

# PROBABILISTIC <sup>COMP6118</sup> & UNSUPERVISED LEARNING

## Introduction & Foundations

- Unsupervised: observe data, find pattern & underlying structure
- Supervised: observe input/output pairs, learn to predict for novel input
- Reinforcement: find policy for action choice to maximise reward

- Probabilistic approach:

$$p(x|\theta) \text{ or } p(y|x, \theta)$$

- generative model / likelihood
- Beliefs: isomorphic to probabilities
- Basic rules of probability:

$$P(x) \geq 0$$

$$\sum_x P(x) = 1$$

$$P(x) = \sum_y P(x, y)$$

$$P(x|y) = P(x, y) / P(y)$$

$$P(y|x) = \frac{P(x, y) P(y)}{P(x)}$$

Dutch book theorem:

if I bet according to my beliefs (ie if  $b(x=T) = 0.9$ , accept 1:9 bet)

• exist set of bets where I am guaranteed to lose money if my beliefs are not consistent

- Bayesian learning:

• Data  $D = \{x_1, \dots, x_n\}$

• Models  $M_1, M_2, \dots$

• Parameters  $\theta_i$  (per model)

Parameter learning

$$\rightarrow P(\theta_i | D, M_j) = \frac{P(D | \theta_i, M_j) P(\theta_i | M_j)}{P(D | M_j)}$$

Model selection

$$\rightarrow P(M_i | D) = \frac{P(D | M_i) P(M_i)}{P(D)} \quad (P(D | M_i) = \int d\theta p(D | M_i, \theta) p(\theta | M_i))$$

- Conjugate priors:

• for a given likelihood, conjugate prior gives posterior in same family as prior

• Exponential family likelihood:

$$P(x|\theta) = g(\theta) \xi(x) e^{\phi(\theta)^T T(x)} \quad \left( \begin{array}{l} \phi(\theta) = \text{natural parameters} \\ T(x) = \text{sufficient statistic} \end{array} \right)$$

$$\Rightarrow p(\theta) \propto g(\theta)^v e^{\phi(\theta)^T \tau}$$

$$(\Rightarrow p(\theta|x) \propto g(\theta)^{v+n} e^{\phi(\theta)^T (y + \sum T(x_i))})$$

$\tau$  = pseudo observations  
 $v$  = 'scale' of prior



- ML learning for a Gaussian:  
 $p(\underline{x}|\underline{\mu}, \underline{\Sigma}) = N(\underline{\mu}, \underline{\Sigma})$

log likelihood  $\rightarrow L = \log \prod_i p(\underline{x}_i|\underline{\mu}, \underline{\Sigma}) = \sum \log p(\underline{x}_i|\underline{\mu}, \underline{\Sigma})$

$$\frac{\partial L}{\partial \underline{\mu}} = 0 \Rightarrow \hat{\underline{\mu}} = \frac{1}{N} \sum \underline{x}_i, \quad \frac{\partial L}{\partial \underline{\Sigma}} = 0 \Rightarrow \hat{\underline{\Sigma}} = \frac{1}{N} \sum (\underline{x}_i - \hat{\underline{\mu}})(\underline{x}_i - \hat{\underline{\mu}})^T$$

- Multivariate linear regression:  
 $p(y|\underline{x}, \underline{W}, \underline{\Sigma}_y) = N(\underline{W}\underline{x}, \underline{\Sigma}_y)$

$$\frac{\partial L}{\partial \underline{W}} \Rightarrow \hat{\underline{W}}_{ML} = \left( \sum_i y_i \underline{x}_i^T \right) \left( \sum_i \underline{x}_i \underline{x}_i^T \right)^{-1}$$

- MAP learning:

$$p(\underline{w}) = N(\underline{0}, \underline{\Theta}^{-1}), \quad y_i \text{ scalar}, \quad y_i|\underline{x}_i, \underline{w}, \sigma_y^2 \sim N(\underline{w}^T \underline{x}_i, \sigma_y^2)$$

$$\Rightarrow \log p(\underline{w}|D, \underline{\Theta}, \sigma_y^2) = L + \log p(\underline{w}|\underline{\Theta})$$

$$= \log N(\underline{w}_{MAP}, \underline{\Sigma}_w)$$

$$[\underline{w}_{MAP} = \underline{\Sigma}_w^{-1} \frac{1}{\sigma_y^2} \sum y_i \underline{x}_i, \quad \underline{\Sigma}_w = (\underline{\Theta} + \frac{1}{\sigma_y^2} \sum \underline{x}_i \underline{x}_i^T)]$$

- Problems with multivariate Gaussian model:
  - no higher order structure

- predicts very few outliers; not robust

- $\frac{D(D+1)}{2}$  parameters

## 2: Latent Variable Models

- Explain correlations in  $\underline{x}$  by assuming dependence on latent variables  $\underline{y}$ 
  - can reduce parameters needed, capture an underlying generative process
  - $p(\underline{y})$ ,  $p(\underline{x}|\underline{y})$ , maybe  $p(\underline{x}, \underline{y})$  in exponential family ( $p(\underline{x})$  rarely is)

### Probabilistic PCA:

$$D = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}, \quad \underline{x}_i \in \mathbb{R}^D \quad (K < D)$$

$$Y = \{y_1, \dots, y_N\}, \quad y_i \in \mathbb{R}^K$$

Model:

$$\underline{x}_i = \underline{\Lambda} y_i + \underline{e}_i, \quad y \sim N(\underline{0}, \underline{I}), \quad \underline{e} \sim N(\underline{0}, \psi \underline{I})$$

$$\Rightarrow \underline{x}|\underline{y} \sim N(\underline{\Lambda} \underline{y}, \psi \underline{I}) \Rightarrow \underline{x} \sim N(\underline{0}, \underline{\Lambda} \underline{\Lambda}^T + \psi \underline{I})$$

$DK+1$  free parameters (vs  $\frac{D(D+1)}{2}$  if  $\underline{x} \sim N(\underline{0}, \underline{\Sigma})$ )  
(more interpretable, but harder to do inference)

noise  $\perp$  to projection  $\rightarrow$  PCA:

- as  $\psi \rightarrow 0$ , can only capture  $K$  dimensions of variance
- $\underline{\Lambda}$  columns = eigenvectors of  $\underline{\Sigma}$ , ordered by eigenvalue
- $(\underline{\Sigma} = \underline{U} \underline{W} \underline{U}^T, \quad \underline{U}$  columns = eigenvectors,  
 $\underline{W} = \text{diag}(\lambda_1, \dots, \lambda_D)$ )

### ML Learning

$$\ell = \log N(\underline{0}, p(\underline{x}|\underline{\Lambda}, \psi))$$

$$\frac{\partial \ell}{\partial \underline{\Lambda}} \Rightarrow \underline{\Lambda} = \underline{0}$$

$$\text{or, } \underline{\Lambda} \underline{\Lambda}^T + \psi \underline{I} = \underline{\Sigma}$$

$$\text{or, } \underline{\Lambda} = \underline{U} \underline{L} \underline{U}^T, \text{ and } \underline{\Sigma} \underline{U} = \underline{U} (\underline{L}^2 + \psi \underline{I})$$

### Projection:

$$\tilde{\underline{x}}_n = \underline{\Lambda} \bar{y}_n, \text{ where } \bar{y}_n = E(y_n | \underline{x}_n)$$

find  $E(y_n | \underline{x}_n)$ :

write  $p(y_n | \underline{x}_n) = p(y_n) p(\underline{x}_n | y_n)$ , consider  $\underline{x}_n$  fixed  
 $\Rightarrow y_n | \underline{x}_n \sim N(\beta \underline{x}_n, \underline{I} - \beta \underline{\Lambda}), \quad \beta = \underline{\Sigma} \underline{\Lambda}^T \psi^{-1}$



## Factor Analysis:

- don't assume equal noise in all dimensions;  $\underline{e} \sim N(\underline{0}, \underline{\Psi})$
- $DK + D - \frac{K(K-1)}{2}$  parameters (not identifiable if  $> \frac{D(D+1)}{2}$ )
- no closed form solution for  $\underline{x} \sim N(\underline{0}, \underline{\Lambda}^T + \underline{\Psi})$   
 $\Rightarrow$  need to use gradient ascent on log likelihood (or EM)
- projection: same as for PPCA, except  $\underline{\Psi}$  <sup>(diagonal)</sup> non matrix

- PCA & PPCA rotationally invariant, FA not
- FA measurement scale invariant, PCA/PPCA not
- FA & PPCA probabilistic model, PCA not
- all can only model very restricted class of densities (too many assumptions)

## Mixture Models:

$$S_i \stackrel{iid}{\sim} \text{Discrete}(\underline{\pi}), \quad \underline{x}_i | S_i \sim P_{S_i}[\underline{\theta}_{S_i}]$$

$$\Rightarrow p(\underline{x}_i) = \sum_m \pi_m p_m(\underline{x}_i | \underline{\theta}_m) \quad (\text{convex combination of component densities})$$

Mixture of Gaussians:

$$\underline{x}_i | S_i \sim N(\underline{\mu}_{S_i}, \underline{\Sigma}_{S_i})$$

$$p(D | \{\underline{\mu}_m\}, \{\underline{\Sigma}_m\}, \underline{\pi}) = \prod_{i=1}^n \sum_{m=1}^K \pi_m N(\underline{x}_i | \underline{\mu}_m, \underline{\Sigma}_m)$$

ML: 
$$L = \sum_{i=1}^n \log \sum_{m=1}^K \pi_m p_m(\underline{x}_i | \underline{\theta}_m)$$

$$\left( \text{For MoG, } \underline{\theta}_m = \{\underline{\mu}_m, \underline{\Sigma}_m\} \right) \frac{\partial L}{\partial \underline{\theta}_m} \stackrel{*}{=} \sum_{i=1}^n r_{im} \frac{\partial}{\partial \underline{\theta}_m} \log p_m(\underline{x}_i | \underline{\theta}_m) \quad \left( r_{im} = P(S_i = m | \underline{x}_i) = \frac{p_m(\underline{x}_i) \pi_m}{\sum_R p_R(\underline{x}_i) \pi_R} \right)$$

$$\frac{\partial L}{\partial \pi_m} = \sum_{i=1}^n \frac{r_{im}}{\pi_m}$$

• k-means:

$$\pi_m = 1/m, \quad \underline{\Sigma}_m = \sigma^2 \underline{I}, \quad \sigma^2 \rightarrow 0$$

$$\Rightarrow \text{MoG approaches k-means } [r_{im} \rightarrow \delta(m, \arg \min_i \|\underline{x}_i - \underline{\mu}_i\|^2)]$$



### 3: Expectation Maximisation

(\* Jensen's inequality)

- Maximisation hard when  $y$  is latent
- Idea:
  - find expectation of  $y$ , then update  $\theta$  w.r.t. this, & alternate
- E-step:
  - find expected values for latent variables
- M-step:
  - maximise likelihood as if latents were not hidden

• Prove we always never decrease likelihood:

$$\begin{aligned} \lambda(\theta) &= \log p(X|\theta) = \log \int p(Y, X|\theta) dY \\ &= \log \int q(Y) \frac{p(Y, X|\theta)}{q(Y)} dY \stackrel{*}{\geq} \int q(Y) \log \frac{p(Y, X|\theta)}{q(Y)} dY \triangleq F(q, \theta) \end{aligned}$$

Free energy

$$\rightarrow F(q, \theta) = \int q(Y) \log p(Y, X|\theta) dY - \int q(Y) \log q(Y) dY$$

$$F(q, \theta) = E_{q(Y)} [\log p(Y, X|\theta)] + H(q) \leftarrow \text{entropy}$$

• E-step:  $q^{(k)}(Y) = \arg \max_{q(Y)} F(q(Y), \theta^{(k-1)}) = p(Y|X, \theta^{(k-1)})$

• M-step:  $\theta^{(k)} = \arg \max_{\theta} F(q^{(k)}(Y), \theta) = \arg \max_{\theta} E_{q^{(k)}(Y)} [\log p(Y, X|\theta)]$

$$F(q, \theta) = \int q(Y) \log p(X|\theta) dY + \int q(Y) \log \frac{p(Y|X, \theta)}{q(Y)} dY$$

$$F(q, \theta) = \lambda(\theta) - KL[q(Y) \| p(Y|X, \theta)]$$

$$\Rightarrow \text{E-step sets } q^{(k)}(Y) = p(Y|X, \theta^{(k-1)}) \text{ [minimize KL-divergence]}$$

$$\lambda(\theta^{(k-1)}) \underset{\substack{\uparrow \\ \text{E-step}}}{=} F(q^{(k)}, \theta^{(k-1)}) \underset{\substack{\uparrow \\ \text{M-step}}}{\leq} F(q^{(k)}, \theta^{(k)}) \underset{\substack{\uparrow \\ \text{Jensen}}}{\leq} \lambda(\theta^{(k)})$$



• EM for MoG:

• E-step:

$$q(s_i) = p(s_i | x_i, \Theta) \propto \pi_{s_i} \times N(\mu_{s_i}, \sigma_{s_i}^2)$$

• M-step:

$$\Theta = \arg \max_{\Theta} E_{q(s)} [\log p(x, s | \Theta)]$$

$$E = \sum q(s) \log [p(s | \Theta) p(x | s, \Theta)]$$

(find  $\mu_m, \sigma_m, \pi_m$  by differentiating E)

$$E = \sum_{i,m} r_{im} [\log \pi_m - \log \sigma_m - \frac{1}{2\sigma_m^2} (x_i - \mu_m)^2]$$

• EM for FA:

• E-step:

$$\begin{aligned} q_m(y_n) &= p(y_n | x_n, \Theta) \propto p(y_n, x_n | \Theta) = p(y_n | \Theta) p(x_n | y_n, \Theta) \\ &= N(\beta x_n, I - \beta \Lambda), \quad \beta = \Sigma \Lambda^T \Psi^{-1} \end{aligned}$$

• M-step:

$$\Theta = \arg \max_{\Lambda, \Sigma, \Psi} E_{q(y_n)} [\log p(x_n, y_n | \Theta)]$$

$$\log p(x_n, y_n | \Theta) = \log p(y_n | \Theta) + \log p(x_n | y_n, \Theta)$$

do maths, take expectation by replacing  $y_n$  w/  $\mu_n$ ,  $y_n y_n^T$  w/  $\mu_n \mu_n^T + \Sigma$

(as  $\Sigma \rightarrow 0$  these tend to equations for ML linear regression)

$$\frac{\partial F}{\partial \Lambda} = 0 \Rightarrow \hat{\Lambda} = \left( \sum_n x_n \mu_n^T \right) \left( N \Sigma + \sum_n \mu_n \mu_n^T \right)^{-1}$$

$$\frac{\partial F}{\partial \Psi^{-1}} = 0 \Rightarrow \hat{\Psi} = \hat{\Lambda} \Sigma \hat{\Lambda} + \frac{1}{N} \sum_n (x_n - \hat{\Lambda} \mu_n)(x_n - \hat{\Lambda} \mu_n)^T$$

• EM for Exponential families:

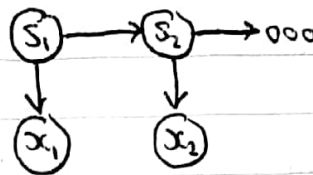
$p(y, x | \Theta)$  has exp-family form

$$\begin{aligned} F(q, \Theta) &= \int q(y) \log p(y, x | \Theta) dy + H(q) \\ &= \Theta^T E_{q(y)} [T(y, x)] - \log Z(\Theta) + c \end{aligned}$$

$\Rightarrow$  E-step: compute expected sufficient stats under  $q$

$$\text{M-step: } \frac{\partial F}{\partial \Theta} = E_{q(y)} [T(y, x)] - E[T(y, x | \Theta)] = 0$$

#### 4: Latent Chain Models



##### • Hidden Markov Models:

•  $S_t \in \{1, \dots, K\}$ ,  $x_t$  discrete or continuous

$$\phi_{ij} = P(S_{t+1}=j | S_t=i)$$

$$\pi_j = P(S_1=j)$$

$$A_{jk} = P(x_t=k | S_t=j) \quad [\text{or } A_j(x) = P(x_t=x | S_t=j) \text{ for continuous } x_t]$$

$$p(x_{1:T}, S_{1:T}) = \pi_{S_1} A_{S_1 x_1} \prod_{t=2}^T \phi_{S_{t-1} S_t} A_{S_t x_t}$$

• Learning: use EM

• E-step: find  $p(S_t | x_{1:T})$

$$p(S_t=i | x_{1:T}) = \frac{P(S_t=i | x_{1:t}) P(x_{t+1:T} | S_t=i)}{P(x_{1:T})} = \frac{\alpha_t(i) \beta_t(i)}{\sum_j \alpha_t(j) \beta_t(j)}$$

• Forward:

in practice, normalise as we go

$$\alpha_t(i) = \pi_i A_i(x_1), \quad \alpha_{t+1}(i) = \left( \sum_j \alpha_t(j) \phi_{ji} \right) A_i(x_{t+1})$$

• Backward:

$$\beta_t(i) = \sum_j \phi_{ij} A_j(x_{t+1}) \beta_{t+1}(j)$$

$$q(S_{1:T}) = p(S_{1:T} | x_{1:T}, \Theta)$$

• M-step:

$$\hat{\pi}_i = q(S_1=i)$$

$$\hat{\phi}_{ij} = \frac{\sum_{t=1}^{T-1} \mathbb{E}_t(i \rightarrow j)}{\sum_{t=1}^{T-1} q(S_t=i)}$$

$$\mathbb{E}_t(i \rightarrow j) = P(S_t=i, S_{t+1}=j | x_{1:T}) = \frac{\alpha_t(i) \phi_{ij} A_j(x_{t+1}) \beta_{t+1}(j)}{P(x_{1:T})}$$

$$\hat{A}_{iR} = \frac{\sum_{t=1}^T q(S_t=i) x_t}{\sum_{t=1}^T q(S_t=i)}$$



$$y_1 \sim N(\mu_0, Q_0)$$

$$y_t, y_{t+1} \sim N(Ay_{t-1}, Q)$$

$$x_t, y_t \sim N(Cy_t, R)$$

• Linear Gaussian State-Space Model:

$$x_t = Cy_t + v_t$$

$$y_t = Ay_{t-1} + w_t$$

$v_t, w_t$  0-mean uncorrelated Gaussian noise

$y_1$  Gaussian,

$p(x_{1:T}, y_{1:T})$  multivariate Gaussian

• Learning:

• E-step:

Kalman Smoothing

$$p(y_t | x_{1:T}) = \int p(y_t | y_{t+1}, x_{1:T}) p(y_{t+1} | x_{1:T}) dy_{t+1}$$

• backwards recursion

• M-step:

$$C_{\text{new}} = \arg \max_C \mathbb{E}_q \left[ \sum_t \ln p(x_t | y_t) \right] = \left( \sum_t x_t \mathbb{E}[y_t] \right) / \left( \sum_t \mathbb{E}[y_t y_t^T] \right)$$

$$A_{\text{new}} = \arg \max_A \mathbb{E}_q \left[ \sum_t \ln p(y_{t+1} | y_t) \right] = \left( \sum_t \mathbb{E}[y_{t+1} y_t^T] \right) / \left( \sum_t \mathbb{E}[y_t y_t^T] \right)$$

• Online:

$$1 = \sum_{t=1}^T \ln p(x_t | x_{1:t-1}) = \sum_{t=1}^T \lambda_t \quad \left( p(x_t | x_{1:t-1}) \text{ determined by } \right)$$

Kalman filtering

• use gradient rules to update  $A, C, Q, R$  as we go  
(learning rate = expectation about non-stationarity)

• Slow Feature Analysis:

$$A \rightarrow I$$

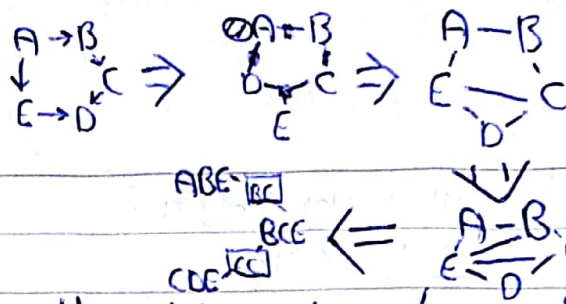
$$Q = I - AA^T \rightarrow 0$$

$$R \rightarrow 0$$

$\Rightarrow$  given series  $x_{1:T}$ , find  $C$  such that  $y_{1:T}$  changes as slowly as possible



## 5: Graphical Models



### Junction Tree Algorithm:

• DAG  $\xrightarrow{1}$  Factor Graph  $\xrightarrow{2}$  Undirected Graph

Message Passing  $\xleftarrow{5}$  Junction Tree  $\xleftarrow{4}$  Chordal Graph

(every step removes conditional independencies,  $\Rightarrow$  ~~not~~ family of possible distributions)

merging DAG does 1 & 2 in 1 step

- 1: Conditional distributions in DAG  $\rightarrow$  Factors in factor graph
- 2: Replace factors by undirected clique
- 3: Add edges to graph so every loop of size  $\geq 4$  has 1+ chords

### Variable elimination:

• eliminate variables in turn, connecting all nodes that variable is connected to

heuristic can significantly impact cost of future message passing

• order of removal heuristic:

• choose variable with where removal would add fewest new edges each step

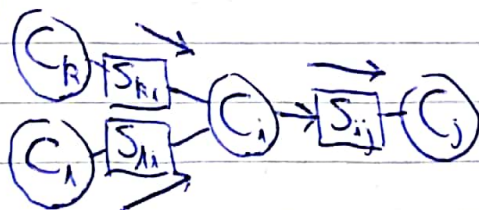
4: Each node in junction tree is a clique, edges are intersections between cliques

### Running intersection property:

if 2 cliques contain X, all cliques & separators between them do so too

• first, find graph where each node is a maximal clique  
• Then, find maximal spanning tree (edge weight = no. variables in separator)

most expensive stage:  
 $O(\sum_i |X_{C_i}|)$



$$M_{i \rightarrow j}(X_{S_{ij}}) = \sum_{X_{C_i} \setminus S_{ij}} \gamma_i(X_{C_i}) \prod_{R \in \text{ch}(i)} M_{R \rightarrow i}(X_{S_{Ri}})$$

$$P(X_G) = \gamma_i(C_i) \prod_{R \in \text{ch}(i)} M_{R \rightarrow i}(X_{S_{Ri}})$$

$$P(X_{S_{ij}}) = M_{i \rightarrow j}(X_{S_{ij}}) M_{j \rightarrow i}(X_{S_{ij}})$$



• Conditional independence:

•  $A \perp\!\!\!\perp B \mid X$  if

• moralise graph of  $A, B, X$  & ancestors

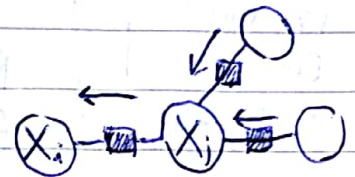
• if  $X$  separates  $A$  &  $B$ , then true (else false)

• Belief propagation in trees:

• undirected graph with no cycles

• factors;  $p(X) \propto \prod_{ij} \xi_{ij}(X_i, X_j)$

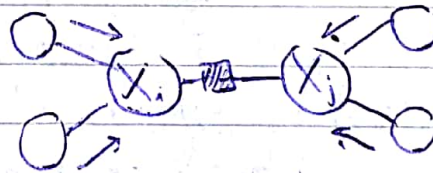
(Start from any leaf)  $M_{j \rightarrow i}(X_i) = \sum_{X_j} \xi_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus X_i} M_{k \rightarrow j}(X_j)$



• Inference:

$$p(X_i) =$$

$$p(X_i, X_j) = \xi_{ij}(X_i, X_j) \prod_{X_k \in \text{ne}(X_j) \setminus X_i} M_{k \rightarrow j}(X_j) \prod_{X_k \in \text{ne}(X_i) \setminus X_j} M_{k \rightarrow i}(X_i)$$



If leaf node is observed, messages are conditioned, not marginalised;

$$M_{a \rightarrow i}(X_i) = \xi_{ai}(X_a = a, X_i)$$



## 6: Model Selection & GPs

### • Model Selection:

- cannot use ML/MAP; more complex models always have higher likelihood
- can compare nested models, stopping when adding complexity doesn't improve some score

• or cross-validation

• or Bayesian model selection:

$$P(m|D) = \frac{p(D|m)p(m)}{p(D)} \quad , \quad p(D|m) = \int p(D|\theta_m, m) p(\theta_m|m) d\theta_m$$

(can ignore)  $\swarrow$  prob of random parameter values generating data

• gives Bayesian Occam's razor;

- too simple models unlikely to generate data
- too complex can generate too many datasets, so unlikely to generate this specific one

• Computing  $p(D|m)$ :

- if likelihood  $p(D|\theta_m, m)$  in exp fam, & using conjugate prior, then joint,  $p(D, \theta_m|m)$  in exp family too
- thus can be tractable, but in general not

• Approximations:

• Laplace: approximate posterior by Gaussian centred at  $\hat{\theta}_{MAP}$

• Bayesian info criterion: take Laplace & let  $N \rightarrow \infty$   
 $\log p(D|m) \approx \log p(D|\theta_m^*, m) - \frac{d}{2} \log N$   
+ quick & easy to compute  
+ can use  $\theta_{ML}$  (as  $N \rightarrow \infty$ )

• Hyperparameters:

- models may be continuously parameterised by hyperparameters  $\eta$ ; need  $p(D|\eta)$
- can try \* exact evidence, approximated, or sample  $\eta$  from  $p(\eta|D) \propto p(D|\eta)p(\eta)$  (by MCMC)

\* [gradient ascent]  
in



⇒ non-parametric model

## Gaussian Process Regression:

- use prediction averaging; integrate out parameters, to predict conditional density at new data point

$$p(y|x, D, m) = \int p(y|x, \theta, m) p(\theta|D, m) d\theta$$

- eg, linear regression:

$$\begin{aligned} \underline{w}|D &\sim N(\hat{\underline{w}}, \Sigma_{\underline{w}}) \\ \Rightarrow y|x, D &\sim N(\hat{\underline{w}}^T \underline{x}, \underline{x}^T \Sigma_{\underline{w}} \underline{x} + \sigma^2) \end{aligned} \quad \begin{pmatrix} \underline{w} \sim N(\underline{0}, \sigma^2 \underline{I}) \\ y|x \sim N(\underline{w}^T \underline{x}, \sigma^2) \end{pmatrix}$$

- marginalised linear regression:

$$\begin{aligned} \underline{Y} &\sim N(\underline{0}, \sigma^2 \underline{X}^T \underline{X} + \sigma^2 \underline{I}) \\ \underline{Y}|\underline{X} &\sim N(\underline{0}, \begin{pmatrix} E[y_i] = E[\underline{w}^T \underline{x}_i] = \underline{0}^T \underline{x}_i = 0 \\ E[(y_i - \bar{y}_i)^2] = \sigma^2 \underline{x}_i^T \underline{x}_i + \sigma^2 \\ E[(y_i - \bar{y}_i)(y_j - \bar{y}_j)] = \sigma^2 \underline{x}_i^T \underline{x}_j \end{pmatrix}) \end{aligned}$$

- predicting:

$$\begin{aligned} \{\underline{y}, y|x\}|\underline{X}, \underline{x} &\sim N\left(\begin{bmatrix} \underline{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 \underline{X}^T \underline{X} + \sigma^2 \underline{I} & \sigma \underline{X}^T \underline{x} \\ \sigma \underline{x}^T \underline{X} & \sigma \underline{x}^T \underline{x} + \sigma^2 \end{bmatrix}\right) \\ &= N\left(\begin{bmatrix} \underline{0} \\ 0 \end{bmatrix}, \begin{bmatrix} K_{\underline{X}\underline{X}} & K_{\underline{X}\underline{x}} \\ K_{\underline{x}\underline{X}} & K_{\underline{x}\underline{x}} \end{bmatrix}\right) \end{aligned}$$

(standard result)  $\Rightarrow y|x, \underline{X}, \underline{Y} \sim N(K_{\underline{x}\underline{X}} K_{\underline{X}\underline{X}}^{-1} \underline{Y}^T, K_{\underline{x}\underline{x}} - K_{\underline{x}\underline{X}} K_{\underline{X}\underline{X}}^{-1} K_{\underline{X}\underline{x}})$

- can replace  $\underline{x}$  with  $\phi(\underline{x})$  & get non-linear regression

- as only inner products, use kernel trick;

$$K(\underline{x}_i, \underline{x}_j) = \phi(\underline{x}_i)^T \phi(\underline{x}_j)$$

$$([K_{\underline{X}\underline{X}}]_{ij} = K(\underline{x}_i, \underline{x}_j), [K_{\underline{X}\underline{x}}]_i = K(\underline{x}_i, \underline{x}), K_{\underline{x}\underline{x}} = K(\underline{x}, \underline{x}))$$

- can also sample vector from GP (parameterised by mean & covariance function)
- can optimise hyperparameters in kernel by gradient ascent in  $\log p(\underline{Y}|\underline{X}, K)$