#### Example: Hidden CRF for Gesture Recognition

 Sy Bor Wang, A. Quattoni, L. . -P. Morency, D. Demirdjian and T. Darrell, "Hidden Conditional Random Fields for Gesture Recognition," 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), 2006, pp. 1521-1527, doi: 10.1109/CVPR.2006.132.



Figure 1. Illustrations of the six gesture classes for the experiments. Below each image is the abbreviation for the gesture class. These gesture classes are: FB - Flip Back, SV - Shrink Vertically, EV - Expand Vertically, DB - Double Back, PB - Point and Back, EH - Expand Horizontally. The green arrows are the motion trajectory of the fingertip and the numbers next to the arrows symbolize the order of these arrows.



Figure 5. Sample image sequence with the estimated body pose superimposed on the user in each frame.

### HCRF Example



Figure 2. HMM model



Figure 3. CRF Model



Figure 4. HCRF Model

Models	Accuracy (%)
HCRF $\omega = 0$	86.44
HCRF $\omega = 1$	96.81
HCRF $\omega = 2$	97.75

top Hidden Markov Model for each gesture.

middle Conditional Random Field: For each frame, the gesture type Y is estimated, the most frequent one is used for classification.

• Different window size  $\omega$  is used.

bottom 12 hidden states for each frame are learned; the overall classification is base on these estimated states.

#### Potentials

Ψ

$$(y, \mathbf{s}, \mathbf{x}; \theta, \omega) = \sum_{j=1} \varphi(\mathbf{x}, j, \omega) \cdot \theta_s[s_j] + \sum_{j=1} \theta_y[y, s_j] + \sum_{(j,k)\in E} \theta_e[y, s_j, s_k]$$
(3)

Models	Accuracy (%)
HMM $\omega = 0$	84.22
$CRF \omega = 0$	86.03
$CRF \omega = 1$	81.75
HCRF (one-vs-all) $\omega = 0$	87.49
HCRF (multi-class) $\omega = 0$	91.64
HCRF (multi-class) $\omega = 1$	93.81

### Undirected (Pairwise, Continuous) Graphical Models

- The generative model represents the full probability distribution P(X).
- Missing edges represent conditional independence of the variables.
- Chapter 17 Elements of Statistical Learning
  - Gaussian graphical models
  - Markov random fields
  - Ising model, (restricted) Boltzmann machine other not mentioned

Undirected Graphical Models 11

- Bayesian networks, Mixed interaction models, ...
- Cytometry dataset (ESLII)
- *N* = 7466 cells
- p = 11 proteins
- We ame to model protein co-occurence probability.

Machine Learning

sklearn.covariance.GraphicalLasso # basics

gRbase # the recommended R package



### Sparse Conditional Gaussian Graphical Model Application

Yin, Jianxin & Li, Hongzhe. (2011). A sparse conditional Gaussian graphical model for analysis of genetical genomics data. The annals of applied statistics. 5. 2630-2650.

- Cytometry dataset (ESLII)
- $p_Y = 54$  gene level expressions
- $p_X = 188$  markers (discrete)
- $Y^{p_Y}|X^{p_X} \sim \mathcal{N}(M^{p_Y \times p_X} X^{p_X}, \Sigma^{p_Y \times p_Y})$  conditional Gaussian distribution
- Top: Black color indicates significant association p value < 0.01 in the linear regression.
- Bottom: The undirected graph of 43 genes constructed on the cGGM.



Data: carcass #Source: Soren Hojsgaard, David Edwards, Steffen Lauritzen: *Graphical Models with R*, Springer.

		mean.					
	Fat11	16.00	_				
	Meat11	52.00					
	Fat12	14.00					
	Meat12	52.00					
	Fat13	13.00					
	Meat13	56.00					
	LeanMeat	59.00					
:	Σ	Fat11	Meat11	Fat12	Meat12	Fat13	Meat13
	Fat11	11.34	0.74	8.42	2.06	7.66	-0.76
	Meat11	0.74	32.97	0.67	35 94	2 01	31.97
		0.1.1	02.01	0.01	55.51	2.01	01.01
	Fat12	8.42	0.67	8.91	0.31	6.84	-0.60
	Fat12 Meat12	8.42 2.06	0.67 35.94	8.91 0.31	0.31 51.79	6.84 2.18	-0.60 41.47
	Fat12 Meat12 Fat13	8.42 2.06 7.66	0.67 35.94 2.01	8.91 0.31 6.84	0.31 51.79 2.18	6.84 2.18 7.62	-0.60 41.47 0.38
	Fat12 Meat12 Fat13 Meat13	8.42 2.06 7.66 -0.76	0.67 35.94 2.01 31.97	8.91 0.31 6.84 -0.60	0.31 51.79 2.18 41.47	6.84 2.18 7.62 0.38	-0.60 41.47 0.38 41.44

LeanMeat -9.08 5.33 -7.95 6.03 -6.93 7.23 12.90

### Gaussian Graphical Models (Undirected Graphs)

• Multivariate Gaussian Distribution on variables  $X = (X_1, \ldots, X_p)$ 

• 
$$\phi(\mathbf{x}) = \frac{1}{\sqrt{|2\pi\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\mu)\Sigma^{-1}(\mathbf{x}-\mu)}$$

- |.| is the determinant. we denote *p* the number of components in **x**. Then  $|2\pi\Sigma| = (2\pi)^p |\Sigma|$ .
- If Σ is not invertible it has dependent columns. It means that the variables x<sub>j</sub> are lineary dependent.
  - If the rank of  $\Sigma$  is  $\ell$  then there exists a matrix A and a vector  $\nu$  so:
  - $x = Az + \nu$  for new coordinates z with  $\ell$  dimensions
  - $\bullet\,$  We just consider the new coordinates and assume  $\Sigma$  has a full rank.



#### • Concentration (Precision, koncentrační) matrix

 $K = \Sigma^{-1}$ 

#### Lemma

For  $u \neq v$ ,  $k_{uv} = 0$  if and only if  $y_u$  and  $y_v$  are conditionally independent given all other variables.

k*100	Fat11	Meat11	Fat12	Meat12	Fat13	Meat13	LeanMeat
Fat11	44	3	-20	-7	-16	4	10
Meat11	3	16	-3	-6	-6	-6	-3
Fat12	-20	-3	54	6	-21	-5	9
Meat12	-7	-6	6	14	-1	-9	-0
Fat13	-16	-6	-21	-1	56	3	7
Meat13	4	-6	-5	-9	3	16	-1
LeanMeat	10	-3	9	-0	7	-1	26

• If looking for small values better to 'scale' the entries into Partial Correlation matrix.

#### Definition (Partial correlation matrix)

Partial correlation matrix is defined from K by

$$\rho_{uv|V\setminus\{uv\}} = \frac{-k_{uv}}{\sqrt{k_{uu}k_{vv}}}.$$

Compare to Pearson correlation  $\frac{cov(u,v)}{\sigma_u\sigma_v}$ 

#### Lemma

In contrast to concentrations, the partial correlations are invariant under a change of scale and origin in the sense that if  $X_j^* = a_j X_j + b_j$ , j = 1, ..., p then  $a_v a_u k_{uv}^* = k_{uv}$  and  $\rho_{uv|V \setminus \{uv\}}^* = \rho_{uv|V \setminus \{uv\}}$ .

$\rho * 100$	Fat11	Meat11	Fat12	Meat12	Fat13	Meat13	LeanMeat
Fat11	-	-11	41	30	32	-16	-29
Meat11	-11	-	9	41	19	35	16
Fat12	41	9	-	-24	38	18	-24
Meat12	30	41	-24	-	2	61	2
Fat13	32	19	38	2	-	-9	-18
Meat13	-16	35	18	61	-9	-	7
LeanMeat	-29	16	-24	2	-18	7	-

#### Models

• The simplest model just removes edges with small  $|\rho_{uv|V\setminus\{uv\}}|$ . Penalized criteria will be introduced later.



### Undirected Gaussian graphical model

#### Definition (Undirected Gaussian graphical model)

An **undirected Gaussian graphical model** is represented by an undirected graph  $\mathcal{G} = (X, E), X = \{X_1, \dots, X_p\}$  represent the set of variables and E is a set of undirected edges.

When a random vector **x** follows a Gaussian distribution  $N_p(\mu, \Sigma)$ , the graph *G* represents the model where  $K = \Sigma^{-1}$  is a positive definite matrix with  $k_{u,v} = 0$  whenever there is no edge between vertices u, v in *G*. This graph is called the **dependence graph** of the model.

#### Lemma

For any non adjacent vertices  $u, v \in \mathcal{G}$  it holds:  $u \perp \!\!\!\perp v | \mathbf{X} \setminus \{u, v\}$ .

#### Definition (Generating class)

Let  $C = \{C_1, \ldots, C_k\}$  be the set of cliques of the dependence graph G. A set of functions  $g_1(), g_2(), \ldots, g_k()$  defined on  $g_i(\mathbf{x}_{C_i})$  is called a **generating class** for the distribution

$$f(\mathbf{x}) = \prod_{i=1,\ldots,k} g_i(\mathbf{x}_{C_i}).$$

### Marginalization

• We have 
$$\frac{1}{\sqrt{|2\pi\Sigma|}}e^{-\frac{1}{2}(\mathsf{x}-\mu)\Sigma^{-1}(\mathsf{x}-\mu)}$$

• We want the distribution over variables  $\{x_3, x_5, x_7\} \subset \{x_1, \dots, x_p\}$ 

#### Marginal of a Gaussian Distribution

The marginal of a Gaussian distribution is calculated by removing appropriate dimensions from the mean and covariance matrix.

• 
$$\mu_{3,5,7} = (\mu_3, \mu_5, \mu_7)$$
 and  
 $\Sigma_{3,5,7} = \begin{bmatrix} \Sigma_{33} & \Sigma_{35} & \Sigma_{37} \\ \Sigma_{53} & \Sigma_{55} & \Sigma_{57} \\ \Sigma_{73} & \Sigma_{75} & \Sigma_{77} \end{bmatrix}$   
•  $\phi_{27} = \phi_{27} = 0$ 

$$\begin{array}{l} \phi_{x_{3,x_{5,x_{7}}}} = \\ \frac{1}{\sqrt{|2\pi\Sigma_{3,5,7}|}} e^{-\frac{1}{2}(x_{3,5,7}-\mu_{3,5,7})\Sigma_{3,5,7}^{-1}(x_{3,5,7}-\mu_{3,5,7})} \\ \end{array}$$

Histogram of s1[, 2]



### Conditioning

• We ame for  $\phi(A|B)$  where

•  $A \subset \{x_1, \ldots, x_p\}$  having q elements,

• the rest  $B = \{x_1, \ldots, x_p\} \setminus A$  has (p - q) elements.

• We rearrange the rows and columns to have A together. Then we get

$$egin{aligned} & x = egin{bmatrix} x_A \ x_B \end{bmatrix}$$
 (one column),  $\mu = egin{bmatrix} \mu_A \ \mu_B \end{bmatrix}$  (one column),  $\Sigma = egin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}$  with dimensions  $egin{bmatrix} q imes q & q imes (p-q) \ (p-q) imes q & (p-q) imes (p-q) \end{bmatrix}$ .

#### Conditional Gaussian

The parameters of the conditional Gaussian distribution  $\phi(A|B = b) = N(\mu_{A|B=b}, \Sigma_{A|B=b})$  are:

$$\begin{aligned} \mu_{A|B=b} &= \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (b - \mu_B) \\ \Sigma_{A|B=b} &= \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}. \end{aligned}$$

Covariance matrix differs but does not depend on the observation b. It depends on the fact B was observed.

### Conditional Gaussian Example

• 
$$\mu^{T} = (1, 2, 3, 4)$$
  
•  $\Sigma = \begin{bmatrix} 10 & 1 & 5 & 4 \\ 1 & 10 & 2 & 6 \\ 5 & 2 & 10 & 3 \\ 4 & 6 & 3 & 10 \end{bmatrix}$ 

- We observed (*X*<sub>3</sub>, *X*<sub>4</sub>) to be (2.8, 4.1)
- We ask for  $\phi(A|B) = \phi(\{X_1, X_2\}|\{X_3, X_4\})$ •  $\Sigma_{AB} = \begin{bmatrix} 5 & 4 \\ 2 & 6 \end{bmatrix}$ •  $\Sigma_{BB} = \begin{bmatrix} 10 & 3 \\ 3 & 10 \end{bmatrix}$ •  $\Sigma_{BB}^{-1} \doteq \begin{bmatrix} 0.11 & -0.033 \\ -0.033 & 0.11 \end{bmatrix}$

• 
$$\Sigma_{AB}\Sigma_{BB}^{-1} \doteq \begin{bmatrix} 0.418 & 0.275 \\ 0.0220 & 0.593 \end{bmatrix}$$
  
•  $\mu_{A|B=b} = \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(b - \mu_B)$   
•  $\mu_{A|B} \doteq \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.418 & 0.275 \\ 0.0220 & 0.593 \end{bmatrix} \begin{bmatrix} (2.8 - 3) \\ (4.1 - 4) \end{bmatrix}$   
•  $\mu_{A|B} \doteq \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} -0.056 \\ 0.055 \end{bmatrix} = \begin{bmatrix} 0.944 \\ 2.055 \end{bmatrix}$   
•  $\Sigma_{A|B=b} = \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$   
•  $\Sigma_{A|B=b} \doteq \begin{bmatrix} 10 & 1 \\ 1 & 10 \end{bmatrix} - \begin{bmatrix} 2.53 & 2.26 \\ 2.26 & 4.13 \end{bmatrix}$   
•  $\Sigma_{A|B=b} \doteq \begin{bmatrix} 7.47 & -1.26 \\ -1.26 & 3.65 \end{bmatrix}$ 

#### Partition Matrix Inverse Properties

• The concentration matrix  $K = \Sigma^{-1}$  is the inverse of the correlation matrix, therefore:

$$\begin{pmatrix} K_{AA} & K_{AB} \\ K_{BA} & K_{BB} \end{pmatrix} \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} = \begin{pmatrix} I_{AA} & \mathbf{0} \\ \mathbf{0} & I_{BB} \end{pmatrix}$$

• From the top right part we get:

$$\begin{aligned}
& \mathcal{K}_{AA}\Sigma_{AB} + \mathcal{K}_{AB}\Sigma_{BB} = \mathbf{0} \\
& -\mathcal{K}_{AA}\Sigma_{AB}\Sigma_{BB}^{-1} = \mathcal{K}_{AB}(1) \\
& \Sigma_{AB}\Sigma_{BB}^{-1} = -\mathcal{K}_{AA}^{-1}\mathcal{K}_{AB}(2).
\end{aligned} \tag{1}$$

• Take the top left part and substitute (1):

$$\begin{array}{rcl} \mathcal{K}_{AA}\Sigma_{AA} & + & \mathcal{K}_{AB}\Sigma_{BA} = I_{AA} \\ \mathcal{K}_{AA}\Sigma_{AA} & + & \left(-\mathcal{K}_{AA}\Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}\right) = I_{AA} \\ \mathcal{K}_{AA}^{-1} & = & \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}. \end{array}$$

#### **Regression Coefficients**

$$\mu_{A|B=b} = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (b - \mu_B)$$
  
 
$$\Sigma_{A|B=b} = \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

• Consider  $x_1$  to be a linear function of others with the noise  $\epsilon_1 \sim N(0, \sigma_1^2)$ :  $x_{1|2...p} = \beta_1 + \beta_{12}x_2 + \beta_{13}x_3 + \ldots + \beta_{1p}x_p + \epsilon_1$ 

Set A the first dimension, B the remaining (p − 1) × (p − 1) matrix:

$$x_{1|B=(x_2,...,x_p)}\tau = \mu_{A|B} + \sum_{AB} \sum_{BB}^{-1} \left( \begin{bmatrix} x_2 \\ \cdots \\ x_p \end{bmatrix} - \mu_B \right) + \epsilon$$

• Recall (2): 
$$\Sigma_{AB}\Sigma_{BB}^{-1} = -K_{AA}^{-1}K_{AB}$$
  
• then  $\sigma_1^2 = \frac{1}{k_{11}}$  with coefficients  $\beta$   
 $(\beta_{12}, \dots, \beta_{1p}) = -\frac{(k_{12}, \dots, k_{1p})}{k_{11}}.$  (3)

#### Fit Linear Gaussian CPD

- To fit ML model of a linear gaussian CPD,
  - you fit the linear regression.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p + \epsilon_1$$
  

$$\hat{\beta} = (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{y}$$
  

$$\hat{\sigma}_Y^2 = Cov(Y, Y) - \sum_i \sum_j \beta_i \beta_j Cov[X_i; X_j]$$
  

$$Cov(X_i; X_j) = \mathbb{E} [(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$$
  

$$\mathbb{E}[X_j] = \frac{1}{N_{rows}} \sum_{i \in rows} x_{ij}$$

from pgmpy.factors.continuous import LinearGaussianCPD ml=maximum\_likelihood\_estimator(data, states) cpdY.fit(data, states, estimator=ml, complete\_samples\_only=True)

 $https://cedar.buffalo.edu/\sim srihari/CSE674/Chap7/7.2-GaussBNs.pdf$ 

#### Parameter Learning for a Gaussian Graphical Model

- Let us have the data  $\mathbf{x}_1^T, \dots, \mathbf{x}_N^T$  over variables  $\mathbf{x} \sim N_p(\mu, \Sigma)$ .
- $S = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i \bar{x}) (\mathbf{x}_i \bar{x})^T$  is the empirical covariance matrix.
- Our model is represented by the concentration matrix  $\Theta = \Sigma^{-1}$  and mean  $\mu$ .
- Log-likelihood of the data is

$$loglik(\Theta,\mu) = \frac{N}{2} \log |\Theta| - \frac{N}{2} tr(\Theta S) - \frac{N}{2} (\bar{x} - \mu)^T \Theta(\bar{x} - \mu).$$

- for a fixed  $\Theta$  is the maximum for  $\mu$ :  $\mu=\bar{x}$  and the last term is 0. We get
- $\textit{loglik}(\Theta,\mu) \propto \textit{log}|\Theta| \textit{tr}(\Theta S)$
- where  $tr(\Theta S) = \sum_{u} \sum_{v} \theta_{uv} s_{uv}$ , therefore only  $s_{uv}$  corresponding to non-zero  $\theta_{uv}$  are considered by the sum.
- We replace the equality conditions by Lagrange multiplyers:  $\ell_C(\Theta) = \log |\Theta| - tr(\Theta S) - \sum_{(j,k) \notin E} \gamma_{jk} \theta_{jk}$
- We maximize. The derivative Θ should be zero (Γ is a matrix with non-zero for missing edges):

$$\Theta^{-1} - S - \Gamma = 0$$

#### Towards the Algorithm

- We iterate one row/column after another.
- We start with the sample covariance matrix

$$W_0 \leftarrow S$$

• We derive the formula for the last row/column: the derivative

$$\begin{pmatrix} W_{11} & w_{12} \\ w_{12}^T & w_{22} \end{pmatrix} - \begin{pmatrix} S_{11} & s_{12} \\ s_{12}^T & s_{22} \end{pmatrix} - \begin{pmatrix} \Gamma_{11} & \gamma_{12} \\ \gamma_{12}^T & \gamma_{22} \end{pmatrix} = 0$$

The upper right block can be written as w<sub>12</sub> - s<sub>12</sub> - γ<sub>12</sub> = 0.
W is inverse of Θ

$$\begin{pmatrix} W_{11} & w_{12} \\ w_{12}^T & w_{22} \end{pmatrix} \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^T & \theta_{22} \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0^T & 1 \end{pmatrix}$$

• therefore the last column without last row is:

$$w_{12} = -W_{11}\theta_{12}/\theta_{22} = W_{11}\beta$$

• Substitute into the derivative  $W_{11}\beta - s_{12} - \gamma_{12} = 0$ 

- we solve for the rows with zero  $\gamma$ :  $\hat{\beta}^* = (W_{11}^*)^{-1} s_{12}^*$ .
- The diagonal  $\theta_{22}$  is (1 bottom right):  $\frac{1}{\theta_{22}} = w_{22} w_{12}^T \beta$ .

#### Estimation of an Undirected Graphical Model Parameters

1: **procedure** GRAPHICAL REGRESSION:( *S* sample covariance )  $W \leftarrow S$  initialize 2: 3. repeat for i = 1, 2, ..., p do 4. Partition W; *i*th row and column,  $W_{11}$  the rest 5 solve  $W_{11}^*\beta^* - s_{12}^* = 0$  for reduced system 6.  $\hat{\beta} \leftarrow \hat{\beta}^*$  by padding with zeros 7: update  $w_{12} \leftarrow W_{11}\hat{\beta}$ 8: end for g٠ until convergence 10: for j = 1, 2, ..., p do 11. lines 5:-8: above and set 12.  $\hat{\theta}_{22} \leftarrow \frac{1}{w_{22} - w_{12}^T \hat{\beta}}$ ▷ the last row on previous slide 13.  $\hat{\theta}_{12} \leftarrow -\hat{\beta} \cdot \hat{\theta}_{22}$ ⊳ (3) 14: end for 15: 16: end procedure

# Example (ESLII)

$(X_1)$ -		$(X_4)$									
$\checkmark$				<b>[</b> 10	.00	1.00	)	5.00	4.(	7 OC	
			$M_{a} = S$	_ 1.	00	10.0	0	2.00	6.0	00	
			$vv_0 = 3$	_ 5.	00	2.00	) [	10.00	3.0	00	
$(X_2)$ -		$-(X_3)$		L 4.	00	6.00	)	3.00	10.	00]	
		[10.00	2.00	6.00				<b>[</b> 10.	00	1.16	4.00 ]
$W_{11}$	=	2.00	10.00	3.00	I	$W_{22}$	=	1.1	16	10.00	3.00
		6.00	3.00	10.00				4.0	)0	3.00	10.00
14/*	_	[10.00	6.00 ]			1/*	_	[10.0	00	1.16	
<i>vv</i> <sub>11</sub>	—	6.00	10.00			vv <sub>22</sub>	_	1.1	6	10.00	
14/*,-1		0.156	-0.09	94]	1478	×,−1		[ 0.1	L01	-0.02	12]
$vv_{11}$	=		4 0.15	6	<i>vv</i> <sub>2</sub>	2	=		012	0.10	1
$\beta^*$	=	[-0.22,	0.53] <sup><i>T</i></sup>		,	β <b>2</b> *	=	[0.08	,0.1	9] <sup>T</sup>	
β	=	[-0.22,	$[0, 0.53]^{T}$			β2	=	- [0.08	,0.1	$[9,0]^T$	
14/1 0	_		16 4 00l	т		14/0	_	[1 00	20		
<i>w</i> <sub>12</sub>	`_	[1.00, 1.	<b>10</b> , <b>4</b> .00]			vv2r	` _	[1.00	', ∠, (		

### Structure Learning

- $\bullet$  We add a lasso penalty  $||\Theta||_1$  which denotes the  $L_1$  norm
  - $\bullet$  the sum of the absolute values of the elements of  $\Theta$  and we ignore the diagonal.
  - The negative penalized log-likelihood is a convex function of Θ.
- we maximize penalized log-likelihood

$$\log|\Theta| - tr(\Theta S) - \lambda ||\Theta||_1 \tag{4}$$

• the gradient equation is now

$$\Theta^{-1} - S - \lambda Sign(\Theta) = 0 \tag{5}$$

sub-gradient notation

• 
$$Sign(\theta_{jk}) = sign(\theta_{jk})$$
 for  $\theta_{jk} \neq 0$ 

- $Sign(\theta_{jk}) \in [-1, 1]$  for  $\theta_{jk} = 0$
- the update for the first row and column will be

$$W_{11}eta - s_{12} + \lambda Sign(eta) = 0$$

• since  $\beta$  and  $\theta_{12}$  have opposite signs.

(6)

1: procedure GRAPHICAL LASSO:( S sample covariance,  $\lambda$  penalty )

2:  $W \leftarrow S + \lambda I$  initialize

3: repeat

4: **for** j = 1, 2, ..., p **do** 5: Partition W; *j*th row and column,  $W_{11}$  the rest

6: solve  $W_{11}\beta - s_{12} + \lambda Sign(\beta) = 0$  using the cyclical 7: ... coordinate-descent algorithm for the modified lasso 8: update  $w_{12}$  by  $W_{11}\hat{\beta}$ 

9: end for

10: **until** convergence

11: **for** 
$$j = 1, 2, \dots, p$$
 **do**  
12: solve  $\hat{\theta}_{22} \leftarrow \frac{1}{1}$ 

12: Solve 
$$\hat{v}_{22} \leftarrow \frac{1}{s_{22} - w_{12}^T \hat{\beta}}$$
  
13: Solve  $\hat{\theta}_{12} \leftarrow -\hat{\beta} \cdot \hat{\theta}_{22}$ 

14: end for

15: end procedure

16: procedure CoordinateDescent:(  $V \leftarrow W_{11}$  )

17: **repeat** 
$$j = 1, 2, \dots, p-1$$

18: 
$$\hat{\beta}_j \leftarrow S(s_{12j} - \sum_{k \neq j} V_{kj}\hat{\beta}_k, \lambda) / V_{jj}$$

19: **until** convergence

20: end procedure

$$\#S(x,t) = sign(x)(|x|-t)_+$$

## Example (glasso)

• 
$$\lambda \leftarrow 1$$
  $W_0 = S + \lambda I = \begin{bmatrix} 11.00 & 1.00 & 5.00 & 4.00 \\ 1.00 & 11.00 & 2.00 & 6.00 \\ 5.00 & 2.00 & 11.00 & 3.00 \\ 4.00 & 6.00 & 3.00 & 11.00 \end{bmatrix}$   
 $W_{11} = \begin{bmatrix} 11.00 & 2.00 & 6.00 \\ 2.00 & 11.00 & 3.00 \\ 6.00 & 3.00 & 11.00 \end{bmatrix}$   $\beta_2^{(2)} = S(1 - \frac{2 \cdot 4}{11} - \frac{6 \cdot 21}{121}, 1)/11 \approx -0.16$   
 $\beta_3^{(2)} = S(1 - \frac{3 \cdot 21}{121}, 1)/11 \approx 0.35$   
 $\beta_4^{(2)} = S(1 - \frac{3 \cdot 4}{11}, 1)/11 = 0$   
 $\psi \leftarrow W_{11}$   $\hat{\beta}_1 \approx [-0.22; 0.32; 0.30]$   
 $\beta_4^{(1)} = S(4 - \frac{3 \cdot 4}{11}, 1)/11 = \frac{21}{121}$   $W_1 \approx \begin{bmatrix} 11.00 & 0.05 & 4.03 & 3.01 \\ 0.05 & 11.00 & 2.00 & 6.00 \\ 4.03 & 2.00 & 11.00 & 3.00 \\ 3.01 & 6.00 & 3.00 & 11.00 \end{bmatrix}$ 

- Computational speed
  - The graphical lasso algorithm is extremely fast
  - can solve a moderately sparse problem with 1000 nodes in less than a minute.
  - It can be modified to have edge–specific penalty parameters  $\lambda_{jk}$
  - setting  $\lambda_{jk} = \infty$  will force  $\hat{ heta}_{jk}$  to be zero
  - graphical lasso subsumes the parameter learning algorithm.
- Missing data
  - some missing observations may be imputed by EM algorithm from the model
  - latent fully unobserved variables do not bring more power in Gaussian graphical model
  - latent variables are very important in discrete distributions.

sklearn.covariance.graphical\_lasso

### Model Quality (Model Selection)

#### Definition (Saturated model, GGM Deviance, iDeviance, Likelihood Ratio Test)

- saturated model full model with all edges, it has maximal loglikelihood
- Deviance

$$D = dev = 2 \cdot (\hat{\ell}_{sat} - \hat{\ell}) = N \log rac{|S^{-1}|}{|\hat{K}|} = -N \log |S\hat{K}|$$

independent model - no edges, it has minimal likelihood
iDeviance

$$iD = idev = 2 \cdot (\hat{\ell} - \hat{\ell}_{ind}) = N\left(\log|\hat{K}| + \sum_{i=1}^{p}\log s_{ii}\right)$$

 $\bullet$  Irt likelihood ratio test for models  $\mathcal{M}_1 \subseteq \mathcal{M}_0$ 

$$Irt = 2 \cdot (\hat{\ell}_0 - \hat{\ell}_1) = N \log \frac{|\hat{K}_0|}{|\hat{K}_1|}.$$

### Undirected Graphical Models and Their Properties

#### Definition (Undirected Graphical Model, Markov Graph)

An **Undirected Graphical Model** (Markov graph, Markov network) is a graph  $\mathcal{G} = (V, E)$ , where nodes V represent random variables and the absence of an edge (A, B) denoted  $A \perp_{\mathcal{G}} B$  implies that the corresponding random variables are conditionally independent given the rest in the probability distribution P(V).

$$A \perp\!\!\!\perp_{\mathcal{G}} B \Longrightarrow A \perp\!\!\!\perp_{P} B | V \setminus \{A, B\}.$$

$$(7)$$

is known as the **pairwise Markov independencies** of  $\mathcal{G}$ .

#### Definition (Separators)

- If *A*, *B* and *C* are subgraphs, then *C* is said to separate *A* and *B* if every path between *A* and *B* intersects a node in *C*.
- C is called a **separator**.

• Separators break the graph into conditionally independent pieces.

#### Definition (Global Markov Property)

A probability measure P over V is (globally) Markov with respect to an undirected graph G iff for any subgraphs A, B and C holds:

• if C separates A and B then the conditional independence  $A \perp _P B | C$  holds, that is

$$A \perp\!\!\!\perp_{\mathcal{G}} B | C \Longrightarrow P(A|C) \cdot P(B|C) = P(A, B|C).$$

#### Theorem

The pairwise and global Markov properties of a graph are equivalent for graphs with strictly positive distributions.

- Gaussian distribution is always positive.
- We may infer global independence relations from simple pairwise properties.
- The global Markov property allows us to decompose graphs into smaller more manageable pieces.

(8)

### Markov Random Fields (Markovská náhodná pole)

• A probability density function *f* over a Markov graph *G* with the set of maximal cliques  $\{C_1, \ldots, C_k\}$  can be represented as

$$f(x) = \prod_{i=1,\ldots,k} \psi_i(x_{C_i}) = \psi_1(x_{C_1}) \cdot \ldots \cdot \psi_k(x_{C_k})$$
(9)

- where  $\psi_i$  are positive functions called **clique potentials**.
- they capture the dependence in  $X_{C_i}$  by scoring certain instances  $x_{C_i}$  higher than others.
- with the normalizing constant (partition function) Z

$$Z = \int_X exp\left(\sum_{i=1,\ldots,k} \log g_i(x_{C_i})\right)$$

- Such set of random variables is called Markov Random Field or Markov graph. If the potentials represent conditional probabilities with respect to some observation, it is called a Conditional Markov Random Field.
- For Markov networks with positive distributions the probability density function (9) implies a graph with independence properties defined by the cliques in the product.

• A graphical model does not always uniquely specify the higher-order dependence structure of ta joint probability distribution.

$$f^{(2)}(x, y, z) = \frac{1}{Z}\psi_1(x, y)\psi_2(x, z)\psi_3(y, z)$$
  
$$f^{(3)}(x, y, z) = \frac{1}{Z}\psi(x, y, z)$$



- For Gaussian distribution, pairwise interactions fully specify the model.
- We focus on pairwise Markov Graphs
  - where at most second order interactions are represented (like  $f^{(2)}$ ).

### MRF for Image Denoising

- Given a noisy image v, perhaps with missing pixels, recover an image u,  $u_{n,m} \in \mathbb{R}$  that is both smooth and close to v.
- Let each pixel be a node in a graph  $\mathcal{G} = (V, E)$ , with 4-connected neighborhood. Only pairwise interactions are present.
- We minimize the energy function (add missing margin on your own)

$$E(u) = \sum_{(m,n)\in P} (u_{m,n} - v_{m,n})^2 + \lambda \sum_{(m,n)\in P} \left[ (u_{n+1,m} - u_{n,m})^2 + (u_{n,m+1} - u_{n,m})^2 \right]$$

We can solve u iteratively

• 
$$s_{m,n} = u_{n-1,m} + u_{n+1,m} + u_{n,m-1} + u_{n,m+1},$$
  
•  $u_{n,m}^{(t+1)} = \begin{cases} \frac{1}{1+4\lambda} (v_{n,m} + \lambda s_{n,m}^{(t)}) \text{ for } (n,m) \in v \\ \frac{1}{4} s_{n,m}^{(t)} & \text{ for missing } v \end{cases}$ 



- The goal is to find the signal u that minimizes the energy E(u).
- Estimation by Max-flow/Min-Cut in a specific graph or Gibbs sampling.
- https://www.cs.toronto.edu/~fleet/courses/2503/fall11/Handouts/mrf.pdf

#### Definition (Ising Model, Boltzmann Machine (in ESLII))

• The **Ising model** is defined by a graph  $\mathcal{G} = (\mathcal{X}, E)$  of binary variables  $X_i \in \mathcal{X}$  and a set of parameters  $\Theta$ . The joint probabilities are given by:

$$p(X,\Theta) = e^{\sum_{(j,k)\in E} \theta_{jk} X_j X_k - \Phi(\Theta)} \text{ for } X \in \mathcal{X}$$
  
$$\Phi(\Theta) = \log \sum_{x \in \mathcal{X}} \left[ e^{\sum_{(j,k)\in E} \theta_{jk} x_j x_k} \right].$$

- it models only binary interactions (and unary)
- for technical reasons requires *constant* node  $X_0 \equiv 1$  to be included.
- Originally from statistical mechanics.
- This model is equivalent to a first-order-interaction Poisson log-linear model for multiway tables of counts (Bishop et al., 1975).
- it implies a logistic form for each node conditional on the others

$$P(X_{j} = 1 | X_{-j} = x_{-j}) = \frac{1}{1 - \exp(-\theta_{j0} - \sum_{(j,k) \in E} \theta_{jk} x_{k})}$$

• Θ is fitted iteratively (Iterative proportional fitting, gradient descend, Poisson log-linear modeling, Mean field approximation, Gibbs sampling).

108 - 145

### Restricted Boltzmann Machines (RBM)

- We have visible  ${\cal V}$  and hidden  ${\cal H}$  variables.
- 'restricted' means the variables are organized in two layers:
  - the hidden layer
  - the visible layer that is split to input V<sub>1</sub> and output variables V<sub>2</sub>. (pixels of image/digit label).
  - no edges inside any layer
  - an edge between each hidden and visible variable.
- The structure enables faster parameter estimation.





$$E(v, h) = -\sum_{i=1}^{n} \sum_{j=1}^{m} w_{ij}h_iv_j - \sum_{j=1}^{m} b_jv_j - \sum_{i=1}^{n} c_ih_iv_j$$

#### Contrastive divergence ×

1: **procedure** CONTRASTIVE DIVERGENCE:(*S* batch for update)  $\Delta w_{ii}, \Delta b_i, \Delta c_i \leftarrow 0$ 2: repeat 3: for each training sample in a batch S do 4: Sample  $\mathcal{H}$  given  $\mathcal{V}_1$ ,  $\mathcal{V}_2$ 5: Sample  $\mathcal{V}_1^{(last)}, \mathcal{V}_2^{(last)}$  given  $\mathcal{H}$ 6: Sample  $\mathcal{H}$  given  $\mathcal{V}_1^{(last)}$ ,  $\mathcal{V}_2^{(last)}$ 7:  $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1|v^{(0)}) \cdot v_i^{(0)} - p(H_i = 1|v^{(last)}) \cdot v_i^{(last)}$ 8:  $\Delta b_i \leftarrow \Delta b_i + v_i^{(0)} - v_i^{(last)}$ 9:  $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1|v^{(0)}) - p(H_i = 1|v^{(last)})$ 10: end for 11: until convergence 12: **return**  $\Delta w_{ii}$ ,  $\Delta b_i$ ,  $\Delta c_i$  to adjust the parameters. 13: 14: end procedure

Fischer, A., Igel, C. (2012). An Introduction to Restricted Boltzmann Machines. In: Alvarez, L., Mejail, M., Gomez, L., Jacobo, J. (eds) Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications. CIARP 2012. Locture Nachine Computered Graphical Modes 111 Springer 108-145

#### Mixed interaction models

- Discrete and Gaussian variables together.
- Conditional Gaussian density for x = (i, y), i is the list of discrete variables, y is the list of continuous variables.
- Directed graphs: a discrete child of a continuous parent is not allowed.
- **Undirected graphs**: If there is a path between two discrete variables *A*, *B*, then they are connected by a path without any continuous variable.
- $f(i,y) = p(i)(2\pi)^{-\frac{q}{2}}|\Sigma|^{-\frac{1}{2}}\exp(-\frac{1}{2}(y-\mu(i))\Sigma^{-1}(y-\mu(i)))$
- The parameters  $p(i), \mu(i), i \in \mathcal{I}, \Sigma$  are called moment parameters
- in the exponential form we get

$$f(i,y) = \exp\left\{g(i) + h(i)^{T}y - \frac{1}{2}y^{T}Ky\right\}$$
$$= \exp\left\{g(i) + \Sigma_{u}h_{u}(i)y_{u} - \frac{1}{2}\Sigma_{u,v}K_{u,v}y_{u}y_{v}\right\}$$

• parameters  $g(i), h(i), i \in \mathcal{I}, K$  are called **canonical parameters**.

- A marginal distribution is not necessarily a conditional Gaussian distribution
- it is a mixture of conditional Gaussians
- it is still tractable for evaluation
- and learning.

• Is this possible for other kind of distributions?

### Markov Properties (Zeros are dangerous)

#### Definition (Markov properties: Global, Local, Pairwise)

Let G be an undirected graph over V, let P be a probability measure P over V. (GM) P is (globally) Markov with respect to G iff

 $\forall (\mathcal{A}, \mathcal{B} \in V, \mathcal{C} \subseteq V) \ \mathcal{A} \perp\!\!\!\perp_{\mathcal{G}} \mathcal{B} | \mathcal{C} \Rightarrow \mathcal{A} \perp\!\!\!\perp_{P} \mathcal{B} | \mathcal{C} \text{ in } \mathsf{P}.$ 

(LM) A probability measure has the local Markov property iff  $(\forall A \in V) : A \perp P V \setminus Fa_A | N_A$ 

(PM) P has the **pairwise Markov property** iff  $\forall A, B \in V, A \neq B$  not connected by an edge holds  $A \perp P B | V \setminus \{A, B\}$ .

#### Theorem

These properties are equivalent for strictly positive measures.

Counterexamples for measures with zero probability everywhere except (0,0,0) and (1,1,1).

See [Milan Studený: Struktury podmíněné nezávislosti, Matfyzpress 2014].

### Examples

Example (*P* has the pairvise but not the local property)

$$\begin{split} V &= \{A, B, C\}, E = \{(b, c)\}. \text{ Let us have a} \\ \text{binary probability measure } V \text{ nonzero at points} \\ (0, 0, 0) \text{ and } (1, 1, 1) \text{ [Studený p.101].} \\ A &\perp B | \{C\} \\ A &\perp C | \{B\} \\ \end{split}$$





# Example (P has the local but not the global property)

$$V = \{A, B, C, D\}, E = \{(a, b), (c, d)\}.$$
 Let  

$$P(V) \text{ be nonzero only at points } (0, 0, 0, 0) \text{ and} (1, 1, 1, 1) [Studený p.101].$$

$$A \perp CD|\{B\}$$

$$B \perp CD|\{B\}$$

$$B \perp CD|\{A\} \& \text{ does not imply } A \perp C|\{\}.$$

$$D \perp AB|\{D\}$$



### Nonparanormal Graphical Models

- A continuous pairwise interaction model.
- We model marginal distributions,
- and the most important relations by gaussian copula.
- https://www.stat.cmu.edu/~larry/=sml/GraphicalModels.pdf





### List of topics

- Linear, ridge, lasso regression, k-neares neighbours,(formulas) overfitting, curse of dimensionality, (LARS)
- Splines the base, natural splines, smoothing splines; kernel smoothing: kernel average, Epanechnikov kernel.
- **O** Logistic regression, Linear discriminant analysis, generalized additive models
- Train/test error and data split, square error, 0-1, crossentropy, AIC, BIC,(formulas) crossvalidation, one-leave-out CV, wrong estimate example
- Idecision trees, information gain/entropy/gini, CART prunning,(formulas)
- random forest (+bagging), OOB error, Variable importance, boosting (Adaboost(formulas) and gradient boosting), stacking, MARS,
- Bayesian learning: MAP, ML hypothesis (formulas), Bayesian optimal prediction, EM algorithm
- Sclustering: k-means, Silhouette, k-medoids, hierarchical
- Apriori algorithm, Association rules, support, confidence, lift
- Inductive logic programming basic: hypothesis space search, background knowledge, necessity, sufficiency and consistency of a hypothesis, Aleph
- Undirected graphical models, Graphical Lasso procedure, deviance, MRF
- Gaussian processes: estimation of the function and its variance (figures, ideas).

### List of topics

- Linear, ridge, lasso regression, k-neares neighbours,(formulas) overfitting, curse of dimensionality, (LARS)
- Splines the base, natural splines, smoothing splines; kernel smoothing: kernel average, Epanechnikov kernel.
- **O** Logistic regression, Linear discriminant analysis, generalized additive models
- Train/test error and data split, square error, 0-1, crossentropy, AIC, BIC,(formulas) crossvalidation, one-leave-out CV, wrong estimate example
- Idecision trees, information gain/entropy/gini, CART prunning,(formulas)
- random forest (+bagging), OOB error, Variable importance, boosting (Adaboost(formulas) and gradient boosting), stacking, MARS,
- Bayesian learning: MAP, ML hypothesis (formulas), Bayesian optimal prediction, EM algorithm
- Sclustering: k-means, Silhouette, k-medoids, hierarchical
- Apriori algorithm, Association rules, support, confidence, lift
- Inductive logic programming basic: hypothesis space search, background knowledge, necessity, sufficiency and consistency of a hypothesis, Aleph
- Undirected graphical models, Graphical Lasso procedure, deviance, MRF
- Gaussian processes: estimation of the function and its variance (figures, ideas).

### Table of Contens

- Overview of Supervised Learning
- Kernel Methods, Basis Expansion and regularization
- 3 Linear Methods for Classification
- 4 Model Assessment and Selection
- 5 Additive Models, Trees, and Related Methods
- 6 Ensamble Methods
- 🕖 Bayesian learning, EM algorithm
- 8 Clustering
- 9 Association Rules, Apriori
- Inductive Logic Programming
- 1 Undirected Graphical Models
- 12 Gaussian Processes
- 13 Support Vector Machines
- (PCA Extensions, Independent CA)