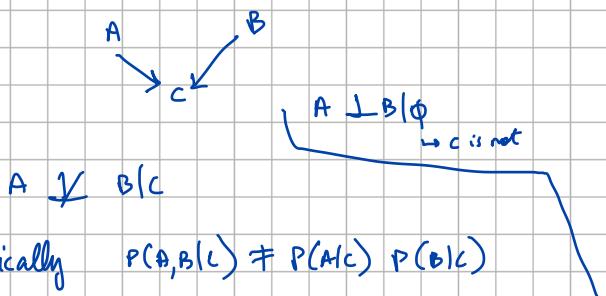
$A \not\perp\!\!\!\perp C$ so, not $p(A, C) = p(A)p(C)$

↑
if we don't
condition on B

proved

∴

Eg:



$$P(A, B) = \prod_{i=1}^n p(A_i, B_i)$$

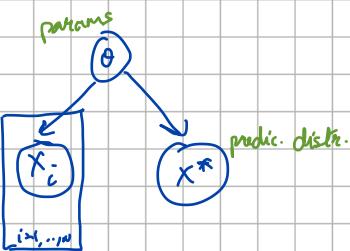
$$= \prod_{i=1}^n p(A_i)p(B_i) p(C_i | A_i, B_i)$$

$$= p(A)p(B) \prod_{i=1}^n p(C_i | A_i, B_i) = 1$$

$$= p(A)p(B)$$

$\therefore A \text{ indep. of } B.$

Eg:



$$p(\theta) \cdot \prod_{i=1}^n p(x_i^* | \theta) p(x_i | \theta)$$

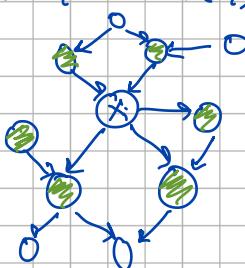
typically, $x^* \neq x_i$
 $x^* \perp\!\!\!\perp x_i$

But $x^* \perp\!\!\!\perp x_i | \theta$
 $\forall x_1..x_n \quad x^* \perp\!\!\!\perp x_i | \theta$

(Once we learn θ & fix it...)

Markov Blanket :

In Bayesian Netw., MB of node X_i : $MB(X_i) = Pa_i \cup Ch_i \cup (Pa_{Ch_i} \setminus i)$



$$P(X_i | X_{MB_i}, X_{rest}) = P(X_i | X_{MB_i})$$

↳ useful during inference.
 Need to know only X_{MB_i} of X_i .
 Not the rest.

Eg:



Undirected graph. Models (MRFs) :

MRF = $\langle G_i, \{\psi_a\}_{a \in \text{cliques}(G_i)} \rangle$
 undirected
 \downarrow
 clique potential

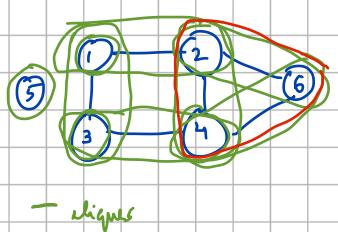
$$P(x_1 \dots x_n) = \frac{1}{Z} \prod_{a \in \text{cliques}} \psi_a(x_a)$$

\downarrow
 partition sum
 (normalizing const.)

nodes appearing
 in clique.

where, a "clique" in undir. 'G' (with nodes $\{1, \dots, N\}$)
 is a fully connected subset of nodes.

Eg:



A clique is maximal if there is no clique that strictly contains it.

Here, max. cliques: $\{5\}, \{1,2\}, \{1,3\}, \{3,4\}, \{2,4,6\}$

$$P(x_1, \dots, x_n) = \frac{1}{Z} \psi_1(x_1) \psi_{1,2}(x_1, x_2) \psi_{1,3}(x_1, x_3) \psi_{1,4}(x_1, x_4) \psi_{1,5}(x_1, x_5) \psi_{1,6}(x_1, x_6)$$

where,

$$Z = \sum_{x_1 \in D_{x_1}} \sum_{x_2 \in D_{x_2}} \dots \sum_{x_6 \in D_{x_6}} (\psi_1 \dots \psi_{1,6})$$

sum over domains

thus, $\sum_{x_1} \dots \sum_{x_6} P(x_1, \dots, x_6) = 1$

★ Calcul. of Z a bottleneck in MRF. \rightarrow use approx.

Separation:

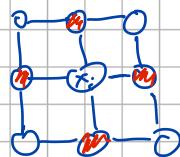
In undir. graph 2 subset of nodes A & B are "separated" given a third subset C, if each path betw. a node in A to a node in B passes through (at least) one node in C.

$$A \perp B | C$$

markov prop.: MRF $\langle G, \{\psi_i\} \rangle$ if $A \perp B | C$ in G , then $x_A \perp\!\!\!\perp x_B | x_C$ in $P(x_1, \dots, x_n) = \frac{1}{Z} \prod \psi_i$

Markov blanket: $MB(i) = \text{nbr}_i$

\downarrow
nbr. = nodes adj. to i / edge connec. with $i(i-j)$.



Hammersley-Cliff. theor.:

Given $p(\vec{x}) > 0$, then this can be represented as MRF.

Eg:

$$\vec{x} \rightarrow (x_1) \rightarrow (x_2) \rightarrow (x_3) \rightarrow (x_4) \dots$$

same cond. indep. as the directed one:

Earlier, $P(x_1) \cdot P(x_2|x_1) P(x_3|x_2) \dots$

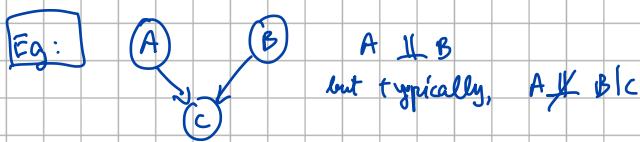
Here, $\frac{1}{Z} \underbrace{\psi_{1,2}(x_1, x_2)}_{\text{more flexible}} \underbrace{\psi_{2,3}(x_2, x_3)}_{\dots} \dots$

no big. prob.
are not uniquely

$$\tilde{\psi}_{1,2}(x_1, x_2) = \psi_{1,2}(x_1, x_2) \circ (\psi(x_2))$$

given we write,

$$\tilde{\psi}_{2,3}() = \frac{\psi_{2,3}(x_2, x_3)}{(\psi(x_2))}$$



$A \perp\!\!\!\perp B$
but typically, $A \not\perp\!\!\!\perp B | C$

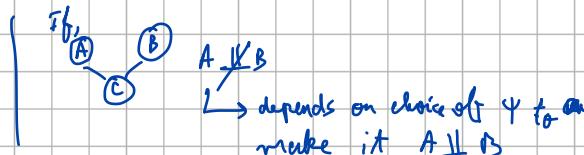
if we model this as MRF:



$A \perp\!\!\!\perp B | C$

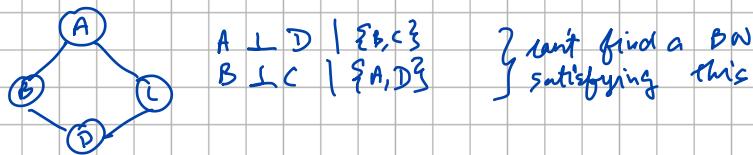
$A \perp\!\!\!\perp B | C$

no paths from A to B that pass through C
(no cross paths)

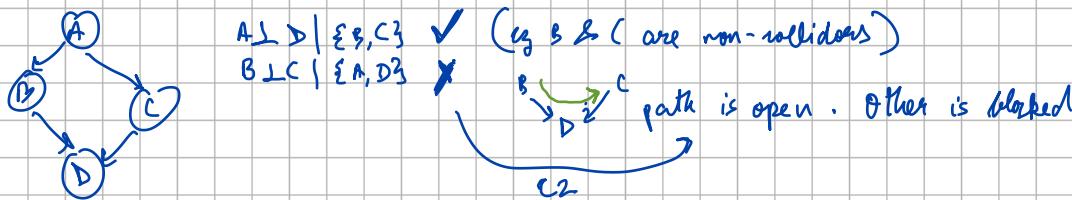


$A \perp\!\!\!\perp B$
depends on choice of ψ to make it $A \perp\!\!\!\perp B$

Eg: MRF but not BN:

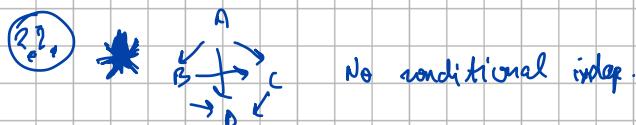


} can't find a BN satisfying this



(eg B & C are non-colliders)

B → D → C path is open. Other is blocked



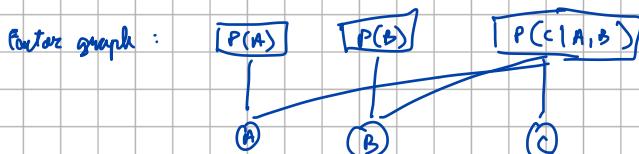
No conditional indep.

Factor graphs:

represent factorization struct. of a prob. dist.

BN: $p(A, B, C) = p(A) \cdot p(B) \cdot p(C | A, B)$ can be different factors of this cond. prob. val.

Eg: $P(B | C, A)$



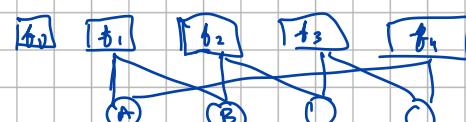
\square = factors
 \circ = variables

} bi-partite graph
(i.e., 2 types of nodes)

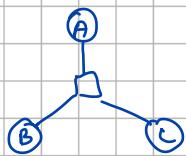
undirected edges betw. 2 node types only.

MRF: $p(A, B, C, D) = \frac{1}{Z} \cdot \psi_{A, B}(A, B) \cdot \psi_{B, C}(B, C) \cdot \psi_{C, D}(C, D) \cdot \psi_{D, A}(D, A)$

$$\underbrace{\psi_1}_{t_1}, \underbrace{\psi_2}_{t_2}, \underbrace{\psi_3}_{t_3}, \underbrace{\psi_4}_{t_4} \xrightarrow{\text{Z}} \frac{1}{Z} \cup \cup \cup \cup$$

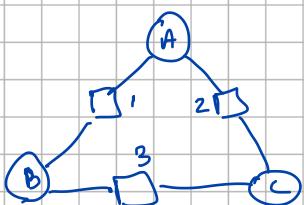


Eg:



$$P(A, B, C) = f(A, B, C)$$

[uninformative]



$$P(A, B, C) = f_1(A, B) \cdot f_2(B, C) \cdot f_3(B, C)$$

↳ Factor grp.

same thing in MRF will be in terms of clique potential.

$$P(A, B, C) = \frac{1}{Z} \Psi_{A, B, C}(A, B, C)$$

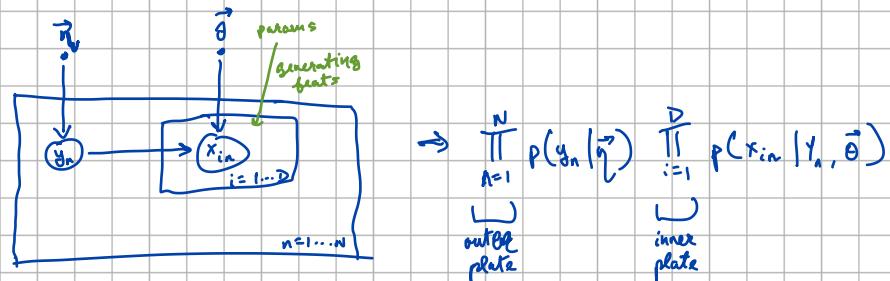
uninform. compared

to factor graph. Less expressive.

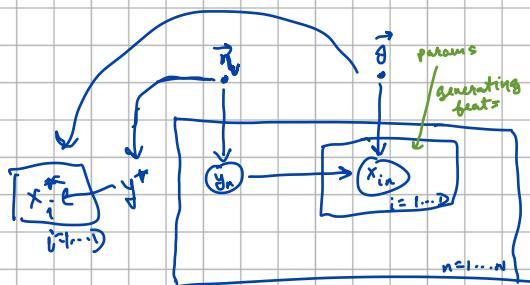
Naive Bayes Classif. :

labels : y_n

feats : x_{in}, \dots, x_{Dn} $n=1 \dots N$

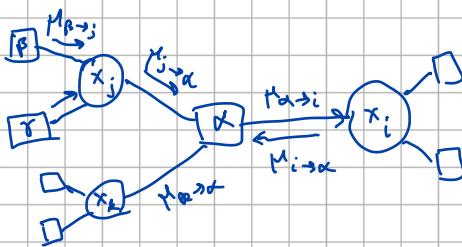


In N.B., $x_{in} \perp\!\!\!\perp x_{in'} | y_n$



19th Sep

Bishop 8.4.4 Sum product Algo / Belief propog.



Factor \rightarrow variable messages :

$$f_{\alpha \rightarrow i}(x_i) = \sum_{f_{\alpha}(x_\alpha)} \prod_{j \in \text{neigh}(i)} M_{j \rightarrow i}(x_j)$$

$$x_\alpha = (x_i, x_j, x_k)$$

$$x_\alpha := (x_i)_{i \in \text{neigh}(\alpha)}$$

variable \rightarrow factor messages :

$$M_{j \rightarrow \alpha}(x_j) = \prod_{i \in \text{neigh}(j) \setminus \alpha} \mu_{i \rightarrow j}(x_i)$$

messages = $2 \times$ no. of factors in graph.

variable beliefs :

$$p(x_\alpha) = \sum_{x_\alpha} p(x) = \frac{1}{Z} \prod_{i \in \text{neigh}(\alpha)} M_{\alpha \rightarrow i}(x_i) = \frac{\prod_{i \in \text{neigh}(\alpha)} M_{\alpha \rightarrow i}(x_i)}{\sum_{x_\alpha} \prod_{i \in \text{neigh}(\alpha)} M_{\alpha \rightarrow i}(x_i)}$$

either calc.
global normaliz.

$$p(x) = \frac{1}{Z} \prod_{\alpha} p(x_\alpha)$$

factor beliefs :

$$p(x_\alpha) = \frac{1}{Z} f_\alpha(x_\alpha) \prod_{i \in \text{neigh}(\alpha)} M_{i \rightarrow \alpha}(x_i)$$

PRO :
So from exp. in # nodes
to # edges

complexity : $O(2EK^m)$ where, $E =$ edges

$|D_{x_i}| \leq K$ (domains repeat bounded by K)

$M =$ max. variables in the = $\max_i |x_i| \leq M$
factors

(for eg, a tree has 3. if no other has more than 3 then $M=3$)

CONS:

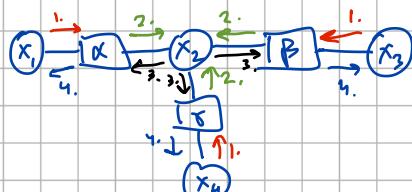
Loops of factors OK \rightarrow can be absorbed

Loops of vars. not OK in this algo. (if loops then its a brute approx.)

if trees/forests
(i.e., no loops)
exact result
loop belief propog.

* alternate for exact inference : var. elimination.

Eg:



$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} f_L(x_1, x_2) f_P(x_2, x_3) f_R(x_2, x_4)$$

Step 1 : Start at leaf node

$$\mu_{i \rightarrow \alpha}(x_1) = 1 \quad (\text{eg no neighbor})$$

$$\mu_{\beta \rightarrow \beta}(x_2) = 1$$

$$\mu_{4 \rightarrow 8}(x_3) = 1$$

Step 2 :

$$\mu_{x \rightarrow 2}(x_2) = \sum_{x_1} f(x_1, x_2) \mu_{i \rightarrow \alpha}(x_1)$$

$$\mu_{p \rightarrow 2}(x_2) = \sum_{x_3} f_p(x_2, x_3) \mu_{3 \rightarrow p}(x_3)$$

$$\mu_{8 \rightarrow 2}(x_2) = \sum_{x_4} f_8(x_2, x_4) \mu_{4 \rightarrow 8}(x_4)$$

Incoming done (x_2) ... Now start with outgoing :

Step 3 :

$$\mu_{2 \rightarrow \alpha}(x_2) = \mu_{p \rightarrow 2}(x_2) \cdot \mu_{8 \rightarrow 2}(x_2) \quad (\text{product of messages})$$

$$\mu_{2 \rightarrow p}(x_2) = \mu_{\alpha \rightarrow 2}(x_2) \cdot \mu_{8 \rightarrow 2}(x_2)$$

$$\mu_{2 \rightarrow 8}(x_2) = \mu_{\alpha \rightarrow 2}(x_2) \cdot \mu_{p \rightarrow 2}(x_2)$$

$$\mu_{\alpha \rightarrow 1}(x_1) = \sum_z f_\alpha(x_1, x_2) \mu_{2 \rightarrow \alpha}(x_2)$$

$$\mu_{p \rightarrow 3}(x_3) = \sum_z f_p(x_3, x_2) \mu_{2 \rightarrow p}(x_2)$$

$$\mu_{8 \rightarrow 4}(x_4) = \sum_z f_8(x_4, x_2) \mu_{2 \rightarrow 8}(x_2)$$

(checking)

$$p(x_2) = \frac{1}{2} \underbrace{\mu_{\alpha \rightarrow 2}(x_2)}_{\substack{\text{substitute}}} \underbrace{\mu_{p \rightarrow 2}(x_2)}_{\substack{\text{substitute}}} \underbrace{\mu_{8 \rightarrow 2}(x_2)}_{\substack{\text{substitute}}}$$

$$= \frac{1}{2} \sum_z \sum_z \sum_z f_\alpha \cdot f_p \cdot f_8$$

$$p(x_2 | x_3 = s) = ?$$

Introduce evidence factors :

$$\text{Give F.G. } p(x_1 \dots x_d) = \frac{1}{2} \prod_a f_a(x_a)$$

$$\text{then, } p(x_i | x_j = \xi_j) = ?$$

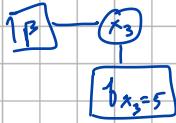
$$\text{Evidence factor } f_{\xi_j}(x_j) = S_{x_j, \xi_j} = \begin{cases} 1, & x_j = \xi_j \\ 0, & x_j \neq \xi_j \end{cases}$$

$$p(x_i | x_j = \xi_j) = \frac{p(x_i, x_j = \xi_j)}{\sum_{x_i} p(\cdot)} \propto p(x_i, x_j = \xi_j) = \sum_{x_j} p(x_i, x_j) S_{x_j, \xi_j}$$

$$= \sum_{x_j} p(x_i, x_j) \cdot f_{\xi_j}$$

$$\text{Soln: extend F.G. } \tilde{p}(x_1 \dots x_d) = \frac{1}{2} \left[\prod_a f_a(x_a) \cdot f_{\xi_j}(x_j) \right]$$

In our eg : extend the graph as :



} run the same alg. again on this F.G.

(intuitively same as not summing over all ξ_j)

($E \# P_2$)

Max-Sum algo :

want to calcul. $x^* = \arg \max_x \frac{1}{Z} \prod_a f_a(x_a)$

State that gives max. prob.

drop c const.
not affects arg.max.

[joint assignment of x_i, x_j & x_k that maximizes values]

$$p(x^*) = \frac{1}{Z} \max_x \prod_a f_a(x_a)$$

value of max. prob.

↓ in log domain

$$\log p(x^*) = -\log Z + \max_x \log \left[\prod_a f_a(x_a) \right]$$

$$= -\log Z + \max_x \sum_a \log f_a(x_a) \quad \left. \begin{array}{l} a + \max(a, c) = \max\{a+b, a+c\} \\ \text{i.e., replacing prod. by sums} \end{array} \right\}$$

sum-prod. algo → max-sum algo.

helps us to sum log more factors → log factors

Factor → variable messages ?

$$\mu_{a \rightarrow i}(x_i) = \sum_{x_{a \setminus i} \in D_{a \setminus i}} f_a(x_a) \prod_{j \neq i} \mu_{j \rightarrow a}(x_j)$$

$$\nu_{a \rightarrow i}(x_i) = \max_{x_{a \setminus i}} \log f_a + \sum_{j \neq i} \nu_{j \rightarrow a}(x_j)$$

$$\text{likewise, } \nu_{j \rightarrow a}(x_j) = \sum_{p \in \text{neigh}(j) \setminus a} \nu_{p \rightarrow j}(x_j)$$

max beliefs / max marginals :

$$q_i(x_i) = \sum_{a \ni i} \nu_{a \rightarrow i}(x_i) - \log Z$$

↓ used for

given max-marginals ($q_i(x_i)$)
 $i = 1, \dots, d$),
run Viterbi alg.

(Bishop 8.4.5)

to get global optimum x^*

$$x^* = \arg \max_x \prod_a f_a(x_a)$$

if $q_i(x_i)$ has unique max. then can use:

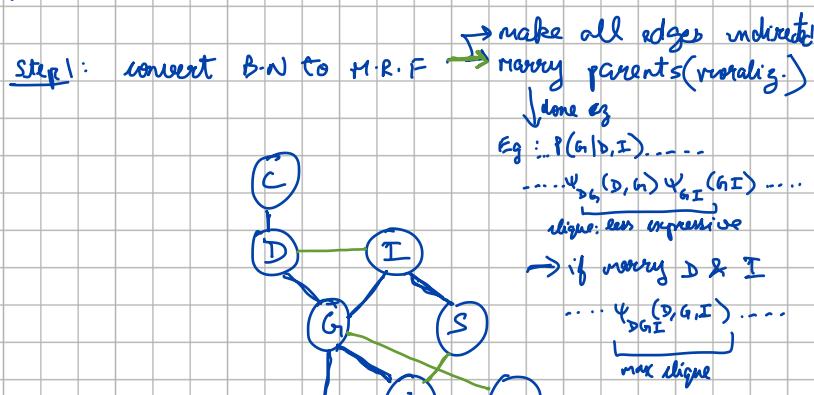
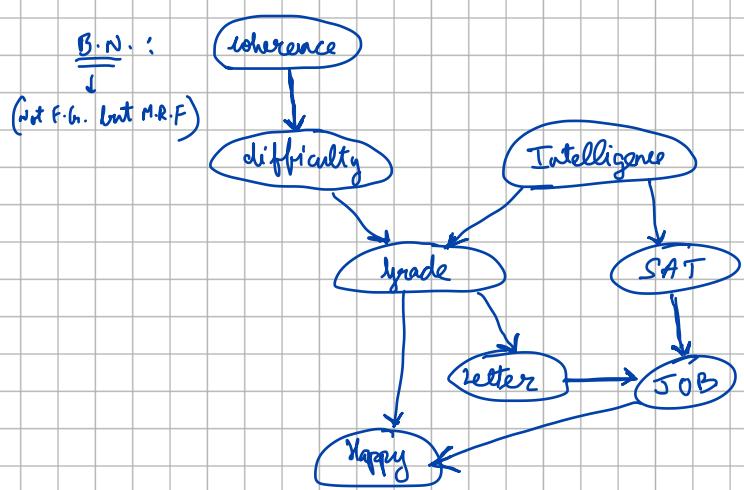
don't have to run Viterbi

$$x_i^* = \arg \max_{x_i} q_i(x_i)$$

$$\text{For eg : } |D_{x_i}| = 5 : \quad p_i(x_i) = (6, 8, 1, 2, \pi) \\ \downarrow \\ = x^* = 2$$

Variable elimin. algot : (b. 8.4.6 exact inference on general graphs)

Ch-20, Murphy for var. elimin. & tree junc. algos.



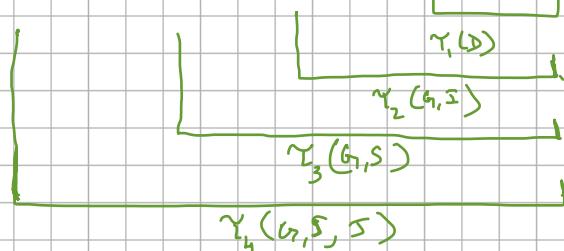
$$\text{Step 2} : P(J) = \sum_L \sum_S \sum_{G,H} \sum_I \sum_D \sum_C p(L, S, G, H, I, D, C)$$

$$= \sum_L \sum_S \sum_{G,H} \sum_I \sum_D \sum_C \underbrace{\psi_c(C)}_{\text{distrb.}} \underbrace{\psi_d(D)}_{\tilde{\psi}_d(C, D)} \underbrace{\psi_i(I)}_{\tilde{\psi}_i(G, I, D)} \underbrace{\psi_g(G, I, D)}_{\tilde{\psi}_g(G, I, D)} \underbrace{\psi_s(S, I)}_{\tilde{\psi}_s(S, I)} \underbrace{\psi_h(H, G)}_{\tilde{\psi}_h(S, L, S)} \underbrace{\psi_l(L, G)}_{\tilde{\psi}_l(S, L, S)} \underbrace{\psi_j(J, L, S)}_{\tilde{\psi}_j(L, S, J)} \underbrace{\psi_n(H, G, J)}_{\tilde{\psi}_n(H, G, J)}$$

using distributive law:

Choose an order for elimin. : C, D, I, H, G, S, L } put them to right as far as possible

$$= \sum_L \sum_S \underbrace{\psi_j(S, L, S)}_{\tilde{\psi}_j(S, L, S)} \underbrace{\psi_h(H, G)}_{\tilde{\psi}_h(H, G, S)} \underbrace{\psi_n(H, G, S)}_{\tilde{\psi}_n(H, G, S)} \underbrace{\psi_s(S, I)}_{\tilde{\psi}_s(S, I)} \underbrace{\psi_g(G, I, D)}_{\tilde{\psi}_g(G, I, D)} \underbrace{\psi_d(D)}_{\tilde{\psi}_d(C, D)}$$



→ only depends on D on summed over C.

$$\tilde{\psi}_5(L, S, J)$$

$$\tilde{\psi}_6(\gamma_5)$$

$$\tilde{\psi}_7(J)$$

= what we wanted

<u>Step</u>	<u>Elim.</u>	<u>Factors Used</u>	<u>Vari. induced</u>	<u>New factor</u>
1	C	$\tilde{\psi}_c = \psi_c \cdot \psi_d$	C, D	$\tilde{\psi}_1(D)$
2	D	$\tilde{\psi}_1, \tilde{\psi}_6 = \psi_6 \cdot \psi_i$	G, I, D	$\tilde{\psi}_2(L, I)$
3	I	$\tilde{\psi}_2, \tilde{\psi}_3$		

$$P(J) \propto \tilde{\psi}_1(J)$$

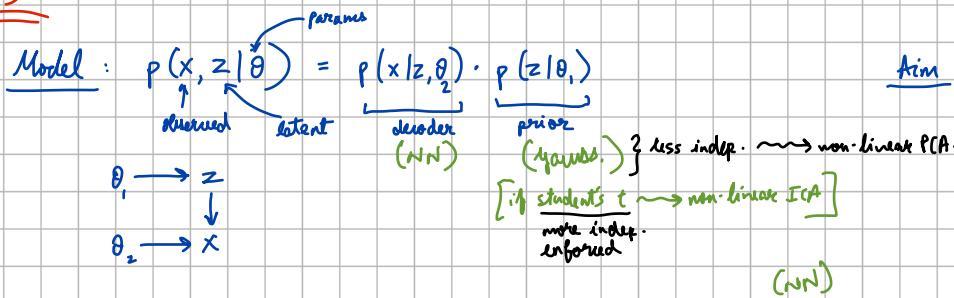
$$= \tilde{\psi}_7(J)$$

$$\uparrow$$

$$z=1$$

26th Sep

VAE's :



Use variational inference. Hence have approx. posterior $q(z|x, \lambda) \approx p(z|x, \theta_1, \theta_2)$

encoder

Thus, we get an exact repres. of decoder but only an approx. of encoder.

Log-likel. : $L(\theta) = \ln p(x|\theta) = \ln \int p(x|z, \theta) p(z|\theta) dz$

hard to approx. \therefore var. infer.

ELBO : $L(\theta, \lambda) = \mathbb{E}_{q(z|x, \lambda)} [\ln p(x|z, \theta_2) + \ln p(z|\theta_1)] + H(q(z|x, \lambda)) \leq L(\theta)$

diffic. integrat. so approx. using sampling.
so optimize this hard to optimize

optimize $\max_{\theta, \lambda} L(\theta, \lambda)$ using SGD

what we actually optimize is $\max_{\theta, \lambda} \hat{L}(\theta, \lambda) = \frac{1}{K} \sum_{k=1}^K [\ln p(x|z^k, \theta_2) + \ln p(z^k|\theta_1)]$

where, $z^k \sim q(z|x, \lambda)$
(k times)

Use: $\mathbb{E}_{q(z|x, \lambda)} f(z) \approx \frac{1}{K} \sum_{k=1}^K f(z^k)$

$z^k \sim q(z|x, \lambda)$
sample

reparam:
 $\begin{cases} z = g_\lambda(x, \varepsilon) \\ \varepsilon \sim p(\varepsilon) \end{cases}$
 like change of variable

$\approx \frac{1}{K} \sum_{k=1}^K f(g_\lambda(x, \varepsilon^k)), \quad \varepsilon^k \sim p(\varepsilon)$

decoder : $p(x|z, \theta_2) = \mathcal{N}(x | \mu_{\theta_2}(z), \Sigma_{\theta_2}(z))$

where, μ & Σ are opf of MLP
with i/p z & wts. θ_2 .

encoder :
The approx. also assumed gaussian:
 $q(z|x, \lambda) = \mathcal{N}(z | \mu_\lambda(x), \Sigma_\lambda(x))$

where, μ & Σ are opf of a / MLP
with i/p x & wts. λ .

$g_\lambda(x, \varepsilon) = \mu_\lambda(x) + \sigma_\lambda(x) \odot \varepsilon$

fixed (constr.)
multi-variate

also gives less variance in grads.
So faster convergence

$\begin{bmatrix} \sigma_{\lambda,1}^2(x) & \dots & \phi \\ \vdots & \ddots & \vdots \\ \phi & \dots & \sigma_{\lambda,D}(x) \end{bmatrix}$

Training procedure :

- sample minibatch (x_1, \dots, x_m) from data
- for $m=1, \dots, M$ sample $z_n^k \sim q(z|x_n, \theta^t)$ (use reparam. trick)
actually, $z_n^k = g_{\theta^t}(x_n, e^k)$, $e^k \sim p(e)$ $\overset{K \text{ samples}}{\uparrow}$
- compute $\hat{L}(\theta^t, \gamma^t)$
- compute $\nabla_{\theta, \gamma} \hat{L}(\theta^t, \gamma^t)$ on all x_m, z_n^k $m=1 \dots M$ $n=1 \dots K$
- SGD update:

$$\theta^{t+1} \leftarrow \theta^t + \alpha \nabla_{\theta} \hat{L}(\theta^t, \gamma^t)$$

$$\gamma^{t+1} \leftarrow \gamma^t + \alpha \nabla_{\gamma} \hat{L}(\theta^t, \gamma^t)$$

30th Sep

Sampling Methods :

Monte-Carlo :

Goal: compute $E_{p(x)}[f(x)] = \int f(x) p(x) dx / \int f(x) dx$

apprx. $\begin{cases} \text{prediction: } p(y^*|x^*) = \int p(y^*|x^*, \theta) p(\theta|x^*) d\theta \\ \text{Estim. evidence: } p(y|x) = \end{cases}$

① Regular sampling : (\wedge or $<$ means estim. this value)

- 1) draw N samples from $p(z)$ $z_i \sim p(z)$ $i=1, \dots, N$
- 2) Calcul. $E[f] \approx \hat{f} = \frac{1}{N} \sum_i f(z_i)$ as $N \rightarrow \infty$ $\hat{f} \rightarrow E[f]$

properties:

\rightarrow Unbiased estimator: $E[\hat{f}] = E[f]$
 $E[\hat{f}] = E\left[\frac{1}{N} \sum_i f(z_i)\right] = \frac{1}{N} \sum_i E[f(z_i)] = E[f]$

\rightarrow Variance:

$$\text{var}[\hat{f}] = E[(\hat{f} - E[\hat{f}])^2]$$

$$\hookrightarrow = \text{var}\left[\frac{1}{N} \sum_i f(z_i)\right]$$

$$= \frac{1}{N^2} \sum_i \text{var}[f(z_i)]$$

$$= \frac{1}{N^2} N \sum_i \text{var}[f(z_i)]$$

$$= \frac{1}{N} \text{var}[f(z)]$$

more samples, less var.

$(\text{var}(a+b) = \text{var}(a) + \text{var}(b) \text{ if } a, b \text{ indep.})$
 Here we have indep. samples!

$$(\text{var}(X) = \sigma^2 \text{var}(x))$$

\Rightarrow RMS error = $\sqrt{\text{var}[\hat{f}]} = \frac{1}{N} \sqrt{\text{var}[f(z)]} \quad \downarrow 0 \text{ as } N \rightarrow \infty$

Discrete Random Variables Sampling

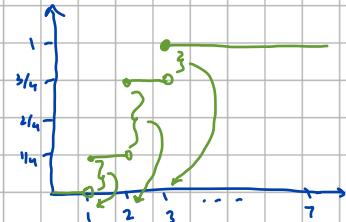
Say $z \in \{1, 2, \dots, K\}$. Given $p(z)$, how to sample from $p(z)$?

Steps:

1) Calcul. CDF/CMF

$$\begin{aligned} K &= 3 \\ p(z=1) &= 1/4 \\ p(z=2) &= 2/4 \\ p(z=3) &= 1/4 \end{aligned}$$

$$\begin{aligned} p(z \leq 0) &= 0 \\ p(z \leq 1) &= 1/4 \\ p(z \leq 2) &= 3/4 \\ p(z \leq 3) &= 1 \\ p(z \leq 4) &= 1 \end{aligned}$$



2) Sample from uniform $u_i \sim U(0, 1)$

$$\text{Eg: } u = 0.35 \rightarrow z_i = 2$$

3) Look up corrsp. value of z_i

$$\begin{aligned} 0 \leq u_i \leq 1/4 &\quad z_i = 1 \\ 1/4 \leq u_i \leq 3/4 &\quad z_i = 2 \\ 3/4 \leq u_i \leq 1 &\quad z_i = 3 \end{aligned}$$

Cont. random sampling

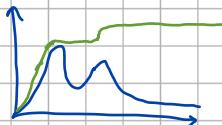
variables in \mathbb{R} with density $p(z)$

1) Calcul. CDF $p(z \leq \varepsilon) = \int_{-\infty}^{\varepsilon} p(z) dz = F(\varepsilon)$ $F: \mathbb{R} \rightarrow [0, 1]$

2) Calcul. F^{-1}

3) sample $u_i \sim U(0, 1)$

4) $z_i \sim F^{-1}(u_i)$



$$F^{-1}: [0, 1] \rightarrow \mathbb{R}$$

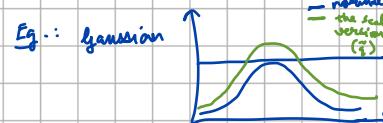
(2) Rejection Sampling

Above was for low dimen. If high dimens. then \int can't calcul.

given: unnormalized density $\tilde{p}(z) \propto p(z)$

1) Find proposal distribs. \tilde{q} (have unnormalized density $\tilde{q}(z) \propto q(z)$ that we can sample from such that:)

$$\begin{aligned} \rightarrow \int \tilde{q}(z) dz &< \infty \\ \rightarrow \tilde{p}(z) &\leq \tilde{q}(z) + \epsilon \end{aligned}$$

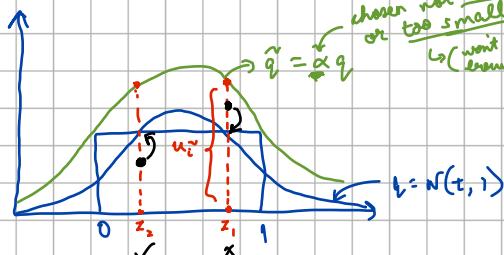


2) For $i = 1, \dots, N$: \rightarrow sample $z_i \sim q$ $\leftarrow (\propto \tilde{q})$
 \rightarrow sample $u_i \sim U(0, \tilde{q}(z_i))$ \leftarrow has to be \tilde{q}
 \rightarrow if $\tilde{p}(z_i) > u_i$ keep sample z_i
 else throw away z_i & re-sample

(Rejected many)
 chosen not too big
 or too small
 (so won't upper bound of \tilde{q})

$$3) \cdot g(z_i) \cdot \frac{\tilde{p}(z_i)}{\tilde{q}(z_i)} \propto p(z_i)$$

* in high dimens., probab. to reject samples is high.



3 Importance Sampling

In above we have to find the upper bound. Here, note -

- Need:
 - unnorm. density $\tilde{p}(z)$ to sample from
 - unnorm. proposal density $\tilde{q}(z)$ that we can sample from.
↳ (does not need to upper bound \tilde{p})

$$p(z) = \frac{\tilde{p}(z)}{Z_p}, \quad q(z) = \frac{\tilde{q}(z)}{Z_q}$$

for $i=1, \dots, n$:

1) Sample $z_i \sim q$

2) Calcul. weight $\frac{\tilde{p}(z_i)}{\tilde{q}(z_i)} = w_i \quad \left\{ \begin{array}{l} \text{if } \tilde{p} > \tilde{q} \text{ (high var.)} \\ \text{if } \tilde{q} \approx 0 \text{ but } \tilde{p} > 0 \Rightarrow \text{dividing by small no.} \end{array} \right.$

3) Calcul. $E[f] = \frac{\sum w_i f(z_i)}{\sum w_i} \quad (\text{in 1, } w_i=1)$

$$\begin{aligned} E[f] &= \int p(z) f(z) dz = \int q(z) \frac{p(z)}{q(z)} f(z) dz = \frac{\int q(z) \frac{\tilde{p}(z)}{\tilde{q}(z)} \cdot \frac{z}{Z_p} f(z) dz}{\int q(z) \frac{\tilde{p}(z)}{\tilde{q}(z)} dz} \approx \frac{\frac{1}{n} \sum_i f(z_i) w_i}{\frac{1}{n} \sum_i w_i} \end{aligned}$$

↳ unnormalized densities
can cancel these off

Ans.

* Imp. Sampling faster than MCMC. (no throwing away of samples, use all)

* But imp. samp. has higher var. than MCMC samp.

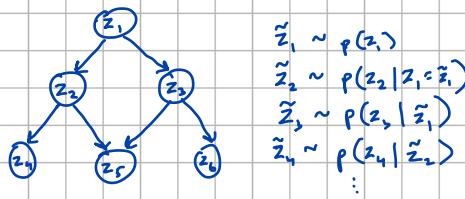
* Imp. sampling also not so good in high dimens. (diffic. to find dens. q , s.t. both p & q are large together)

4 Ancestral Sampling

Given Bayesian NN, if want to sample from it topological ordering z_1, \dots, z_d (i.e., $z_j \leq z_i$ if $j \in \text{pa}(i)$)

For $k=1, \dots, d$ (#variables)

1) Sample $z_k \sim p(z_k | z_{\text{pa}(k)})$



5 MCMC Sampling



Given target distn. $p(x)$ that we want to sample from, setup a Markov chain s.t. $p(x_n) \rightarrow p(x)$ as $n \rightarrow \infty$ i.e., s.t. $p(x)$ is its equilibrium distrib..

1) $x_i \sim q(x_i)$ (q = proposal distrib.)

2) $x_i \sim T(x_i | x_i)$

← Transition kernel T or trans. prob. $T(x_i | x_i)$

3) $q(x_i, x_2) = q(x_i) T(x_2 | x_i)$

$$q(x_i) = \int q(x_i) T(x_2 | x_i) dx_2$$

$$x_i \sim T(x_i | x_{i-1})$$

$$q(x_1, \dots, x_n) = q(x_1) \prod_{i=2}^n q(x_i | x_{i-1})$$

$$= q(x_1) \prod_{i=2}^n T(x_i | x_{i-1})$$

Suppose after convergence: if $x_n \sim p(x)$ then like to have $x_{n+1} \sim p(x)$ ("invariance")

$$p(x_{n+1}) = \int T(x_{n+1} | x_n) p(x_n) dx_n$$

this gives a constraint
on T .

eigen value λ_1 with $\lambda = 1$

Detailed balance implies invariance:

$$\begin{aligned} T(x_{n+1} | x_n) p(x_n) &= T(x_n | x_{n-1}) p(x_{n-1}) \\ \Rightarrow T(y|x) p(x) &= T(x|y) p(y) \quad \forall x, y \end{aligned}$$

detailed balance property (in equilib. state)

$$\int T(x_{n+1} | x_n) p(x_n) dx_n = \int T(x_n | x_{n+1}) p(x_{n+1}) dx_n = 1 \cdot p(x_{n+1})$$

One can show: given Markov chain with initial distrib. q & trans. kernel T

then p is unique equilib. distrib. of Markov chain.

$$\left\{ \begin{array}{l} p \text{ is invariant} \\ \text{ergodicity holds} \end{array} \right. \quad (T(y|x) > 0 \quad \forall x, y)$$

3rd Oct

MCMC Sampling: either in high dim. or running along many $Q(x_{t+1}|x_t)$ to estim. target dist. without relying too much on its know. to come up with 1 single good proposal dist.

Detailed balance:

$$p^*(x_{t+1}) = \int \underbrace{T(x_{t+1} | x_t)}_{\sum} p^*(x_t) dx_t \quad \forall x_{t+1}$$

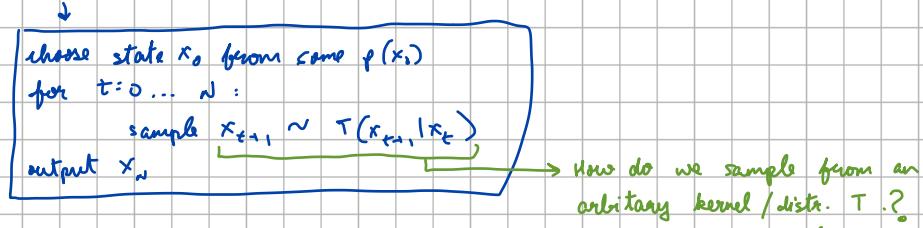
$$p^*(x_t) T(x_{t+1} | x_t) = p^*(x_{t+1}) T(x_t | x_{t+1})$$

Valid kernels:

If T_1 & T_2 valid:

- 1) $\alpha T_1 + (1-\alpha) T_2$ ($\alpha \in [0, 1]$) \Rightarrow with prob. α choose T_1 at 't' & $1-\alpha$ choose T_2 .
- 2) $T_2 \circ T_1$ (composition)
i.e., $T_3 = T_2 \circ T_1$: $T_3(x_{t+1} | x_t) = \int T_2(x_{t+1} | x_t') T_1(x' | x_t) dx'$

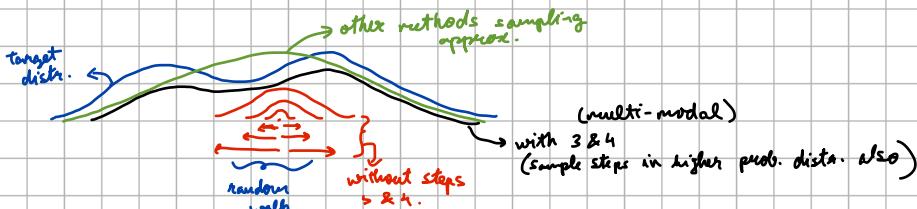
The MCMC sampler uses those transition prob.:



Metropolis Hastings algo : (inner loop of MCMC sampler)

Given a proposal trans. kernel $\theta(x_{t+1}|x_t)$ (Something that we can sample from)
 To sample $x_{t+1} \sim T(\cdot)$:

- 1) sample $\tilde{x}_{t+1} \sim Q(\cdot)$
 - 2) compute acceptance prob. $\alpha(x_{t+1} | x_t) = \min\left(1, \frac{p(\tilde{x}_{t+1}) \& (x_t | \tilde{x}_{t+1})}{p(x_t) Q(\tilde{x}_{t+1} | x_t)}\right)$
 - 3) sample $u_t \sim U(0,1)$
 - 4) Σ_f $u \leq \alpha$: accept the sample $\rightarrow x_{t+1} = \tilde{x}_{t+1}$
 else : reject it but don't throw it away, use it as new state : $x_{t+1} = \tilde{x}_t$



Sketch that is satisfies detailed balance:

* if ^{new} sample is accepted - ($x_{t+1} = \tilde{x}_{t+1}$)

$$\begin{aligned}
 p^*(x_t) & \otimes (x_{t+1}|x_t) \min \left(\underbrace{\dots}_{\text{start in } x_t} \underbrace{\dots}_{\text{happens with prob. } \alpha} \right) = \min \left(p^*(x_t) \otimes (x_{t+1}|x_t), p^*(x_{t+1}) \otimes (x_t|x_{t+1}) \right) \\
 & = \min \left(p^*(x_t) \otimes (x_t|x_{t+1}), p^*(x_{t+1}) \otimes (x_{t+1}|x_t) \right) \quad \xrightarrow{x_t \leftrightarrow x_{t+1}} \text{to prove detail level.} \\
 & = p^*(x_{t+1}) \otimes (x_t|x_{t+1}) \min \left(1, \underbrace{p^*(x_t) \dots}_{\text{proved}} \right)
 \end{aligned}$$

3) combine random walk with another algo that can take longer steps.

6 Gibbs Sampler

(good in high dim.)

$$\text{D-dimen. : } \vec{x}^t = (x_1^t, \dots, x_D^t) \quad (x^t = \text{Tageszeit-Niveau})$$

* can change ordering. If ordering is fixed detailed balance does not hold.

* Here we are sampling from conditionals (not joint target dist.) & this sometimes possible. E.g.:

$$\begin{aligned} l_{step-in} &= \left\{ \begin{array}{l} x_1^{t+1} \sim p(x_1 | x_1^{t+1}, x_2^t, \dots, x_D^t) \\ x_2^{t+1} \sim p(x_2 | x_1^{t+1}, x_2^{t+1}, x_3^t, \dots, x_D^t) \\ x_3^{t+1} \sim p(x_3 | x_1^{t+1}, x_2^{t+1}, x_3^{t+1}, \dots, x_D^t) \\ \vdots \\ x_D^{t+1} \sim p(x_D | x_1^{t+1}, \dots, x_{D-1}^{t+1}) \end{array} \right. \\ &\text{MC} \\ &\downarrow \\ &\text{repeat} \end{aligned}$$

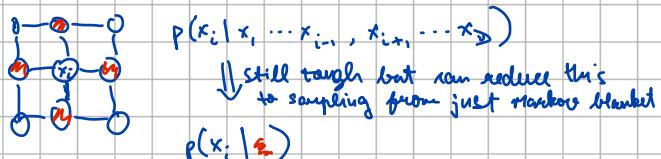
NOTE :

$$x \sim N(\mu, \sigma^2) \quad , \quad p(x) = \underbrace{N(x | \mu, \sigma^2)}_{\text{density}}$$

↓

$$\text{not } p(x) = N(\mu, \sigma^2)$$

func. of x but no x
in the R.H.S.



10th Oct

EM for HMMs

1. Calcul. initial values θ^{old} with $\theta = (\pi, A, \phi)$

2. Iterate until convergence:

E-step : (calcul. $\alpha(\theta, \theta^{\text{old}})$)

Baum-Welch algo (message passing step) {

- 1) Run forward recursion to calcul. $\alpha(z_1), \dots, \alpha(z_n)$
- 2) Run backward pass recursion $\beta(z_n), \dots, \beta(z_1)$
- 3) (calcul. suff. stats) $\tau(z_n), \varepsilon(z_{n+1}, z_n)_{n=1, \dots, N}$, $p(x | \theta^{\text{old}})$

M-step : Calcul. $\theta^{\text{new}} = \arg\max_{\theta} Q(\theta, \theta^{\text{old}})$

$$\pi^{\text{new}} : (13 \cdot 18)$$

$$A^{\text{new}} : (13 \cdot 19)$$

$$\text{ratig. emissions} : (13 \cdot 20, 13 \cdot 21)$$

$$\text{lyric. emissions} : (13 \cdot 23)$$

$$\text{other emissions} : Q(\theta, \theta^{\text{old}}) = \dots$$

Viterbi algo (max-product for HMMs)

Calcul. most prob. seq. of latent states given x .



if observed all the x then
can get the most probable latent seqs.



$$h(z_1) = p(z_1 | \pi) p(x_1 | z_1, \phi)$$

$$f_n(z_1, \dots, z_n) = p(z_n | z_{n-1}, A) p(x_n | z_n, \phi)$$

def $\omega(z_n) = \max_{z_{n+1} \rightarrow z_n} z_{n+1} \quad z_{n+1} \in \mathcal{Z} \quad \left. \begin{array}{l} \text{do forward pass } (\omega \text{ messages}) \text{ but no backwd. pass.} \\ \text{Instead do dynamic programming.} \end{array} \right\}$

$$\omega(z_1) = \ln p(z_1) + \ln p(x_1 | z_1)$$

for $n = 1 \dots N-1$

$$\omega(z_{n+1}) = \ln p(x_{n+1} | z_{n+1}) + \max_{z_n} [\ln p(z_{n+1} | z_n) + \omega(z_n)]$$

fn

the z_n that gave this max
 incoming message

$$\dashrightarrow \Psi_n(z_{n+1}) = \arg\max_{z_n} [\ln p(z_{n+1} | z_n) + \omega(z_n)] \quad \left. \begin{array}{l} \text{keep track of this.} \end{array} \right\}$$

each $z_n^{\text{max}} = \arg\max_{z_n} \omega(z_n)$

(reason backw.) {
 $z_n^{\text{max}} = \Psi_n(z_{n+1})$
 dynamic prog.

linear Dynamical Sys. (13.3)

Same as HMM's but var. continuous & prob. = Gaussian.
(same graph)

Now : $\{x_n, z_n\}$ are continuous $x_n \in \mathbb{R}^d$, $z_n \in \mathbb{R}^{d_z}$

linear Gauss. model : all condi. dist. in Bay. Netw. are Gaussian with means that depend linearly on parents.

Viterbi : for Gauss. mean = median (mean is the most probable)
= MAP

// AM.

→ no viterbi necessary.

Forw. passing message seqs : Kalman filter seqs.

Backw. " " " : Kalman smoother seqs.

Transition : $p(z_n | z_{n-1}) = N(z_n | Az_{n-1}, T)$ → trans. uncertainty

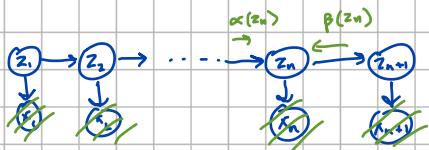
Emission : $p(x_n | z_n) = N(x_n | Cz_n, \Sigma)$ const. in time } (i.e., mean shifts & covar. gets larger & larger.)
= I if good
observ. noise



Initial State : $p(z_0) = N(z_0 | \mu_0, V_0)$

Params : $\theta = (A, T, C, \Sigma, \mu_0, V_0)$

Inference in LDS (for E-step) (13.3.1) (prob. of latent given observed)



Forw. seqs. : $\hat{\alpha}(z_n) = N(z_n | \mu_n, V_n)$

not normalized
(messages)

update these in
message update step.

U
S
E

T
R
I
C
K
S

$$\begin{aligned} p(x) &= N(x | \mu, \Lambda^{-1}) \\ p(y|x) &= N(y | Ax + b, L^{-1}) \\ \text{(integrate over } x \text{ to get msg. } \hat{\alpha} \text{)} \\ p(y) &= N(y | A\mu + b, \Sigma + A\Lambda^{-1}A^T) \\ p(x|y) &= N(x | \Sigma(A^T L(y - b) + \Lambda\mu), \Sigma) \quad , \quad \Sigma = (I + A^T L A)^{-1} \end{aligned}$$

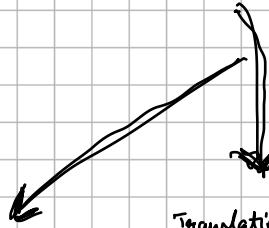
Basic idea : $p(x)p(y|x) = p(y)p(x|y)$

} 2.115
} 2.116

Forw. eqs :

$$\hat{\alpha}(z_n) = \frac{1}{C_n} p(x_n | z_n) \int \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1}) dz_{n-1}$$

$p(x)$



$$C_n \cdot N(z_n | \mu_n, V_n) = \underbrace{N(x_n | c z_n, \Sigma)}_{p(y|x)} \underbrace{\int N(z_{n-1} | \mu_{n-1}, V_{n-1}) N(z_n | A z_{n-1}, T) dz_{n-1}}_{p(y|x)}$$

$p(y|x)$

$N(z_n | A\mu_{n-1}, T + AV_{n-1}A^T)$

P_{n-1}

now just product of 2 Gaussians.

Translation Table :

x	$\rightarrow z_n$
y	$\rightarrow x_n$
M	$\rightarrow A\mu_{n-1}$
Σ	$\rightarrow P_{n-1}$
A	$\rightarrow C$
b	$\rightarrow 0$
ϵ	$\rightarrow \xi$
Σ	$\rightarrow (P_{n-1}^{-1} + C^T \xi^{-1} C)^{-1}$

$$\begin{aligned} C_n &= N(x_n | CA\mu_{n-1}, \Sigma + CP_{n-1}C^T) \\ \mu_n &= (P_{n-1}^{-1} + C^T \xi^{-1} C)^{-1} (C^T \xi^{-1} x_n + P_{n-1}^{-1} A\mu_{n-1}) \\ V_n &= (P_{n-1}^{-1} + C^T \xi^{-1} C)^{-1} \end{aligned}$$

Apply Bishop C.7 (Woodbury identity)

$$(A + B)^{-1} = A^{-1} - A^{-1}B(I + CA^{-1}B)^{-1}CA^{-1}$$

$$\underline{V_n} = P_{n-1} - \underbrace{P_{n-1} C^T (\Sigma + CP_{n-1} C^T)^{-1} C P_{n-1}}_{K_n} = \underline{(I - K_n C)P_{n-1}}$$

(Kalman gain matrix)

less inversions

Apply Bishop C.5 :

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$$

$$\begin{aligned} \underline{\mu_n} &= K_n x_n + (I - K_n C) A \mu_{n-1} \\ &= A \mu_{n-1} + K_n (x_n - A \mu_{n-1}) \end{aligned}$$

new pred.
observ.

Backw. eqs :

not in β .

formulated in terms of $\gamma(z_n) = \hat{\alpha}(z_n) \hat{\beta}(z_n)$

variable beliefs (prob. of incoming msg)

$$= N(z_n | \hat{\mu}_n, \hat{V}_n)$$

factor beliefs $\varepsilon(z_{n-1}, z_n) = N((z_{n-1}) | (\hat{\mu}_{n-1}) \begin{pmatrix} \hat{V}_{n-1} & J_{n-1} \hat{V}_n \\ \hat{V}_n J_{n-1}^T & \hat{V}_n \end{pmatrix})$

} used for EM

Learning in LDS using EM : (3.3.2)

complete data log-likelihood :

$$\ln p(x, z | \theta) = \ln(z_1 | \mu_0, V_0) + \sum_{n=2}^N \ln p(z_n | z_{n-1}, A, T) + \sum_{n=1}^N p(x_n | z_n, C, \xi)$$

$$Q(\theta, \theta^{\text{old}}) = E_{z| \theta^{\text{old}}} [\ln p(x, z | \theta)]$$

Use inference results : $E[z_n | \theta^{\text{old}}] = \hat{\mu}_n$ (mean of $\gamma(z_n)$)

$(\beta_3 \cdot 10^5 - \beta_3 \cdot 10^7)$
slipping

$$E[z_n z_{n-1}^T | \theta^{\text{old}}] \approx J_{n,n} \hat{\mu}_n + \hat{\mu}_n \hat{\mu}_n^T \quad (\text{variance of } \Sigma(z_n))$$

$$E[z_n z_m^T | \theta^{\text{old}}] = \hat{\mu}_n + \hat{\mu}_m \hat{\mu}_m^T$$

$$\left[\begin{array}{l} \Sigma \begin{pmatrix} x \\ y \end{pmatrix}, \begin{pmatrix} \varepsilon_{xx} & \varepsilon_{xy} \\ \varepsilon_{yx} & \varepsilon_{yy} \end{pmatrix} \\ E[xx^T] = E[(x - \mu_x)(x - \mu_x)^T] + \dots \\ = \Sigma_{xx} + \mu_x \mu_x^T \end{array} \right]$$

M-step :

$$\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{\text{old}}) \quad xy^T$$

quadratic in μ_0 , optimum in z_i .

$$\text{Eq: } \mu_0, V_0 : Q(\theta, \theta^{\text{old}}) = \frac{-1}{2} \ln |V_0| - \frac{1}{2} E_{z| \theta^{\text{old}}} [(z_i - \mu_0^T) V_0^{-1} (z_i - \mu_0)] + \text{const}$$

$$\mu_0^{\text{new}} = E[z | \theta^{\text{old}}]$$

$$V_0^{\text{new}} = E[z z^T | \theta^{\text{old}}] - E[z | \theta^{\text{old}}] E[z^T | \theta^{\text{old}}]$$

\vdots
updates for A, Γ, C, Σ

14th Oct

Causal Relations : A causes B if change in A \rightarrow change in B.

$x_1 \rightarrow x_2$ x_1 causes x_2