

## Chapter 1

# Neural Network and Learning Algorithm



## Chapter 2

# Independent Component Analysis(ICA)

ICA tries to find original signals  $\hat{S}$ . from mixed signal  $X$ , called Blind Source Separation problem.

$$X = AS \quad \rightarrow \quad \hat{S} = WX$$

where

- $X \in \mathbb{R}^{N \times p}$
- $A, W \in \mathbb{R}^{N \times N}$
- $S, \hat{S} \in \mathbb{R}^{N \times p}$
- $N$  is no. sources
- $p$  is no. observations in the source

One could say that ICA tries to find  $W \approx A^{-1}$ .

## 2.1 ICA vs PCA

The goal of PCA is to find directions that don't correlate with each other( orthogonal projections ). However, ICA finds sources that are independent to each other. The right most figure of Fig 2.1 shows that even correlation  $C_{12}, C_{21}$  is zero, but  $x_1, x_2$  are dependent to each other( knowing one of them can tell position of the other ).

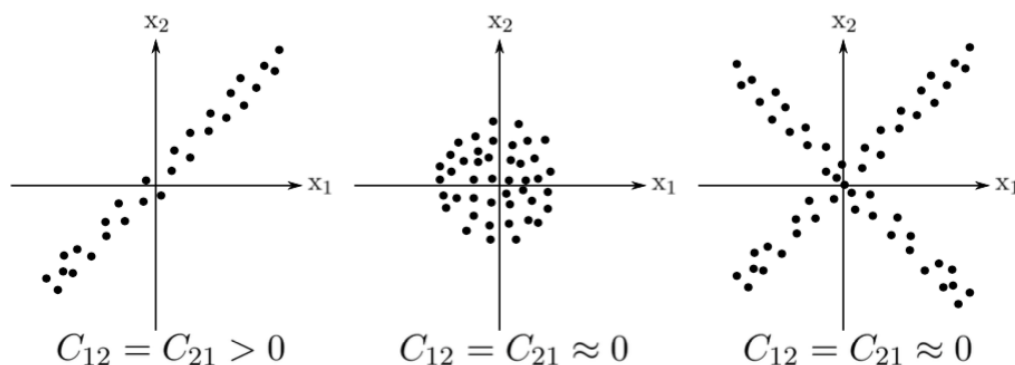


Figure 2.1: Decorrelation vs Independent (Drawn from MI2 Slides)

## 2.2 ICA Limitations

### 2.2.1 Permutation of Sources

$\hat{s}_i$  that found that ICA is not necessary the same as the original source  $s_i$ .

## 2.2.2 Discovered Amplitude

## 2.2.3 Gaussian Sources

The underly theory behind ICA is Central Limit Theorem in which it implies that sum of random variables will approach Gaussian distribution. Thus, if the original sources come from Gaussian distribution, we won't have any clue to separate the sources.

## 2.3 ICA Approaches

An ICA algorithm can be derived using different cost functions using assumption that

$$\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}}) = \prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)$$

and  $P_{\mathbf{s}}(\hat{\mathbf{s}}) = P_{\mathbf{s}}(W\mathbf{x})$

### 2.3.1 Informax

From Kullback-Leibler divergence, we know that

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Thus, we can formulate the problem as

$$\min_W D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}}))$$

However, deriving cost from  $D_{KL}$  directly would not work as it require probability density estimation. Hence, we shall use a transformation :

$$\hat{u}_i = \hat{f}(\hat{s}_i) \quad \text{s.t.} \quad \hat{P}_{u_i}(\hat{u}_i) = \text{const.}$$

From Conservation of Probability, we know that

$$\hat{P}_{u_i}(\hat{u}_i) d\hat{u}_i = \hat{P}_{s_i}(\hat{s}_i) d\hat{s}_i$$

Then

$$\begin{aligned} \hat{P}_{u_i} &= \left| \frac{d\hat{s}_i}{d\hat{u}_i} \right| \hat{P}_{s_i}(\hat{s}_i) \\ &= \left| \frac{1}{\frac{d\hat{u}_i}{d\hat{s}_i}} \right| \hat{P}_{s_i}(\hat{s}_i) \\ &= \left| \frac{1}{\hat{f}'(\hat{s}_i)} \right| \hat{P}_{s_i}(\hat{s}_i) \\ &\stackrel{!}{=} 1 \end{aligned}$$

This implies that  $\hat{f}(\hat{s}_i)$  is CDF of  $\hat{P}_{s_i}(\hat{s}_i)$ . From this knowledge, we can derive  $D_{KL}$  :

$$\begin{aligned} D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}})) &= \int_{-\infty}^{\infty} d\hat{\mathbf{s}} P_{\mathbf{s}}(\hat{\mathbf{s}}) \log \frac{P_{\mathbf{s}}(\hat{\mathbf{s}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{s}} P_{\mathbf{s}}(\hat{\mathbf{s}}) \log \frac{P_{\mathbf{s}}(\hat{\mathbf{s}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \frac{\prod_{i=1}^N 1/\hat{f}'}{\prod_{i=1}^N 1/\hat{f}'} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log \frac{P_{\mathbf{u}}(\hat{\mathbf{u}})}{\prod_{i=1}^N \hat{P}_{u_i}(\hat{u}_i)} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) - \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log \prod_{i=1}^N \hat{P}_{u_i}(\hat{u}_i) \end{aligned}$$

Hence,

$$\begin{aligned}
\min D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}})) &\approx \min \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) \\
&\approx \min -H(u) \\
&\approx \max H(u)
\end{aligned}
\tag{Informax Principle}$$

### Learning via Neural Network and Empirical Risk Minimization

From Conservation of Probability, we know that

$$\begin{aligned}
P(\hat{\mathbf{u}})d\hat{\mathbf{u}} &= P(\hat{\mathbf{s}})d\hat{\mathbf{s}} \\
&= P(\mathbf{x})d\mathbf{x}
\end{aligned}$$

and,

$$\begin{aligned}
\hat{\mathbf{u}} &= \hat{f}(\hat{\mathbf{s}}) \\
&= \hat{f}(W\mathbf{x})
\end{aligned}$$

Thus,

$$\begin{aligned}
P(\hat{\mathbf{u}}) &= P(\mathbf{x})d\mathbf{x} \\
&= \left| \frac{d\mathbf{x}}{d\hat{\mathbf{u}}} \right| P(\mathbf{x}) \\
&= \frac{1}{|M|} P(\mathbf{x})
\end{aligned}$$

where  $|M|$  is functional determinant and

$$\begin{aligned}
|M| &= \left| \frac{\partial \hat{\mathbf{u}}}{\partial \mathbf{x}} \right| \\
&= |W| \prod_{l=1}^N \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right)
\end{aligned}$$

Hence,

$$\begin{aligned}
H(\mathbf{u}) &= - \int d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) \\
&= - \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log \frac{P_{\mathbf{x}}(\hat{\mathbf{x}})}{|M|} \\
&= - \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log P_{\mathbf{x}} + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |M|
\end{aligned}$$

Because  $-\int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log P_{\mathbf{x}}$  is a constant if we optimize over  $w$ , therefore, the cost function  $E^G$  is reduced to

$$\begin{aligned}
E^G &= \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |M| \\
&= \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |W| + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log \prod_{l=1}^N \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right) \\
&= \log |W| + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \left\{ \sum_{l=1}^N \log \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right) \right\}
\end{aligned}$$

Using Empirical Risk Minimization, we convert expectation to average thus

$$\begin{aligned}
E^T &= \log |W| + \frac{1}{p} \sum_{\alpha=1}^p \sum_{l=1}^N \log \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right) \\
&\stackrel{!}{=} \max_W
\end{aligned}$$

We can learn  $W$  using gradient ascent :

$$W \leftarrow W + \eta \nabla_W E^T$$

For individual cost of each observation  $e^{(\alpha)}$ ,

$$\frac{\partial e^{(\alpha)}}{\partial w_{ij}} = W_{ji}^{-1} + \frac{\hat{f}''\left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)}\right)}{\hat{f}'\left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)}\right)} x_j^{(\alpha)}$$

Or

$$\frac{\partial}{\partial W} e^{(\alpha)} = (W^{-1})^T + \psi(Wx^{(\alpha)})(x^{(\alpha)})^T$$

One could observe that there is  $W^{-1}$  in the computation which causes performance.

### Learning via Natural Gradient

Due to the fact that  $W$  space is not orthogonal coordinate, hence gradient doest point to the direction of steepest ascent. Thus, we need to transform the space back to comparable space before computing gradient. This yields

$$dZ = dW \cdot W^{-1}$$

Using Taylor expansion of  $e^{(\alpha)}$  or  $e$  for short, around  $W$ , we have :

$$\begin{aligned} e(W + dW) &= e(W) + \nabla_W e^T dW \\ &= e(W) + \eta \nabla_W e^T z_w \end{aligned}$$

where  $z_w = \eta dw$

To find the direction of steepest ascent, we have to

$$\nabla_W e^T z_w \stackrel{!}{=} \max \quad \text{s.t.} \quad (z_w \cdot W^{-1})^2 \stackrel{!}{=} 1$$

Using Lagrange multiplier we have

$$\mathcal{L}(z_w, \lambda) = \nabla_W e^T z_w - \lambda((z_w \cdot W^{-1})^2 - 1)$$

Taking derivative wrt to  $z_w$  and set to zero, we have

$$\begin{aligned} 0 &= \nabla_W e - 2\lambda(z_w \cdot W^{-1})(W^{-1})^T \\ \nabla_W e &= 2\lambda(z_w \cdot W^{-1})(W^{-1})^T \\ z_w &= \frac{1}{2\lambda} \nabla_W e W^T W \end{aligned}$$

Hence the direction for **Natural Gradient** is

$$\Delta W = \eta \nabla_W e W^T W$$

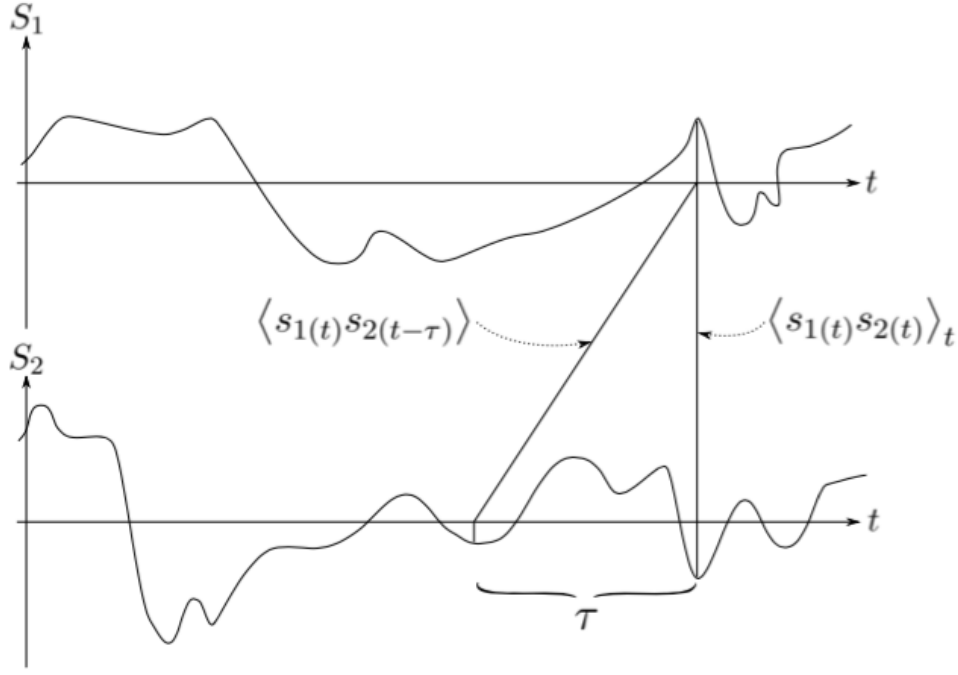
where  $\nabla_W e$  is original gradient.

Putting things together, we have

$$\Delta W = \eta \left( \mathbb{I} + \psi(Wx)(Wx)^T \right) W$$

Commonly, people use the sigmoid function as CDF :

$$f(x) = \frac{1}{1 + \exp(-x)} \quad \rightarrow \quad \psi(x) = \frac{f''(x)}{f'(x)} = 1 - 2f(x)$$

Figure 2.2: Autocorrelation with time  $\tau$  shifted (Drawn from MI2 Slides)

### 2.3.2 Second-order Approach

Due to the fact that, higher moment measurements are sensitive to noise, especially when we have small set of samples, hence it might yield unreliable result. Moreover, signals, such as EEG, fMRI data, might correlate to each other in temporal pattern. Because most of the time we know the length of the signal, we can exploit this fact by using autocorrelation.

Formally, the motivation of 2nd order approach is to find unmixing matrix  $W$  that cross-correlation vanishes.

#### Source Separation with 2 shifts

1. Whitening data

$$\mathbf{x}_c = \mathbf{x} - \frac{1}{p} \sum_{\alpha} \mathbf{x}^{(\alpha)}$$

Then compute eigen value decomposition of  $C_x^{(0)}$  where  $(C_x)_{ij} = \langle x_i, x_j \rangle$ . Hence,

$$\begin{aligned} \mathbf{u} &= \Lambda_0^{-1/2} \cdot E_0 \mathbf{x} \\ &= M_0 \mathbf{x} \end{aligned}$$

where  $\Lambda_0^{-1/2}$  is a diagonal of inverse eigen values and  $E_0$  is matrix of eigen vectors.

2. Orthogonal Transformation

Assume :  $\mathbf{s} = B\mathbf{u}$ . We know that

$$\begin{aligned} \langle s_i s_j \rangle &= \delta_{ij} \\ &= \sum_{k,l=1}^N B_{ik} \langle u_k u_l \rangle B_{l,j}^T \\ &= \sum_{k,l=1}^N B_{ik} B_{l,j}^T \quad (\langle u_k u_l \rangle = 1 \text{ from Whitening}) \end{aligned}$$

Hence,

$$BB^T = B^T B = \mathbb{I} \rightsquigarrow \text{Orthogonal Transformation}$$

3. Find a candidate of  $B$  using time-shifted cross-correlation matrix. Apply eigen value decomposition on  $C_u^{(\tau)}$  where

$$(C_u^{(\tau)})_{ij} = \langle u_i^{(t)} u_j^{(t-\tau)} \rangle$$

This results in a matrix of eigen vectors  $E_\tau$  whose  $E_\tau^T E_\tau = \mathbb{I}$ . Putting things together, we have

$$\hat{\mathbf{s}} = E_\tau \Lambda_0^{-1} E_0 \mathbf{x}$$

### Noise robust algorithms in general case

Given a set of  $T$  zero mean observations :  $\mathbf{x}^{(t)}$ . We compute joint cross-correlation metric  $C_x^{(\tau)}$  as :

$$(C_x^{(\tau)})_{ij} = \frac{1}{T} \sum_{t=0}^{T-1} x_i^{(t)} x_j^{(t-\tau)}$$

To find  $N \times N$  matrix  $W$  that diagonalizes  $C^{(\tau)}$ ,  $\tau = 0, 1, \dots$ , we use sum of off-diagonal entries of  $W C^{(\tau)} W^T$  as a cost function.

1. QDIAG Algorithm

$$E_W^T = \sum_{\tau} \alpha_{\tau} \sum_{i \neq j} \left( W C^{(\tau)} W^T \right)_{ij}^2$$

where  $\alpha_{\tau}$  is a weighting factor and the optimization problem is  $E^T \stackrel{!}{=} \min$  s.t.

$$\forall i, \quad \left( W C^{(0)} W^T \right)_{ii} = 1$$

2. FFDIAG Algorithm

$$E_W^T = \sum_{\tau} \sum_{i \neq j} \left( W C^{(\tau)} W^T \right)_{ij}^2$$

with a constraint that  $W$  is a non-singular matrix.

### 2.3.3 Maximize Non-Gaussianity

According to Central Limit Theorem, the sum of random variables is more Gaussian than the original sources. Hence, we can leverage this knowledge by finding a unmixing matrix that minimizes Gaussianity of the unmixed sources yielding interesting directions (projection pursuit). There are 2 ways to measure Gaussianity, namely Kurtosis (4th moment) and Negentropy.

#### Kurtosis

$$\text{kurt}(x) = \langle x^4 \rangle - 3 \underbrace{(\langle x^2 \rangle)^2}_{\substack{1 \text{ for} \\ \text{sphered data}}}$$

- $\text{kurt}(x) = 0$  for Gaussian distribution
- $\text{kurt}(x) > 0$  for **Super**-Gaussian : long-tail, e.g. Laplace Distribution
- $\text{kurt}(x) < 0$  for **Sub**-Gaussian : constant distribution.

For independent variables  $x$  and  $y$  we have :

$$\begin{aligned} \text{kurt}(x + y) &= \text{kurt}(x) + \text{kurt}(y) \\ \text{kurt}(ax) &= a^4 \text{kurt}(x) \end{aligned}$$

Because,

$$\hat{\mathbf{s}} = W \mathbf{x} = W A \mathbf{s} = b^T \tilde{A} \mathbf{s} = b^T \mathbf{u}$$

Hence, the optimization problem can be written as



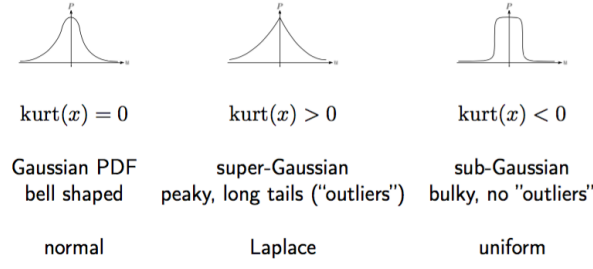


Figure 2.3: Distribution and Kurtosis (Drawn from MI2 Slides)

$$\max_{\mathbf{b}} |\text{kurt}(\mathbf{b}^T \mathbf{u})| \quad \text{s.t.} \quad |\tilde{A}^T \mathbf{b}| = |\mathbf{b}| = 1$$

We have :

$$\begin{aligned} \frac{\partial}{\partial \mathbf{b}} |\text{kurt}(\mathbf{b}^T \mathbf{u})| &= 4 \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \left( \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \langle (\mathbf{b}^T \mathbf{u})^2 \rangle \right) \\ &\stackrel{!}{=} 0 \quad \text{s.t.} \quad |\mathbf{b}| = 1 \end{aligned}$$

Due to the fact that

$$\begin{aligned} 3\mathbf{b} \langle (\mathbf{b}^T \mathbf{u})^2 \rangle &= 3\mathbf{b} \langle ((\mathbf{b}^T \mathbf{u})(\mathbf{b}^T \mathbf{u})^T) \rangle \\ &= 3\mathbf{b} \langle ((\mathbf{b}^T \mathbf{u})(\mathbf{u}^T \mathbf{b})) \rangle \\ &= 3\mathbf{b} \mathbf{b}^T \langle \mathbf{u} \mathbf{u}^T \rangle \mathbf{b} \\ &= 3\mathbf{b} \mathbf{b}^T \mathbb{I} \mathbf{b} \\ &= 3\mathbf{b} |\mathbf{b}|^2 \\ &= 3\mathbf{b} \end{aligned} \quad (|\mathbf{b}|^2 = 1)$$

which changes only direction of  $\mathbf{b}$ , we can drop it and hence the problem is reduced to

$$\begin{aligned} 4 \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle &\stackrel{!}{=} 0 \\ \epsilon \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle &\approx 0 \end{aligned}$$

Therefore for **batch** learning, we have

$$\begin{aligned} \Delta \mathbf{b} &= \epsilon \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

where  $\text{kurt}$  and  $\langle \cdot \rangle$  are replaced with corresponding empirical method and  $\mathbf{b}_0$  is a arbitrary vector with unit length.

For **online** approach, we use running average of kurt  $\gamma$  where  $\gamma_0 = 0$ .

$$\begin{aligned} \Delta \mathbf{b} &= \epsilon \text{sgn}[\gamma] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle \\ \Delta \gamma &= \eta [(\mathbf{b}^T \mathbf{u})^4 - 3 - \gamma] \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

At the equilibrium point of gradient ascent, one could observe that

$$\begin{aligned} \mathbf{b} &\propto \Delta \mathbf{b} \\ &\rightsquigarrow \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \end{aligned}$$

This yields fixed-point algorithm ( **Kurtosis-based fastICA** )

$$\begin{aligned} \mathbf{b} &\leftarrow \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

For whitened data  $\mathbf{u}^{(\alpha)}$ ,  $\alpha = 1, \dots, p$ , we have

$$\begin{aligned}\mathbf{b} &\leftarrow \frac{1}{p} \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} (\mathbf{b}^T \mathbf{u}^{(\alpha)})^3 - 3\mathbf{b} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}|\end{aligned}$$

One remark is that **kurtosis** tends to be very sensitive to outliers.

### 2.3.4 Negentropy

$$J(\hat{\mathbf{s}}) = H_{(\hat{\mathbf{s}})}^{\text{Gauss}} - H(\hat{\mathbf{s}})$$

Properties of  $J(\hat{\mathbf{s}})$ , non-negative and scale-invariant :  $J(\alpha \hat{\mathbf{s}}) = J(\hat{\mathbf{s}})$ ,  $\forall \alpha \neq 0$

Computing  $J(\hat{\mathbf{s}})$  directly requires estimation of density  $p(\hat{\mathbf{s}})$ . Hence, we shall use **Nonpolynomial moment** contrast function  $G$ :

$$J(\hat{\mathbf{s}}) \approx (\langle G(\hat{\mathbf{s}}) \rangle_{\mathbb{P}(\mathbf{s})} - \langle G(\hat{\mathbf{s}}) \rangle_{\text{Gauss}})^2$$

#### Choices of Contrast Function

##### General Purpose

$$G_1(\hat{\mathbf{s}}) = \frac{1}{a} \log \cosh(a\hat{\mathbf{s}}) \quad G_1'(\hat{\mathbf{s}}) = \tanh(a\hat{\mathbf{s}}) \quad G_1''(\hat{\mathbf{s}}) = a(1 - \tanh^2(a\hat{\mathbf{s}}))$$

##### Good for super Gaussian sources with many outliers

$$G_2(\hat{\mathbf{s}}) = -\exp\left(-\frac{\hat{\mathbf{s}}^2}{2}\right) \quad G_2'(\hat{\mathbf{s}}) = \hat{\mathbf{s}} \exp\left(-\frac{\hat{\mathbf{s}}^2}{2}\right) \quad G_2''(\hat{\mathbf{s}}) = (1 - \hat{\mathbf{s}}^2) \exp\left(-\frac{\hat{\mathbf{s}}^2}{2}\right)$$

##### Good for sub-Gaussian sources with few outliers

$$G_3(\hat{\mathbf{s}}) = \frac{1}{4} \hat{\mathbf{s}}^4 \quad G_3'(\hat{\mathbf{s}}) = \hat{\mathbf{s}}^3 \quad G_3''(\hat{\mathbf{s}}) = 3\hat{\mathbf{s}}^2$$

The optimization problem can be written as :

$$J(\mathbf{b}^T \mathbf{u}) = \left( \langle G(\mathbf{b}^T \mathbf{u}) \rangle_{\mathbb{P}(\mathbf{u})} - \langle G(\hat{\mathbf{u}}_{\text{Gauss}}) \rangle_{\mathcal{N}(0,1)} \right)^2$$

#### Batch Learning

$$\begin{aligned}\Delta \mathbf{b} &= \epsilon \left\{ \langle G(\mathbf{b}^T \mathbf{u}) \rangle_{\mathbb{P}(\mathbf{u})} - \langle G(\hat{\mathbf{u}}_{\text{Gauss}}) \rangle_{\mathcal{N}(0,1)} \right\} \langle \mathbf{u} G'(\mathbf{b}^T \mathbf{u}) \rangle_{\mathbb{P}(\mathbf{u})} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}|\end{aligned}$$

where  $\mathbf{b}_0$  is a random unit vector.

#### Online Learning

$$\begin{aligned}\Delta \mathbf{b} &= \epsilon \gamma \langle \mathbf{u} G'(\mathbf{b}^T \mathbf{u}) \rangle_{\mathbb{P}(\mathbf{u})} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}| \\ \Delta \gamma &= \eta (\langle G(\mathbf{b}^T \mathbf{u}) \rangle_{\mathbb{P}(\mathbf{u})} - \langle G(\hat{\mathbf{u}}_{\text{Gauss}}) \rangle_{\mathcal{N}(0,1)} - \gamma)\end{aligned}$$

where  $\gamma_0 = 0$ .

**Fixed point algorithm (fastICA)**

$$\begin{aligned}\mathbf{b} &= \langle \mathbf{u} G'(\mathbf{b}^T \mathbf{u}) \rangle - \langle G''(\mathbf{b}^T \mathbf{u}) \rangle \mathbf{b} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}|\end{aligned}$$

**One Unit Neural Network Implementation**

1. whiten data  $\mathbf{u}^{(\alpha)} = \Lambda_0^{-1/2} \mathbf{E}^T x_{\text{centered}}^{(\alpha)}$
2. choose a random unit vector  $\mathbf{b}$
3. Loop

$$\begin{aligned}\mathbf{b} &\leftarrow \frac{1}{p} \left\{ \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} G'(\mathbf{b}^T \mathbf{u}^{(\alpha)}) - \mathbf{b} \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} G''(\mathbf{b}^T \mathbf{u}^{(\alpha)}) \right\} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}|\end{aligned}$$

and  $\mathbf{w}$  for original centered data is  $\mathbf{w} = \mathbf{E} \Lambda^{-1/2} \mathbf{b}$ .

**Multiple Units Neural Network Implementation**

1. whiten data  $\mathbf{u}^{(\alpha)} = \Lambda_0^{-1/2} \mathbf{E}^T x_{\text{centered}}^{(\alpha)}$
2. choose a random orthogonal matrix  $B = (\mathbf{b}_1, \dots, \mathbf{b}_N)$ .
3. Outer Loop
4. Inner Loop for all  $i$

$$\mathbf{b}_i \leftarrow \frac{1}{p} \left\{ \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} G'(\mathbf{b}_i^T \mathbf{u}^{(\alpha)}) - \mathbf{b}_i \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} G''(\mathbf{b}_i^T \mathbf{u}^{(\alpha)}) \right\}$$

5. End inner loop and do orthogonalization

$$\begin{aligned}B &\leftarrow B / \max_i |\mathbf{b}_i| \\ B &\leftarrow (BB^T)^{-\frac{1}{2}} B\end{aligned}$$

and  $W$  for original centered data is  $W = \mathbf{E} \Lambda^{-1/2} B$ .



# Chapter 3

## Stochastic Optimization

### 3.1 Simulated Annealing

The goal of Simulated Annealing(SA) is to find a state  $\mathbf{s}$  that minimize a cost function  $E_{T(\mathbf{s})}$  where  $s_i \in \mathbf{s}$  is a discrete variable.

$$\mathbf{s}^* = \underset{\mathbf{s}}{\operatorname{argmin}} E_{T(\mathbf{s})}$$

```
Data:  $\mathbf{s}_0, \tau, M, \beta_0$ 
for annealing loop,  $t = 1, 2, \dots$  do
     $\mathbf{s}_t = \mathbf{s}_{t-1}$ 
    for Miterations do
         $\mathbf{s}' = \text{do local permutation from } \mathbf{s}_t$  ;
         $\Delta E = E_{\mathbf{s}'} - E_{\mathbf{s}_t}$  ;
        switch to  $\mathbf{s}_t$  to  $\mathbf{s}'$  with probability  $W = \frac{1}{1 + \exp(\beta_t \Delta E)}$ 
    end
     $\beta_{t+1} = \tau \beta_t \quad \tau \in [1.01 - \dots]$ 
end
```

**Algorithm 1:** Simulated Annealing

where  $\beta = 1/\tau$  and  $M \approx 2000$ . Ideally, we want  $\beta_t = \ln t$  but this is too slow in practice.

Influence of  $T$  and  $\beta$  to transition probability  $W$ .

- $\downarrow T \uparrow \beta$  :  $W$  changes **dramatically** when  $\Delta E$  has a slight change.
- $\uparrow T \downarrow \beta$  :  $W$  **barely** changes when  $\Delta E$  has a huge change.

### 3.2 The Gibbs Distribution

If we fix  $T$  or  $\beta$  at a specific value, this results in a stochastic process with Markov property<sup>1</sup>. Hence  $P_{(\mathbf{s})}$  converges to a stationary distribution as  $t \rightarrow \infty$ .

$$\lim_{t' \rightarrow \infty} \prod (\mathbf{s}, t') = P_{(\mathbf{s})}$$

With **Detailed Balance** assumption which states as follow :

$$P_{(\mathbf{s})} W_{\mathbf{s} \rightarrow \mathbf{s}'} = P_{(\mathbf{s}')} W_{\mathbf{s}' \rightarrow \mathbf{s}}$$

We can find  $P_{(\mathbf{s})}$  analytically :

---

<sup>1</sup>what is this property?

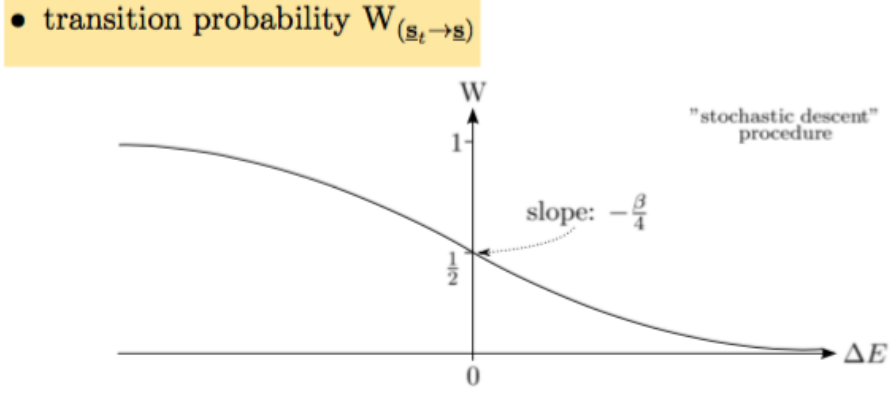


Figure 19: Transition probabilities

- limiting cases

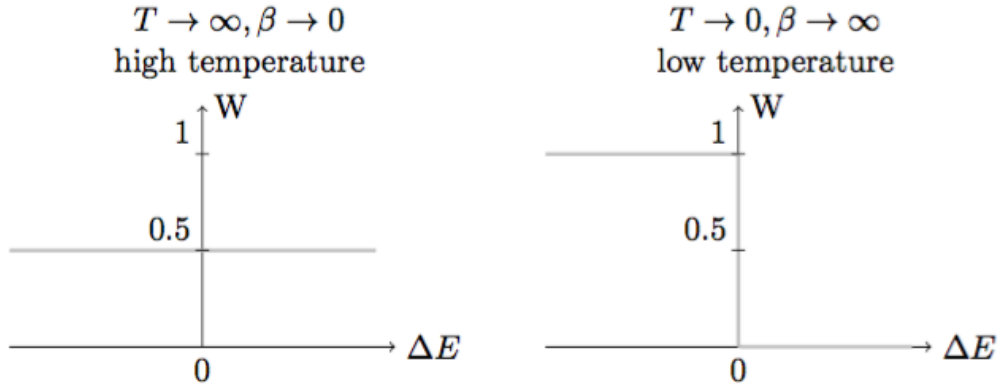


Figure 3.1: Transition Probability

$$\begin{aligned}
 \frac{P_{(s)}}{P_{(s')}} &= \frac{W_{s' \rightarrow s}}{W_{s \rightarrow s'}} \\
 &= \frac{1 + \exp(\beta(E_{s'} - E_s))}{1 + \exp(\beta(E_s - E_{s'}))} \\
 &= \frac{1 + \exp(\beta\Delta E)}{1 + \exp(-\beta\Delta E)} \\
 &= \exp(\beta\Delta E)
 \end{aligned}$$

To fulfill this condition, we need  $P_{(s)}$  to be a Gibbs-Boltzmann Distribution

$$P_{(s)} = \frac{1}{Z} \exp(-\beta E_s)$$

where  $Z = \sum_s \exp(-\beta E_s)$ .

### 3.3 Mean-field Annealing

As we know that  $P_{(s)}$  has an analytic solution, we can exploit this knowledge to make the optimization faster, but estimate  $P_{(s)}$  is difficult and its moment doesn't have analytic solution.

To alleviate the problem, we will choose a distribution  $Q_{(s)}$  that approximates  $P_{(s)}$ . Formally,  $Q_{(s)}$  is a factorizing distribution that  $E_Q$  is a linear combination of  $s_k$ .

$$\begin{aligned}
Q_{\mathbf{s}} &= \frac{1}{Z_Q} \exp \left( -\beta E_Q \right) \\
&= \frac{1}{Z_Q} \exp \left( -\beta \sum_k e_k s_k \right)
\end{aligned}$$

$e_k$  is the mean field that parameterizes family of distributions. As  $\beta \rightarrow \infty$ , the first moment  $\langle \mathbf{s} \rangle_Q$  well characterizes the distributions. For example, it will peak at the location of optimum value.

For factorized distribution,

$$Q(\mathbf{s}) = \prod_k Q(s_k)$$

and the moment of 2 independent variables  $f(\mathbf{s}/s_l)$  and  $g(s_l)$  of  $\mathbf{s}$  can be computed as follows :

$$\begin{aligned}
\langle f(\mathbf{s}/s_l) g(s_l) \rangle &= \frac{1}{Z_Q} \sum_{\mathbf{s}} f(\mathbf{s}/s_l) g(s_l) \exp \left( -\beta \sum_{s_k \in \mathbf{s}} e_k s_k \right) \\
&= \frac{1}{Z_Q} \left[ f(\mathbf{s}/s_l) g(s_{l,i}) \exp \left( -\beta \sum_{s_k \in \mathbf{s}/s_l} e_k s_k \right) \exp \left( -\beta e_l s_{l,i} \right) + \dots \right] \\
&= \frac{1}{Z_Q} \left[ f(\mathbf{s}/s_l) \exp \left( -\beta \sum_{s_k \in \mathbf{s}/s_l} e_k s_k \right) \left( \sum_{s_{l,i} \in s_l} g(s_{l,i}) \exp \left( -\beta e_l s_{l,i} \right) \right) + \dots \right] \\
&= \frac{1}{Z_Q} \left[ \sum_{\mathbf{s}/s_l} f(\mathbf{s}/s_l) \exp \left( -\beta \sum_{s_k \in \mathbf{s}/s_l} e_k s_k \right) \right] \left[ \sum_{s_{l,i} \in s_l} g(s_{l,i}) \exp \left( -\beta e_l s_{l,i} \right) \right] \\
&= \frac{1}{Z_Q} \left[ \sum_{\mathbf{s}/s_l} f(\mathbf{s}/s_l) \exp \left( -\beta \sum_{s_k \in \mathbf{s}/s_l} e_k s_k \right) \right] \frac{\sum_{s_{l,i} \in s_l} \exp \left( -\beta e_l s_{l,i} \right)}{\sum_{s_{l,i} \in s_l} \exp \left( -\beta e_l s_{l,i} \right)} \left[ \sum_{s_{l,i} \in s_l} g(s_{l,i}) \exp \left( -\beta e_l s_{l,i} \right) \right] \\
&= \frac{1}{Z_Q} \underbrace{\left[ \sum_{\mathbf{s}/s_l} f(\mathbf{s}/s_l) \exp \left( -\beta \sum_{s_k \in \mathbf{s}} e_k s_k \right) \right]}_{\langle f(\mathbf{s}/s_l) \rangle_Q} \underbrace{\frac{1}{\sum_{s_{l,i} \in s_l} \exp \left( -\beta e_l s_{l,i} \right)} \left[ \sum_{s_{l,i} \in s_l} g(s_{l,i}) \exp \left( -\beta e_l s_{l,i} \right) \right]}_{\langle g(s_l) \rangle_Q}
\end{aligned}$$

Therefore,

$$\langle s_l \rangle_Q = \frac{\sum_{s_{l,i} \in s_l} s_{l,i} \exp \left( -\beta e_l s_{l,i} \right)}{\sum_{s_{l,i} \in s_l} \exp \left( -\beta e_l s_{l,i} \right)}$$

### The mean-field approximation

- True Distribution :  $P(\mathbf{s}) = \frac{1}{Z} \exp(-\beta E_p)$
- Approximation Distribution :  $Q(\mathbf{s}) = \frac{1}{Z} \exp(-\beta \sum_k e_k s_k)$  and  $e_k$  are mean fields which we shall find.

To find  $e_k$ , we minimize Kullback-Leibler Divergence as follow :

$$\begin{aligned}
D_{KL}(Q||P) &= \sum_{\mathbf{s}} Q(\mathbf{s}) \ln \frac{Q(\mathbf{s})}{P(\mathbf{s})} \\
&\stackrel{!}{=} \min
\end{aligned}$$

To find  $e_l$ , we do  $\frac{\partial D_{KL}}{\partial e_l} \stackrel{!}{=} 0$ .

$$\begin{aligned}
D_{KL} &= \sum_{\mathbf{s}} Q_{(\mathbf{s})} \ln Q_{(\mathbf{s})} - Q_{(\mathbf{s})} \ln P_{(\mathbf{s})} \\
&= \sum_{\mathbf{s}} Q_{(\mathbf{s})} \ln \frac{1}{Z_Q} \exp(-\beta E_Q) - Q_{(\mathbf{s})} \ln \frac{1}{Z_P} \exp(-\beta E_P) \\
&= \sum_{\mathbf{s}} Q_{(\mathbf{s})} \left( -\beta E_Q - \ln Z_Q \right) - Q_{(\mathbf{s})} \left( -\beta E_P - \ln Z_P \right) \\
&= \sum_{\mathbf{s}} -\beta E_Q Q_{(\mathbf{s})} - Q_{(\mathbf{s})} \ln Z_Q + \beta E_P Q_{(\mathbf{s})} + Q_{(\mathbf{s})} \ln Z_P \\
&= \sum_{\mathbf{s}} -\beta E_Q Q_{(\mathbf{s})} - \sum_{\mathbf{s}} Q_{(\mathbf{s})} \ln Z_Q + \sum_{\mathbf{s}} \beta E_P Q_{(\mathbf{s})} + \sum_{\mathbf{s}} Q_{(\mathbf{s})} \ln Z_P \\
&= \sum_{\mathbf{s}} -\beta E_Q Q_{(\mathbf{s})} - \ln Z_Q + \sum_{\mathbf{s}} \beta E_P Q_{(\mathbf{s})} + \ln Z_P \\
&= \sum_{\mathbf{s}} -\beta E_Q Q_{(\mathbf{s})} - \ln Z_Q + \beta \langle E_P \rangle_Q + \ln Z_P
\end{aligned}$$

Therefore,

$$\begin{aligned}
\frac{\partial}{\partial e_l} D_{KL} &= \frac{\partial}{\partial e_l} \sum_{\mathbf{s}} -\beta E_Q Q_{(\mathbf{s})} - \frac{\partial}{\partial e_l} \ln Z_Q + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q + \frac{\partial}{\partial e_l} \ln Z_P \\
&= -\beta \sum_{\mathbf{s}} \frac{\partial}{\partial e_l} E_Q Q_{(\mathbf{s})} - \frac{1}{Z_Q} \frac{\partial}{\partial e_l} Z_Q + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&= -\beta \sum_{\mathbf{s}} \left[ E_Q \frac{\partial}{\partial e_l} Q_{(\mathbf{s})} + Q_{(\mathbf{s})} \frac{\partial}{\partial e_l} E_Q \right] - \frac{1}{Z_Q} \frac{\partial}{\partial e_l} \sum_{\mathbf{s}} \exp \left( -\beta \sum_k e_k s_k \right) + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&= -\beta \sum_{\mathbf{s}} \left[ \sum_k e_k s_k \frac{\partial}{\partial e_l} Q_{(\mathbf{s})} + Q_{(\mathbf{s})} s_l \right] - \frac{1}{Z_Q} \sum_{\mathbf{s}} \exp \left( -\beta \sum_k e_k s_k \right) (-\beta s_l) + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&= -\beta \sum_{\mathbf{s}} \left[ \sum_k e_k s_k \frac{\partial}{\partial e_l} Q_{(\mathbf{s})} \right] - \beta \left[ \sum_{\mathbf{s}} Q_{(\mathbf{s})} s_l \right] + \frac{\beta}{Z_Q} \sum_{\mathbf{s}} s_l \exp \left( -\beta \sum_k e_k s_k \right) + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&= -\beta \sum_k \left[ e_k \frac{\partial}{\partial e_l} \sum_{\mathbf{s}} s_k Q_{(\mathbf{s})} \right] - \beta \left[ \sum_{\mathbf{s}} Q_{(\mathbf{s})} s_l \right] + \beta \sum_{\mathbf{s}} Q_{\mathbf{s}} s_l + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&= -\beta \sum_k \left[ e_k \frac{\partial}{\partial e_l} \langle s_k \rangle_Q \right] + \beta \frac{\partial}{\partial e_l} \langle E_P \rangle_Q \\
&\stackrel{!}{=} 0
\end{aligned}$$

Hence,

$$\sum_k \left[ e_k \frac{\partial}{\partial e_l} \langle s_k \rangle_Q \right] = \frac{\partial}{\partial e_l} \langle E_P \rangle_Q$$

Due to the independent assumption that  $s_i \perp s_j$ , yields

$$\frac{\partial}{\partial e_l} \langle E_P \rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q = 0$$

### Mean Field for Ising Model

$\mathbf{s} = \{s_i\}_{i=1,\dots,N}$  where  $s_i \in \mathcal{S} = -1, 1$  and

$$E_{(\mathbf{s})} = -\frac{1}{2} \sum_{i,j=1, i \neq j}^N W_{ij} s_i s_j$$

where  $W$  is a real symmetric matrix with zero diagonal entries.



**Algorithm 6:** Mean Field Annealing

---

```

initialization:  $\langle \mathbf{s} \rangle_0, \beta_0, t = 0$ 
begin Annealing loop
  repeat
    calculate mean-fields:  $e_k, k = 1, \dots, N$ 
    calculate moments:  $\langle s_k \rangle_Q, k = 1, \dots, N$ 
  until  $|e_k^{\text{old}} - e_k^{\text{new}}| < \varepsilon$ 
  increase  $\beta$ 
end

```

---

$\Rightarrow$  inner loop: fixed-point iteration for the mean-fields  $e_k$

$\Rightarrow$  the moments  $\langle s_k \rangle$  for not too small temperatures are in general not from the discrete set  $\mathcal{S}$  but range continuously between the elements of  $\mathcal{S}$

$\Rightarrow$  however:  $\beta \rightarrow \infty (T \rightarrow 0) : \langle s_k \rangle \rightarrow s_k^*$   
because  $P(\mathbf{s})$  becomes singular at the state  $\mathbf{s}^*$  of minimal cost!

$\Rightarrow$  deterministic (fast) rather than stochastic (slow) optimization method  
(given that mean-field equations can be easily evaluated)

Figure 3.2: Drawn from MI2 course note

1. The first moment  $\langle s_l \rangle_Q$  can be computed follows :

$$\begin{aligned}
 \langle s_l \rangle_Q &= \frac{\sum_{s_{l,i} \in \mathcal{S}_l} s_{l,i} \exp \left( -\beta e_l s_{l,i} \right)}{\sum_{s_{l,i} \in \mathcal{S}_l} \exp \left( -\beta e_l s_{l,i} \right)} \\
 &= \frac{\exp(-\beta e_l) - \exp(\beta e_l)}{\exp(-\beta e_l) + \exp(\beta e_l)} \\
 &= \tanh(-\beta e_l)
 \end{aligned}$$

2.  $e_l$  is computed by

$$\begin{aligned}
 0 &= \frac{\partial}{\partial e_l} \langle E_P \rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q \\
 &= \frac{\partial}{\partial e_l} \left\langle -\frac{1}{2} \sum_{i,j=1, i \neq j}^N W_{ij} s_i s_j \right\rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q \\
 &= -\frac{1}{2} \frac{\partial}{\partial e_l} \sum_{i,j=1, i \neq j}^N W_{ij} \langle s_i \rangle_Q \langle s_j \rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q \\
 &= \sum_{i=1, i \neq l}^N W_{il} \langle s_i \rangle_Q \frac{\partial}{\partial e_l} \langle s_l \rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q
 \end{aligned}$$

Hence,

$$e_l = \sum_{i=1, i \neq l}^N W_{il} \langle s_i \rangle_Q$$



## Chapter 4

# Clustering and Embedding

Instead of finding interesting projections from data, we can use clustering algorithms to find groups or categories of the data. Each cluster is represented by a prototype  $\mathbf{w}$ .

Given set of  $\mathbf{x}^{(\alpha)}, \alpha = 1, \dots, p; \mathbf{x}^{(\alpha)} \in \mathbb{R}^N$ , we want to assign them into  $M$  clusters where the indicator

$$m_q^{(\alpha)} = \begin{cases} 1, & \text{if } \mathbf{x}^{(\alpha)} \in \text{cluster } q \\ 0, & \text{otherwise} \end{cases}$$

and  $\sum_q m_q^{(\alpha)} = 1$ .

### 4.1 K-Means Clustering

#### 4.1.1 Cost function

$$E^T_{[\{m_q^{(\alpha)}\}, \{\mathbf{w}_q\}]} = \frac{1}{p} \sum_{\alpha, q} m_q^{(\alpha)} \left( \mathbf{x}^{(\alpha)} - \mathbf{w}_q \right)^2$$

- Cluster center is defined by continuous variable
- Cluster assign is binary
- Dissimilarity is measured by Euclidean distance.

**Data:**  $\mathbf{w}_q = \langle \mathbf{x} \rangle + \boldsymbol{\eta}_q$  where  $\boldsymbol{\eta}_q$  is a random vector

**while** *not converge* **do**

$$\left| \begin{array}{l} \forall \alpha \in p : m_q^{(\alpha)} = 1 \left\{ q = \underset{\gamma}{\operatorname{argmin}} \left( \mathbf{x}^{(\alpha)} - \mathbf{w}_\gamma \right)^2 \right\} \\ \forall q \in M : \mathbf{w}_q = \frac{\sum_{\alpha} m_q^{(\alpha)} \mathbf{x}^{(\alpha)}}{\sum_{\alpha} m_q^{(\alpha)}} \end{array} \right.$$

**end**

**Algorithm 2:** Batch K-means

- $E^T$  is monotonic decreasing.
- $\mathbf{w}_q$  is center of mass  $\rightarrow E^T$  is total variance.

**Data:**  $\mathbf{w}_q = \langle \mathbf{x} \rangle + \boldsymbol{\eta}_q$  where  $\boldsymbol{\eta}_q$  is a random vector

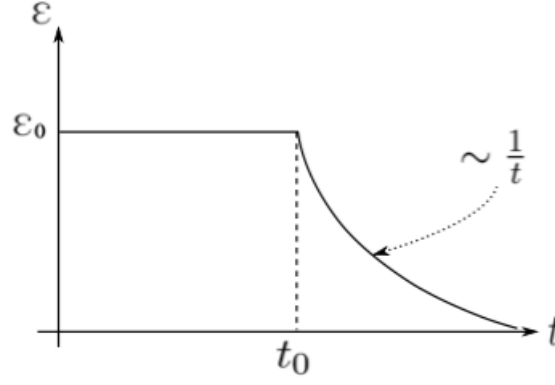
**Data:**  $\epsilon \in [0, 1]$

**while** *not converge* **do**

$$\left| \begin{array}{l} \text{Select a point } \mathbf{x}^{(\alpha)} \\ q \leftarrow \underset{\gamma}{\operatorname{argmin}} \left( \mathbf{x}^{(\alpha)} - \mathbf{w}_\gamma \right)^2 \\ \mathbf{w}_q \leftarrow \mathbf{w}_q + \epsilon (\mathbf{x}^{(\alpha)} - \mathbf{w}_q) \end{array} \right.$$

**end**

**Algorithm 3:** Online K-means

Figure 4.1:  $\epsilon$  Annealing Schedule

### 4.1.2 Number of prototypes

- $E^T \propto \frac{1}{M}$
- $E^T = 0$  when  $M = p$

**Data:**  $w_q = \frac{1}{p} \sum_{\alpha} x^{(\alpha)}$

**Data:**  $(E^T)^* = \text{some constant}$

**Data:**  $M = 1$

**while**  $E^T > (E^T)^*$  **do**

$$q \leftarrow \underset{\gamma \in M}{\operatorname{argmax}} \frac{\sum_{\alpha} m_{\gamma}^{(\alpha)} (x^{(\alpha)} - w_q)^2}{\sum_{\alpha} m_{\gamma}^{(\alpha)}}$$

$w_{M+1} \leftarrow w_q + n_q$

$M \leftarrow M + 1$

do regular K-Means algorithm with  $M$  prototypes

**end**

**Algorithm 4:** Iterative refinement

### 4.1.3 Validation Measure

Given  $Y = (y_1, \dots, y_p)$ ,  $y_i \in \{1, \dots, M\}$  solution of the clustering algorithm. We can measure the solutions  $Y_1, Y_2$  by

$$d = \frac{1}{|Y_1|} \sum_{\alpha} \mathbf{1}\{Y_{1,\alpha} \neq Y_{2,\alpha}\}$$

---

**Algorithm 4: Validation measure**

---

```

begin for each  $M \in \{M_{min}, \dots, M_{max}\}$ 
  begin loop for  $r$  splits of data
    Split data  $\mathbf{X}$  into  $\mathbf{X}_1^{(i)}$  and  $\mathbf{X}_2^{(i)}$  and find corresponding clustering solutions  $\mathbf{Y}_1^{(i)}$ 
    and  $\mathbf{Y}_2^{(i)}$ 
    Compute dissimilarity  $d_i := \frac{1}{|\mathbf{X}_2|} \sum_{\alpha} \mathbf{1}\{\mathbf{Y}_2^{(\alpha)} \neq \phi_1[\mathbf{X}_2^{(\alpha)}]\}$ 
    Where  $\phi_1$  denotes the label of the nearest prototype taking into account
    permutations
  end
  Compute average dissimilarity:  $\hat{S}_{clustering} = \frac{1}{r} \sum_r d_i$ 

  Sample  $s$  random clustering assignments and compute estimate average of the
  distances to estimate  $\hat{S}_{random}$ 

  Calculate stability index:  $\bar{S}_M = \frac{\hat{S}_{clustering}}{\hat{S}_{random}}$ 
end
Return  $\hat{M} = \operatorname{argmin}_M(\bar{S}_M)$ 

```

---

Figure 4.2: Validation measure

## 4.2 Pairwise Clustering

distance matrix  $\{d_{\alpha\alpha'}\}$  represents pairwise data and it is symmetric and its on diagonal elements are zero.

### 4.2.1 Cost Function

$$\begin{aligned}
E_{[\{m_q^{(\alpha)}\}]} &= \frac{1}{2p} \sum_q \frac{\sum_{\alpha\alpha'} m_q^{(\alpha)} m_q^{(\alpha')} (\mathbf{x}^{(\alpha)} - \mathbf{x}^{(\alpha')})^2}{\sum_{\alpha} m_q^{(\alpha)}} \\
&= \frac{1}{2p} \sum_q \frac{\sum_{\alpha\alpha'} m_q^{(\alpha)} m_q^{(\alpha')} \left( (\mathbf{x}^{(\alpha)})^2 - 2(\mathbf{x}^{(\alpha)})^T \mathbf{x}^{(\alpha')} + (\mathbf{x}^{(\alpha')})^2 \right)}{\sum_{\alpha} m_q^{(\alpha)}} \\
&= \frac{1}{2p} \sum_q \frac{\sum_{\alpha\alpha'} m_q^{(\alpha)} m_q^{(\alpha')} (\mathbf{x}^{(\alpha)})^2}{\sum_{\alpha} m_q^{(\alpha)}} - 2 \frac{\sum_{\alpha\alpha'} m_q^{(\alpha)} m_q^{(\alpha')} \mathbf{x}^{(\alpha)} \mathbf{x}^{(\alpha')}}{\sum_{\alpha} m_q^{(\alpha)}} + \frac{\sum_{\alpha\alpha'} m_q^{(\alpha)} m_q^{(\alpha')} (\mathbf{x}^{(\alpha')})^2}{\sum_{\alpha} m_q^{(\alpha)}} \\
&= \frac{1}{2p} \sum_q \frac{\sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 \sum_{\alpha'} m_q^{(\alpha')}}{\sum_{\alpha} m_q^{(\alpha)}} - 2 \frac{\sum_{\alpha} m_q^{(\alpha)} \mathbf{x}^{(\alpha)} \sum_{\alpha'} m_q^{(\alpha')} \mathbf{x}^{(\alpha')}}{\sum_{\alpha} m_q^{(\alpha)}} + \frac{\sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 \sum_{\alpha'} m_q^{(\alpha')}}{\sum_{\alpha} m_q^{(\alpha)}} \\
&= \frac{1}{2p} \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - 2 \sum_{\alpha} m_q^{(\alpha)} \mathbf{x}^{(\alpha)} \frac{\sum_{\alpha'} m_q^{(\alpha')} \mathbf{x}^{(\alpha')}}{\sum_{\alpha'} m_q^{(\alpha')}} + \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 \\
&= \frac{1}{p} \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q - \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q + \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q - \sum_q \sum_{\alpha} m_q^{(\alpha)} \frac{\sum_{\alpha'} m_q^{(\alpha')} (\mathbf{x}^{(\alpha')})^T}{\sum_{\alpha'} m_q^{(\alpha')}} \mathbf{w}_q + \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q - \sum_q \frac{\sum_{\alpha} m_q^{(\alpha)}}{\sum_{\alpha'} m_q^{(\alpha')}} \sum_{\alpha'} m_q^{(\alpha')} (\mathbf{x}^{(\alpha')})^T \mathbf{w}_q + \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q - \sum_q \sum_{\alpha'} m_q^{(\alpha')} (\mathbf{x}^{(\alpha')})^T \mathbf{w}_q + \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \left\{ \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^2 - 2 \sum_q \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)})^T \mathbf{w}_q + \sum_q \sum_{\alpha} m_q^{(\alpha)} \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \sum_{q,\alpha} m_q^{(\alpha)} \left\{ (\mathbf{x}^{(\alpha)})^2 - 2(\mathbf{x}^{(\alpha)})^T \mathbf{w}_q + \mathbf{w}_q^T \mathbf{w}_q \right\} \\
&= \frac{1}{p} \sum_{q,\alpha} m_q^{(\alpha)} \left\{ \mathbf{x}^{(\alpha)} - \mathbf{w}_q \right\}^2 \\
&\implies \text{K-means cost function}
\end{aligned}$$

### 4.2.2 Mean-field approximation for pairwise-clustering

Define

- $\mathcal{M}$  : Cartesian product of all assignments for  $\forall \alpha$
- $\mathcal{M}_{\gamma} = \mathcal{M} / \{\mathbf{m}^{(\gamma)}\}$

From Gibbs Distribution, we know that

$$\begin{aligned}
\langle m_q^{(\gamma)} \rangle &= \frac{1}{\sum_{\mathcal{M}} \exp \left( -\beta \sum_{r,\delta} m_r^{(\delta)} e_r^{(\delta)} \right)} \sum_{\mathcal{M}} m_q^{(\gamma)} \exp \left( -\beta \sum_{r,\delta} m_r^{(\delta)} e_r^{(\delta)} \right) \\
&= \frac{\left[ \sum_{\mathcal{M}_\gamma} \exp \left( -\beta \sum_{r,\delta \neq \gamma} m_r^{(\delta)} e_r^{(\delta)} \right) \right] \left[ \sum_{\{\mathbf{m}^{(\gamma)}\}} m_q^{(\gamma)} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right) \right]}{\sum_{\mathcal{M}} \exp \left( -\beta \sum_{r,\delta} m_r^{(\delta)} e_r^{(\delta)} \right)} \\
&= \frac{\left[ \sum_{\mathcal{M}_\gamma} \exp \left( -\beta \sum_{r,\delta \neq \gamma} m_r^{(\delta)} e_r^{(\delta)} \right) \right] \left[ \sum_{\{\mathbf{m}^{(\gamma)}\}} m_q^{(\gamma)} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right) \right]}{\left[ \sum_{\mathcal{M}_\gamma} \exp \left( -\beta \sum_{r,\delta \neq \gamma} m_r^{(\delta)} e_r^{(\delta)} \right) \right] \left[ \sum_{\{\mathbf{m}^{(\gamma)}\}} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right) \right]} \\
&= \frac{\sum_{\{\mathbf{m}^{(\gamma)}\}} m_q^{(\gamma)} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right)}{\sum_{\{\mathbf{m}^{(\gamma)}\}} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right)} \\
&= \frac{\exp \left( -\beta e_q^{(\gamma)} \right)}{\sum_{\{\mathbf{m}^{(\gamma)}\}} \exp \left( -\beta \sum_r m_r^{(\gamma)} e_r^{(\gamma)} \right)} \quad (\text{only remain when } m_q^{(\alpha)} = 1) \\
&= \frac{\exp \left( -\beta e_q^{(\gamma)} \right)}{\sum_{\{\mathbf{m}^{(\gamma)}\}} \exp \left( -\beta e_r^{(\gamma)} \right)} \quad (\text{only remain when } m_r^{(\alpha)} = 1) \\
&= \frac{\exp \left( -\beta e_q^{(\gamma)} \right)}{\sum_r \exp \left( -\beta e_r^{(\gamma)} \right)} \quad (r \in \mathbf{m}^{(\gamma)}) \\
&\approx \text{soft-max of the mean fields}
\end{aligned}$$

$$e_q^{(\alpha)} = \underbrace{\frac{2}{p} \frac{1}{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q} \sum_\delta \langle m_q^{(\delta)} \rangle_Q \left\{ d_{\delta\alpha} - \underbrace{\frac{1}{2} \frac{1}{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q} \sum_\varepsilon \langle m_q^{(\varepsilon)} \rangle_Q d_{\varepsilon\delta}}_{\substack{\text{distance between data objects } \alpha \text{ and } \delta, \\ \text{corrected by the average distance between} \\ \text{objects of the cluster, to which } \delta \text{ belongs (here: } q\text{)}}} \right\}}_{\substack{\text{average corrected distance between data objects} \\ \alpha \text{ and all objects } \delta \text{ of cluster } q}} \quad (4.27)$$

interpretation of the “mean fields”:

$\Rightarrow$  ”metric visualization” of the  $e_q^{(\alpha)}$ :

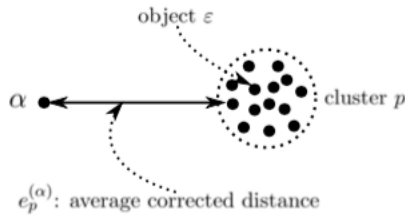


Figure 4.3:

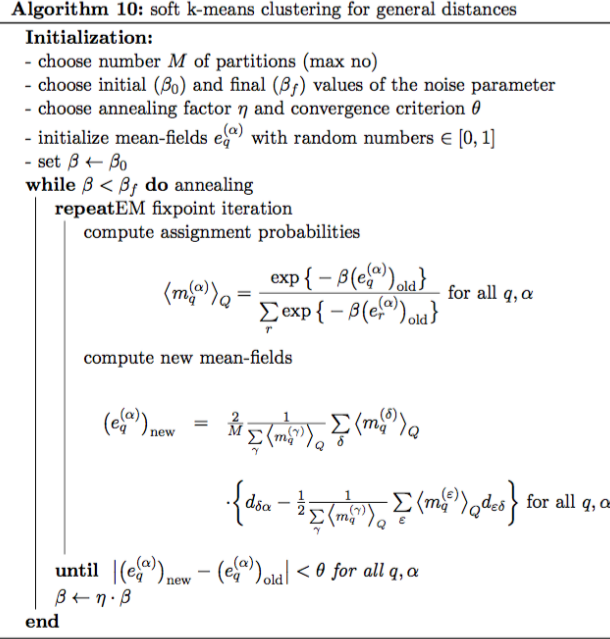


Figure 4.4:

### 4.2.3 Missing data

Sometimes, we don't have access to all  $d_{q\alpha}$  pairs. We can derive as follow :

$$\bar{d}_{q\alpha} = \frac{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q d_{\gamma\alpha}}{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q}$$

$$e_q^{(\alpha)} = \left\{ \bar{d}_{q\alpha} - \frac{1}{2} \sum_\delta \langle m_q^{(\delta)} \rangle_Q \bar{d}_{\gamma\delta} \right\}$$

Therefore, we can compute  $\bar{d}_{\gamma\delta}$  using only measured data that we have.

### 4.2.4 Soft K-means for Euclidean distances

$$d_{\alpha\alpha'} = \frac{1}{2} \left( \mathbf{x}^{(\alpha)} - \mathbf{x}^{(\alpha')} \right)^2$$

$$e_q^{(\alpha)} = \frac{1}{2} \left( \mathbf{x}^{(\alpha)} - \mathbf{w}_q \right)^2$$

and

$$\mathbf{w}_q = \frac{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q \mathbf{x}^{(\gamma)}}{\sum_\gamma \langle m_q^{(\gamma)} \rangle_Q}$$

For online version, we compute  $\langle m_q^{(\alpha)} \rangle_Q, \forall q \in M$  and then update  $\mathbf{w}_q$  with the rule :

$$\mathbf{w}_q \leftarrow \mathbf{w}_q + \eta \langle m_q^{(\alpha)} \rangle_Q \left( \mathbf{x}^{(\alpha)} - \mathbf{w}_q \right)$$

### 4.2.5 Influence of $\beta$ to cluster structure

Average cost is

$$\langle E \rangle = \frac{1}{Z} \sum_{\{m_p^{(\alpha)}\}} E \exp(-\beta E)$$

$\uparrow \beta$  means  $\downarrow$  average cost,  $\downarrow$  cluster size and  $\uparrow$  hierachical clustering.

**Algorithm 11:** soft k-means clustering for Euclidean distances

---

**Initialization:**

- choose no.  $M$  of partitions (max no)
- choose initial ( $\beta_0$ ) and final ( $\beta_f$ ) values of the noise parameter
- initialize prototypes:  $\mathbf{w}_q = \frac{1}{p} \sum_{\alpha} \mathbf{x}^{(\alpha)} + \underline{\eta}_q$  (small random vector)
- choose annealing factor  $\eta$
- choose convergence criterion  $\theta$
- $\beta \leftarrow \beta_0$

**while**  $\beta < \beta_f$  (*annealing*) **do**

**repeat** EM

compute assignment probabilities

$$\langle m_q^{(\alpha)} \rangle_Q = \frac{\exp \left\{ -\frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q^{\text{old}})^2 \right\}}{\sum_r \exp \left\{ -\frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_r^{\text{old}})^2 \right\}} \text{ for all } \alpha, q$$

compute new prototypes

$$\mathbf{w}_q^{\text{new}} = \frac{\sum_{\alpha} \langle m_q^{(\alpha)} \rangle_Q \mathbf{x}^{(\alpha)}}{\sum_{\alpha} \langle m_q^{(\alpha)} \rangle_Q} \text{ for all } q$$

center of mass of the data points  
which belong to cluster  $q$  - weighted  
by assignment probability

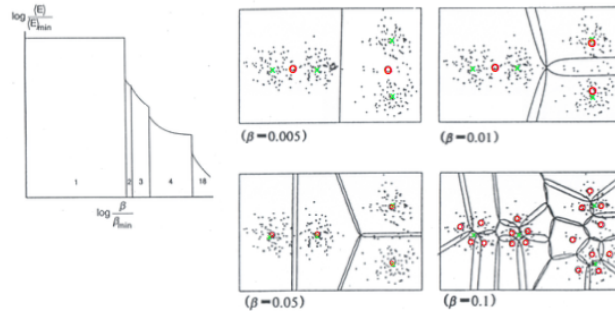
**until**  $|\mathbf{w}_q^{\text{new}} - \mathbf{w}_q^{\text{old}}| < \theta$  for all  $q$

$\beta \leftarrow \eta\beta$

**end**

---

Figure 4.5: Soft K-means with Euclidean distance

Figure 4.6: Influence of  $\beta$  to cluster structure

### 4.3 Self-Organizing Maps(SOM)

SOM is a clustering algorithm that is based on similarity, e.g. Euclidean distance. It's a low-dimensional embedding technique that preverves neighborhood.

**Data:** choose of no. partitions ( clusters )  $M$

**Data:** choose of annealing schedule  $\epsilon$  and  $\sigma$

**for**  $\alpha = \text{shuffle}(\{1, \dots, p\})$  **do**

Find the closest prototype :  $p = \underset{r}{\operatorname{argmin}} |\mathbf{x}^{(\alpha)} - \mathbf{w}_r|$

$\forall q$  : update  $\mathbf{w}_q$  respect to :

$$h_{qp} = \exp \left( -\frac{(p-q)^2}{2\sigma^2} \right) \quad (\text{Neighborhood function})$$

$$\Delta \mathbf{w}_q = \epsilon h_{qp} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)$$

**end**

**Algorithm 5:** On-line learning for SOM

Noting that, if  $\sigma = 0$  the algorithm becomes an standard on-line K-means.



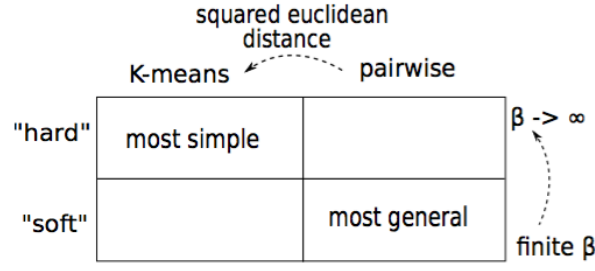


Figure 4.7: Summary of clustering algorithms

#### 4.3.1 Neighborhood function $h_{qp}$

$$h_{qp} = \exp\left(-\frac{(p-q)^2}{2\sigma^2}\right)$$

- $\uparrow \sigma$  gives us smooth cost function  $\Rightarrow$  lower-cost optimization but results might not well represent the data.
- $\downarrow \sigma$  gives us rough cost function  $\Rightarrow$  higher-cost optimization but has more representation capabilities.

#### 4.3.2 Self-organizing maps for pairwise data

- For each pair of data, we compute dissimilarity  $d_{\alpha\alpha'}$ .
- Construct  $M$  clusters with geometrical structure ( 1d line, 2d grid )
- $m_q^{(\alpha)}$  binary assignment variable ( normalize to 1 for each  $\alpha$  )

$$m_q^{(\alpha)} = 1\{\mathbf{x}^{(\alpha)} \text{ belongs to } \mathbf{w}_q\}$$

#### Cost Function

$$E_{[\{m_q^{(\alpha)}\}]} = \frac{1}{p} \sum_r \frac{\sum_{\alpha, \alpha'} \left( \sum_q h_{rq} m_q^{(\alpha)} \right) \left( \sum_q h_{rq} m_q^{(\alpha')} \right) d_{\alpha\alpha'}}{\sum_{\alpha} \left( \sum_q h_{rq} m_q^{(\alpha)} \right)}$$

$$\stackrel{!}{=} \min$$

This cost function is similar to pairwise clustering but  $\mathbf{w}_q$  is replaced by  $\sum_q h_{rq} m_q^{(\alpha)}$ .

From Figure 4.8, if we replace  $h_{qp}$  with Direc's  $\delta_{qp} = 1\{q = p\}$ , the computation becomes standard pairwise clustering. This is called "Kohonen-approximation", hence no need  $\sigma$  for annealing.

#### Euclidean distances

$$d_{\alpha\alpha'} = \frac{1}{2} \left( \mathbf{x}^{(\alpha)} - \mathbf{x}^{(\alpha')} \right)^2$$

$$E_{[\{m_q^{(\alpha)}\}]} = \frac{1}{p} \sum_{q, \alpha} \left( \sum_p h_{qp} m_p^{(\alpha)} \right) (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2$$

$$= \frac{1}{p} \sum_{p, \alpha} m_p^{(\alpha)} \left( \sum_q h_{qp} \right) (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2$$

$$\mathbf{w}_q = \frac{\sum_{\alpha'} \left( \sum_p h_{qp} m_p^{(\alpha')} \mathbf{x}^{(\alpha')} \right)}{\sum_{\alpha'} \left( \sum_p h_{qp} m_p^{(\alpha')} \right)}$$

**Initialization:**

- choose no.  $M$  of partitions, initial  $(\beta_0)$  and final  $(\beta_f)$  noise parameters, annealing factor  $\eta$ , width  $\sigma$  of neighborhood function  $h_{sq}$ , tolerance  $\theta$
- initialize mean-fields  $e_{\underline{q}}^{(\alpha)}$  : random numbers  $\in [0, 1]$

$\beta \leftarrow \beta_0$

**while**  $\beta < \beta_f$  **do** annealing

**repeat** EM

        compute assignment probabilities  $\langle m_{\underline{q}}^{(\alpha)} \rangle_Q = \frac{\exp \{-\beta (e_{\underline{q}}^{(\alpha)})_{\text{old}}\}}{\sum_{\underline{r}} \exp \{-\beta (e_{\underline{r}}^{(\alpha)})_{\text{old}}\}} \quad \forall \underline{q}, \alpha$

        compute new mean-fields

$$(e_{\underline{q}}^{(\alpha)})_{\text{new}} = \frac{1}{p} \sum_{\underline{s}} h_{sq} \left[ \frac{1}{\sum_{\gamma} \left( \sum_{\underline{r}} h_{sr} \langle m_{\underline{r}}^{(\gamma)} \rangle_Q \right)} \sum_{\delta} \left( \sum_{\underline{r}} h_{sr} \langle m_{\underline{r}}^{(\delta)} \rangle_Q \right) \right]$$

$$\cdot \left\{ d_{\delta\alpha} - \frac{1}{2} \frac{1}{\sum_{\gamma} \left( \sum_{\underline{r}} h_{sr} \langle m_{\underline{r}}^{(\gamma)} \rangle_Q \right)} \sum_{\epsilon} \left( \sum_{\underline{r}} h_{sr} \langle m_{\underline{r}}^{(\epsilon)} \rangle_Q \right) d_{\epsilon\delta} \right\} \quad \forall \underline{q}, \alpha$$

**until**  $|(e_{\underline{q}}^{(\alpha)})_{\text{new}} - (e_{\underline{q}}^{(\alpha)})_{\text{old}}| < \theta \quad \forall \underline{q}, \alpha$

$\beta \leftarrow \eta\beta$

**end**

Figure 4.8: SOM for Pairwise data with mean-field approximation

$w_q$  is now a center of mass of the cluster weighted by neighborhood function  $h_{qp}$ .

**Data:** choose of no. partitions ( clusters )  $M$

**Data:** choose of annealing schedule  $\epsilon$

**for**  $\alpha = \text{shuffle}(\{1, \dots, p\})$  **do**

    Find the closest prototype :  $p = \underset{r}{\operatorname{argmin}} \sum_q h_{rq} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2$

$\forall q$  : update  $w_q$  respect to :

$$h_{qp} = \exp \left( - \frac{(p - q)^2}{2\sigma^2} \right) \quad (\text{Neighborhood function})$$

$$\Delta \mathbf{w}_q = \epsilon h_{qp} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)$$

**end**

**Algorithm 6:** On-line learning for SOM with Euclidean distances

## 4.4 Locally Linear Embedding

Locally linear embedding (LLE) seeks a lower-dimensional projection of the data which preserves distances within local neighborhoods - quoted from scikit-learn<sup>1</sup>.

Given data  $\mathbf{x}^{(\alpha)} \in \mathbb{R}^N, \alpha = \{1, \dots, p\}$ .

1. Find  $K$ -nearest neighbors of  $\mathbf{x}^{(\alpha)}$
2. Calculate reconstruction weights  $W$  of  $\mathbf{x}^{(\alpha)}$  from its  $K$ -neighbors.
3. Obtain embedding  $\mathbf{u}^{(\alpha)} \in \mathbb{R}^M, M \ll N$  from  $W$  and its neighbors' embedding.

### 4.4.1 Find $K$ -nearest neighbors

Traditional  $K$ -nearest neighbors approach has  $O(Np^2)$  complexity, but k-d tree reduces the complexity to  $O(Np \log p)$ .

<sup>1</sup><http://scikit-learn.org/stable/modules/manifold.html#locally-linear-embedding>

### 4.4.2 Calculate reconstruction weights

We want to find  $W$  that minimizes reconstruction cost

$$E(W) = \sum_{\alpha=1}^p \left( \mathbf{x}^{(\alpha)} - \sum_{\beta=1}^p W_{\alpha\beta} \mathbf{x}^{(\beta)} \right)^2$$

$$\stackrel{!}{=} \min$$

s.t.

$$W_{\alpha\beta} = 0 \text{ if } \mathbf{x}^{(\beta)} \notin \text{KNN}(\mathbf{x}^{(\alpha)})$$

$$\sum_{\beta} W_{\alpha\beta} = 1$$

Nothing that  $W$  is asymmetric sparse matrix that have up to  $K$  nonempty elements for each row. It is also invariant to linear transformations e.g. rotation, scaling and translation.

For each  $\mathbf{x}^{(\alpha)}$ ,

- compute  $C^{(\alpha)} \in \mathbb{R}^{K,K}$

$$C_{ij}^{(\alpha)} = (\mathbf{x}^{(\alpha)} - \mathbf{x}^{(\beta_i)})^T (\mathbf{x}^{(\alpha)} - \mathbf{x}^{(\beta_j)})$$

where  $\mathbf{x}^{(\beta_i)}, \mathbf{x}^{(\beta_j)} \in \text{KNN}(\mathbf{x}^{(\alpha)})$

- solve linear system

$$C^{(\alpha)} \tilde{\mathbf{w}}^{(\alpha)} = (1, \dots, 1)^T$$

- rescale weight

$$W_{\alpha\beta_i} = \frac{\tilde{\mathbf{w}}_i^{(\alpha)}}{\sum_{j=1}^K \tilde{\mathbf{w}}_j^{(\alpha)}}$$

### 4.4.3 Obtain Embedding

We want to find embedding  $\forall \alpha \in \{1, \dots, p\}, \mathbf{u}^{(\alpha)} \in \mathbb{R}^M$  s.t.

$$F(U) = \sum_{\alpha=1}^p \left( \mathbf{u}^{(\alpha)} - \sum_{\beta=1}^p W_{\alpha\beta} \mathbf{u}^{(\beta)} \right)^2$$

$$= \sum_{\alpha, \beta=1}^p g_{\alpha\beta} (\mathbf{u}^{(\alpha)})^T \mathbf{u}^{(\beta)}$$

$$\stackrel{!}{=} \min$$

s.t.

$$\sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} = 0$$

(remove translation)

$$\frac{1}{p} \sum_{\alpha=1}^p (\mathbf{u}^{(\alpha)})^T \mathbf{u}^{(\alpha)} = I$$

(prevent trivial solutions e.g.  $\mathbf{u}^{(\alpha)} = 0$ )

where  $g_{\alpha\beta} = \delta_{\alpha\beta} - W_{\alpha\beta} - W_{\beta\alpha} + \sum_{\gamma=1}^p W_{\gamma\alpha} W_{\gamma\beta}$ .

Equivalently,

$$\mathbf{G} \in \mathbb{R}^{p,p}$$

$$= \{g_{\alpha\beta}\}, \forall \alpha, \beta = \{1, \dots, p\}$$

$$= (I - W^T)(I - W)$$

The solution to  $F(U)$  is to compute eigen value decomposition of  $\mathbf{G}$  with  $M + 1$  from the lowest eigen values. We discard eigen vector  $e_p$  whose eigen value is zero.

$$U = \begin{pmatrix} e_{p-M}^T \\ \vdots \\ e_{p-1}^T \end{pmatrix} = (u^{(1)}, \dots, u^{(p)}) \in \mathbb{R}^{M,p}$$

parameters:  $K, M$

- ① find  $K$  nearest neighbors  $\text{KNN}(\underline{\mathbf{x}}^{(\alpha)}) = \{\beta_1^{(\alpha)}, \dots, \beta_K^{(\alpha)}\} \forall \alpha = 1, \dots, p$
- ② calculate (locally invariant) reconstruction weights  $\underline{\mathbf{W}}$ :

$$\underline{\mathbf{C}}^{(\alpha)} \underline{\tilde{\mathbf{w}}}^{(\alpha)} = (1, \dots, 1)^T, \quad \forall \alpha = 1, \dots, p$$

$$W_{\alpha\beta_j^{(\alpha)}} = \frac{\tilde{w}_j^{(\alpha)}}{\sum_{k=1}^K \tilde{w}_k^{(\alpha)}}$$

- ③ calculate the embedding coordinates  $\underline{\mathbf{U}}$ :

compute the  $M + 1$  eigenvectors  $(\mathbf{e}_p, \dots, \mathbf{e}_{p-M})$  of  $\underline{\mathbf{G}}$   
with the smallest eigenvalues

$$g_{\alpha\beta} = \delta_{\alpha\beta} - W_{\alpha\beta} - W_{\beta\alpha} + \sum_{\gamma=1}^p W_{\gamma\alpha} W_{\gamma\beta}$$

$$\underline{\mathbf{G}} \cdot \mathbf{e}_j = \lambda_j \mathbf{e}_j \quad \underline{\mathbf{U}} = \begin{pmatrix} \mathbf{e}_{p-M}^T \\ \vdots \\ \mathbf{e}_{p-1}^T \end{pmatrix} = (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(p)})$$

Figure 4.9: Locally Linear Embedding Algorithm

## Chapter 5

# Probability Density Estimation

Density Estimation is important because if  $P(x)$  is known, we can easily compute moments of  $x$ .

There are 2 approaches to estimate density, namely

1. **Nonparametric** : histogram-based or kernel-based
2. **Parametric** : Gaussian densities

### 5.1 Nonparametric Methods

#### 5.1.1 Gliding Histogram

This approach finds  $\hat{P}(\mathbf{x})$  by counting number of points in volume centered by  $\mathbf{x}$ .

$$H(\mathbf{u}) = \begin{cases} 1, & |u_j| < \frac{1}{2} \quad \forall j \in 1, \dots, n \\ 0, & \text{otherwise} \end{cases}$$

where  $n$  is no. of dimensions of  $\mathbf{u}$ .

Given observations  $\mathbf{x}^{(\alpha)} \in \mathbb{R}^n, \alpha \in 1, \dots, p$ . We can derive :

$$\hat{P}(\mathbf{x}) = \frac{1}{h^n} \cdot \frac{1}{p} \sum_{\alpha=1}^p H\left(\frac{\mathbf{x} - \mathbf{x}^{(\alpha)}}{h}\right)$$

where  $h$  is width of the counting volume and  $\frac{1}{h^n}$  plays a normalization role.

One remark of Gliding Histogram is the **discontinuity** of the estimated density.

#### 5.1.2 Gaussian Kernel

Instead of relying on count of observations in volume ( Histogram Kernel), we can apply a kernel computation, for example Gaussian Kernel.

$$\begin{aligned} H(\mathbf{u}) &= \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(-\frac{\mathbf{u}^2}{2}\right) \\ \hat{P}(\mathbf{x}) &= \frac{1}{h^n} \cdot \frac{1}{p} \sum_{\alpha=1}^p H\left(\frac{\mathbf{x} - \mathbf{x}^{(\alpha)}}{h}\right) \\ &= \frac{1}{h^n} \cdot \frac{1}{p} \sum_{\alpha=1}^p \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(-\frac{(\mathbf{x} - \mathbf{x}^{(\alpha)})^2}{2h^2}\right) \\ &= \frac{1}{p} \sum_{\alpha=1}^p \frac{1}{(2\pi h^2)^{\frac{n}{2}}} \exp\left(-\frac{(\mathbf{x} - \mathbf{x}^{(\alpha)})^2}{2h^2}\right) \end{aligned}$$

where  $h$  is width of the kernel.

smoothness of estimated density  $\propto h$

small  $h \rightarrow$  non smooth estimation ( kind of overfitting )

large  $h \rightarrow$  smooth estimation ( could be underfitting )

## 5.2 Parametric Methods

We assume that observations  $x^{(\alpha)} \in \mathbb{R}^n, \alpha = 1, \dots, p$  are generated from a model. The goal here is to find a parameterized model via  $\mathbf{w}$  parameters ( $\hat{P}(\mathbf{x}; \mathbf{w})$ ) s.t.

$$P(\mathbf{x}) = \hat{P}(\mathbf{x}; \mathbf{w})$$

Using Kullback-Leibler Divergence

$$D_{KL}[P(\mathbf{x}), \hat{P}(\mathbf{x}; \mathbf{w})] = \int d\mathbf{x} P(\mathbf{x}) \ln \frac{P(\mathbf{x})}{\hat{P}(\mathbf{x}; \mathbf{w})}$$

$$\stackrel{!}{=} \min$$

Nothing that  $D_{KL} \geq 0$  and  $D_{KL} = 0$  iff  $P(\mathbf{x}) = \hat{P}(\mathbf{x}; \mathbf{w})$ .

Therefore,

$$\begin{aligned} \mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \int d\mathbf{x} P(\mathbf{x}) \ln \frac{P(\mathbf{x})}{\hat{P}(\mathbf{x}; \mathbf{w})} \right\} \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \int d\mathbf{x} P(\mathbf{x}) \ln P(\mathbf{x}) - \int d\mathbf{x} P(\mathbf{x}) \ln \hat{P}(\mathbf{x}; \mathbf{w}) \right\} \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ - \int d\mathbf{x} P(\mathbf{x}) \ln \hat{P}(\mathbf{x}; \mathbf{w}) \right\} \\ &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ E^G \right\} \end{aligned}$$

Using Empirical Risk Minimization (ERM), we can approximate expectation  $E^G$  using empirical average. As a result, the training cost  $E^T$  is

$$\begin{aligned} E^T &= -\frac{1}{p} \sum_{\alpha=1}^p \ln \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \\ &= \frac{1}{p} \sum_{\alpha=1}^p -\ln \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \\ &= \frac{1}{p} \sum_{\alpha=1}^p e_{[\mathbf{w}]^{(\alpha)}} \\ &\stackrel{!}{=} \min \end{aligned}$$

Using SGD, we can update  $\mathbf{w}$  as follow

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} E^T && \text{(batch-learning)} \\ \mathbf{w} &\leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} e_{[\mathbf{w}]}^{(\alpha)} && \text{(online-learning)} \end{aligned}$$

There are several approaches to validate quality of model, e.g. test-set, cross-validation set. Overfitting,  $E^T$  low but  $E^G$  high.

Another way to derive this learning procedure is to look at **likelihood function**. The likelihood of a model is the probability that the model generates the given observations with iid assumption.

$$\hat{P}(\{\mathbf{x}^{(\alpha)}\}; \mathbf{w}) = \prod_{\alpha=1}^p \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w})$$

$$\stackrel{!}{=} \max$$

$$\begin{aligned}
\mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w}} \left\{ \prod_{\alpha=1}^p \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \right\} \\
&= \operatorname{argmax}_{\mathbf{w}} \left\{ \ln \prod_{\alpha=1}^p \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \right\} \\
&= \operatorname{argmax}_{\mathbf{w}} \left\{ \sum_{\alpha=1}^p \ln \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \right\} \\
&= \operatorname{argmin}_{\mathbf{w}} \left\{ - \sum_{\alpha=1}^p \ln \hat{P}(\mathbf{x}^{(\alpha)}; \mathbf{w}) \right\}
\end{aligned}$$

### 5.2.1 Multivariate Gaussian

$$\hat{P}(\{\mathbf{x}^{(\alpha)}\}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \left( \frac{1}{\sqrt{(2\pi)^n \det \boldsymbol{\Sigma}}} \right)^p \cdot \prod_{\alpha=1}^p \exp \left( -\frac{1}{2} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}) \right)$$

Hence,

$$\begin{aligned}
E_{[\boldsymbol{\mu}, \boldsymbol{\Sigma}]}^T &= -\ln \hat{P}(\{\mathbf{x}^{(\alpha)}\}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
&= -\ln \left( \frac{1}{\sqrt{(2\pi)^n \det \boldsymbol{\Sigma}}} \right)^p - \ln \prod_{\alpha=1}^p \exp \left( -\frac{1}{2} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}) \right) \\
&= -\ln \left( (2\pi)^n \det \boldsymbol{\Sigma} \right)^{-\frac{p}{2}} - \sum_{\alpha=1}^p \left( -\frac{1}{2} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}) \right) \\
&= \frac{p}{2} \ln \left( (2\pi)^n \det \boldsymbol{\Sigma} \right) + \frac{1}{2} \sum_{\alpha=1}^p \left( (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}) \right) \\
&= \frac{pn}{2} \ln(2\pi) + \frac{p}{2} \ln(\det \boldsymbol{\Sigma}) + \frac{1}{2} \sum_{\alpha=1}^p \left( (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}) \right)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial E_{[\boldsymbol{\mu}, \boldsymbol{\Sigma}]}^T}{\partial \boldsymbol{\mu}} &= 0 \Rightarrow \boldsymbol{\mu}^* = \frac{1}{p} \sum_{\alpha=1}^p \mathbf{x}^{(\alpha)} \\
\frac{\partial E_{[\boldsymbol{\mu}, \boldsymbol{\Sigma}]}^T}{\partial \boldsymbol{\Sigma}} &= 0 \Rightarrow \boldsymbol{\Sigma}^* = \frac{1}{p} \sum_{\alpha=1}^p (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}^*)(\mathbf{x}^{(\alpha)} - \boldsymbol{\mu}^*)^T
\end{aligned}$$

Noting that  $\boldsymbol{\mu}^*$  is unbiased but  $\boldsymbol{\Sigma}^*$  is biased estimator.

### 5.2.2 Gaussian Mixture Model

Here we assume that data is generated from various sources. Formally,

$$P(\mathbf{x}) = \sum_{q=1}^M P(\mathbf{x}|q)P(q)$$

where

1.  $\mathbf{x} \in \mathbb{R}^d$ .
2.  $P(q)$  is a mixture parameter that component  $q$  creates a data point

3.  $P(\mathbf{x}|q)$  is probability that  $\mathbf{x}$  is created by component  $q$ .

There are several choices of model that we can use for  $P(\mathbf{x}|q)$ . One popular choice is a Gaussian.

$$\begin{aligned} P(\mathbf{x}|q) &= \mathcal{N}(\mathbf{x}; \mathbf{w}_q, \sigma_q^2) \\ &= \frac{1}{(2\pi\sigma_q^2)^{d/2}} \exp\left(-\frac{(\mathbf{x} - \mathbf{w}_q)^2}{2\sigma_q^2}\right) \end{aligned}$$

Hence, parameters to determine for this kind of mixture model are  $P(q), \mathbf{w}_q, \sigma_q^2$ .

Using principle of maximum likelihood, we have

$$\begin{aligned} P(\{\mathbf{x}^{(\alpha)}\}) &= \prod_{\alpha=1}^p \sum_{q=1}^M P(\mathbf{x}|q)P(q) \\ &= \prod_{\alpha=1}^p \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2)P(q) \\ &\stackrel{!}{=} \max \end{aligned}$$

Hence,

$$\begin{aligned} \{\mathbf{w}, \sigma\} &= \operatorname{argmax}_{\{\mathbf{w}, \sigma\}} \prod_{\alpha=1}^p \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2)P(q) \\ &= \operatorname{argmax}_{\{\mathbf{w}, \sigma\}} \sum_{\alpha=1}^p \ln \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2)P(q) \\ &= -\operatorname{argmin}_{\{\mathbf{w}, \sigma\}} \sum_{\alpha=1}^p \ln \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2)P(q) \end{aligned}$$

### Relation of Soft-Clustering method

Assumptions:

1. We have  $M$  Gaussian components.
2.  $\sigma_q^2 = \sigma^2 := \frac{1}{\beta}$  for all components.
3. Constant mixture parameter  $P(q) = \frac{1}{M}$ .

We can compute posterior that  $\mathbf{x}$  is created by component  $q$  via :

$$\begin{aligned} P(q|\mathbf{x}) &= \frac{P(\mathbf{x}|q)P(q)}{P(\mathbf{x})} \\ &= \frac{\frac{1}{M} \cdot \mathcal{N}(\mathbf{x}; \mathbf{w}_q, \sigma_q^2)}{\sum_{\gamma=1}^M \frac{1}{M} \cdot \mathcal{N}(\mathbf{x}; \mathbf{w}_\gamma, \sigma_\gamma^2)} \\ &= \frac{\mathcal{N}(\mathbf{x}; \mathbf{w}_q, \sigma_q^2)}{\sum_{\gamma=1}^M \mathcal{N}(\mathbf{x}; \mathbf{w}_\gamma, \sigma_\gamma^2)} \end{aligned}$$



Hence, the objective function is reduced to

$$\begin{aligned}
& - \operatorname{argmin} \sum_{\alpha=1}^p \ln \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2) P(q) \\
& - \operatorname{argmin} \sum_{\alpha=1}^p \ln \frac{1}{M} \sum_{q=1}^M \mathcal{N}(\mathbf{x}^{(\alpha)}; \mathbf{w}_q, \sigma_q^2) \\
& - \operatorname{argmin} \sum_{\alpha=1}^p \ln \left\{ \frac{1}{M} \sum_{q=1}^M \left( \frac{\beta}{2\pi} \right)^{d/2} \exp \left( - \frac{\beta(\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2}{2} \right) \right\} \\
& - \operatorname{argmin} \sum_{\alpha=1}^p \ln \left\{ \frac{1}{M} \left( \frac{\beta}{2\pi} \right)^{d/2} \sum_{q=1}^M \exp \left( - \frac{\beta(\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2}{2} \right) \right\} \\
& - \operatorname{argmin} \sum_{\alpha=1}^p \ln \left\{ \sum_{q=1}^M \exp \left( - \frac{\beta(\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2}{2} \right) \right\} + \operatorname{const}(\beta, M)
\end{aligned}$$

Given

$$E^T = - \sum_{\alpha=1}^p \ln \left\{ \sum_{q=1}^M \exp \left( - \frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2 \right) \right\}$$

We have :

$$\begin{aligned}
\frac{\partial E^T}{\partial \mathbf{w}_r} &= - \sum_{\alpha=1}^p \ln \left\{ \sum_{q=1}^M \exp \left( - \frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2 \right) \right\} \\
&= - \sum_{\alpha=1}^p \frac{\exp \left( - \frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_r)^2 \right)}{\sum_{q=1}^M \exp \left( - \frac{\beta}{2} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2 \right)} \beta (\mathbf{x}^{(\alpha)} - \mathbf{w}_r) \\
&= - \sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)}) \beta (\mathbf{x}^{(\alpha)} - \mathbf{w}_r) \\
&\stackrel{!}{=} 0
\end{aligned}$$

Hence,

$$\mathbf{w}_r = \frac{\sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)}) \mathbf{x}^{(\alpha)}}{\sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})}$$

### New interpretation of Soft-Clustering

1. All clusters have the same width  $\sigma^2 = \frac{1}{\beta}$ .
2. and same number of data points  $P(q) = \frac{1}{M}$ .

### Mixture Models ( an extension of soft-clustering )

1. each cluster has different width  $\sigma_q^2$ . and number of points  $P(q)$ .

### Optimization of Gaussian Mixture Model

We know that

$$E^T = - \sum_{\alpha=1}^p \ln \left( \sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q) P(q) \right) \stackrel{!}{=} \min$$

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{w}_r} P(\mathbf{x}^{(\alpha)}|r) &= \frac{(\mathbf{x} - \mathbf{w}_r)}{\sigma^2} P(\mathbf{x}^{(\alpha)}|r) \\
\frac{\partial}{\partial \sigma_r} P(\mathbf{x}^{(\alpha)}|r) &= \left\{ - \frac{d}{\sigma_r} + \frac{(\mathbf{x} - \mathbf{w}_r)^2}{\sigma_r^3} \right\} P(\mathbf{x}^{(\alpha)}|r)
\end{aligned}$$

$$\begin{aligned}
\frac{\partial E^T}{\partial \mathbf{w}_r} &= - \sum_{\alpha=1}^p \frac{\partial}{\partial \mathbf{w}_r} \ln \left( \sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q) \right) \\
&= - \sum_{\alpha=1}^p \frac{P(r) \frac{\partial}{\partial \mathbf{w}_r} P(\mathbf{x}^{(\alpha)}|r)}{\sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q)} \\
&= - \sum_{\alpha=1}^p \frac{P(r) \frac{(\mathbf{x} - \mathbf{w}_r)}{\sigma^2} P(\mathbf{x}^{(\alpha)}|r)}{\sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q)} \\
&= - \sum_{\alpha=1}^p \frac{(\mathbf{x} - \mathbf{w}_r)}{\sigma^2} \frac{P(r)P(\mathbf{x}^{(\alpha)}|r)}{P(\mathbf{x}^{(\alpha)})} \\
&= - \sum_{\alpha=1}^p \frac{(\mathbf{x} - \mathbf{w}_r)}{\sigma^2} P(r|\mathbf{x}^{(\alpha)}) \\
&\stackrel{!}{=} 0
\end{aligned}$$

Hence,

$$\mathbf{w}_r = \frac{\sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})\mathbf{x}^{(\alpha)}}{\sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})}$$

Similarly,

$$\begin{aligned}
\frac{\partial E^T}{\partial \sigma_r} &= - \sum_{\alpha=1}^p \frac{\partial}{\partial \sigma_r} \ln \left( \sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q) \right) \\
&= - \sum_{\alpha=1}^p \frac{P(r) \frac{\partial}{\partial \sigma_r} P(\mathbf{x}^{(\alpha)}|r)}{\sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q)} \\
&= - \sum_{\alpha=1}^p \frac{P(r) \left\{ -\frac{d}{\sigma_r} + \frac{(\mathbf{x} - \mathbf{w}_r)^2}{\sigma_r^3} \right\} P(\mathbf{x}^{(\alpha)}|r)}{\sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q)P(q)} \\
&= - \sum_{\alpha=1}^p \frac{\left\{ -\frac{d}{\sigma_r} + \frac{(\mathbf{x} - \mathbf{w}_r)^2}{\sigma_r^3} \right\} P(\mathbf{x}^{(\alpha)}|r)P(r)}{P(\mathbf{x}^{(\alpha)})} \\
&= - \sum_{\alpha=1}^p \left\{ -\frac{d}{\sigma_r} + \frac{(\mathbf{x} - \mathbf{w}_r)^2}{\sigma_r^3} \right\} P(r|\mathbf{x}^{(\alpha)}) \\
&= -\frac{1}{\sigma_r} \sum_{\alpha=1}^p \left\{ -d + \frac{(\mathbf{x} - \mathbf{w}_r)^2}{\sigma_r^2} \right\} P(r|\mathbf{x}^{(\alpha)}) \\
&= -\frac{1}{\sigma_r^3} \sum_{\alpha=1}^p \left\{ -d\sigma_r^2 P(r|\mathbf{x}^{(\alpha)}) + (\mathbf{x} - \mathbf{w}_r)^2 P(r|\mathbf{x}^{(\alpha)}) \right\} \\
&\stackrel{!}{=} 0
\end{aligned}$$

Hence,

$$\sigma_r^2 = \frac{1}{d} \frac{\sum_{\alpha=1}^p (\mathbf{x} - \mathbf{w}_r)^2 P(r|\mathbf{x}^{(\alpha)})}{\sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})}$$

For  $P(q)$ , we use Lagrange Multiplier as follow:

$$\frac{\partial}{\partial P(r)} \left\{ E^T + \lambda \left( \sum_{q=1}^M P(q) - 1 \right) \right\} \stackrel{!}{=} 0$$

$$\begin{aligned}
\frac{\partial}{\partial P(r)} \left\{ E^T + \lambda \left( \sum_{q=1}^M P(q) - 1 \right) \right\} &= \frac{\partial}{\partial P(r)} \left\{ - \sum_{\alpha=1}^p \ln \sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q) P(q) + \lambda \left( \sum_{q=1}^M P(q) - 1 \right) \right\} \\
&= - \sum_{\alpha=1}^p \frac{P(\mathbf{x}^{(\alpha)}|r)}{\sum_{q=1}^M P(\mathbf{x}^{(\alpha)}|q) P(q)} + \lambda \\
&= - \sum_{\alpha=1}^p \frac{P(\mathbf{x}^{(\alpha)}|r)}{P(\mathbf{x}^{(\alpha)})} + \lambda \\
&= - \sum_{\alpha=1}^p \frac{P(\mathbf{x}^{(\alpha)}, r)}{P(r) P(\mathbf{x}^{(\alpha)})} + \lambda \\
&= - \sum_{\alpha=1}^p \frac{P(r|\mathbf{x}^{(\alpha)})}{P(r)} + \lambda
\end{aligned}$$

Hence,

$$P(r) = \frac{1}{\lambda} \sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})$$

and because

$$\begin{aligned}
1 &= \sum_{q=1}^M P(q) \\
&= \sum_{q=1}^M \frac{1}{\lambda} \sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)}) \\
&= \frac{1}{\lambda} \sum_{\alpha=1}^p \sum_{q=1}^M P(q|\mathbf{x}^{(\alpha)}) \\
&= \frac{1}{\lambda} \sum_{\alpha=1}^p 1 \\
&= \frac{1}{\lambda} p
\end{aligned}$$

Therefore,

$$P(r) = \frac{1}{p} \sum_{\alpha=1}^p P(r|\mathbf{x}^{(\alpha)})$$

**Expectation Maximization Algorithm**

**Data:** choose of no. Gaussian ( clusters )  $M$

**Data:** Initialize  $P(q) = \frac{1}{M}$

**Data:** Initialize  $\mathbf{w}_q = \frac{1}{p} \sum_{\alpha=1}^p \mathbf{x}^{(\alpha)} + \eta_q$ ,  $\sigma_q^2 = \frac{1}{p} \sum_{\alpha=1}^p (\mathbf{x}^{(\alpha)} - \boldsymbol{\mu})^2 + \epsilon_q$ ,  $\eta_q, \epsilon_q$  is random vectors.

**while** *until convergence* **do**

    E-step :  $\forall q = \{1, \dots, M\}$

$$P(q|\mathbf{x}^{(\alpha)}) \leftarrow \frac{P(\mathbf{x}^{(\alpha)}|q)P(q)}{\sum_{r=1}^M P(\mathbf{x}^{(\alpha)}|r)P(r)}$$

    M-step :  $\forall q = \{1, \dots, M\}$

$$\begin{aligned} \mathbf{w}_q^{new} &\leftarrow \frac{\sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)})\mathbf{x}^{(\alpha)}}{\sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)})} \\ (\sigma_q^2)^{new} &\leftarrow \frac{\sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)})(\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2}{\sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)})} \\ P(q) &\leftarrow \frac{1}{p} \sum_{\alpha=1}^p P(q|\mathbf{x}^{(\alpha)}) \end{aligned}$$

**end**

**Algorithm 7:** EM Algorithm

**Proof of EM algorithm**

Given  $\mathbf{x}^{(\alpha)} \in \mathbb{R}^N, \alpha \in \{1, \dots, p\}$  and  $\mathbf{m}^{(\alpha)} = \{m_1^{(\alpha)}, \dots, m_M^{(\alpha)}\}^T \in \{0, 1\}^M$  where  $m_q^{(\alpha)} = 1$  if component  $q$  generates  $\mathbf{x}^{(\alpha)}$ , otherwise 0. We call these  $\mathbf{m}^{(\alpha)}$  latent variables.

Latent variables and maximum likelihood,

$$\begin{aligned} P(\{\mathbf{x}^{(\alpha)}\}|\boldsymbol{\theta}) &= \prod_{\alpha=1}^p P(\mathbf{x}^{(\alpha)}|\boldsymbol{\theta}) \\ &= \prod_{\alpha=1}^p P(\mathbf{x}^{(\alpha)}, \mathbf{m}^{(\alpha)}) \\ &\stackrel{!}{=} \max \end{aligned}$$

We know that :

$$0 \leq P(q) \leq 1 \text{ and } \sum_{q=1}^M P(q) = 1$$

$$\begin{aligned} P(\mathbf{x}|\boldsymbol{\theta}) &= \sum_{q=1}^M P(q)\mathcal{N}(\mathbf{x}|\mathbf{w}_q, \sigma_q^2) \\ &= \sum_{\mathbf{m}} P(\mathbf{x}, \mathbf{m}) \\ &= \sum_{\mathbf{m}} P(\mathbf{x}|\mathbf{m})P(\mathbf{m}) \end{aligned}$$

and the prior distribution :

$$P(\mathbf{m}) = \prod_{q=1}^M P(q)^{m_q}$$

similarly,

$$P(\mathbf{x}|\mathbf{m}) = \prod_{q=1}^M \mathcal{N}(\mathbf{x}|\mathbf{w}_q, \sigma_q^2)^{m_q}$$

Hence,

$$P(\mathbf{x}, \mathbf{m}) = \prod_{q=1}^M P(q)^{m_q} \mathcal{N}(\mathbf{x} | \mathbf{w}_q, \sigma_q^2)^{m_q}$$

Therefore,

$$\begin{aligned} \boldsymbol{\theta}^* &= \underset{\{P(q)\}, \{\mathbf{w}_q\}, \{\sigma_q^2\}}{\operatorname{argmax}} \prod_{\alpha=1}^p \prod_{q=1}^M P(q)^{m_q^{(\alpha)}} \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q, \sigma_q^2)^{m_q^{(\alpha)}} \\ &= \underset{\alpha=1}{\operatorname{argmax}} \sum_{\alpha=1}^p \sum_{q=1}^M \ln P(q)^{m_q^{(\alpha)}} \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q, \sigma_q^2)^{m_q^{(\alpha)}} \\ &= \underset{\alpha=1}{\operatorname{argmax}} \sum_{\alpha=1}^p \sum_{q=1}^M m_q^{(\alpha)} \left( \ln P(q) + \ln \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q, \sigma_q^2) \right) \\ &= \underset{\alpha=1}{\operatorname{argmin}} - \sum_{\alpha=1}^p \sum_{q=1}^M m_q^{(\alpha)} \left( \ln P(q) + \ln \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q, \sigma_q^2) \right) \end{aligned}$$

Posterior distribution :

$$P(\mathbf{m} | \mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{m})}{P(\mathbf{x})} = \frac{\prod_{q=1}^M [P(q) \mathcal{N}(\mathbf{x} | \mathbf{w}_q, \sigma_q^2)]^{m_q}}{\sum_{q=1}^M P(q) \mathcal{N}(\mathbf{x} | \mathbf{w}_q, \sigma_q^2)}$$

Hence, posterior distribution of all hidden variables

$$P(\{\mathbf{m}^{(\alpha)}\} | \{\mathbf{x}^{(\alpha)}\}, \boldsymbol{\theta}) = \prod_{\alpha=1}^p \frac{\prod_{q=1}^M [P(q) \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q^{(\alpha)}, \sigma_q^2)]^{m_q^{(\alpha)}}}{\sum_{q=1}^M P(q) \mathcal{N}(\mathbf{x}^{(\alpha)} | \mathbf{w}_q^{(\alpha)}, \sigma_q^2)}$$

Therefore,

$$\langle m_q^{(\alpha)} \rangle_{P(\{\mathbf{m}^{(\alpha)}\} | \{\mathbf{x}^{(\alpha)}\}, \boldsymbol{\theta})} = P(q | \mathbf{x}^{(\alpha)})$$