

K-medoids algorithm:

Fin: $\{x_i\}_{i=1}^N$, distance measure $D(x, \mu)$

randomly initialise μ

• Heate until c is no longer changing:-

1. For each c_i ; set $c_i = \arg\min_{\mu} D(x_i, \mu_K)$

2. For each μ_K ; set $\mu_K = \arg\min_{\mu} \sum_{i:c_i=k} D(x_i, \mu)$ (*)

• Extension using general distance metric

• K-means $D(x, \mu) = \|x - \mu\|_2^2$ (Euclidean); $c_i: D(x_i, \mu) = \|x_i - \mu\|_1$ (robust to outliers)

LIS - Key equations

Probabilistic data modelling: - Bayes classifier, logistic regression, least-squares, ridge-reg (ML/MAP interp.), Bayesian LR

Data modelling: - Perceptron, decision trees, SVM, K-means

(non-probabilistic) 1) Model param. 2) Data $\{(x_i)\}_{i=1}^N$ 3) Distri family $p(x|\theta)$

ML estimation: 4) ingredients (situation I) a) iid $x_i \sim p(x|\theta)$ b) $\arg\max_{\theta} \prod_{i=1}^N p(x_i|\theta) = \arg\max_{\theta} \sum_{i=1}^N \ln p(x_i|\theta)$

$$\textcircled{m} \quad \hat{\theta}_{ML} = \arg\max_{\theta} p(x_1, \dots, x_N|\theta) \stackrel{\text{a)}{=} \arg\max_{\theta} \prod_{i=1}^N p(x_i|\theta) \stackrel{\text{b)}{=} \arg\max_{\theta} \sum_{i=1}^N \ln p(x_i|\theta)$$

$$\text{a) iid } \text{b) } f(y) > f(x) \Rightarrow \ln f(y) > \ln f(x); \arg\max_y \ln g(y) = \arg\max_y g(y)$$

LSLR + Bayes class: $x_i \stackrel{\text{iid}}{\sim} N(\bar{x}|\mu, \Sigma)$ $\theta = (\mu, \Sigma) \quad \text{J}_\theta \ln \left(\prod_{i=1}^N p(x_i|\theta) \right) = 0$ yields:-

$$(\text{Analytic sol.}) \quad \hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N x_i \quad \hat{\Sigma}_{ML} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{ML})(x_i - \hat{\mu}_{ML})^\top \quad (\text{LSLR}) \quad (\text{BC})$$

$$\hat{\mu}_R = \frac{1}{NR} \sum_{i=1}^N \mathbb{I}(y_i=R)x_i \quad \hat{\Sigma}_R = \frac{1}{NR} \sum_{i=1}^N \mathbb{I}(y_i=R)(x_i - \hat{\mu}_R)(x_i - \hat{\mu}_R)^\top \quad \hat{\Sigma}_K = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i=K)$$

ML estimation: split parameters into two groups θ_1, θ_2 , maximise likelihood

(situation II)

$$\textcircled{m} \quad \hat{\theta}_{1,ML}, \hat{\theta}_{2,ML} = \arg\max_{\theta_1, \theta_2} \sum_{i=1}^N p(x_i|\theta_1, \theta_2)$$

can solve one given
the other, cannot
solve simultaneously

Co-ordinate ascent (?): For $t=1, 2, \dots$

(general, probabilistic) 1. optimise $\hat{\theta}_1^{(t)} = \arg\max_{\theta_1} \sum_{i=1}^N \ln p(x_i|\theta_1, \theta_2^{(t-1)})$

2. optimise $\hat{\theta}_2^{(t)} = \arg\max_{\theta_2} \sum_{i=1}^N \ln p(x_i|\theta_1^{(t)}, \theta_2)$

ML situation: \textcircled{m} : find $\hat{\theta}_{1,ML} = \arg\max_{\theta_1} \sum_{i=1}^N \ln p(x_i|\theta_1)$ via:

(situation III)

$$\text{E-M} \quad \hat{\theta}_{1,ML} = \arg\max_{\theta_1} \sum_{i=1}^N \ln p(x_i|\theta_1, \theta_2|\theta_1) \quad (*)$$

(add a second variable θ_2)
lt side condition. \rightarrow prior on θ_2

θ_2 is a 'latent variable'; as marginal distri. $p(z_i | \theta)$ is tricky to optimise, we integrate over the joint distribution of z_i and latent θ_2 .
E-M algorithm finds $\hat{\theta}_{1, \text{ML}}$ via this method.

General Expectation-Maximisation:

Setup: Two parameter sets θ_1, θ_2 : $p(z | \theta_1) = \int p(z, \theta_2 | \theta_1) d\theta_2$ (marginal distri.)

objective criteria: Define objective
criteria
1) optimise marginal $p(z_i | \theta_1)$ wrt θ_1 ,
2) $p(z, \theta_2 | \theta_1)$ used for comp. convenience
3) we know ^{and} posterior $p(\theta_2 | z, \theta_1)$ given data z, θ_1

EM-objective function:

$$\ln p(z | \theta_1) = \underbrace{\int q(\theta_2) \ln \frac{p(z, \theta_2 | \theta_1)}{q(\theta_2)} d\theta_2}_{1. L(z, \theta_1)} + \underbrace{\int q(\theta_2) \ln \frac{q(\theta_2)}{p(\theta_2 | z, \theta_1)} d\theta_2}_{2. KL(q || p)}$$

1. $L(z, \theta_1)$ can be calculated (evaluated as integral), leaving only function of θ_1 for a particular setting $q(\theta_2)$

2. $KL(q || p)$: $KL \geq 0$, $KL=0$ when $q=p$

EM-objective: product rule of prob.; conditional prob., $\ln p(z | \theta_1) \neq f(\theta_1)$ (factor out)
(identity?)

EM algorithm output:

- Heavily optimise $\ln p(z | \theta_1)$ wrt $\theta_1 \Rightarrow$ generate a sequence of values $\theta_1^{(1)}, \dots, \theta_1^{(T)}$: 1. sequence $\{\theta_1^{(t)}\}_{t=1}^T$: $\ln p(z | \theta_1^{(t)}) > \ln p(z | \theta_1^{(t-1)})$
- we want $\theta_1^{(t)}$ converge to a local maximum in $\ln p(z | \theta_1)$
- EM generates #1, satisfies #2

EM-Algorithm:

Given: $\theta_1^{(t)}$, find $\theta_1^{(t+1)}$:

E-step: set $q_t(\theta_2) = p(\theta_2 | z, \theta_1^{(t)})$, evaluate

$$L_t(z, \theta_1) = \int q_t(\theta_2) \ln \frac{p(z, \theta_2 | \theta_1)}{q_t(\theta_2)} d\theta_2 = \int q_t(\theta_2) \ln p(z, \theta_2 | \theta_1) d\theta_2 - \int q_t(\theta_2) \ln \underbrace{q_t(\theta_2)}_{\text{ignore}} d\theta_2$$

M-step: set $\theta_1^{(t+1)} = \underset{\theta_1}{\operatorname{argmax}} L_t(z, \theta_1)$

EM results in monotonic increases in observed data log-likelihood $\ln p(z | \theta_1)$

• In other presentations, E-step is framed as taking expected value of the complete-data log likelihood with respect to the conditional distribution of $\underline{\theta}_2$ (or \underline{z}) given \underline{x} and current estimates of parameters $\underline{\theta}_1^{(t)}$
 i.e. $E_{\underline{\theta}_2 | \underline{x}, \underline{\theta}_1^{(t)}} [\log L(\underline{\theta}_1; \underline{x}, \underline{\theta}_2)]$ (?) come back (*) harmonise with other sources

EM-missing data: - $\underline{x}_i \in \mathbb{R}^d$ - a vector with missing data
 $\underline{x}_i = \begin{cases} 1. \underline{x}_i^o \text{ (observed sub-vector of } \underline{x}_i) \\ 2. \underline{x}_i^m \text{ (unobserved sub-vector of } \underline{x}_i) \end{cases}$ - missing dim. can be different for different i

- Assuming $\underline{x}_i \stackrel{iid}{\sim} N(\underline{\mu}, \Sigma)$, we require $\hat{\underline{\mu}}_{ML}, \hat{\Sigma}_{ML} = \arg \max_{\underline{\mu}, \Sigma} \sum_{i=1}^N \ln p(\underline{x}_i^o | \underline{\mu}, \Sigma)$

- Intractable: as \underline{x}_i^o (each) has different dimensionality; each dim corresponds to different subsets of dimensions in $\underline{\mu}$ and Σ

Tractable optimisation: If we knew \underline{x}_i^m (and Σ_i), following is tractable:-
 $\hat{\underline{\mu}}_{ML}, \hat{\Sigma}_{ML} = \arg \max_{\underline{\mu}, \Sigma} \sum_{i=1}^N \ln p(\underline{x}_i^o, \underline{x}_i^m | \underline{\mu}, \Sigma)$ easy to optimise

MLEstimation + imputation of missing values: A method for optimising $\sum_{i=1}^N \ln p(\underline{x}_i^o | \underline{\mu}, \Sigma)$ and imputing $\{\underline{x}_1^m, \dots, \underline{x}_N^m\}$
Setup: $p(\underline{x}_i^o | \underline{\mu}, \Sigma) = \int p(\underline{x}_i^o, \underline{x}_i^m | \underline{\mu}, \Sigma) d\underline{x}_i^m = N(\underline{\mu}_i^o, \Sigma_i^o)$ $\left\{ \begin{array}{l} \underline{\mu}_i^o - \text{subvector} \\ \text{of } \underline{\mu} \\ \Sigma_i^o - \text{submatrix} \\ \text{of } \Sigma \end{array} \right.$

Single Gaussian EM:

$$\sum_{i=1}^N \ln p(\underline{x}_i^o | \underline{\mu}, \Sigma) = \sum_{i=1}^N \int q(\underline{x}_i^m) \ln \frac{p(\underline{x}_i^o, \underline{x}_i^m | \underline{\mu}, \Sigma)}{q(\underline{x}_i^m)} d\underline{x}_i^m \quad \text{from } \underline{x}_i^o \quad \left\{ \begin{array}{l} \text{from } \underline{x}_i^o \\ \Sigma_i^o - \text{submatrix} \\ \text{of } \Sigma \end{array} \right.$$

E-step (I): set $q(\underline{x}_i^m) = p(\underline{x}_i^m | \underline{x}_i^o, \underline{\mu}, \Sigma)$ under current $\underline{\mu}, \Sigma$
 $\cdot \underline{x}_i = \begin{bmatrix} \underline{x}_i^o \\ \underline{x}_i^m \end{bmatrix} \sim N \left(\begin{bmatrix} \underline{\mu}_i^o \\ \underline{\mu}_i^m \end{bmatrix}, \begin{bmatrix} \Sigma_i^{oo} & \Sigma_i^{om} \\ \Sigma_i^{mo} & \Sigma_i^{mm} \end{bmatrix} \right) \rightarrow p(\underline{x}_i^m | \underline{x}_i^o, \underline{\mu}, \Sigma) = N(\hat{\underline{\mu}}_i, \hat{\Sigma}_i)$
 via cond. Gaussian results. $\left\{ \begin{array}{l} \hat{\underline{\mu}}_i = \underline{\mu}_i^o + \Sigma_i^{mo} (\Sigma_i^{oo})^{-1} (\underline{x}_i^o - \underline{\mu}_i^o) \\ \hat{\Sigma}_i = \Sigma_i^{mm} - \Sigma_i^{mo} (\Sigma_i^{oo})^{-1} \Sigma_i^{om} \end{array} \right.$

E-step (II): $E_q(\underline{x}_i^m) [\ln p(\underline{x}_i^o, \underline{x}_i^m | \underline{\mu}, \Sigma)]$

• For each i , evaluate:-

$$E_q[(\underline{x}_i - \underline{\mu})^T \Sigma^{-1} (\underline{x}_i - \underline{\mu})] = E_q[\text{trace} \{ \Sigma^{-1} (\underline{x}_i - \underline{\mu})(\underline{x}_i - \underline{\mu})^T \}] \\ = \text{trace} \{ \Sigma^{-1} E_q[(\underline{x}_i - \underline{\mu})(\underline{x}_i - \underline{\mu})^T] \}$$

• Calculated using $q(\underline{x}_i^m) = p(\underline{x}_i^m | \underline{x}_i^o, \underline{\mu}, \Sigma)$ - only \underline{x}_i^m portion of \underline{x}_i integr.
 $q(\underline{x}_i^m) = N(\hat{\underline{\mu}}_i, \hat{\Sigma}_i)$

M-step: Define \hat{x}_i^o -vector where we replace missing values in x_i with \hat{x}_i^m
 \hat{V}_i -matrix of 0s plus sub-matrix $\hat{\Sigma}_i$ in missing dimensions

$$\text{maximise } \sum_{i=1}^N \mathbb{E}_q[\ln p(x_i^o, x_i^m | \mu, \Sigma)]$$

$$\cdot \hat{\mu}_{\text{up}}, \hat{\Sigma}_{\text{up}} = \underset{\mu, \Sigma}{\operatorname{argmax}} \sum_{i=1}^N \mathbb{E}_q[\ln p(x_i^o, x_i^m | \mu, \Sigma)]$$

$$\cdot \hat{\mu}_{\text{up}} = \frac{1}{N} \sum_{i=1}^N \hat{x}_i^o \quad \hat{\Sigma}_{\text{up}} = \frac{1}{N} \sum_{i=1}^N \{(\hat{x}_i^o - \hat{\mu}_{\text{up}})(\hat{x}_i^o - \hat{\mu}_{\text{up}})^T + \hat{V}_i\}$$

Return to E-step: calculate new $p(x_i^m | x_i^o, \hat{\mu}_{\text{up}}, \hat{\Sigma}_{\text{up}})$

High-level: - Hoover mechanics again (1), (2), not fully clear.

(Strategy): 1. maximise log-likelihood only considering observed point x_i^o over unknown μ, Σ

2. use q-distribution (conditional posterior on missing x_i^m $p(x_i^m | x, \theta)$) given data and parameter estimates to make probabilistic state about conjectured values for x_i^m

Implementation: - Each EM update \rightarrow monotonic increase in log likelihood of incomplete dataset ~~p(x|θ)~~ $\log p(x^o | \theta)$

Convergence: - marginal improvements in $\ln p(x_i^o | \theta)$

- initialise μ and Σ (set x_i^m (missing) = 0 and use initial $\hat{\mu}_{\text{ME}}$, $\hat{\Sigma}_{\text{ME}}$ or random initialisation)

- EM objective after each iteration update $\overset{\text{obs.}}{\ln p(x^o | \theta)}$ to μ, Σ $\overset{\text{opt.}}{\ln p(x^o | \theta)}$



① 1) Output $\hat{\mu}_{\text{ME}}, \hat{\Sigma}_{\text{ME}}$ that maximises $\sum_{i=1}^N \ln p(x_i^o | \mu, \Sigma)$

② 2) Output conditional posterior $q(x_i^m) = p(x_i^m | x_i^o, \hat{\mu}_{\text{ME}}, \hat{\Sigma}_{\text{ME}})$ on missing dimensions which gives mean and uncertainty of missing data x_i^m
 i.e. probability distri on $x_i^m \rightarrow \text{mixture } \{x_1^m, \dots, x_N^m\}$

16 - Key equations

K-means: Hard clustering model, each observation x_i assigned to only one cluster
 $c_i = k$ for some $k \in \{1, \dots, K\}$, no accounting for boundary cases probabilist.

soft clustering: Data can be clustered ^{semi-}probabilistically (K-means is non-prob.)
weighted K-means (algorithm) (probabilistic semantics, but no formal distributions defined)

Input: data $\{x_i\}_{i=1}^N, x_i \in \mathbb{R}^d$
 GOAL: Minimise $L = \sum_{i=1}^N \sum_{k=1}^K \phi_{ik}(k) \frac{\|x_i - \mu_k\|_2^2}{\beta}$ over all ϕ_{ik} $i=1, \dots, N$
 $\mu_k \quad k=1, \dots, K$

conditions: $\phi_{ik} > 0, \sum_{k=1}^K \phi_{ik} = 1, \beta > 0$

Heuristic: 1. update ϕ : For each i , update word allocation weights

$$\phi_{ik} = \frac{\exp\left\{-\frac{1}{\beta} \|x_i - \mu_k\|_2^2\right\}}{\sum_j \exp\left\{-\frac{1}{\beta} \|x_i - \mu_j\|_2^2\right\}} \quad \text{for } k=1, \dots, K$$

2. update μ : For each k , update μ_k with weighted average:

$$\mu_k = \frac{\sum_{i=1}^N x_i \phi_{ik}}{\sum_{i=1}^N \phi_{ik}}$$

• substitute $\prod_{i=1}^N c_i = k$ (K-means) $\leftarrow \frac{\phi_{ik}}{\beta}$ into L and μ_k update with $\phi_{ik} = \frac{\exp\left\{-\frac{1}{\beta} \|x_i - \mu_k\|_2^2\right\}}{\sum_j \exp\left\{-\frac{1}{\beta} \|x_i - \mu_j\|_2^2\right\}}$
 (normalised RBF kernel)

• $\phi_{ik} = \{0, 1\} \rightarrow$ hard clustering

• ϕ_{ik} large $\rightarrow x_i$ close to μ_k ; ϕ_{ik} far from μ_k

Mixture models: Above: weight vector ϕ_i is akin to probability x_i assigned to cluster

Mixture models: A probabilistic model where ϕ_i is actually a probability distribution according to model

Mixture model: 1. Prior distribution on cluster assignment indicator c_i
 (high-level) 2. likelihood \rightarrow on observation x_i , given assign. c_i

Analog with SL: 1. cluster prior \rightarrow class prior π_L
 (Bayes class.) 2. cluster-cond. likelihood \rightarrow class-cond. likelihood π_{iL} $\begin{cases} \text{but cluster} \\ i \text{ is} \\ \text{latent/auxiliary} \end{cases}$

Mixture models: 1. generative model (derives p.d. on data x)
 (features) 2. A weighted combination of simple distributions
 (see diagram $\begin{cases} \text{discrete} \\ - \text{complex distn} \end{cases}$ p.d.) (same distn family)

generic mixture model - data: $\{\underline{x}_i\}_{i=1}^N \in \mathbb{R}^d$
(generative process) model param: $\pi \in \mathbb{R}^K \quad \underline{\pi} = [\pi_1 \dots \pi_K] \quad \underline{\theta}_1, \dots, \underline{\theta}_K$

generative process: For obs $i=1, \dots, N$:

1. generate cluster assign. $c_i \sim \text{discrete}(\pi)$ $p(c_i=k|\pi) = \pi_k$
 2. generate obs. $\underline{x}_i \sim p(\underline{x}|c_i)$
- each \underline{x}_i randomly assigned to cluster c_i via distri π
 - c_i indexes cluster assignment for $\underline{x}_i \rightarrow$ i) selects index of $\underline{\theta}$ 2) \underline{x}_i since param \rightarrow clustered together $\underline{\theta}$ to gen. $\underline{x}_i (\theta_1, \dots, \theta_K)$

EM-generic mixture algorithm:

TIN: $\{\underline{x}_i\}_{i=1}^N \in \mathbb{R}^d$

GOAL: Maximize $L = \ln p(\underline{x}|\pi, \underline{\theta})$ where $p(\underline{x}|\underline{\theta}_K)$ is problem spec.

Update until incremental improvement to L is small:-

• HMM until incremental improvement to L is small:- (posterior prob.)

1. E-step: For $i=1, \dots, n$ set:-

$$\phi_i(k) = \frac{\pi_k p(\underline{x}_i | \underline{\theta}_k)}{\sum_j p(\underline{x}_i | \underline{\theta}_j) \pi_j} \quad \text{for } k=1, \dots, K$$

2. M-step: For $k=1, \dots, K$ define $n_k = \sum_{i=1}^N \phi_i(k)$ and set

$$\pi_k = \frac{n_k}{n} \quad \underline{\theta}_k = \underset{\underline{\theta}}{\operatorname{argmax}} \sum_{i=1}^N \phi_i(k) \ln p(\underline{x}_i | \underline{\theta})$$

(*) (*)

]

- Similar to generalisation of Bayes class for any $p(\underline{x}|\underline{\theta}_K)$

- comments: mixture model \rightarrow specify cluster prior and cluster cond. likelihood;

(*) (*) (?)

Gaussian Mixture Models: parameters: $\pi \in \mathbb{R}^K \quad \underline{\pi} = [\pi_1, \dots, \pi_K]$
 (μ_R, Σ_R) - mean and covariance of k^{th} Gaussian
 $\mu_R \in \mathbb{R}^d, \Sigma_R \in \mathbb{R}^{d \times d}; k=1, \dots, K$

generate data: For i^{th} obs \underline{x}_i :

1. Assign i^{th} obs \rightarrow to cluster $c_i \sim \text{discrete}(\pi)$
2. Generate $\underline{x}_i \sim N(\mu_{c_i}, \Sigma_{c_i})$

$\underline{\mu} = [\mu_1, \dots, \mu_K] \in \mathbb{R}^{d \times K}$ $\underline{\Sigma} = \{\Sigma_1, \dots, \Sigma_K\} \in \mathbb{R}^{d \times d \times K}$ goal: estimate $\pi, \underline{\mu}, \underline{\Sigma}$

(numpy)

Generic EM objective:

$$\ln p(\underline{x} | \theta_1) = \int q(\theta_2) \ln \frac{p(\underline{x}, \theta_2 | \theta_1)}{q(\theta_2)} + \int q(\theta_2) \ln \frac{q(\theta_2)}{p(\theta_2 | \underline{x}, \theta_1)} d\theta_2$$

EM objective for GMM: (Joint)

$$\sum_{i=1}^N \ln p(\underline{x}_i | \underline{\pi}, \underline{\mu}, \underline{\Sigma}) = \sum_{i=1}^N \sum_{k=1}^K q(c_i=k) \ln \frac{p(\underline{x}_i, c_i=k | \underline{\pi}, \underline{\mu}, \underline{\Sigma})}{q(c_i=k)} + \sum_{i=1}^N \sum_{k=1}^K q(c_i=k) \ln \frac{q(c_i=k)}{p(c_i | \underline{x}_i, \underline{\pi}, \underline{\mu}, \underline{\Sigma})}$$

Maximum likelihood EM for GMM:

FIN: $\{\underline{x}_i\}_{i=1}^N$ $\underline{x}_i \in \mathbb{R}^d$

GOAL: Maximise $L = \sum_{i=1}^N \ln p(\underline{x}_i | \underline{\pi}, \underline{\mu}, \underline{\Sigma})$

• Heave until incremental improvement to L is small:-

1. E-step: For $i=1, \dots, N$, set:-

$$\phi_i(k) = \frac{\pi_k N(\underline{x}_i | \underline{\mu}_k, \underline{\Sigma}_k)}{\sum_j \pi_j N(\underline{x}_i | \underline{\mu}_j, \underline{\Sigma}_j)} \quad k=1, \dots, K$$

2. M-step: For $k=1, \dots, K$,

a) Define $\alpha_k = \sum_{i=1}^N \phi_i(k)$

$$b) \text{ Update values: } \hat{\pi}_k = \frac{\alpha_k}{N} \quad \hat{\underline{\mu}}_k = \frac{1}{\alpha_k} \sum_{i=1}^N \phi_i(k) \underline{x}_i \quad \hat{\underline{\Sigma}}_k = \frac{1}{\alpha_k} \sum_{i=1}^N \phi_i(k) (\underline{x}_i - \hat{\underline{\mu}}_k)(\underline{x}_i - \hat{\underline{\mu}}_k)^T$$

EM-GMM: (one it.)

PRELIM: Set N q -distributions over auxiliary variables c_1, \dots, c_N via Bayes

$$q(c_i=k) \leftarrow p(c_i=k | \underline{x}_i, \underline{\pi}, \underline{\mu}, \underline{\Sigma}) \quad \text{as } N \text{ condit. post } p(c_i=k | \underline{x}_i, \underline{\pi}, \underline{\mu}, \underline{\Sigma}) \text{ indexed by } \underline{x}_i$$

$$p(c_i=k | \underline{x}_i, \underline{\pi}, \underline{\mu}, \underline{\Sigma}) \propto p(c_i=k | \underline{\pi}) p(\underline{x}_i | c_i=k, \underline{\mu}, \underline{\Sigma})$$

$$q(c_i=k) = \frac{\pi_k N(\underline{x}_i | \underline{\mu}_k, \underline{\Sigma}_k)}{\sum_j \pi_j N(\underline{x}_i | \underline{\mu}_j, \underline{\Sigma}_j)} \rightarrow \phi_i(k) \quad \text{④ solve for posterior of } c_i \text{ given } \underline{\pi}, \underline{\mu}, \underline{\Sigma}$$

E-step: Take expectations using updated q s

$$Q = \sum_{i=1}^N \sum_{k=1}^K \phi_i(k) \ln p(\underline{x}_i, c_i=k | \underline{\pi}, \underline{\mu}_k, \underline{\Sigma}_k) + \text{constant wrt } \underline{\pi}, \underline{\mu}, \underline{\Sigma}$$

M-step: Maximise Q wrt $\underline{\pi}$ and each $\underline{\mu}_k, \underline{\Sigma}_k$

Reasons for using EM:

original objective: $L = \sum_{i=1}^N \ln \left(\sum_{k=1}^K p(x_i, c_i=k | \pi, \mu_k, \Sigma_k) \right) = \sum_{i=1}^N \ln \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$

log-sum \Rightarrow optimising π, μ_k, Σ_k hard.

EM-M-step: $Q = \sum_{n=1}^N \sum_{k=1}^K d_{ik} \{ \ln \pi_k + \ln N(x_i | \mu_k, \Sigma_k) \} + \text{constant w.r.t. } \pi, \mu, \Sigma$

$$\ln p(x_i, c_i=k | \pi, \mu_k, \Sigma_k)$$

cluster no. selection: ML for GMM \rightarrow overfitting. use Dirichlet priors on π (mixing weight vector) which encourages many Gaussians to disappear

U7 - Key equations:

object recommendation: content filtering
collaborative filtering

content filtering: use known information about products, uses to make recommended.
(e.g. movie info, prod info, demograp, questionnaire w.r.t.) (a prior.)

collaborative filtering: use previous user's input/behav. to make recommended.
(use feedback, domain fee)

collaborative filtering: item-based (evaluate user pref on 'neighbouring' items by some user)

(neighborhood based) \hookrightarrow user-based (like-minded users who complement ratings, similarity score between user and them; let them vote on recommendation)

collaborative-fitting: ratings \rightarrow characterise items and users from factors inferred from ratings
(latent/latent factor)

factors measure how much user likes movies that score high on that factor \rightarrow movie or user.

①: location-based / latent factor models embed objects and users into points in \mathbb{R}^d

②: recommendations based on proximity of users and objects in latent embedding space inferred from data

embeddings estimated via matrix factorisation, used to produce estimated ratings

Matrix factorisation: form ratings matrix M
(technique, prob., goal)

Many $M_{ij} = \text{NaN}$ (sparse matrix)

goal: probabilistically fill in missing values

solution: hope to have prediction of $\text{missing}_{\text{prior}}$ ratings for recommended highly rated obj j amongst pred.

③: MF - estimate a low-rank factorisation using observed data, ignoring unobserved.

low-rank matrix factorisation

- 1. Account for sparsity

$m \geq M$

- 2. low rank $\Rightarrow d \ll \min\{m, M\}$

- 3. Estimate $\{N_1, N_2\}$

N_1 : user-location vectors $u_i \in \mathbb{R}^d$

N_2 : object-location vectors $v_j \in \mathbb{R}^d$

4. d - no. of factors (latent things considered)

$$\left(\begin{array}{c|cc|c} & \vdots & \vdots & \\ \text{M}_{ij} & | & | & \\ & \vdots & \vdots & \\ \hline & \vdots & \vdots & \end{array} \right) = \left(\begin{array}{c|c} & \vdots \\ u_i \in \mathbb{R}^{N_1 \times d} & | \\ & \vdots \end{array} \right) \left(\begin{array}{c|c} & \vdots \\ v_j \in \mathbb{R}^{d \times N_2} & | \\ & \vdots \end{array} \right)$$

low-rank factorisation of M : (i) MF model maps users, items to a joint latent factor space of dimensionality d , determined here by low rank ($\text{rank}(M)$)

user-item interactions (i.e. ratings) modelled as inner products in that space
user i , object j associated with $u_i \in \mathbb{R}^d, v_j \in \mathbb{R}^d$ respectively
elements of u_i and v_j measure extent of interest user has in items high on those factors / extent to which item possesses those factors

$$\Rightarrow M_{ij} \approx u_i^T v_j \quad (\hat{M}_{ij} = u_i^T v_j) \quad (ii)$$

challenge \rightarrow mapping $i, j \rightarrow u_i, v_j \in \mathbb{R}^d$ (factor-vectors)

conventional SVD fails due to sparsity; missing data imputation - comp. expensive

low-rank factors - intuition:

sparsity, subset of N_1 users rated movie ad , another subset 24 (think columns of M)

some ratings will overlap

similar movies (e.g. spy thriller / any factor) \rightarrow rating across columns of M highly correlated

low-rank factorisation enforce this correlation by requiring users and objects to be in a smaller latent embedding space $\mathbb{R}^d \rightarrow$ restrictions on u and v

low-rank $\Rightarrow N_1$ dim columns of M do not 'fill up' \mathbb{R}^{N_1} , span region of \mathbb{R}^{N_1} that is smaller (i.e. \mathbb{R}^d)

low-rank restriction / high correlation \Rightarrow 'borrow' information across movies to learn factors

some agents for user-oriented (rows of M) (users with similar tastes across many)

Probabilistic Matrix Factorisation: Model for estimating low-rank factorisation (setup / gen. model)

observation/index sets: define observation set Ω containing observed (i, j)

$$\Omega = \{(i, j) : M_{ij} \text{ is measured } (\neq \text{NaN})\} \quad (i, j) \in \Omega \text{ if measured}$$

define index sets: index set of objects rated by user i Ω_{ui}

$$\Omega_{ui} = \{j : M_{ij} \text{ is measured } (\neq \text{NaN})\} \quad j \in \Omega_{ui} \text{ if measured}$$

index set of users who rated object j Ω_{vj}

$$\Omega_{vj} = \{i : M_{ij} \text{ is measured } (\neq \text{NaN})\}$$

- PMF: • For N_1 users, N_2 objects
(generative model) $u_i \sim N(0, \lambda^{-1} I) \quad i=1, \dots, N_1$ } priors with
 $v_j \sim N(0, \lambda^{-1} I) \quad j=1, \dots, N_2$ } zero-mean isotropic covariance
- Data distribution: $M_{ij} \sim N(u_i^T v_j, \sigma^2)$ for each $(i, j) \in \Omega$
- Criteria for EM: M_0 - observed M_m missing (include, EM not necessary)
(discussion) $p(M_0 | U, V) = \int p(M_0, M_m | U, V) dM_m \rightarrow$ EM tool for $\max_{U, V} p(M_0 | U, V)$
- Criteria: 1. $p(M_0 | U, V)$ - diffi to max 2. $p(M_0, M_m | U, V)$ - easy to use
(marginal likel.) (joint likel.)
3. $p(M_m | M_0, U, V)$ - known
- If $p(M_0 | U, V)$ no issues for inference \rightarrow NO (i.e. closed form updates for ML/MAP estimates)
- Try $\max_{U, V} p(M_0, U, V)$ over U, V : $p(M_0, U, V) = \prod_{(i, j) \in \Omega} p(M_{ij} | u_i, v_j) \times \left[\prod_{i=1}^{N_1} p(u_i) \right] \times \left[\prod_{j=1}^{N_2} p(v_j) \right]$

PMF: $U_{MAP}, V_{MAP} = \underset{U, V}{\operatorname{argmax}} \ln p(M_0, U, V)$
(log-joint like and MAP objective)

$$= \underset{U, V}{\operatorname{argmax}} \sum_{(i, j) \in \Omega} \ln p(M_{ij} | u_i, v_j) + \sum_{i=1}^{N_1} \ln p(u_i) + \sum_{j=1}^{N_2} \ln p(v_j)$$

MAP objective: $L = - \sum_{(i, j) \in \Omega} \frac{1}{\sigma^2} \|M_{ij} - u_i^T v_j\|_2^2 - \sum_{i=1}^{N_1} \frac{\lambda}{2} \|u_i\|_2^2 - \sum_{j=1}^{N_2} \frac{\lambda}{2} \|v_j\|_2^2$

• $\nabla_{u_i} L$ and $\nabla_{v_j} L$ yields solutions for u_i and v_j ; but they are dependent on each other \rightarrow coordinate ascent (k-means / GMM)

PMF (MAP inference co-ordinate ascent): Algorithm:

IN: Incomplete ratings matrix M , indexed by Ω , rank d
OUT: N_1 user locations, N_2 object locations, $u_i, v_j \in \mathbb{R}^d$
INIT: Each v_j e.g. $v_j \sim N(0, \lambda^{-1} I)$ $\forall j = 1, \dots, N_2$

For each iteration do:

For $i = 1, \dots, N_1$ update user location:

$$u_i = (\lambda \sigma^2 I + \sum_{j \in \Omega_{u_i}} v_j v_j^T)^{-1} \left(\sum_{j \in \Omega_{u_i}} M_{ij} v_j \right)$$

For $j = 1, \dots, N_2$

$$v_j = (\lambda \sigma^2 I + \sum_{i \in \Omega_{v_j}} u_i u_i^T)^{-1} \left(\sum_{i \in \Omega_{v_j}} M_{ij} u_i \right)$$

PREDICT: that $\hat{M}_{ij} = \hat{u}_{i, \text{MAP}}^T \hat{v}_{j, \text{MAP}}$ (i.e. with estimated v_j, u_i after algo converge.)

- PMF: - consider updating for object j for covariate vector \mathbf{v}_j
- Analogies with SL: - "minimise sum of squared error"
- (RR) $\sum_{(i,j) \in Q} \frac{1}{\delta^2} (M_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2$ with penalty $\lambda \|\mathbf{v}_j\|_2^2$
- use location (\mathbf{u}_i)
is covariate vector
for rating M_{ij}
- Ridge regression for \mathbf{v}_j : $\mathbf{v}_j = (\lambda \mathbf{I} + \sum_{i \in \mathcal{S}_{\mathbf{v}_j}} \mathbf{u}_i \mathbf{u}_i^\top)^{-1} \left(\sum_{i \in \mathcal{S}_{\mathbf{v}_j}} M_{ij} \mathbf{u}_i \right)$
- M_{ij} analogous to SL output we want to predict
- Analogy- $N_1 + N_2$ coupled ridge regression problems
- No covariates \rightarrow re estimate them from the distribution family defining our model
- $\hat{w}_{RR} = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \Leftrightarrow \mathbf{v}_j = (\lambda \mathbf{I} + \sum_{i \in \mathcal{S}_{\mathbf{v}_j}} \mathbf{u}_i \mathbf{u}_i^\top)^{-1} \left(\sum_{i \in \mathcal{S}_{\mathbf{v}_j}} M_{ij} \mathbf{u}_i \right)$ By symmetry, holds for \mathbf{u}_i
- PMF: - remove Gaussian priors on \mathbf{u}_i and \mathbf{v}_j
- Analogies with SL: $\mathbf{v}_j = \left(\sum_{i \in \mathcal{S}_{\mathbf{v}_j}} \mathbf{u}_i \mathbf{u}_i^\top \right)^{-1} \left(\sum_{i \in \mathcal{S}_{\mathbf{v}_j}} M_{ij} \mathbf{u}_i \right) \Leftrightarrow \hat{w}_{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$
- LS
- sequence of updates for columns and rows \rightarrow maximises condit. likeli. $p(M|U, V)$
- ML only max conditional likelihood
- MAP accounts for priors, estimate free mode of posterior \rightarrow maximise joint likelihood $p(M, U, V)$
- Additional restriction: invertibility of $\mathbf{u}_i \mathbf{u}_i^\top$ and $\mathbf{v}_j \mathbf{v}_j^\top \Rightarrow$ all users rated > d objects
all objects rated > d times

8. Key equations

Topic Models: used in info retrieval - discover themes, annotate documents, organise etc.

- Probabilistic topic model:
1. Topics - estimate probability distributions on words
(topic-specific word distributions) on a fixed vocabulary. Each p.d. assigns most probability to subsets of co-occurring words in a 'theme'
 2. Topic proportions - estimates a probability distribution on topics for each document. Documents contain variable proportions of no. of topics
 - 3) Topic assignment - Assigns every word in a document to a topic

1), 2), 3) must be estimated

Latent Dirichlet Allocation - most well-known topic model intra-document

use of bag-of-words representation for a document (sequential dependencies between words binned)

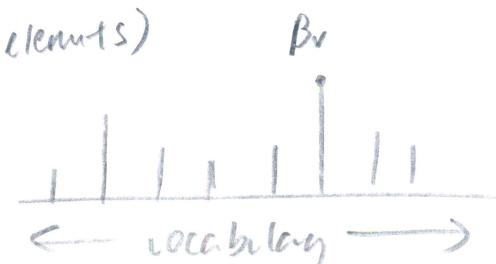
LDA (ingredients): 1. collection of topics 2. collection of topic proportions

topics: - represented by a probability-distribution on a global, corpus-wide set of vocabulary words (V , words indexed by $v=1, \dots, V$)

(representation) some set of vocabulary words (V , words indexed by $v=1, \dots, V$)

r^{th} topic β_R of K topics - $\beta_R = [\beta_{R1}, \dots, \beta_{RV}]_R$ (V elements)

each element β_{Rv} is probability of word v occurring in topic R



document / topic proportion:

• Represent a document as a document-specific distribution on topics / signature topic proportion.

• for d th document out of D documents:

d th topic proportion vector $\underline{\theta}_d$ $\underline{\theta}_d = [\underline{\theta}_d^1, \underline{\theta}_d^2, \dots, \underline{\theta}_d^K]$

• $\underline{\theta}_d = [\theta_1, \dots, \theta_K]_d$ (K -elements) \leftarrow topics \rightarrow

• element θ_K is probability of K th topic occurrence in document d .

LDA (Generative process):

1. generate each topic, a distribution over words:- (Topic)

$$\underline{\beta}_K \sim \text{Dirichlet}(\gamma) \quad k=1, \dots, K$$

2. For each document $d=1, \dots, D$, generate distribution over topics:- (Topic prop.)

$$\underline{\theta}_d \sim \text{Dirichlet}(\alpha) \quad d=1, \dots, D$$

3. For n th word in d th document, (Topic assignment) (*) @

a) Allocate word to a topic $c_{dn} \sim \text{Discrete}(\underline{\theta}_d)$

b) Generate word from selected topic $x_{dn} \sim \text{Discrete}(\underline{\beta}_{c_{dn}})$

• 1./2. Priors \Rightarrow topic and topic proportions drawn iid from a Dirichlet distribution over K topics and D documents

• 3. For d th document, n th word, decide which topic word comes from i.e. a topic indicator from a discrete distri. (Multinomial) using distribution on topics for d th document ($\underline{\theta}_d$) (i.e. $\underline{\theta}_d$ acts as a parameter of discrete distri.)

$$c_{dn} = \{1, \dots, K\}$$

topic chosen \Rightarrow at data level, n th word in d th document is generated from a discrete distri. using topic selected by topic indicator c_{dn} .

⑦ Do not know c_{dn} (topic indicators), $\underline{\theta}_d$ (topic proportions), $\underline{\beta}_K$ (Topics)

Dirichlet distribution: A continuous distribution over discrete probability vectors

• $\underline{\beta}_K$ - probability vector

• $\underline{\gamma}$ - positive parameter vector

$$\underline{\gamma} \in \mathbb{R}^V \quad \beta_{k,v} \in \mathbb{R}^V$$

↳ Shaded indexes
Additionality

concentration param. (γ):

$$p(\underline{\beta}_K | \underline{\gamma}) = \frac{P\left(\sum_{v=1}^V \gamma_v\right)}{\prod_{v=1}^V P(\gamma_v)} \prod_{v=1}^V \frac{\gamma_{v-1}}{\prod_{r=1}^V \beta_{r,v}} \quad V \geq 2 \text{ (no. of words)}$$

→ As K -topics, there will be K draws from Dirichlet distri.
• γ controls the extent to which probability mass is concentrated across each element $\beta_{k,v}$ for a particular β_K vector.

- $\gamma \rightarrow \infty$: each element $\beta_{k,v}$ of topic vector β_k becomes uniform for all K topics
- $\gamma \rightarrow 0$: probability mass concentrated on subset of dimensions V (β_k element sparsity)
- $\gamma < 1$: A priori belief that topics place probability mass on small subset of total words in vocab (V).

- Prior distribution on topic proportions parameterised by $\alpha_{ik} < 1 \Rightarrow$ a priori, document only uses subset of topics available

Inference algorithm: variational inference
LDA outputs: 1) A set of K topics specific word distributions $\{\hat{\beta}_1, \dots, \hat{\beta}_K\} = \hat{\beta}$
2) A set of D document-specific topic distributions $\{\hat{\theta}_1, \dots, \hat{\theta}_D\}$

LDA and MF @ For a document, what is $p(x_{dn}=i | \beta, \theta_d)$
(Probability of n^{th} word in d^{th} doc equal to in vocab list, given set of topics and topic proportion for d^{th} doc?)

① Marginal distri on x_{dn} with c_{dn} integrated out:-

$$p(x_{dn}=i | \beta, \theta_d) = \sum_{R=1}^K p(x_{dn}=i, c_{dn}=R | \beta, \theta_d) = \sum_{R=1}^K p(x_{dn}=i | \beta, c_{dn}=R) p(c_{dn}=R | \theta_d)$$

• Let $B = [\beta_1, \dots, \beta_K]$, $B \in \mathbb{R}^{V \times K}$ $\Theta = [\theta_1, \dots, \theta_D] \Theta \in \mathbb{R}^{K \times D}$

• $p(x_{dn}=i | \beta, \theta) = (B\Theta)_{id}$ = product of 2 matrices with non-negative entries

• LDA is an example of NMF

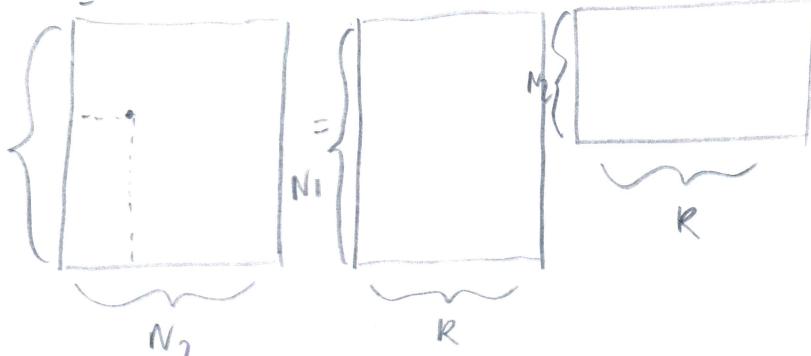
Non-negative matrix factorisation:

• Data X - none missing $(X)_{ij} \quad N_1$

• $X_{ij} \geq 0 \quad \forall i, j$

• Factors $\rightarrow W_{ik}, H_{kj} \geq 0 \quad \forall k, j, i$

⑦: $X_{ij} \geq \sum_{k=1}^K W_{ik} H_{kj}$



NMF application domains

1) Text data - X_{ij} word term freq (occurrence word i in doc j)

$X_{i,j}$ - word histogram word term freq in doc j

$X_{i,:}$ - histo of word i occurrences

2) Images - vectorise $N \times M$ image into $X_{:,j}$ (column)

3) Discrete grouped data - Quantise continuous set of K means (compress info \rightarrow prototype)
- X_{ij} - count in times group j uses cluster i (cluster medoids)

- X_{ij} - count in times group j uses cluster i in graph, e.g. indexing ft. j index songs

2 NMF objectives: NMF #1: $\|X - WH\|_F^2 = \sum_i \sum_j (X_{ij} - (WH)_{ij})^2$ (squared error)

NMF #2: $D(X||WH) = -\sum_i \sum_j [X_{ij} \ln(WH)_{ij} - (WH)_{ij}] (KL\text{-diverg.})$

s.t. $W_{ik}, H_{kj} \geq 0 \quad \forall i, k, j$

NMF Algorithm - works to minimise objective $\min_{\underline{W}, \underline{H}} F(\underline{h})$ by generating seq $\underline{f}(\underline{h}^1) \geq F(\underline{h}^2) \dots$ and ensuring it converges to a local minimum.

- minimise via an auxiliary function \rightarrow multiplicative update for $\underline{W}, \underline{H}$

NMF #1 (multiplicative update for $\|X - \underline{W}\underline{H}\|_2^2$) Algorithm:

$$\text{minimise}_{\underline{W}, \underline{H}} \sum_{ij} (X_{ij} - (\underline{W}\underline{H})_{ij})^2 \text{ st. } W_{ik} \geq 0, H_{kj} \geq 0$$

- randomly init $\underline{H}, \underline{W}$ with non-negative values

- iterate first for all values in \underline{H} , then all in \underline{W} :

$$H_{kj} \leftarrow H_{kj} \frac{(W^T X)_{kj}}{(W^T W H)_{kj}}$$

$$W_{ik} \leftarrow W_{ik} \frac{(X H^T)_{ik}}{(W H H^T)_{ik}}$$

- until change in $\|X - \underline{W}\underline{H}\|_2^2$ is small

NMF #2 (multiplicative update for $D(X||W\underline{H})$) Algorithm

$$\text{minimise}_{\underline{W}, \underline{H}} \sum_{ij} \left[X_{ij} \ln \frac{1}{(W H)_{ij}} + (W H)_{ij} \right] \text{ st. } W_{ik} \geq 0, H_{kj} \geq 0$$

- randomly initialise \underline{W} and \underline{H} with non-negative values

- iterate first for all values in \underline{H} , then all in \underline{W} :

$$H_{kj} \leftarrow H_{kj} \left(\frac{\sum_i W_{ik} X_{ij} / (W H)_{ij}}{\sum_i W_{ik}} \right)$$

$$W_{ik} \leftarrow W_{ik} \left(\frac{\sum_j H_{kj} X_{ij} / (W H)_{ij}}{\sum_j H_{kj}} \right)$$

until change in $D(X||W\underline{H})$ is small

($X_{ij} \geq 0 \Rightarrow$ not true)
NMF #1 & #2 - NMF #1 - squared error penalty $\Rightarrow X_{ij} \stackrel{iid}{\sim} N(\sum_k W_{ik} H_{kj}, \sigma^2)$ st
 through prob. lens) - NMF #2 - divergence penalty $\Rightarrow X_{ij} \stackrel{iid}{\sim} Po((W H)_{ij})$ $Po(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$ non-negative

#2 - NMF \rightarrow Estimating expected value for data matrix $E[X]$

under generative assumption that likelihood is poisson, param by
 $(i, j)^{\text{th}}$ element of product of NMF factors $(W H)_{ij}$

- LDA-NMF:
1. Form term freq matrix X (x_{ij} = # times word i appears in doc j)
 2. Run NMF for estimates $\underline{W}, \underline{H}$
 3. For $K=1, \dots, K$:
 - i) Set $\underline{w}_k = \sum_i \underline{w}_{ik}$
 - ii) Divide \underline{w}_{ik} by w_k for all i (n^{th} column of \underline{w})
 - iii) Multiply \underline{h}_{kj} by w_k for all i (n^{th} row of \underline{H})
- } 'post processing'
- (*) Each column vector \underline{w}_k of \underline{W} (after post-processing) \rightarrow p.d. on words for topic K (i.e. topic)
- Each column vector \underline{h}_j of \underline{H} ($-\cdot-\cdot-$) \rightarrow p.d. on topics for document j

Application: Put face images along columns of data matrix X

(Face modelling): $X \approx \underline{W}\underline{H}^T$; $x_{ij} \approx \sum_{k=1}^K w_{ik} h_{kj}$

- Columns of \underline{W} - basis images
- Columns of \underline{H} - encodings ($1-1$ correspondence with a face x_{ij})
- Encoding-efficient where face is a linear comb. of basis images

K-means: Face(column of X) represented by prototypical basis image (column of \underline{W}) (n.a.)

- Columns of \underline{H} (n.a.) (m cluster assign.)

SVD: Face(column of X) represented by an eigenface basis image (column of \underline{W})

(Face modelling): and a non-sparse ratings^{m.} of encodings \underline{H}

Face is a linear comb of eigenfaces

• Difficult to interpret

NMF: Face(column of X) represented by interpretable parts based localised features as single basis image (col. of \underline{W}) col sparse encodgs mat. \underline{H}

NMF \Rightarrow non-negativity of word H

• Face is linear comb of parts/localised features

Topic modelling - Each col. of X (doc j) as linear combination of topics

PCA - Key equations

Dimensionality reduction: High dimensional data, where 'info' does not use full d dimensions

Principal Components Analysis: A way of automatically mapping data \underline{x}_i into a low-d dimensional co-ordinate system

1. Captures most info in data \rightarrow few dimensions 2. Extensions allows to handle missing data and unmap data

PCA (2D \rightarrow 1D intuition): Approximate $\underline{x}_i \in \mathbb{R}^2$ via unit length vector q ($\|q\|=1$) - projecting vector onto a scalar

(Minimum error form): Unit length q : $q^T q = 1$

$\cdot x$ (length) - $(q^T x_i)$ after

$\cdot x - (q^T x_i)q$ - projecting x onto line def by q

takes vector q and stretches it to x (scalar mult. of q)

$\underline{x}_i \in \mathbb{R}^d$



PCA (2D → 1D): selecting q : minimise squared approximation error
 (min. error term) $q = \underset{q}{\operatorname{argmin}} \sum_{i=1}^N \|x_i - q q^T x_i\|_2^2 \text{ s.t. } q^T q = 1$

$\cdot q q^T x_i = (q^T x_i) q$
 -approximation of x_i
 by stretching q to X .

PCA (2D → 1D) (eigenvector, eigenvalue, maximise data cov.) $\left[\begin{array}{l} \text{- defining } \underline{X} = [\underline{x}_1, \dots, \underline{x}_N] \text{ so } \underline{X} \in \mathbb{R}^{d \times N} \\ \Leftrightarrow \underset{q}{\operatorname{argmax}} q^T (\underline{X} \underline{X}^T) q \text{ s.t. } q^T q = 1 \end{array} \right]$ (so in a different way to previous stacking)
 Eigenvalue/eigenvector problem $q = 1^{\text{st}} \text{ eigenv. of } \underline{X} \underline{X}^T$
 $\lambda = q^T (\underline{X} \underline{X}^T) q = 1^{\text{st}} \text{ eigenvalue}$

PCA (general): consider K eigenvectors:

$q = \underset{q}{\operatorname{argmin}} \sum_{i=1}^N \|x_i - \sum_{k=1}^K (x_i^T q_k) q_k\|_2^2 \text{ s.t. } q_k^T q_{k'} = \begin{cases} 1 & k=k' \\ 0 & \text{otherwise} \end{cases}$ (maximising projected variance)

$\Leftrightarrow \underset{q}{\operatorname{argmin}} \sum_{i=1}^N x_i^T x_i - \sum_{k=1}^K q_k^T \left(\sum_{i=1}^N x_i x_i^T \right) q_k$

(*) - units to eigenvalues as optimisation in slides 29

K vectors $Q = [q_1 | \dots | q_K]$ define a K dimensional sub-space with which to represent data

$\cdot \underline{x}_{\text{proj}} = \begin{bmatrix} q_1^T \underline{x} \\ \vdots \\ q_K^T \underline{x} \end{bmatrix}$ (scalars that vertically stretch q_k by magnitude / projection of \underline{x} onto line defined by q_k (see earlier))

(*) reconstruction approx: $\underline{x} \approx \sum_{k=1}^K (q_k^T \underline{x}) q_k = Q \underline{x}_{\text{proj}}$

(*) call built-in functions to find K largest eigenvectors and eigenvalues

PCA (Eigenvalue/eigenvectors / SVD)

reformulate problem as: - find (λ, q) : $(\underline{X} \underline{X}^T) q = \lambda q$

since $(\underline{X} \underline{X}^T)$ is P.S.D. $\exists \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$

$q_k^T q_k = 1 \quad q_k^T q_{k'} = 0$

Justifying $\underline{X} \underline{X}^T \in S^r_+$ via SVD:

$\underline{X} \in \mathbb{R}^{d \times n} \rightarrow \underline{U} \underline{S} \underline{V}^T$ with $\underline{U} \in \mathbb{R}^{d \times d}$ $\underline{S} \in \mathbb{R}^{d \times n}$ $\underline{V}^T \in \mathbb{R}^{n \times n}$ (Full-SVD) $s_{ii} > 0$ (square)

(orthog.) $\Rightarrow \underline{U}_{:,i}^T \underline{U}_{:,j} = 0 \quad \|\underline{U}_{:,i}\|^2 = 1 \quad \forall i, j, i \neq j \Rightarrow \underline{U}^T \underline{U} = \underline{U} \underline{U}^T = I$

$(\underline{X} \underline{X}^T) = \underline{U} \underline{S} \underline{V}^T \underline{V}^T \underline{S}^T \underline{U}^T = \underline{U} \underline{S}^2 \underline{U}^T$ with $Q = \underline{U} \quad \lambda_i = (S^2)_{ii} \geq 0 \Leftrightarrow \boxed{(\underline{X} \underline{X}^T) \underline{U} = \underline{U} \underline{S}^2}$

PCA (pre-processing): - PCA sensitive to scaling

- mean centering - subtract mean from every dimension of x_i
- ensures every principal component passes through origin

PCA-digits example: 16×16 handwritten 3s $x_i \in \mathbb{R}^{256}$; stack $\underline{X} \in \mathbb{R}^{256 \times N}$ (images)

- $\underline{X} \underline{X}^T \in \mathbb{R}^{256 \times 256}$ ad demean

$\sum_{R=1}^K \frac{1}{\sqrt{R}} \mathbf{z}_R \mathbf{z}_R^\top$

... or eigenvectors / principal components used to determine degree of compression

Probabilistic PCA: = mean + $\mathbf{Q}\mathbf{z}_{\text{proj}}$ with $(K - (m-1))$ dimensions of \mathbf{Q} and \mathbf{z}_{proj} snatched off

- PCA, SVD \rightarrow probabilistic model on PCA, estimate param via EM (with gen model)
- Allows 1. missing data 2. additional param e.g. noise 3. modularity + point est.
- 4. Distributions \rightarrow characterise uncertainty on precision

Probabilistic PCA (setup): Analogous to matrix factorisation, with gen process and distri exp. defined.

① Approximate $\underline{\mathbf{x}} = \underline{\mathbf{w}}\underline{\mathbf{z}}$ $\mathbf{x} \in \mathbb{R}^{d \times n}$ $\mathbf{w} \in \mathbb{R}^{d \times K}$ - factor loadings (similar to eigenvectors with no orthonorm, unit length const.)
 $\mathbf{z} \in \mathbb{R}^{K \times n}$ - i-th column of $\underline{\mathbf{z}}$, $\underline{\mathbf{z}}_{::i} = \underline{\mathbf{z}}_i \in \mathbb{R}^K$, a low-dim represent. of \mathbf{x}_i
 $\mathbf{z}_i \sim N(0, I)$ (A) better prior elts (Bish), papers

Do not know \mathbf{w}, \mathbf{z}_i - no orthogonality and unit length constraints

PCA (tractability of marginal likelihood): $\underline{\mathbf{w}}_{\text{ML}} = \underset{\mathbf{w}}{\operatorname{argmax}} \ln p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n \ln p(\mathbf{x}_i | \mathbf{w})$ (not consistent with Bishop, lit.)
Intractable as $p(\mathbf{x}_i | \mathbf{w}) = N(\mathbf{x}_i | \mathbf{0}, \sigma^2 I + \mathbf{w}\mathbf{w}^T)$ (log-sum eventually?) and cannot find gradient wrt \mathbf{w} ??

CA (setup): • marginal log-likelihood:-

$$\begin{aligned} \sum_i \ln p(\mathbf{z}_i | \mathbf{w}) &= \sum_{i=1}^n \ln \int p(\mathbf{x}_i, \mathbf{z}_i | \mathbf{w}) d\mathbf{z}_i = \sum_{i=1}^n \int q(\mathbf{z}_i) \ln \frac{p(\mathbf{x}_i, \mathbf{z}_i | \mathbf{w})}{q(\mathbf{z}_i)} d\mathbf{z}_i + \sum_{i=1}^n \int q(\mathbf{z}_i) \ln \frac{q(\mathbf{z}_i)}{p(\mathbf{z}_i | \mathbf{x}_i, \mathbf{w})} d\mathbf{z}_i \\ &+ q(\mathbf{z}_i) = p(\mathbf{z}_i | \mathbf{x}_i, \mathbf{w}) \quad \forall i=1, \dots, N \quad \left. \begin{array}{l} \text{requires calc. of } p(\mathbf{z}_i | \mathbf{x}_i, \mathbf{w}) \text{ and tractable max. of } L \\ \text{maximise } L \text{ wrt } \mathbf{w} \end{array} \right\} \end{aligned}$$

for PPPCA (Algorithm):

$\{\mathbf{z}_i\}_{i=1}^N, \mathbf{z}_i \in \mathbb{R}^d$, model $\mathbf{x}_i \sim N(\mathbf{w}\mathbf{z}_i, \sigma^2 I)$ $\mathbf{z}_i \sim N(0, I)$ $\mathbf{z}_i \in \mathbb{R}^K$

Step: For all $i=1, \dots, N$ set $q(\mathbf{z}_i) = p(\mathbf{z}_i | \mathbf{x}_i, \mathbf{w}) = N(\mathbf{z}_i | \underline{\mu}_i, \Sigma_i)$ where

$$\Sigma_i = ((I + \mathbf{w}\mathbf{w}^T)/\sigma^2)^{-1} \quad \underline{\mu}_i = \Sigma_i \mathbf{w}^T \mathbf{x}_i / \sigma^2$$

Step: Update \mathbf{w} by maximising L from E-step:

$$\mathbf{w} = \left[\sum_{i=1}^N \mathbf{z}_i \underline{\mu}_i^T \right] \left[\sigma^2 I + \sum_{i=1}^N (\underline{\mu}_i \underline{\mu}_i^T + \Sigma_i) \right]^{-1}$$

Point estimate $\hat{\mathbf{w}}$, conditional posterior $p(\mathbf{z}_i | \mathbf{x}_i, \mathbf{w})$ on each \mathbf{z}_i , $\hat{\mathbf{K}}, \hat{\sigma}^2$

Heate until increase in
 $\sum_{i=1}^N \ln p(\mathbf{z}_i | \mathbf{w})$
is small
(marg likel)

]

• Conditional posterior $q(z_i | x_i, w)$ captures belief about what low-dim embedding should be; but also a distribution with uncertainty

PPCA Applications: #1 - Image processing; 8×8 patch in image; vectorise $\rightarrow x_i \in \mathbb{R}^{64}$ overlap

then stack into columns of $X: X \in \mathbb{R}^{256 \times 262,144}$

• Factor model: Approx. $x_i \sim w\mu_i$ where μ_i is posterior mean of z_i ($z_{:,i}$ - i th column of Z)

• Reconstruct image by replacing x_i with $w\mu_i$ (and averaging) $\odot(*)$

• Noise from denoised patches \rightarrow captured in noise variance param σ^2 #2 denoising

$z_i \sim N(wz_i, \sigma^2 I)$

ass 1
denoised
mean

\sim noise variance soaks up noise ; σ^2 estimated

#3 - Missing data

Kernel PCA

• Dot products, generalise with non-linear kernel

• Map to high-dim feature space, then execute dimensionality reduction

• PCA \rightarrow find eigenvalues, eigenvectors of $\underline{X}\underline{X}^T = \sum_{i=1}^N z_i z_i^T$ \rightarrow due to definition of $\underline{X} = \begin{bmatrix} z_1 & \dots & z_N \\ 1 & \dots & 1 \end{bmatrix}$

• Define $\phi(x)$ as feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$

• Solve the eigen decomposition: $-\left[\sum_{i=1}^N \phi(z_i) \phi(z_i)^T \right] q_R = \lambda q_R \quad \forall R = 1, \dots, K$ without using $\phi(\cdot)$ and $q_R \in \mathbb{R}^D$

$$\text{(see derivations)} \Rightarrow q_R = \sum_{i=1}^N a_{R,i} \phi(z_i) \text{ where } a_{R,i} = \frac{\phi(z_i)^T q_R}{\lambda_R} \quad q_R = \begin{bmatrix} a_{R,1} \\ \vdots \\ a_{R,N} \end{bmatrix}$$

KPCA: Perform PCA on kernel matrix K rather than $\underline{X}\underline{X}^T$ -

$K \underline{a}_R = \lambda_R \underline{a}_R$ where K is kernel matrix constructed on data with $K \in \mathbb{R}^{N \times N}, S^N$
(i.e. rewrite earlier eigen decomposition in terms of \underline{a}_R, K)

Kernel PCA (Algorithm):

FIN: $\{x_i\}_{i=1}^N, x_i \in \mathbb{R}^d; K(x_i, x_j)$ - kernel fn

- construct K , kernel matrix on data e.g. $K_{ij} = \exp\left\{-\frac{1}{b} \|x_i - x_j\|_2^2\right\}$

- solve: eigen decomposition

$$K \underline{a}_k = \lambda_k \underline{a}_k \quad \text{for first } r \text{ eigenvalue/vector pairs } (\lambda_1, \underline{a}_1), \dots, (\lambda_r, \underline{a}_r)$$

OUT: new co-ordinate system for x_i by implicitly mapping $\phi(x_i)$ and projecting $(q_R^T \phi(x_i))$

$$x_i \xrightarrow{\text{proj}} \begin{bmatrix} \lambda_1 a_{1,i} \\ \vdots \\ \lambda_r a_{r,i} \end{bmatrix} \text{ with } a_{k,i} - i^{\text{th}} \text{ dimension of } k^{\text{th}} \text{ eigenvector } \underline{a}_k$$