# Probabilistic Graphical Models

## Lecture 10 – Undirected Models

CS/CNS/EE 155

Andreas Krause

# Announcements

- Homework 2 due this Wednesday (Nov 4) in class
- Project milestones due next Monday (Nov 9)
  - About half the work should be done
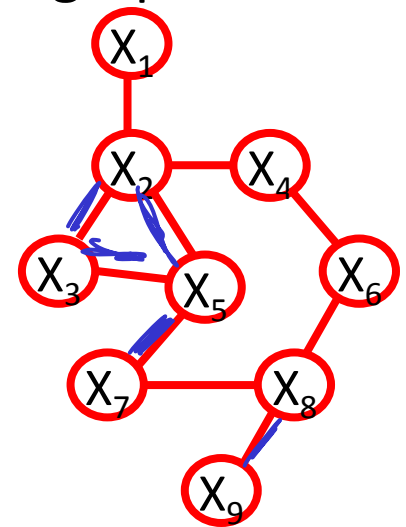  - 4 pages of writeup, NIPS format
  - http://nips.cc/PaperInformation/StyleFiles

# Markov Networks
## (a.k.a., Markov Random Field, Gibbs Distribution, …)

- A Markov Network consists of
  - An undirected graph, where each node represents a RV
  - A collection of factors defined over cliques in the graph
- Joint probability

$$P(x) = \frac{1}{Z} \prod_i \Psi_i(C_i)$$

$P$

- A distribution factorizes over undirected graph G if

$\exists$ factors $\Psi_1 \cdots \Psi_k$ over cliques of G s.t.

$$P(x) = \frac{1}{Z} \prod_i \Psi_i(C_i)$$

# Computing Joint Probabilities

- Computing joint probabilities in BNs

$$P(X_1, \ldots, X_m) = \prod_i P(X_i \mid Pa_i)$$

$$P(X_1 \mid X_m)$$
actually comp. $P(X_1, X_m)$

$$Z = \sum_X \prod_i \psi_i(C_i)$$

- Computing joint probabilities in Markov Nets
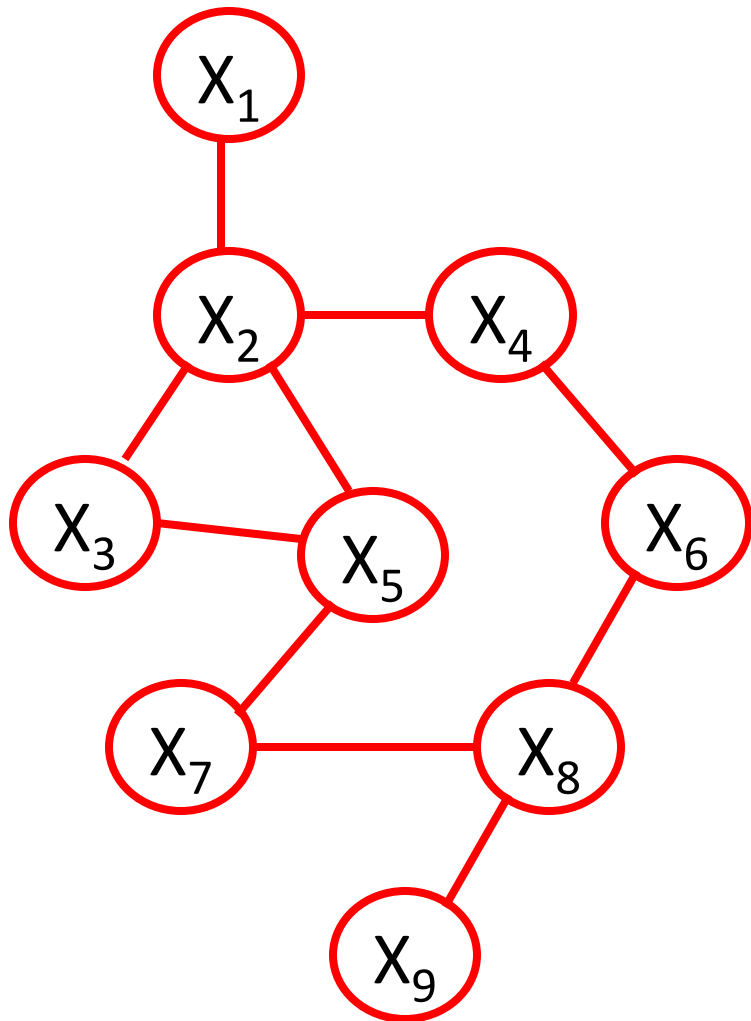
$$P(X_1 \ldots X_m) = \frac{1}{Z} \prod_i \psi_i(C_i)$$

Can do via variable elimination

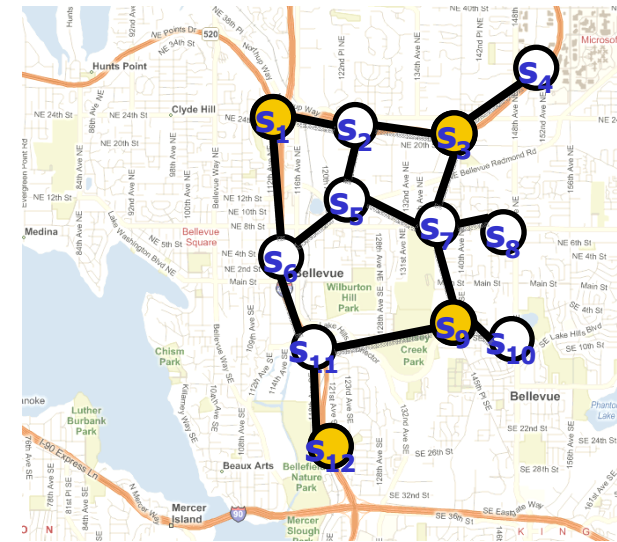$\uparrow$ Need to know partition "function" $Z$

$$\text{Can compute } \frac{P(X_1 \ldots X_m)}{P(X_1' \ldots X_m')} = \frac{\prod_i \psi_i(C_i)}{\prod_i \psi_i(C_i')}$$
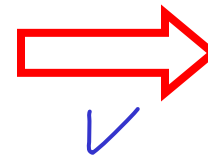
# Local Markov Assumption for MN



- The **Markov Blanket MB(X)** of a node X is the set of neighbors of X

- Local Markov Assumption:
  $X \perp \text{EverythingElse} \mid MB(X)$

- $I_{loc}(G)$ = set of all local independences

- G is called an I-map of distribution P if $I_{loc}(G) \subseteq I(P)$

# Factorization Theorem for Markov Nets "➔"



True distribution P
can be represented exactly as
a Markov net (G,P)

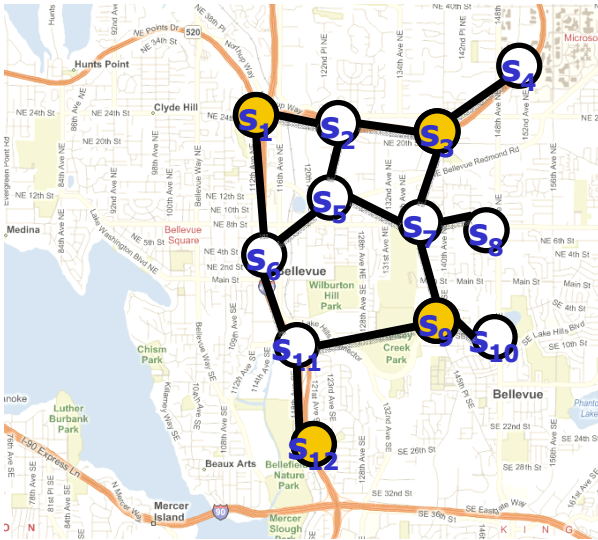$$P(X_1, ..., X_n) = \frac{1}{Z} \prod_i \phi_i(\mathbf{C}_i)$$

$I_{loc}(G) \subseteq I(P)$

G is an **I-map** of P
(independence map)

$$I_{loc}(G) \subseteq I(P)$$

G is an **I-map** of P
(independence map)
**and** P>0

True distribution P
can be represented exactly as

$$P(X_1, ..., X_n) = \prod_i P(X_i \mid \mathbf{Pa}_{X_i})$$

i.e., P can be represented as
a Markov net (G,P)

# Global independencies



- A trail $X—X_1—...—X_m—Y$ is called active for evidence $E$, if none of $X_1,...,X_m \in E$

- Variables X and Y are called **separated** by $E$ if there is no active trail for $E$ connecting X, Y Write sep(X,Y | $E$)

- $I(G) = \{X \perp Y \mid E: sep(X,Y|E)\}$

# Soundness of separation

- Know: For positive distributions P>0

$$I_{loc}(G) \subseteq I(P) \Leftrightarrow P \text{ factorizes over G}$$

- **Theorem**: Soundness of separation

  For positive distributions P>0

$$I_{loc}(G) \subseteq I(P) \Leftrightarrow I(G) \subseteq I(P)$$

- Hence, separation captures only true independences

- How about $I(G) = I(P)$?

# Completeness of separation

**Theorem**: Completeness of separation
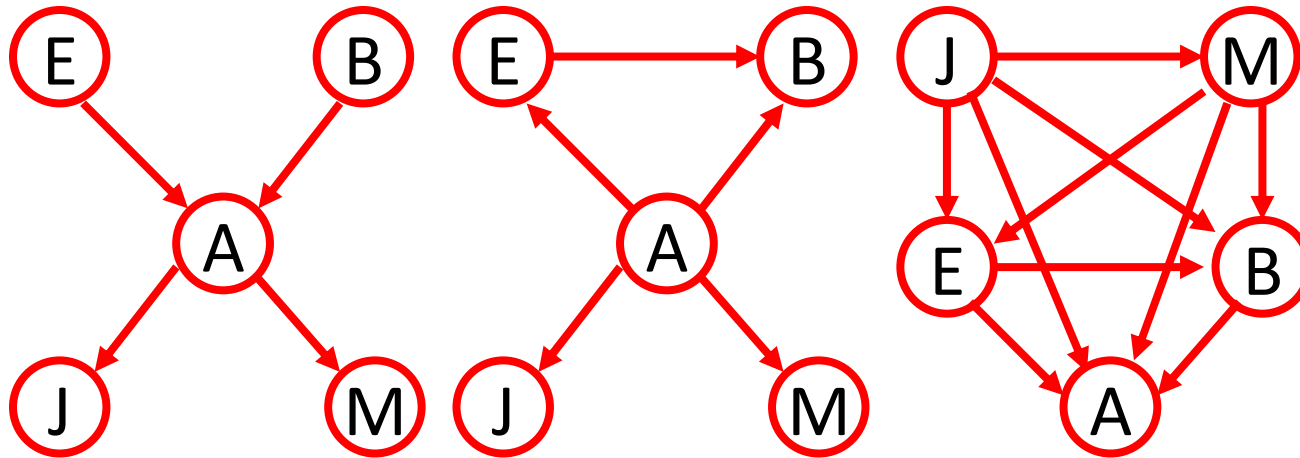
$$I(G) = I(P)$$

for "almost all" distributions P that factorize over G

"almost all": Except for of potential parameterizations of measure 0 (assuming no finite set have positive measure)
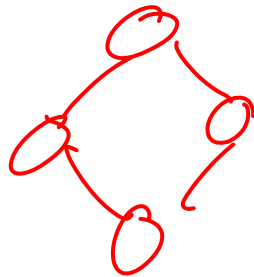
- For BNs: Minimal I-map not unique



- For MNs: For positive P, minimal I-map is unique!!

# P-maps

- Do P-maps always exist?

- For BNs: no

- How about Markov Nets?

does not have
M N P-map !

# Exact inference in MNs

- Variable elimination and junction tree inference work exactly the same way!
  - Need to construct junction trees by obtaining chordal graph through triangulation

# Pairwise MNs

- A pairwise MN is a MN where all factors are defined over single variables or pairs of variables
- Can reduce any MN to pairwise MN!

# Logarithmic representation

- Can represent any positive distribution in log domain

$$P(x) = \frac{1}{Z} \prod_i \psi_i(C_i)$$

$$\log P(x) = \sum_i \underbrace{\log \psi_i(C_i)}_{\varphi_i(C_i)} - \log Z$$

$$P(x) = \frac{1}{Z} \exp\left( \sum_i \varphi_i(C_i) \right)$$

# Log-linear models

- Feature functions $\phi_i(D)$ defined over cliques

$$\phi_i(X_i, X_{i+1}) = \begin{cases} 1 & \text{if } X_i = X_{i+1} \\ 0 & \text{otherwise} \end{cases}$$

- Log linear model over undirected graph G
  - Feature functions $\phi_1(D_1),...,\phi_k(D_k)$
  - Domains $D_i$ can overlap
  - Set of weights $w_i$ learnt from data

$$P(X) = \frac{1}{Z} \exp\left( \sum_i w_i^T \phi_i(C_i) \right)$$

**Theorem**: Moralized Bayes net is minimal Markov I-map

Resulting BN has far fewer
cond. independencies than original
MN

$$I(G') \subseteq I(G)$$

$$\underset{BN}{\uparrow} \qquad \underset{MN}{\uparrow}$$

**Theorem**: Minimal Bayes I-map for MN must be chordal

# So far

- Markov Network **Representation**

  - Local/Global Markov assumptions; Separation

  - Soundness and completeness of separation

- Markov Network **Inference**

  - Variable elimination and Junction Tree inference work exactly as in Bayes Nets

- How about **Learning** Markov Nets?

$$\log P(D \mid \theta) = \log \prod_\ell \prod_i P(X_i^{(\ell)} \mid Pa_i^{(\ell)} ; \theta)$$

Parameter indepent

$$= \sum_\ell \sum_i \log P(X_i^{(\ell)} \mid Pa_i^{(\ell)} ; \theta_{X_i (Pa_i)})$$

$$\frac{\partial}{\partial \theta_{X_i \mid Pa_i}} \log P(D \mid \theta) = \sum_j \sum_\ell \frac{\partial}{\partial \theta_{X_i (Pa_i)}} \log P(X_j^{(\ell)} \mid Pa_j^{(\ell)} ; \theta_{X_j (Pa_j)})$$

$$= \sum_\ell \frac{\partial}{\partial \theta_{X_i (Pa_i)}} \log P(X_i^{(\ell)} \mid Pa_i , \theta_{X_i (Pa_i)}) \stackrel{!}{=} 0$$

Problem breaks down into independent subproblems

Learn every CPD independent of others

# Algorithm for BN MLE

Given BN structure $G$

For each variable $X_i$,

learn $\hat{\Theta}_{X_i | Pa_i} = \dfrac{Count(X_i, Pa_i)}{Count(Pa_i)}$

$\Rightarrow$ globally maximum likelihood estimate for fixed structure $G$

# MLE for Markov Nets

- Log likelihood of the data

$$\log P(D \mid \theta) = \sum_{\ell} \log P(x^{(\ell)} \mid \theta)$$

$$= \sum_{\ell=1}^{m} \log \frac{1}{Z} \prod_i \psi_i \left( C_i^{(\ell)} \right)$$

$$= \sum_{\ell} \sum_i \log \psi_i \left( C_i^{(\ell)} \right) - m \log Z$$

$$= m \sum_i \sum_{C_i} \hat{P}(C_i) \underbrace{\log \psi_i (C_i)}_{\log P(X_i \mid Pa_i)} - m \underbrace{\log Z}_{\text{not in B}}$$
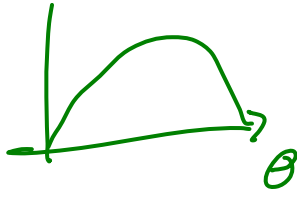
in BN

$$Z = Z(\theta) = \sum_X \prod_i \psi_i (C_i) \qquad \log Z = \log \sum_X \prod_i \psi_i (C_i)$$

avg $\theta$

# Log-likelihood doesn't decompose

- Log likelihood

$$\log P(\mathcal{D} \mid \theta) = m \underbrace{\sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i)}_{\text{decomposes nicely}} - \underbrace{m \log Z(\theta)}_{\text{does \underline{not} decompose}}$$

- l(D | θ) is concave function! *log P(D|θ)*

  No local optima!

  Gradient ascent won't get stuck!  θ

- Log Partition function log Z(θ) doesn't decompose

$$\log P(\mathcal{D} \mid \theta) = m \sum_i \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \log \psi_i(\mathbf{c}_i) - m \log Z(\theta)$$

$$\frac{\partial \log P(D \mid \theta)}{\partial \psi_i(c_i)} = m \sum_j \sum_{c_j} \hat{P}(c_j) \frac{\partial}{\partial \psi_i(c_i)} \log \psi_j(c_j) - m \frac{\partial}{\partial \psi_i(c_i)} \log Z(\theta)$$

$$= m \quad \hat{P}(c_i) \frac{1}{\psi_i(c_i)} - m \frac{\partial}{\partial \psi_i(c_i)} \log Z(\theta)$$

# Derivative of log-likelihood

$$\frac{\partial \log P(\mathcal{D} \mid \theta)}{\partial \psi_i(\mathbf{c}_i)} = m\frac{\hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - m\frac{\partial \log Z(\theta)}{\partial \psi(\mathbf{c}_i)}$$

$\psi_i:$

| A | B | $\psi_i(A,B)$ |
|---|---|---|
| 0 | 0 | $\psi_i(0,0)$ |
| 0 | 1 | $\psi_i(0,1)$ |
| 1 | 0 | |
| 1 | 1 | |

$$\frac{\partial \log Z(\theta)}{\partial \psi_i(c_i)} = \frac{\frac{\partial}{\partial \psi_i(c_i)} Z(\theta)}{Z(\theta)} = \frac{Z\, P(c_i \mid \theta)}{Z\, \psi_i(c_i)} = \frac{P(c_i \mid \theta)}{\psi_i(c_i)}$$

$$\frac{\partial Z(\theta)}{\partial \psi_i(c_i)} = \frac{\partial}{\partial \psi_i} \sum_x \prod_j \psi_j(c_j) = \sum_x \frac{\partial}{\partial \psi_i} \prod_j \psi_j(c_j)$$

$$\frac{}{Z(\theta)}$$

0 if $x \not\sim c_i$
$x$ inconsistent w $c_i$

$$= \sum_{x \sim c_i} \prod_{j \neq i} \psi_j(c_j) \frac{\psi_i(c_i)}{\psi_i(c_i)}$$

$$= \frac{Z\, P(c_i \mid \theta)}{\psi_i(c_i)}$$

25

# Computing the derivative

- Derivative

$$\frac{\partial \log P(\mathcal{D} \mid \theta)}{\partial \psi_i(\mathbf{c}_i)} = m \frac{\hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - m \frac{P(\mathbf{c}_i \mid \theta)}{\psi(\mathbf{c}_i)}$$

$$\frac{\partial \log P(D \mid \theta)}{\partial \psi_G(G=1, D=0, I=1)} = m \frac{\hat{P}(1,0,1)}{\psi_i(1,0,1)} - m \frac{P(1,0,1 \mid \theta)}{\psi_i(1,0,1)}$$



Can do this using VE
Junction tree...

- Computing P(c$_i$ | $\theta$) requires inference!

- Can optimize using conjugate gradient etc.

- At optimum, it must hold that

$$\frac{\partial \log P(\mathcal{D} \mid \theta)}{\partial \psi_i(\mathbf{c}_i)} = m \frac{\hat{P}(\mathbf{c}_i)}{\psi_i(\mathbf{c}_i)} - m \frac{P(\mathbf{c}_i \mid \theta)}{\psi_i(\mathbf{c}_i)} = 0$$

At opt.: $\dfrac{\hat{P}(c_i)}{\psi_i(c_i)} = \dfrac{P(c_i|\theta)}{\psi_i(c_i)}$    "Data agrees with model on marginals"

➜ Solve fixed point equation    $\psi_i^{(0)}(c_i) = 1$

$$\psi_i^{(t+1)}(c_i) = \psi_i^{(t)}(c_i) \cdot \frac{\hat{P}(c_i)}{P(c_i|\theta)}$$

- Must recompute parameters every iteration

$$P(c_i|\theta)$$

# Parameter learning for log-linear models

- Feature functions $\phi_i(C_i)$ defined over cliques

- Log linear model over undirected graph G
  - Feature functions $\phi_1(C_1),\ldots,\phi_k(C_k)$
  - Domains $C_i$ can overlap
- Joint distribution

$$P(X_1,\ldots,X_n) = \frac{1}{Z}\exp\left(\sum_i w_i^T \phi_i(C_i)\right)$$

- How do we get weights $w_i$?

$$\frac{\partial \log P(\mathcal{D} \mid \theta)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \frac{\partial w_i^T \phi_i(\mathbf{c}_i)}{\partial w_i} - m \frac{\partial \log Z(w)}{\partial w_i}$$

$$= \quad m \sum_{c_i} \hat{P}(c_i) \, \phi_i(c_i) \quad - \quad m \frac{\partial \log Z(w)}{\partial w_i}$$

$$\hat{E}[\phi_i]$$

$$\text{If } \hat{\phi}_i(X_i, X_{i+1}) = \begin{cases} 1 & \text{if } X_i = X_{i+1} \\ 0 & \text{otherwise} \end{cases}, \quad \hat{E}[\phi_i] = \frac{Count(X_i = X_{i+1})}{m}$$

$$\frac{\partial \log P(\mathcal{D} \mid \theta)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i)\phi_i(\mathbf{c}_i) - m \frac{\partial \log Z(w)}{\partial w_i}$$

$$\frac{\partial}{\partial w_i} \log Z(w) = \frac{1}{Z(w)} \frac{\partial}{\partial w_i} \sum_x \exp\left( \sum_i w_i^\top \phi_i(c_i) \right)$$

$$= \frac{1}{Z(w)} \sum_x \phi_i(c_i) \exp\left( \sum_i w_i^\top \phi_i(c_i) \right)$$

$$= \sum_x \phi_i(c_i) \underbrace{\frac{1}{Z} \exp\left( \sum_i w_i^\top \phi_i(c_i) \right)}_{P(x)}$$

$$= \sum_{c_i} \phi_i(c_i) P(c_i \mid w)$$

$$= \mathbb{E}_w(\phi_i)$$

# Optimizing parameters

- Gradient of log-likelihood

$$\frac{\partial \log P(\mathcal{D} \mid w)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i)\phi_i(\mathbf{c}_i) - m \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w)\phi_i(\mathbf{c}_i)$$

$$\underbrace{\phantom{m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i)\phi_i(\mathbf{c}_i)}}_{\hat{\mathbb{E}}(\phi_i)} \qquad \underbrace{\phantom{m \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w)\phi_i(\mathbf{c}_i)}}_{\mathbb{E}_w(\phi_i)}$$

- Thus, w is MLE $\Leftrightarrow$ $\hat{\mathbb{E}}[\phi_i] = \mathbb{E}_w[\phi_i]$

# Regularization of parameters

- Put prior on parameters w $\quad P(w)$

$$\frac{\partial \log P(\mathcal{D} \mid w) P(w)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i)\phi_i(\mathbf{c}_i) - m \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w)\phi_i(\mathbf{c}_i) + \frac{\partial \log P(w)}{\partial w_i}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{Last slide}} \qquad \underbrace{\quad}_{(*)}$$

Prior: $\quad P(w) = \mathcal{N}(w; 0, I) \propto \exp\left(-\sum_i w_i^2\right)$

$(*) \quad \log P(w) = \;\sim\; -\sum_i w_i^2$

$\frac{\partial}{\partial w_i} \log P(w) = -2 w_i$

32

# Summary: Parameter learning in MN

- MLE in BN is easy (score decomposes)

- MLE in MN requires inference (score doesn't decompose)

- Can optimize using gradient ascent or IPF

# Tasks

- Read Koller & Friedman Chapters 20.1-20.3, 4.6.1