

- today:
- furnish factor analysis
 - Bayesian non-parametrics
 - . GP
 - . DP

factor analysis continuation :

$$X|Z \sim N(Wz + \mu, D)$$

$p(x)$ is Gaussian

$$\mathbb{E}[x] = \mathbb{E}[\underbrace{\mathbb{E}[x|z]}_{Wz + \mu}] = W\mathbb{E}[z] + \mu = \mu$$

$$\begin{aligned}\text{cov}(x, x) &= \text{cov}(Wz + \mu + \varepsilon, Wz + \mu + \varepsilon) \\ &\quad \text{indep.} \\ &= \text{cov}(Wz, Wz) + \text{cov}(\varepsilon, \varepsilon) \\ &= W \underbrace{\text{cov}(z, z)}_{I_K} W^T + D \\ &= WW^T + D\end{aligned}$$

equivalent model on

$$X \sim N(\mu, \underbrace{WW^T + D}_{\text{rank } K \text{ and diagonal}})$$

but of rank K low rank covariance assumption $\rightarrow d$ degrees of freedom

estimate $W \& D \& \mu$ by MLE

→ do EM (latent variable model)

get $p(z|x) \rightarrow$ Gaussian with mean
 $\mathbb{E}[z|x] = W^T(WW^T + D)^{-1}(x - \mu)$

probabilistic PCA: special case of factor analysis where suppose $D = \sigma^2 I_d$

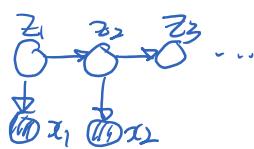
$$\lim_{\sigma \rightarrow 0} W^T(WW^T + \sigma^2 I)^{-1} = W^+ \leftarrow \text{pseudo-inverse}$$

$$= W^T \quad \text{if } W^T W = I_K$$

this suggests that PCA is limit of PFA as $\sigma \rightarrow 0$

Kalman filter

Factor analysis



more to state-space model: unroll in time
 (LMM style)

Kalman filter: $z_t | z_{t-1} \sim N(Az_{t-1}, B)$

→ doing "sum-product" alg. in HMM $p(z_t | x_{1:t})$

get "Kalman filtering" alg.

variational auto-encoder



$$z \sim N(0, I_k)$$

diagonal noise

$$x|z \sim N(\mu_w(z), \sigma_w^2(z))$$

where $\mu_w(z)$ ← output of a NN
"decoder"

MLE → use EM

↳ $p(z|x)$ is intractable \Rightarrow approximate with variational inference

approximate $p(z|x)$ with $q_\phi(z|x)$

$$z|x \sim N(\underbrace{\mu_\phi(x)}, \underbrace{\sigma_\phi^2(x)})$$

in EM $\log p(x) \geq \mathbb{E}_q [\log p(z|x)] + H(q)$ output of a NN → "encoder"
 $= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - KL[q_\phi(z|x) || p(z)]$

allows "reparameterization trick"

$$z|x \rightarrow \mu_\phi(x) + \sigma_\phi(x) \cdot \xi$$

$\xi \sim N(0, 1)$

- VAE innovations:
 - share parameters ϕ among data points for their variational approximation $q_\phi(z|x)$
 - re-parameterization trick to only have parameters appear in simple deterministic transformation, stochasticity is all left in $N(0,1)$ noise variables (no parameters) => allow simple backpropagation of gradient through expectations
 - see also: <https://gregorygundersen.com/blog/2018/04/29/reparameterization/>
 - for more details, see: [Slides on VAE](#) by Aaron Courville - deep learning class Winter 2017

15h57

Bayesian non-parametrics

non-parametric model → infinite # of parameters
(or growing with # of data pts.)

e.g. KNN classifier → boundary complexly grows with # of pts.

\dots, n density on x given x_i

- Kernel density estimation $\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K(x, x_i)$
- Bayesian non-parametric \rightarrow need prior on infinite dim. parameter
 \rightarrow define dist. on ∞ -vector (stochastic process)

Stochastic process

collection of random variables indexed by a (potentially infinite) index set T

$$\{X(t) : t \in T\}$$

Examples:

$$T = \{1, \dots, n\} \quad X(t) = X_t \rightarrow \text{random vector } (X_1, \dots, X_n)$$

but also $T = \mathbb{N} \rightarrow \infty$ -sequence (X_1, X_2, X_3, \dots)

or $T = \mathbb{R} \rightarrow \text{"random function"}$

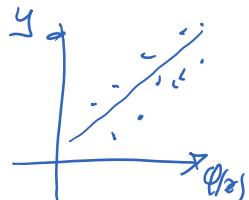
Gaussian process \rightarrow random function

Dirichlet process \rightarrow random measure/distribution
useful non-mixture model

Gaussian process:

motivation from Bayesian linear regression:

consider fixed location x_1, \dots, x_N [conditional model of $Y|X$]



model of $Y|X$: $y = w^\top \phi(x) + \sigma_y \epsilon$

thus $Y|X, w \sim N(w^\top \phi(x), \sigma_y^2)$

Bayesian: prior on $w \sim N(0, \sigma_w^2 I)$

$$y_i \triangleq y(x_i)$$

$$\mathbb{E}[\mathbb{E}[y(x)|w]] = \mathbb{E}[w^\top \phi(x) + \epsilon] = 0 + 0$$

$$\mathbb{E}[y_i y_j] = \mathbb{E}[w^\top \phi(x_i) w^\top \phi(x_j) + 0 + \sigma_y^2 \epsilon^2]$$

$$= \mathbb{E}[\text{tr}(\underbrace{\phi(x_i) \phi(x_j)^\top}_{K(x_i, x_j)} w w^\top)] + \sigma_y^2 = \sigma_w^2 k(x_i, x_j) + \sigma_y^2$$

$k(x_i, x_j) \rightarrow$ similarity

\Rightarrow generally, marginal on y 's [function values] a priori:

$$y_{1:N} \sim N(0, \sigma_w^2 \underbrace{\Phi_N \Phi_N^\top}_{K \text{ (kernel matrix)}} + \sigma_y^2)$$

$\Phi_N = \begin{pmatrix} -\phi(x_1) \\ \vdots \\ -\phi(x_N) \end{pmatrix}$

observation noise

Gaussian prior on function outputs

Gaussian process: generalization of Gaussian to ∞ -dimension

moments: $E[1] \dots E[\cdot \cdot \cdot] \dots E[\cdot \cdot \cdot \cdot \cdot]$

1) Gaussian process: generalization of Gaussian to ∞ -dimension

parameterized by $\mu(x)$ $\kappa(x, x')$

prior mean covariance (kernel)

stochastic process $Y(x)$ where for any x_1, \dots, x_n (and n)

marginal: $(Y(x_1), \dots, Y(x_n))$ follows a Gaussian with mean $\begin{pmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{pmatrix}$ and covariance Σ
s.t. $\Sigma_{ij} = \kappa(x_i, x_j)$

has to be PSD

special case of GP: Bayesian linear regression:

$$\text{use } \kappa(x, x') = \sigma_w^2 \phi(x)\phi(x')^\top + \sigma_y^2 I$$

but more generally, square exponential kernel \propto uncertainty

$$\kappa(x, x') = C \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

Length scale

(like ∞ -dim ϕ)

so Bayesian inference: suppose observed y_1, \dots, y_n (for x_1, \dots, x_n)
what is posterior on $y(x)$?

Simple? condition in Gaussian model?

$$Y_1, \dots, Y_n, Y(x) \sim N(0, \begin{pmatrix} C_n & k \\ k^\top & \kappa(x, x) \end{pmatrix})$$

$$k_i \triangleq \kappa(x_i, x)$$

↳ how output at x covaries with output at x_i

$$\text{get } Y(x) \mid Y_1, \dots, Y_n \sim N\left(0 + k^\top C_n^{-1} \vec{y}_{1:n}, \kappa(x) - k^\top C_n^{-1} k\right)$$

That's it!

note that only need to compute once

demo: <http://chifeng.scripts.mit.edu/stuff/gp-demo/>

⊕ applications to Bayesian optimization

- GPC use $p(y=1|x) = \sigma(f(x))$

- hyperparameter selection → can maximize the marginal likelihood

↳ closed form expression
(model selection)

Dirichlet process:

→ used to model ∞ -mixing model in Bayesian model

Bayesian mixture model (finite)

$$z \sim \text{Mult}(\pi) \quad \pi \in \Delta_k$$

$$x|z \sim p(x|\theta_z) \quad [\text{e.g. } N(x|\mu_z, \Sigma_z)]$$

Bayesian \rightarrow put a prior on π [Dirichlet($\alpha_1, \dots, \alpha_k$)]

and θ_z [say $\theta_z \stackrel{\text{iid}}{\sim} G_0$]

would like $k \rightarrow \infty$ can do $\theta_z \& \pi$ together using DP

Dirichlet process

$$G \sim DP(\alpha, G_0)$$

$\underset{P}{\text{random measure}}$

G_0 : dist. on Θ
 α : concentration parameter

stochastic process
 indexed by measurable set
 of Θ

for every partition of Θ in A_1, \dots, A_n

$$\text{then } (G(A_1), \dots, G(A_n)) \sim \text{Dir}(\alpha G_0(A_1), \dots, \alpha G_0(A_n))$$

stick breaking construction:

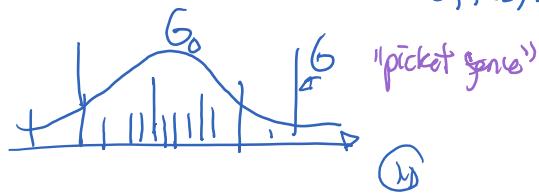
turns out that

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$$

dirac

where $\pi = (\pi_1, \dots) \sim \text{GEM}(\alpha)$

and $\theta_1, \theta_2, \dots \stackrel{\text{iid}}{\sim} G_0$



stick-breaking construction

$$\pi_i \sim \text{Beta}(1, \alpha)$$

$\underset{\Delta \Theta}{\text{all}}$

$$V_i \stackrel{\text{iid}}{\sim} \text{Beta}(1, \alpha)$$

$$\pi_2 = (1 - \pi_1) V_2$$

$$\pi_3 = (1 - \pi_2 - \pi_1) V_3$$

