

Mixture Models and SVD

Monday, August 30, 2021 6:57 AM

gdu

Gaussian Mixture Model: $N(\mu_1, \Sigma_1), N(\mu_2, \Sigma_2) \dots$

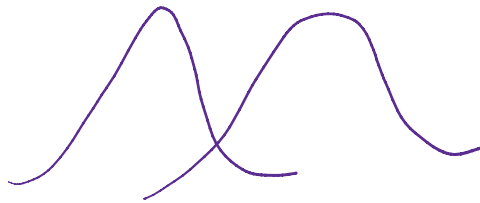
$$\omega_1 \geq 0 \quad \omega_2 \geq 0 \quad \dots \omega_k$$

$$\sum_{i=1}^k \omega_i = 1.$$

Problem 1. Given random samples from an unknown k -GMM estimate its parameters.

$k=1$: $\omega_1 = 1$, $\mu =$ Sample mean $\Sigma =$ sample covariance

$k=2$?



Special case: Separable GMMs.

The components are pairwise separated.
Has to measure separation.

1-dim $|\mu_i - \mu_j| > ?? \max\{\sigma_i, \sigma_j\}$

d-dim (geometric)

(probabilistic) $d_{TV}(F_i, F_j) \geq 1 - \epsilon.$

$$d_{TV}(p, q) = \frac{1}{2} \int |p(x) - q(x)| dx$$

1-dim τ 1-dim d_{TV} "large"

Lemma. In 1-dim, d_{TV} "large"

\Rightarrow either $\|\mu_i - \mu_j\|$ is large

or $\max\left\{\frac{\sigma_i}{\sigma_j}, \frac{\sigma_j}{\sigma_i}\right\}$ is large.

Thm. (concentration)
 $X \sim N(\mu, \sigma^2 I)$ $\exists C > 0$: $\Pr(|X - \mu| > Ct\sigma\sqrt{d}) \leq e^{-\frac{d t^2}{2}}$
 $\Pr(|\|X - \mu\|^2 - d\sigma^2| > t\sigma^2\sqrt{d}) \leq 2e^{-t^2/8}$

How to solve P1?

In the separable case, cluster and estimate each component separately.

P2. Cluster a sample from a K-GMM by component of origin.

How? Suppose mean separated



$$\begin{aligned} X, Y \in F_i \quad \mathbb{E}(\|X - Y\|^2) &= \mathbb{E}(\|X - \mu_i - (Y - \mu_i)\|^2) \\ &= \mathbb{E}(\|X - \mu_i\|^2) + \mathbb{E}(\|Y - \mu_i\|^2) + 0. \end{aligned}$$

$$\begin{aligned}
 x, y \in F_i & \dots \\
 & = E(\|x - \mu_i\|^2) \rightarrow E(\|y - \mu_i\|^2) + 0. \\
 & = 2d\sigma^2.
 \end{aligned}$$

$$\begin{aligned}
 x \in F_i, y \in F_j & \\
 E(\|x - y\|^2) & = E(\|x - \mu_i - (y - \mu_j) + \mu_i - \mu_j\|^2) \\
 & = E(\|x - \mu_i\|^2) + E(\|y - \mu_j\|^2) + \|\mu_i - \mu_j\|^2 \\
 & = 2d\sigma^2 + \|\mu_i - \mu_j\|^2.
 \end{aligned}$$

By conc. with prob $\geq 1 - 2e^{-t^2/8}$

$$\|x - y\|^2 \leq 2d\sigma^2 + 2t\sqrt{d}\sigma^2 \quad x, y \in F_i$$

$$\|x - y\|^2 > 2d\sigma^2 + \|\mu_i - \mu_j\|^2 - 2t\sqrt{d}\sigma^2 \quad \begin{matrix} x \in F_i \\ y \in F_j \end{matrix}$$

\therefore it suffices to have

$$\|\mu_i - \mu_j\|^2 > 4t\sqrt{d}\sigma^2$$

to ensure that pairs from same gaussian are closer.

Cluster using distances

- put nearest pair in same cluster
- Repeat till only k clusters.

Thm. with prob $1 - \delta$, random sample with m points

$$\text{and } \|\mu_i - \mu_j\| > C \left(\log \frac{m}{\delta} \cdot d \right)^{\frac{1}{4}} \max\{\sigma_i, \sigma_j\}$$

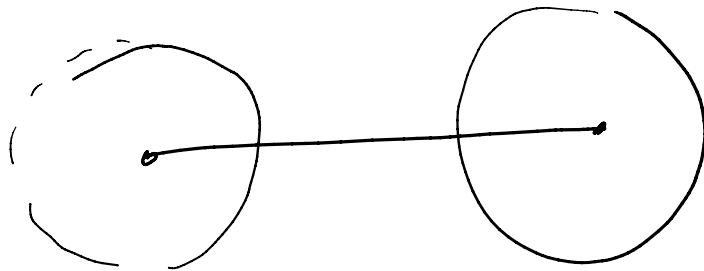
can be clustered using pairwise distances in

Can we cluster using ~~...~~
polynomial time.

Pf. Set $t = \sqrt{C \log \frac{m}{\delta}}$.

Is this the right answer?
Separation grows with $d^{1/4}$.

No!



Project to line joining v_i, v_j .

Separation needed is $O(\sigma)$. Not $d^{1/4}\sigma$.

Q. But how to find line joining means?

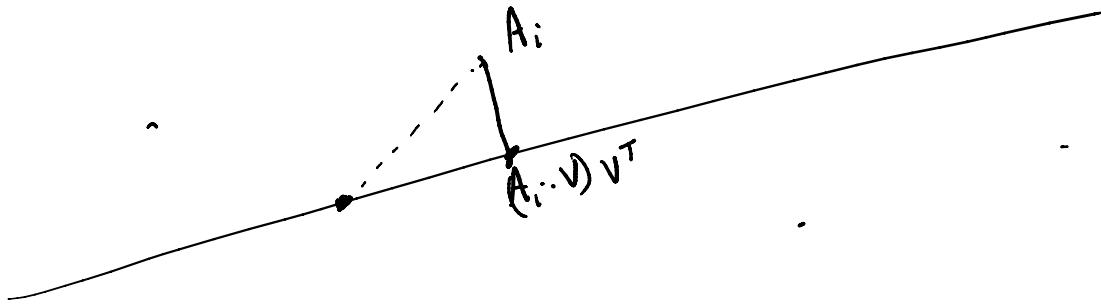
Best fit line: $v = \operatorname{argmax}_{\|v\|=1} \|Av\|^2$

maximizes sum of squared projections of
rows of A .

For each row A_i

$$\|A_i\|^2 = (A_i v)^2 + \|A_i - (A_i v) v^T\|^2$$

$$\|A_i\| = (A_i \cdot v) + \|A_i - (A_i \cdot v)v\|$$



$$\arg \max \|Av\|^2 = \arg \min \|A - (Av)v^T\|^2$$

least squared error.

$$A \in \mathbb{R}^{n \times d}$$

u, v : left, right singular vectors of A

$$Av = \sigma u \quad \sigma \geq 0.$$

$$A^T u = \sigma v$$

Lemma $v = \arg \max_{\|v\|=1} \|Av\|^2$ is a right singular vector of A with largest singular value.

Pf.

$$A^T Av = A^T (\sigma u) = \sigma^2 v$$

v is an eigen vector of $A^T A$.

goes both ways: $A^T Av = \sigma^2 v$

dot \cdot $\|v\|=1$ $Av = \sigma v$

0 \dots v
define $u = \frac{1}{\sigma} Av$

$$\text{Then } A^T u = \frac{1}{\sigma} A^T A v = \sigma v.$$

Now consider $f(v) = \|Av\|^2 = v^T A^T A v$.

$$\nabla_v f(v) = 2A^T A v$$

at any local max/min $\nabla_v f(v) = \lambda v$

$$\Rightarrow A^T A v = \lambda v.$$

Hence the maximizer is an eigenvector.

$$\left. \begin{array}{l} \|Av\|^2 + \lambda (1 - \|v\|^2) \\ 2A^T A v + 2\lambda v = 0 \\ \|v\|^2 = 1 \end{array} \right\}$$

SVD:

$$v_1 = \operatorname{argmax} \|Av_1\| \quad \sigma_1 = \|Av_1\|$$

$$v_2 = \operatorname{argmax} \|Av_2\| \\ v_2 \perp v_1$$

$$\vdots \\ v_k = \operatorname{argmax} \|Av_k\| \quad \sigma_k \\ v_k \perp v_1 \dots v_{k-1}$$

$$\begin{aligned} & \|Aw_1\|^2 + \|Aw_2\|^2 + \dots + \|Aw_k\|^2 \\ & \leq \sum_{i=1}^n d(A_{(i)}, V_{k-1})^2 + \|Aw_k\|^2 \\ & = \|Av_1\|^2 + \dots + \|Av_{k-1}\|^2 + \|Av_k\|^2 \end{aligned}$$

Hence $V_k' = \text{Span } V_{k-1} \cup \{w_k\}$
wlog

But $w_k \perp V_{k-1}$ and must be a maximizer.
∴ $V_k = V_k'$.

$(\sum \sigma_i u_i v_i^T) v_j$ is the same as Av_j
Hence also for any $x (= \sum \alpha_j v_j)$. □

Back to k -GMMs.

Thm (Mean Subspace). V_k for a mixture of spherical Gaussians
 $\supseteq \text{Span } \{v_1, \dots, v_k\}$.

Algorithm. - Project sample to top k -dim SVD
subspace.
- Cluster according to distances in \mathbb{R}^k .

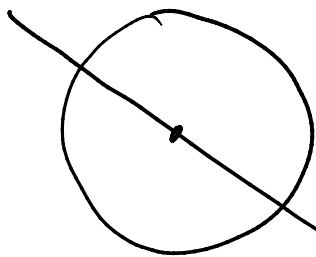
- Unusual notation for v

Thm. (SVD + Cluster) $| \mu_i - \mu_j | > C \left(\log \frac{m}{\delta} \cdot K \right)^{\frac{1}{4}} \max\{\sigma_i, \sigma_j\}$
suffices!

Pf (of Thm. [Mean Subspace]) K instead of d .

Suppose $K=1$. Q. What is the best subspace?

A. line through μ !



$$v = \arg \min \mathbb{E} (\|x - (x \cdot v)v\|^2)$$
$$= \arg \max \mathbb{E} (\|x \cdot v\|^2)$$

$$= \mathbb{E} [\|(x - \mu) \cdot v + \mu \cdot v\|^2]$$
$$= \sigma^2 + (\mu \cdot v)^2 + 0.$$

to maximize set $v = \frac{\mu}{\|\mu\|}$.

For 1 Gaussian

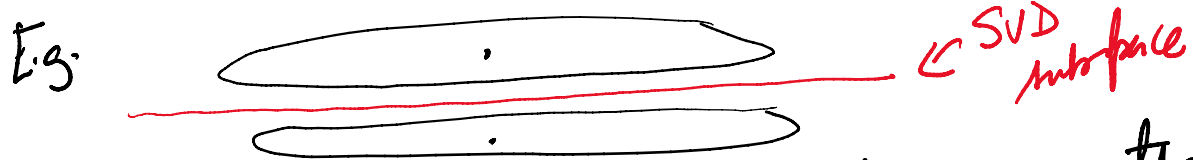
best K -d subspace is any subspace containing μ .

Eigenvalues of $\mathbb{E}(XX^T)$ are $\sigma^2 + \|\mu\|^2$
 σ^2

⋮
 σ^2

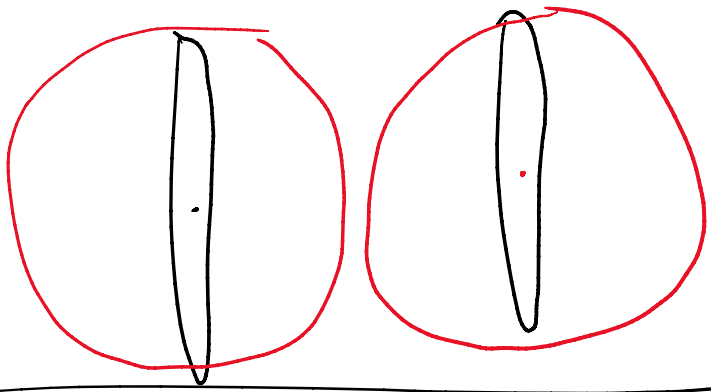
So for K Gaussians, best K -dim subspace is any subspace containing all their means!

Q. Does SVD work for general Gaussians
A. Only if the pairwise separation grows with the largest variance of each component.



Projecting to SVD subspace simply merges the two "pancakes".

Need:



Two questions:

- 1) smaller separation?
- 2) equal Gaussians?

- 1) ~~isotropic~~
- 2) general Gaussians?

Idea: SVD is about second moments.

Q. Would higher moments help? $(E((x \cdot v)^2))$

Mixture of k spherical Gaussians

with $N(\mu_i, \sigma_i^2 I)$.

Assume $\{\mu_1, \dots, \mu_k\}$ are lin. ind.

① Make isotropic i.e. $E_F(x) = 0$

$$E_F(x x^T) = I.$$

Thm [isotropic transformation]

Any distribution with bounded second moments can be made isotropic via an affine transformation.

Pf. Suppose $E(x) = \mu$ $E((x-\mu)(x-\mu)^T) = A$.

Then $Y = x - \mu$ has $E(Y) = 0$.

and $Y = A^{-1/2} (x - \mu)$ is isotropic!

$$\begin{aligned}
 E(Y) &= 0 \\
 E(Y Y^T) &= A^{-1/2} E((x - \mu)(x - \mu)^T) (A^{-1/2})^T \\
 &= A^{-1/2} A (A^{-1/2})^T = I
 \end{aligned}$$

$$= \bar{A}^{-1/2} A (\bar{A}^{-1/2})^T = I$$

What is $\bar{A}^{-1/2}$?

Note $\mathbb{E}((x-\mu)(x-\mu)^T) \succ 0$

Positive definite matrix.

$\therefore A = BB^T$ for some B .

$$\bar{A}^{-1/2} = B^{-1}$$

$$\text{So } \bar{A}^{-1/2} A (\bar{A}^{-1/2})^T = B^{-1} B B^T (B^{-1})^T = I$$

So assume that $F = \sum_i w_i F_i$ is isotropic.

Next consider $\mathbb{E}(x \otimes x \otimes x) = T$

this is a 3-dim array, a tensor of size $d \times d \times d$.

$$T_{ijk} = \mathbb{E}(x_i x_j x_k)$$

$$T = \sum_i w_i \mathbb{E}_{F_i}(x \otimes x \otimes x)$$

So we need

$$\mathbb{E}(x \otimes x \otimes x) \quad x \sim N(\mu, \sigma^2 I)$$

$$= \mathbb{E}((x-\mu+\mu) \otimes (x-\mu+\mu) \otimes (x-\mu+\mu))$$

$$= \mathbb{E}(\otimes^3(x-\mu)) + \mathbb{E}((x-\mu) \otimes (x-\mu) \otimes \mu)$$

$$+ \mathbb{E}((x-\mu) \otimes \mu \otimes (x-\mu))$$

$$+ E((X-\mu) \otimes \mu \otimes (X-\mu))$$

$$+ E(\mu \otimes (X-\mu) \otimes (X-\mu))$$

$$\left. \begin{array}{l} + E((X-\mu) \otimes \mu \otimes \mu) \\ + \\ \vdots \\ + \end{array} \right\} = 0$$

$$+ \mu \otimes \mu \otimes \mu.$$

$$E((X-\mu)_i (X-\mu)_j (X-\mu)_k) = \begin{cases} 0 & \text{if } i, j, k \text{ are not all equal} \\ E((X-\mu)_i^3) & \text{if } i=j=k \\ = 0 & \end{cases}$$

So we are left with

$$E(X \otimes X \otimes X) = \sigma^2 I \otimes \mu + \mu \otimes \sigma^2 I$$

$$+ E((X-\mu) \otimes \mu \otimes (X-\mu))$$

$$+ \mu \otimes \mu \otimes \mu.$$

$$E((X-\mu)_i \mu_j (X-\mu)_k) = \begin{cases} \sigma^2 \mu_j & \text{if } i=k \\ 0 & \text{o.w.} \end{cases}$$

$$= \left(\sigma^2 \sum_{i=1}^d 0 \otimes \mu \otimes e_{ii} \right) + \left(\sigma^2 \mu_j \text{ if } i=k=l \right)$$

$$= \left(\sigma^2 \sum_{l=1}^d e_l \otimes \mu \otimes e_l \right)_{ijk} = \begin{cases} \sigma^2 \mu_j & \text{if } i=k=l \\ 0 & \text{o.w.} \end{cases}$$

Note $I = \sum_{l=1}^d e_l \otimes e_l$.

So,

Lemma (a) $E(X \otimes X \otimes X) = \mu \times \mu \times \mu$

$X \sim N(\mu, \sigma^2 I)$ $+ \sigma^2 \sum_{i=1}^d e_i \times e_i \times \mu + e_i \times \mu \times e_i + \mu \times e_i \times e_i$.

(b) $X \sim \sum w_i N(\mu_i, \sigma_i^2 I)$

$$E(X \otimes X \otimes X) = \sum_i w_i \mu_i \otimes \mu_i \times \mu_i + \sum_i w_i \sigma_i^2 \sum_{j=1}^d e_j \times e_j \times \mu_i + e_j \times \mu_i \times e_j + \mu_i \times e_j \times e_j$$

Can we estimate this ?!

For any tensor $T = (T_{ijk})$ $x, y, z \in \mathbb{R}^d$

$T(x, y, z) = \sum_{i,j,k} T_{ijk} x_i y_j z_k$ ← scalar

$T(\cdot, y, z) = \sum_{j,k} T_{ijk} y_j z_k$ ← vector

matrix

$$T(\cdot, \cdot, \cdot) = \sum_{j,k} T_{ijk} z_j z_k$$

$$T(\cdot, \cdot, z) = \sum_k T_{ijk} z_j \leftarrow \text{matrix.}$$

Consider $E(X \otimes (X-\mu) \otimes (X-\mu)) [\cdot, v, v]$

$$= E(X \otimes X \otimes X) [\cdot, v, v]$$

$$+ E(X \otimes -\mu \otimes (X-\mu)) [\cdot, v, v]$$

We choose $v \perp \mu$
 $\|v\|=1$ $+ E(X \otimes (X-\mu) \otimes -\mu) [\cdot, v, v]$

$$= E(X \otimes X \otimes X) [\cdot, v, v]$$

use Lemma (b):

$$= 0 + \underbrace{\sum_i w_i \sigma_i^2 \mu_i}_{\text{call this vector } u.}$$

Then $E(X \otimes X \otimes X) = \sum w_i \mu_i \otimes \mu_i \otimes \mu_i$

$$+ \sum_{j=1}^d u \otimes e_j \otimes e_j + e_j \otimes u \otimes e_j + e_j \otimes e_j \otimes u$$

We know this!

and this!

So, we have $T = \sum w_i \mu_i \otimes \mu_i \otimes \mu_i$

So, we have $T = \sum w_i \psi_i \otimes \psi_i \otimes \psi_i$

and $E(x x^T) = \sum_i w_i \psi_i \otimes \psi_i + \sum_i w_i \sigma_i^2 I$

$E((x-r) \cdot v)^2 = \sum w_i \sigma_i^2 = \hat{\sigma}^2$

\therefore we also have $\sum w_i \psi_i \times \psi_i = I$.
and by linear transformation we have $= I$.

Claim. $\sum_{i=1}^K w_i \psi_i \times \psi_i = I \Rightarrow \psi_i$ are orthogonal
 $\sqrt{w_i} \psi_i$ are orthonormal.
 $\|\psi_i\|^2 = \frac{1}{w_i}$

$\sum a_i a_i^T = I$

$AA^T = I \Rightarrow \|a_i\|^2 = 1 \quad a_i^T a_j = 0 \quad i \neq j$.

ψ_i are orthogonal and we know $\sum w_i \psi_i \otimes \psi_i \otimes \psi_i$

Is this enough?

[Tensor Decomposition]

Thm. Given $T = \sum_i \alpha_i u_i \otimes u_i \otimes u_i$ $\{u_i\}$ orthogonal

there is a polytime algorithm to recover α_i, u_i

Pf. Consider the iteration

...

Pf. Consider the iteration

$$x = \frac{T(\cdot, x, x)}{\|T(\cdot, x, x)\|}$$

starting with x_0
random.

assume $\|u_i\| = 1$

$$x = \sum \beta_i u_i$$

$$x^{(1)} \propto T(\cdot, x, x) = \sum_i \alpha_i \beta_i^2 u_i$$

$$x^{(2)} \propto T(\cdot, x^{(1)}, x^{(1)}) = \sum_i \alpha_i \beta_i^4 u_i$$

$$= \sum_i \alpha_i \beta_i^8 u_i$$

$$= \sum_i (\alpha_i \beta_i)^{2^{k-1}} \beta_i u_i$$

So i with largest $\alpha_i \beta_i$ will quickly dominate!

$$x^{(k)} \rightarrow u_i$$

Peel off and repeat.

(*) need to be careful about error accumulation.

So now we have an algorithm!

$$F = \sum w_i N(\mu_i, \sigma_i^2 I) \quad \{\mu_1, \dots, \mu_k\} \text{ lin.}$$

$$F = \sum_i w_i N(\mu_i, \sigma_i^2 I) \quad \{\mu_1, \dots, \mu_k\} \text{ lin. ind.}$$

① $M = \mathbb{E}_S(X \otimes X)$ find top k eigenvectors.
 $\hat{\sigma}^2 = (k+1)^{\text{th}}$ eigenvalue v_1, v_2, \dots, v_k .

② $(M - \hat{\sigma}^2 I) = WW^T$
 compute $\hat{S} = W^{-1}S \leftarrow$ (sample)

③ $v \perp \{W^{-1}v_1, \dots, W^{-1}v_k\}$
 $u = \mathbb{E}_S \left(X \frac{(X - \mu) \cdot v}{\hat{S}} \right)^2$

$$T = \mathbb{E}_S (X \otimes X \otimes X) - \left(\sum_j u \otimes e_j \otimes e_j + e_j \otimes u \otimes e_j + e_j \otimes e_j \otimes u \right)$$

④ Decompose T using tensor decomposition
 for vector y set $\hat{\mu}_i = T(y, y, y) y$

and $w_i = \frac{1}{\|\hat{\mu}_i\|^2}$

and finally $\sigma_i^2 = u \cdot \hat{\mu}_i = w_i \sigma_i^2 \|\hat{\mu}_i\|^2 = \sigma_i^2$.

Note: complexity depends on the condition number

Note: complexity depends on the condition number
of $\begin{pmatrix} -M_i & - \\ -M_k & - \end{pmatrix}$.
