

Probabilistic Graphical Models in Computer Vision (IN2329)

Csaba Domokos

Summer Semester 2017

4. Mixture of Gaussians & Graph cut

Agenda for today's lecture

Mixture of Gaussians Graph Cut Flow network

In the **previous lecture** we learnt about the *energy function* corresponding to the problem of **binary image segmentation**:

$$E(\mathbf{y}; \mathbf{x}) = \sum_{i \in \mathcal{V}} -\log \phi_i(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}} w_{ij} \mathbb{I}[y_i \neq y_j]$$

Today we are going to learn about

- the definition of the *unary energies* $E_i(y_i; x_i)$. In fact we estimate *mixture of Gaussians* f_{fg} and f_{bg} for the foreground and background, respectively, by making use of the *EM algorithm*.

$$E_i(y_i; x_i) = \begin{cases} -\log f_{bg}(x_i) & \text{if } y_i = 0 \\ -\log f_{fg}(x_i) & \text{otherwise} \end{cases} = \begin{cases} 0 & \text{if } y_i = 0 \\ -\log \frac{f_{fg}(x_i)}{f_{bg}(x_i)} & \text{otherwise} \end{cases}$$

- Graph cuts*, which will be applied to **minimize** the *energy function*.

Mixture of Gaussians

Multivariate Gaussian distribution

Mixture of Gaussians Graph Cut Flow network

Assume a *D-dimensional random vector* $\mathbf{X} = (X_1, \dots, X_D)$, i.e. a vector whose components are random variables, with the joint density function

$$p(x_1, \dots, x_D) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

\mathbf{X} is said to have **multivariate Gaussian (or Normal) distribution** with parameters $\boldsymbol{\mu} \in \mathbb{R}^D$ and $\Sigma \in \mathbb{R}^{D \times D}$ assuming that Σ is *positive definite*.

$\boldsymbol{\mu}$ is called the **mean vector** and Σ is called the **covariance matrix**. We often use the notation $\mathbf{X} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$ denoting \mathbf{X} has Normal distribution.

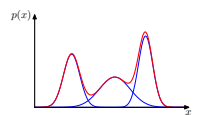
Reminder. A symmetric $\mathbf{A} \in \mathbb{R}^{n \times n}$ matrix is said to be **positive definite**, if $\mathbf{u}^T \mathbf{A} \mathbf{u} > 0$ for all non-zero $\mathbf{u} \in \mathbb{R}^n$.

Mixture of Gaussians

Mixture of Gaussians Graph Cut Flow network

Let us consider a superposition of *K* Gaussian densities

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$$



Mixture of three Gaussians
Source: C. Bishop: PRML, 2006.

which is called a **mixture of Gaussians**.

The parameters π_k are called **mixing coefficients**.

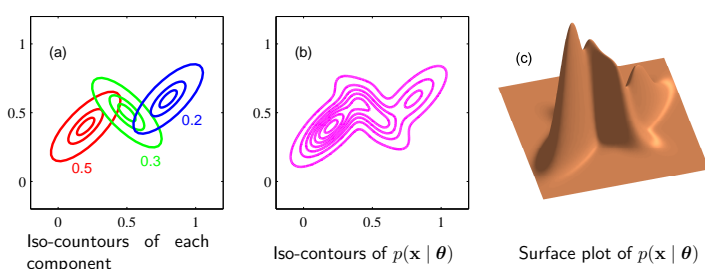
$$1 = \int_{\mathbb{R}^D} p(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^D} \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k) d\mathbf{x} = \sum_{k=1}^K \pi_k$$

All the density functions are non-negative, hence $\pi_k \geq 0$ for $1 \leq k \leq K$, therefore

$$0 \leq \pi_k \leq 1 \quad \text{for all } k = 1, \dots, K$$

Example: Mixture of three 2D Gaussians *

Mixture of Gaussians Graph Cut Flow network



Source: C. Bishop: Pattern Recognition and Machine Learning, 2006.

Parameter estimation *

Mixture of Gaussians Graph Cut Flow network

We are interested in a method to find the **maximum likelihood estimator** of a **parameter** θ of a **probability distribution** $p(\mathbf{x} | \theta)$.

Reminiscent of naming conventions:

$$p(\theta | \mathbf{x}) = \frac{p(\mathbf{x} | \theta)p(\theta)}{p(\mathbf{x})} \propto p(\mathbf{x} | \theta) p(\theta)$$

Posterior probability

Likelihood

Prior probability ↑

We are given finite amount of **measurement** (i.e. observed data) x_1, x_2, \dots , and also know the probability distribution $p(\mathbf{x} | \theta)$. The maximum likelihood estimate of θ is given by

$$\hat{\theta} \in \arg\max_{\theta} p(\mathbf{x} | \theta)$$

A *possible solution*: *Expectation-maximization algorithm*, which iteratively makes guesses about the data \mathbf{x} , and iteratively maximizes $p(\mathbf{x} | \theta)$ over θ .

We introduce a K -dimensional **binary random variable** $z \in \mathbb{B}^K$ having a *1-of- K representation*, i.e. $z_k = 1$ and all other elements are equal to 0. Let us define the *marginal distribution*

$$p(z_k = 1) = \pi_k,$$

which is considered as the *prior probability* of picking the k^{th} component of a mixture of Gaussians. This distribution can be also written as a joint distribution

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}.$$

Moreover, the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} , i.e. the *likelihood*, can be written as

$$p(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \text{thus} \quad p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}.$$

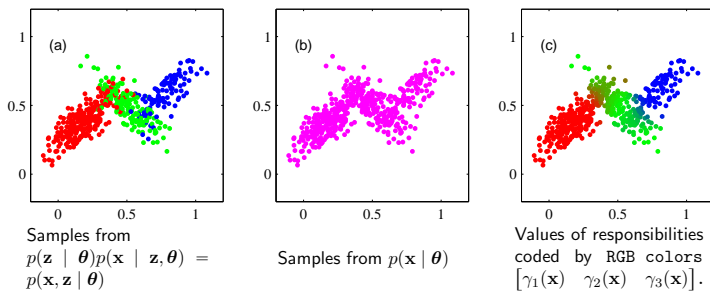
The **distribution of mixture of Gaussian**, specified by the parameter vector $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \in \mathbb{R}^K \times \mathbb{R}^{D \times K} \times \mathbb{R}^{D \times D} \times K$, is given by

$$p(\mathbf{x}) \triangleq p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{z} | \boldsymbol{\theta}) p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}) \\ = \sum_{\mathbf{z}} \prod_{k=1}^K (\pi_k p(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_k} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

The *posterior probabilities* $p(z_k = 1 | \mathbf{x})$, denoted by $\gamma_k(\mathbf{x})$, a.k.a. **responsibilities**, show the probability that a given sample \mathbf{x} belongs to the k^{th} component.

$$\gamma_k(\mathbf{x}) \triangleq p(z_k = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | z_k = 1) p(z_k = 1)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | z_k = 1) p(z_k = 1)}{\sum_{l=1}^K p(z_l = 1) p(\mathbf{x} | z_l = 1)} \\ = \frac{p(z_k = 1) p(\mathbf{x} | z_k = 1)}{\sum_{l=1}^K p(z_l = 1) p(\mathbf{x} | z_l = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}.$$

Example: Mixture of three 2D Gaussians *



Source: C. Bishop: Pattern Recognition and Machine Learning, 2006.

Estimation of a mixture of Gaussians

Suppose we have a set of *i.i.d.* data samples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ drawn from a mixture of Gaussians. The data set is represented by $\mathbf{X} \in \mathbb{R}^{N \times D}$.

The goal is to find the parameter vector $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$, specifying the model from which the samples \mathbf{x}_n have most likely been drawn. We may find the parameters which maximize the *likelihood function* $p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$. To simplify the optimization we use the **log-likelihood function** $\mathcal{L}(\boldsymbol{\theta})$

$$\boldsymbol{\theta}^* \in \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) \stackrel{i.i.d.}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \ln \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta}) \\ = \operatorname{argmax}_{\boldsymbol{\theta}} \ln \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}) p(\mathbf{z}_n | \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \ln \prod_{n=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{nk}} \\ = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)).$$

Note that there is no closed-form solution for this model \Rightarrow iterative solution.

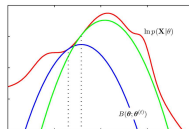
Recall the EM algorithm *

- 1: Choose an initial setting for the parameters $\boldsymbol{\theta}^{(0)}$
- 2: $t \rightarrow 0$
- 3: **repeat**
- 4: $t \rightarrow t + 1$
- 5: **E step.** Evaluate $q^{(t-1)}(\mathbf{Z}) \triangleq p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{(t-1)})$
- 6: **M step.** Evaluate $\boldsymbol{\theta}^{(t)}$ given by

$$\boldsymbol{\theta}^{(t)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}),$$

$$\text{where } Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t-1)}) \triangleq \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) | \mathbf{X}, \boldsymbol{\theta}^{(t-1)}] \\ = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{(t-1)}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$$

- 7: **until** convergence of either the parameters $\boldsymbol{\theta}$ or the log likelihood $\mathcal{L}(\boldsymbol{\theta}; \mathbf{X})$



Source: C. Bishop: PRML, 2006.

E step *

We need to calculate $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$, which is calculated based on $p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}^{\text{old}})$ for all $n = 1, \dots, N$ as follows (see Exercise)

$$p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}^{\text{old}}) = \frac{p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta}^{\text{old}}) p(\mathbf{z}_n | \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{x}_n | \boldsymbol{\theta}^{\text{old}})} \\ = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)} \\ \triangleq \gamma_k(\mathbf{x}_n).$$

Therefore, in the **E step** we need to calculate the *responsibilities* $\gamma_k(\mathbf{x}_n)$ for all data points \mathbf{x}_n and components $k = 1, \dots, K$.

M step for $\boldsymbol{\mu}$ *

We have already known that $z_{nk} = \gamma_k(\mathbf{x}_n)$. Therefore, we may consider

$$\boldsymbol{\theta}^* \in \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \sum_{k=1}^K \gamma_k(\mathbf{x}_n) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad \text{s.t.} \quad \pi_k > 0, \sum_{k=1}^K \pi_k = 1.$$

Setting the derivative of $\mathcal{L}(\boldsymbol{\theta})$ w.r.t. $\boldsymbol{\mu}_k$ to 0, one can obtain that (see Exercise)

$$\frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)} = \boldsymbol{\mu}_k.$$

M step for $\boldsymbol{\Sigma}$ *

$$\boldsymbol{\theta}^* \in \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \sum_{k=1}^K \gamma_k(\mathbf{x}_n) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad \text{s.t.} \quad \pi_k > 0, \sum_{k=1}^K \pi_k = 1.$$

Setting the derivative of $\mathcal{L}(\boldsymbol{\theta})$ w.r.t. $\boldsymbol{\Sigma}_k$ to 0, one can obtain (see Exercise)

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)}.$$

Remark: A $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$ matrix, calculated as

$$\boldsymbol{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T,$$

is called **sample covariance matrix** of data points $\{\mathbf{x}_n \in \mathbb{R}^D\}_{n=1}^N$, where $\boldsymbol{\mu} \in \mathbb{R}^D$ is the **sample mean**.

To integrate the conditions on π we use the **Lagrange multiplier method**

$$\theta^* \in \operatorname{argmax}_{\theta} \sum_{n=1}^N \sum_{k=1}^K \gamma_k(\mathbf{x}_n) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) + \lambda (1 - \sum_{k=1}^K \pi_k).$$

Setting the derivative w.r.t. π_k to 0, we obtain

$$\sum_{n=1}^N \frac{\gamma_k(\mathbf{x}_n)}{\pi_k} - \lambda = 0$$

$$\sum_{n=1}^N \sum_{k=1}^K \gamma_k(\mathbf{x}_n) = \lambda \sum_{k=1}^K \pi_k \Rightarrow N = \lambda$$

therefore

$$\pi_k = \frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)}{N}.$$

- 1: Initialize the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$ and mixing coefficients π_k for all $k = 1, \dots, K$
- 2: **repeat**
- 3: **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma_k(\mathbf{x}_n) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)} \quad \text{for } 1 \leq n \leq N \text{ and } 1 \leq k \leq K.$$

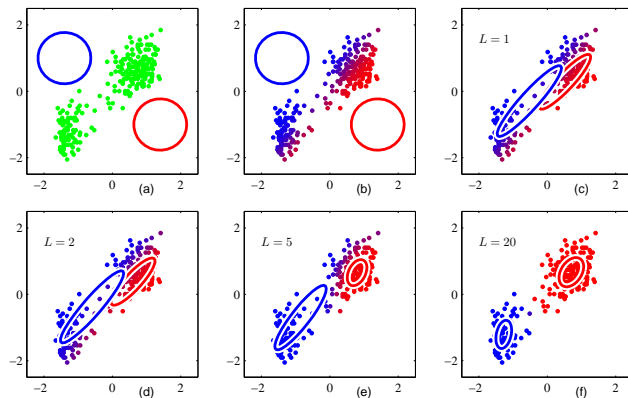
- 4: **M step.** Re-estimate the parameters $(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ for all $k = 1, \dots, K$

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)}, \quad \boldsymbol{\Sigma}_k^{\text{new}} = \frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T}{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)}$$

$$\pi_k^{\text{new}} = \frac{\sum_{n=1}^N \gamma_k(\mathbf{x}_n)}{N}$$

- 5: **until** convergence of either the parameters θ or the log likelihood $\mathcal{L}(\theta)$

Example *



Remarks

- The EM algorithm is **not limited** to mixture of Gaussians, but it can also be applied to *other probability distributions*.
- The algorithm does **not** necessarily yield global maxima. In practice, it is restarted with *different initializations* and after convergence the result with the highest log-likelihood is chosen.
- One can think the EM algorithm as an **alternating minimization** procedure. Considering $f(\theta, q)$ as the objective function, one iteration of the EM algorithm can be reformulated as

$$\text{E-step: } q^{(t+1)} \in \operatorname{argmax}_q f(\theta^{(t)}, q)$$

$$\text{M-step: } \theta^{(t+1)} \in \operatorname{argmax}_{\theta} f(\theta, q^{(t)})$$

Graph Cut

Assume a **weighted directed graph** $G = (\mathcal{V}, \mathcal{E}, c)$

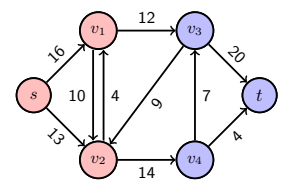
- $\mathcal{V} = \{1, \dots, n\}$ is a finite set of nodes,
- $\mathcal{E} \subseteq \{(i, j) \in \mathcal{V} \times \mathcal{V} \mid i \neq j\}$ is the set of edges,
- $c: \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ is a weight function. (For any $(i, j) \notin \mathcal{E}$, $c(i, j) = 0$.)

A **cut** $(\mathcal{S}, \mathcal{T})$ of G is a *disjoint* partition of \mathcal{V} into \mathcal{S} and $\mathcal{T} = \mathcal{V} \setminus \mathcal{S}$.

The **capacity** of the cut $(\mathcal{S}, \mathcal{T})$ is defined as

$$\text{cut}(\mathcal{S}, \mathcal{T}) = \sum_{(i,j) \in \mathcal{S} \times \mathcal{T}} c(i, j).$$

Assume distinct nodes $s, t \in \mathcal{V}$, a cut $(\mathcal{S}, \mathcal{T})$ is called **$s-t$ cut** if $s \in \mathcal{S}$ and $t \in \mathcal{T}$.



The **minimum $s-t$ cut problem** is to find an $s-t$ cut with the lowest cost.

Example: $\text{cut}(\mathcal{S}, \mathcal{T}) = c(v_1, v_3) + c(v_2, v_4) = 12 + 14 = 26$.

Flow network and flow

Let $G = (\mathcal{V}, \mathcal{E}, c)$ be a *directed weighted graph* with **non-negative** edge weights. Given two distinct nodes, a **source** s and a **sink** t , we call $(\mathcal{V}, \mathcal{E}, c, s, t)$ a **flow network**.

Let $(\mathcal{V}, \mathcal{E}, c, s, t)$ be a *flow network*. A function $f: \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ is called a **flow** if it satisfies the following properties:

1. **Capacity constraint:**

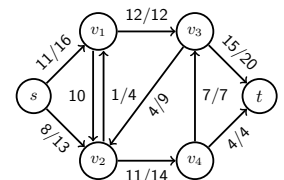
$$f(i, j) \leq c(i, j) \quad \text{for all } i, j \in \mathcal{V}.$$

2. **Skew-symmetry:**

$$f(i, j) = -f(j, i) \quad \text{for all } i, j \in \mathcal{V}.$$

3. **Flow conservation:**

$$\sum_{j \in \mathcal{V}} f(i, j) = 0 \quad \text{for all } i \in \mathcal{V} \setminus \{s, t\}.$$



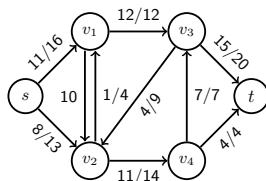
The edges are labeled by $f(i, j)/c(i, j)$.
Only positive $f(i, j)$ are shown.

Flow network

The **value** of a flow f is defined as

$$|f| \triangleq \sum_{(s,i) \in \mathcal{E}} f(s,i) = - \sum_{(i,t) \in \mathcal{E}} f(i,t).$$

The **maximum-flow problem** is to find a flow f with the highest cost for a given flow network G .



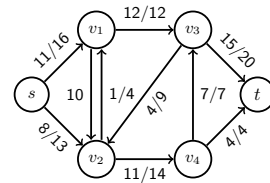
The edges are labeled by $f(i,j)/c(i,j)$.
 $|f| = 19$.

Now we give a *more intuitive definition* of flows. We will see that the previous definition is more helpful for the analysis of the *maximum-flow algorithm*.

Let $(\mathcal{V}, \mathcal{E}, c, s, t)$ be a flow network. A function $f : \mathcal{E} \rightarrow \mathbb{R}^+$ is called a **flow** if it satisfies the following two properties:

1. $f(i,j) \leq c(i,j)$ for all $(i,j) \in \mathcal{E}$.
2. For all $i \in \mathcal{V} \setminus \{s, t\}$

$$\sum_{(i,j) \in \mathcal{E}} f(i,j) = \sum_{(j,i) \in \mathcal{E}} f(j,i).$$



The edges are labeled by $f(i,j)/c(i,j)$.

One can see that the two definitions of the flow are equivalent. (See Exercise)

Working with flows *

Let $G = (\mathcal{V}, \mathcal{E}, c, s, t)$ be a flow network and let f be a flow in G . We will use the following notation for $A, B \subseteq \mathcal{V}$

$$f(A, B) = \sum_{a \in A} \sum_{b \in B} f(a, b).$$

It is easy to see that $|f| = f(\mathcal{V}, \{t\})$, and $f(\{i\}, \mathcal{V}) = 0$ for all $i \in \mathcal{V} \setminus \{s, t\}$ due to *flow conservation*.

Let $G = (\mathcal{V}, \mathcal{E}, c, s, t)$ be a flow network and let f be a flow in G . Then the following equalities hold:

- i) For all $A \subseteq \mathcal{V}$, we have $f(A, A) = 0$.
- ii) For all $A, B \subseteq \mathcal{V}$, we have $f(A, B) = -f(B, A)$.
- iii) For all $A, B, C \subseteq \mathcal{V}$ with $A \cap B = \emptyset$, we have

$$f(A \cup B, C) = f(A, C) + f(B, C) \text{ and } f(C, A \cup B) = f(C, A) + f(C, B).$$

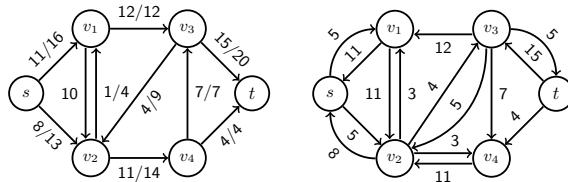
Proof. Exercise. \square

Residual network

Let $G = (\mathcal{V}, \mathcal{E}, c, s, t)$ be a flow network and let f be a flow in G . The *weighted directed graph* $G_f = (\mathcal{V}, \mathcal{E}_f, c_f)$ is called **residual network** of G induced by f , where

$$c_f(i, j) = c(i, j) - f(i, j),$$

$$\mathcal{E}_f = \{(i, j) \in \mathcal{V} \times \mathcal{V} : c_f(i, j) > 0\}.$$



A path p from s to t in G_f is called an **augmenting path**.

Max-flow-min-cut theorem

Let f be a flow in a flow network $G = (\mathcal{V}, \mathcal{E}, c, s, t)$. Then the following conditions are equivalent:

- 1) f is a maximal flow in G .
- 2) The residual graph G_f contains no augmenting paths.
- 3) $|f| = \text{cut}(\mathcal{S}, \mathcal{T})$ for some $s-t$ cut of G .

Proof of Max-flow-min-cut theorem

1) \Rightarrow 2) *

Suppose that f is *maximum flow* in G , but and there exists an *augmenting path* p in the residual graph G_f .

The maximum amount by which we can **increase** the flow in p is the **residual capacity** of p , given by

$$c_f(p) = \min\{c_f(i, j) : (i, j) \text{ is on } p\}.$$

Furthermore, let us define $f_p : \mathcal{E} \rightarrow \mathbb{R}$ as follows:

$$f_p(i, j) = \begin{cases} c_f(p) & \text{if } (i, j) \text{ is on } p \\ -c_f(p) & \text{if } (j, i) \text{ is on } p \\ 0 & \text{otherwise.} \end{cases}$$

One can see that f_p is a flow in G_f with value $|f_p| = c_f(p) > 0$. Therefore the flow $f + f_p$ has the value $|f| + |f_p| > |f|$, which contradicts the optimality of f .

Proof of Max-flow-min-cut theorem

2) \Rightarrow 3) *

Suppose that G_f has no augmenting path, i.e. s and t are disconnected in G_f . Define

$$\mathcal{S} := \{v \in \mathcal{V} : \text{there exists a path from } s \text{ to } v \text{ in } G_f\}.$$

Obviously, $(\mathcal{S}, \mathcal{T})$ is a cut of G , where $\mathcal{T} = \mathcal{V} \setminus \mathcal{S}$.

For each pair of $(i, j) \in \mathcal{S} \times \mathcal{T}$, we have $f(i, j) = c(i, j)$, otherwise $(i, j) \in \mathcal{E}_f$ would be held, which would imply that $j \in \mathcal{S}$.

One can see that the flow across $(\mathcal{S}, \mathcal{T})$ is $|f|$:

$$f(\mathcal{S}, \mathcal{T}) \stackrel{\text{iii)}}{=} f(\mathcal{S}, \mathcal{V}) - f(\mathcal{S}, \mathcal{S}) \stackrel{\text{i)}}{=} f(\mathcal{S}, \mathcal{V}) \stackrel{\text{iii)}}{=} f(\{s\}, \mathcal{V}) + f(\mathcal{S} \setminus \{s\}, \mathcal{V}) = f(\{s\}, \mathcal{V}) = |f|.$$

Therefore $|f| = f(\mathcal{S}, \mathcal{T}) = \text{cut}(\mathcal{S}, \mathcal{T})$.

Proof of Max-flow-min-cut theorem

3) \Rightarrow 1) *

Let f be a flow in G such that $|f| = \text{cut}(\mathcal{S}, \mathcal{T})$. In general, for **any** flow f in G the following holds:

$$|f| = f(\mathcal{S}, \mathcal{T}) = \sum_{i \in \mathcal{S}} \sum_{j \in \mathcal{T}} f(i, j) \leq \sum_{i \in \mathcal{S}} \sum_{j \in \mathcal{T}} c(i, j) = \text{cut}(\mathcal{S}, \mathcal{T}).$$

Hence $|f| = \text{cut}(\mathcal{S}, \mathcal{T})$ is maximal (equivalently $\text{cut}(\mathcal{S}, \mathcal{T})$ is minimal). \square

Ford-Fulkerson algorithm *

Mixture of Gaussians Graph Cut Flow network

Input: A flow network $G = (\mathcal{V}, \mathcal{E}, c, s, t)$

Output: A minimum $s - t$ cut $(\mathcal{S}, \mathcal{T})$ of G

```

1: for all  $(i, j) \in \mathcal{E}$  do
2:    $f(i, j) \leftarrow 0$  and  $f(j, i) \leftarrow 0$ 
3: end for
4: while there exists a path  $p$  from  $s$  to  $t$  in the residual network  $G_f$  do
5:    $c_f(p) \leftarrow \min\{c_f(i, j) : (i, j) \text{ is in } p\}$ 
6:   for all  $(i, j)$  in  $p$  do
7:      $f(i, j) \leftarrow f(i, j) + c_f(p)$ 
8:      $f(j, i) \leftarrow -f(i, j)$ 
9:   end for
10: end while
11:  $\mathcal{S} \leftarrow \{v \in \mathcal{V} : \text{there exists a path from } s \text{ to } v \text{ in } G_f\}$  and  $\mathcal{T} \leftarrow \mathcal{V} \setminus \mathcal{S}$ 

```

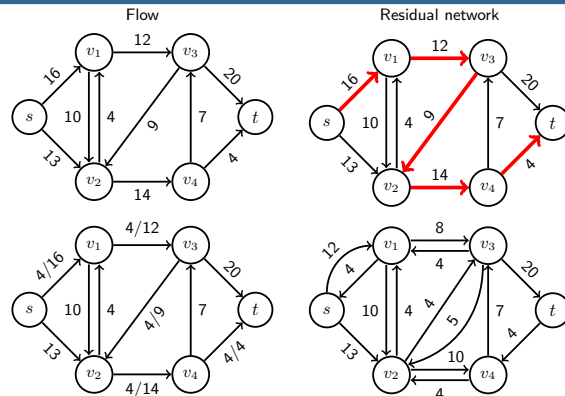
The complexity of this algorithm is $\mathcal{O}(|\mathcal{E}| \cdot \|f^*\|)$, where f^* is the value of the maximal flow.

IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 33 / 41

Example: iteration 1 *

Mixture of Gaussians Graph Cut Flow network

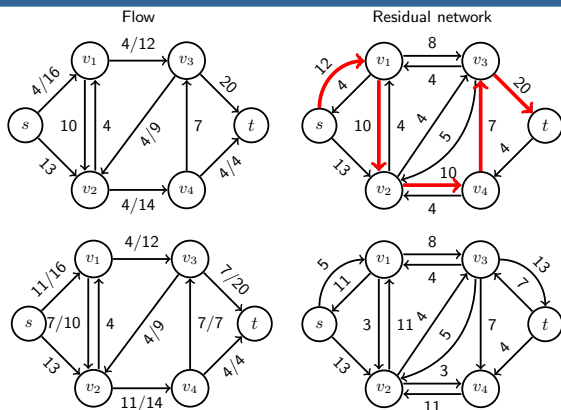


IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 34 / 41

Example: iteration 2 *

Mixture of Gaussians Graph Cut Flow network

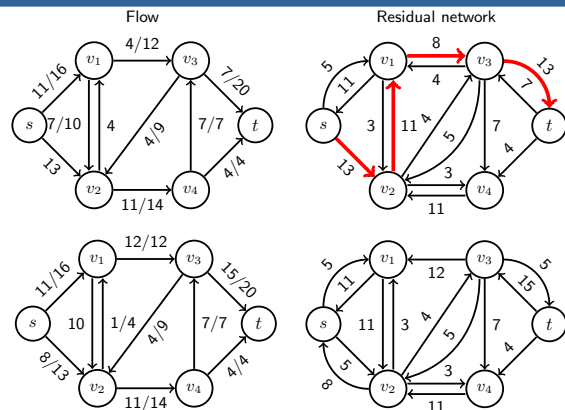


IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 35 / 41

Example: iteration 3 *

Mixture of Gaussians Graph Cut Flow network

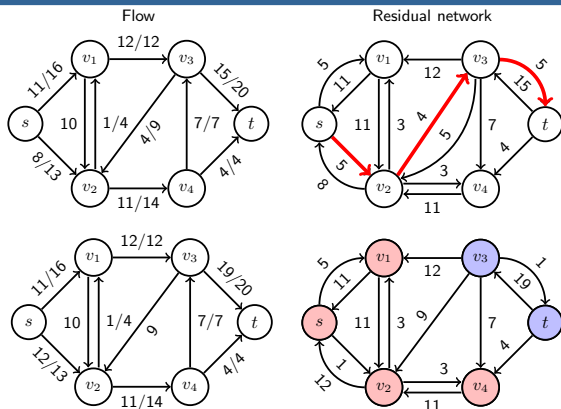


IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 36 / 41

Example: iteration 4 *

Mixture of Gaussians Graph Cut Flow network

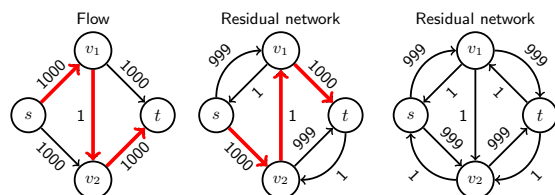


IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 37 / 41

A "bad" example *

Mixture of Gaussians Graph Cut Flow network



Note that there exists an example, where the flow, computed by the Ford-Fulkerson algorithm, does not even converge to the maximum flow.

More precisely, if a flow network has integer (\mathbb{N}_0) or rational (\mathbb{Q}_0^+) capacities, then the Ford-Fulkerson algorithm terminates and it computes a maximum flow.

IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 38 / 41

Edmonds-Karp algorithm

Mixture of Gaussians Graph Cut Flow network

Input: A flow network $G = (\mathcal{V}, \mathcal{E}, c, s, t)$

Output: A minimum $s - t$ cut $(\mathcal{S}, \mathcal{T})$ of G

```

1: for all  $(i, j) \in \mathcal{E}$  do
2:    $f(i, j) \leftarrow 0$  and  $f(j, i) \leftarrow 0$ 
3: end for
4: while there exists a path  $p$  from  $s$  to  $t$  in the residual network  $G_f$  do
5:    $p \leftarrow \text{shortestPath}(G_f, s, t)$ 
6:    $c_f(p) \leftarrow \min\{c_f(i, j) : (i, j) \text{ is in } p\}$ 
7:   for all  $(i, j)$  in  $p$  do
8:      $f(i, j) \leftarrow f(i, j) + c_f(p)$ 
9:      $f(j, i) \leftarrow -f(i, j)$ 
10:  end for
11: end while
12:  $\mathcal{S} \leftarrow \{v \in \mathcal{V} : \text{there exists a path from } s \text{ to } v \text{ in } G_f\}$  and  $\mathcal{T} \leftarrow \mathcal{V} \setminus \mathcal{S}$ 

```

The complexity of this algorithm is $\mathcal{O}(|\mathcal{V}| \cdot |\mathcal{E}|^2)$. There exist more efficient algorithms for maximum flow calculation with complexity $\mathcal{O}(|\mathcal{V}|^2 \cdot |\mathcal{E}|)$ and $\mathcal{O}(|\mathcal{V}|^3)$.

IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 39 / 41

Summary *

Mixture of Gaussians Graph Cut Flow network

- **Max-flow-min-cut theorem** tells us that the minimum cut problem can be solved via maximum flow. These two problems are dual to each other, moreover strong duality holds.
- **Edmonds-Karp algorithm:** The Ford-Fulkerson algorithm becomes polynomial, if the shortest path is used as augmented path.

In the **next lecture** we will learn about

- Exact solution for *binary image segmentation* via *graph cut*
- *Boykov-Kolmogorov algorithm*
- Multi-label problem (e.g., stereo matching)

IN2329 - Probabilistic Graphical Models in Computer Vision

4. Mixture of Gaussians & Graph cut - 40 / 41

The EM algorithm for mixture of Gaussians

1. Frank Dellaert. The expectation maximization algorithm. Technical Report GIT-GVU-02-20, Georgia Institute of Technology, Atlanta, GA, USA, 2002
2. Shane M. Haas. The expectation-maximization and alternating minimization algorithms. Unpublished, 2002
3. Christopher Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006

Graph cut, Maximum flow

4. Ronald L. Rivest Thomas H. Cormen, Charles E. Leiserson and Clifford Stein. *Introduction to Algorithms*. MIT Press, 2009
5. Tanmay Verma and Dhruv Batra. MaxFlow revisited: An empirical comparison of MaxFlow algorithms for dense vision problems. In Richard Bowden, John Collomosse, and Krystian Mikolajczyk, editors, *In Proceedings of British Machine Vision Conference*, pages 61.1–61.12, Surrey, UK, September 2012. BMVA Press