

Implementazione algoritmo PC

A lezione abbiamo visto che dato un insieme di n variabili casuali:

$$\{X_1, X_2, \dots, X_n\}$$

È possibile definire su di esso la matrice varianza covarianza come segue:

$$\Sigma \in \mathbb{R}^{n \times n} = \begin{bmatrix} E[(X_1 - \bar{X}_1)(X_1 - \bar{X}_1)] & E[(X_1 - \bar{X}_1)(X_2 - \bar{X}_2)] & \dots & E[(X_1 - \bar{X}_1)(X_n - \bar{X}_n)] \\ E[(X_2 - \bar{X}_2)(X_1 - \bar{X}_1)] & E[(X_2 - \bar{X}_2)(X_2 - \bar{X}_2)] & \dots & E[(X_2 - \bar{X}_2)(X_n - \bar{X}_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \bar{X}_n)(X_1 - \bar{X}_1)] & E[(X_n - \bar{X}_n)(X_2 - \bar{X}_2)] & \dots & E[(X_n - \bar{X}_n)(X_n - \bar{X}_n)] \end{bmatrix}$$

Questa matrice è simmetrica definita positiva, dunque da un punto di vista puramente matematico è *sempre* non singolare, quindi *sempre* invertibile.

La sua inversa Σ^{-1} è chiamata matrice di precisione, ed è legata alla *correlazione parziale* tra le n variabili casuali secondo la relazione:

$$\rho_{ij.rest} = \frac{-\sigma^{ij}}{\sqrt{\sigma^{ii}\sigma^{jj}}} \quad (1.0)$$

Dove: $\rho_{ij.rest}$ è il coefficiente di correlazione parziale fra la variabile X_i e la variabile X_j condizionato a tutte le rimanenti, e σ^{ij} è l' i, j -esimo elemento di Σ^{-1} .

Vale inoltre la seguente relazione:

$$X_i \perp\!\!\!\perp X_j \mid rest \Leftrightarrow \rho_{ij.rest} = 0$$

Il PC algorithm costruisce lo scheletro del DAG con la seguente logica "grossolana":

- Inizializza il grafo come isomorfo a $K_{n \times n}$
- Per ogni coppia ordinata (i, j) di nodi (variabili) tale per cui esiste un arco non direzionato da i a j :
 - controlla se una certa disequazione di test è soddisfatta
 - se sì, allora le due variabili sono indipendenti dato il resto, quindi l'arco fra i due rispettivi nodi nel grafo non deve esistere

La disequazione che ho chiamato "di test" utilizza la trasformata z di Fisher sulla matrice delle correlazioni parziali (che può essere ottenuta da Σ^{-1}), ed è la seguente:

$$\sqrt{N - |K| - 3} |z(i, j | rest)| \leq \Phi^{-1}(1 - \frac{\alpha}{2})$$

Dove N è il numero di variabili, $|K| = l$ = numero di vicini del nodo i , escluso il nodo j , z è la trasformata di Fisher, $\Phi^{-1}(1 - \frac{\alpha}{2})$ è l'inversa della funzione Φ distribuzione cumulativa della Normale, valutata in $(1 - \frac{\alpha}{2})$ e infine α è detto *tuning parameter*.

Per provare a valutare le performance di questo algoritmo, ho implementato il seguente codice python:

```
1 def pc_algorithm(a, sigma_inverse):
2     l = - 1
3     n = len(sigma_inverse)
4     z = lambda i, j : -sigma_inverse[i][j]/((sigma_inverse[i][i]*sigma_inverse[j][j])**0.5)
5     act_g = complete(n)
6     act = 0
7     while l < n-1:
8         l = l + 1
9         for i in range(0, n):
10             for j in range(i+1, n):
11                 if (act_g[i][j] != 1):
12                     continue
13                 adjacents = adj(i, act_g)
14                 if len(adjacents) == 0:
15                     continue
16                 act_set = setdiff(adjacents, [j])
17                 all_k = findsubsets(act_set, l)
```

```

18     counter = 0
19     while (act_g[i][j]!=0):
20         if (counter >= len(all_k)):
21             break
22     K = all_k[counter];
23     counter = counter+1
24     if test(n,z,i,j,a,K):
25         act_g[i][j] = 0
26 for i in range(0,n):
27     for j in range(0,n):
28         if i>j:
29             act_g[i][j] = act_g[j][i]
30 return act_g;
31

```

Problemi

Non sono ovviamente sicuro di aver implementato correttamente l'algoritmo in se, ma nello specifico un primo problema è stato quello di definire la trasformata z . Infatti la definisco come una funzione anonima che, data la matrice Σ^{-1} mi restituisce la correlazione parziale di i, j dato il resto, secondo la relazione (1.0):

```

1 z(i,j) = sigma.inverse[i][i]/sqrt(sigma.inverse[i][i]*sigma.inverse[j][j])

```

Non sono sicuro che questo sia corretto perchè il coefficiente di correlazione parziale andrebbe calcolato dato K , e non date tutte le altre variabili.

Un altro problema riguarda invece la funzione Φ^{-1} , infatti non posso pensare di calcolarla direttamente in quanto non è esprimibile come combinazione di funzioni analitiche e/o polinomiali, l'unica cosa che posso fare è approssimarla. Infatti:

$$\Phi(x) = \frac{1}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

È un integrale irrisolvibile in forma chiusa ed esplicita.

Ho scelto di utilizzare la seguente approssimazione nota della funzione quantile Φ^{-1} :

$$\Phi^{-1}(p) = \sqrt{2} \operatorname{erf}^{-1}(2p - 1)$$

Dove erf^{-1} è a sua volta l'inversa della *error function*, la quale necessita a sua volta di essere approssimata:

$$\operatorname{erf}^{-1}(x) = \operatorname{sgn}(x) \sqrt{\sqrt{\left(\frac{2}{\pi a} + \frac{\ln(1-x^2)}{2}\right)^2 - \frac{\ln(1-x^2)}{a}} - \left(\frac{2}{\pi a} + \frac{\ln(1-x^2)}{2}\right)}$$

Una buona approssimazione la si ottiene ponendo $a = 0.147$

```

1 def test(n,z,i,j,a,K):
2     root = (n-len(K)-3)**0.5
3     return root*abs(z(i,j)) <= phi(1-a/2)
4
5 def phi(p):
6     return (2**0.5)*erfinv(2*p-1)
7
8 def erfinv(x):
9     sgn = 1
10    a = 0.147
11    PI = numpy.pi
12    if x<0:
13        sgn = -1
14    temp = 2/(PI*a) + numpy.log(1-x**2)/2
15    add_1 = temp**2
16    add_2 = numpy.log(1-x**2)/a
17    add_3 = temp
18    rt1 = (add_1-add_2)**0.5
19    rtarg = rt1 - add_3
20    return sgn*(rtarg**0.5)

```

Un ultimo problema è che, generando dati casualmente, spesso ottengo anche matrici varianza-covarianza singolari in aritmetica finita, ma questo mi sembra che ce lo disse anche lei a lezione, quindi lo considero normale.

Output

La prima cosa che ho notato è che per $n > 10$ variabili, i tempi di esecuzione sono molto lunghi. Allego alcune immagini di grafi prodotti utilizzando il seguente dataset di 7 variabili generato in maniera casuale, con la seguente strategia:

$$X_i \sim N(\mu, \sigma^2) = N[(i-1), i] \quad i = 1, \dots, 7$$

X_1 :	0.1298712897211212	-0.2112987219871212	-0.1749521075950111	0.2156423156412200	-0.0215615646515618	0.3231256154842521	0.1235484556489215
X_2 :	1.9741980022248782	0.8451937696200296	1.1799283254919182	0.6031724223406041	-0.6566008384851423	-1.0064929865840662	3.4676455398519845
X_3 :	3.100078139389854	-0.7838869190170845	5.165622481002579	3.1184726581121383	4.899115147443788	1.870911939491902	1.5530716047378323
X_4 :	5.198030621914989	2.8160245411363602	5.482744040414655	2.17400869757652	8.353322676841408	2.4365539380996406	0.6299492890543488
X_5 :	4.1699869117111215	9.644663893776446	4.189526185481068	-1.0367886406038753	5.959916183783191	5.331092113213319	3.9619019745126614
X_6 :	4.376395614608083	10.209356572618773	-1.201476701612889	-1.7526211220068442	1.2869284076136944	13.388500730751154	2.7765624095474357
X_7 :	-0.430623316687897	1.6433833402269427	-10.033331689878722	8.137353672070207	-2.3771004289730833	-0.55908452236672	6.661604133285603

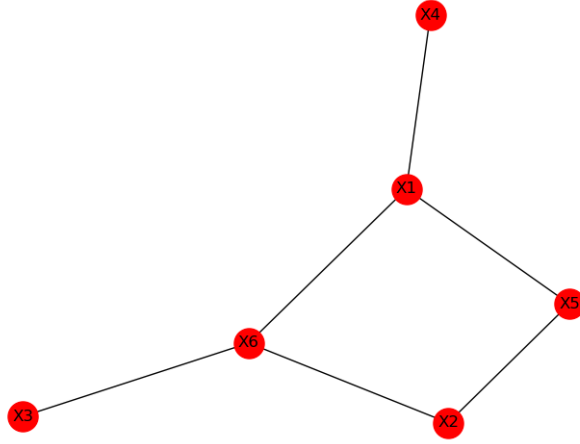


Figure 1: Dataset casuale, alpha = 0.20

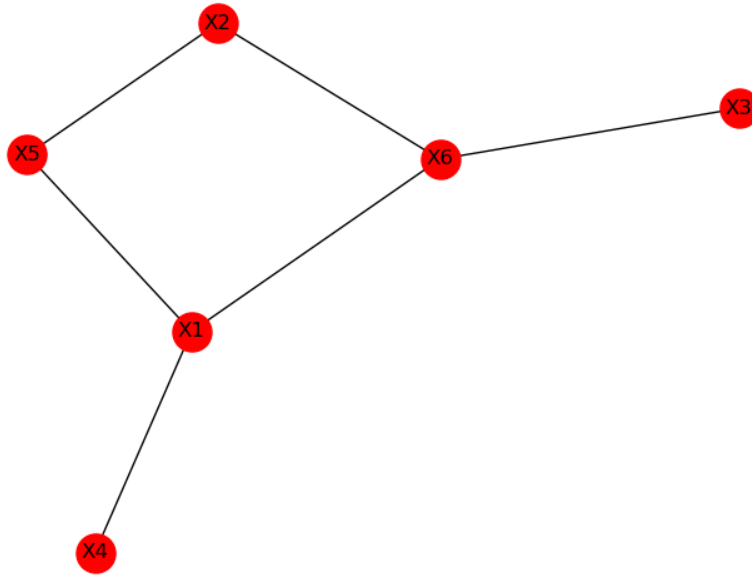


Figure 2: Dataset casuale, alpha = 0.50

Non ritengo i grafici attendibili in questo esempio perchè le variabili si distribuiscono come normali univariate indipendenti fra di loro; sebbene questo sia vero solo in teoria (in pratica i numeri generati sono pseudo-casuali, quindi non perfettamente indipendenti).

Utilizzando invece direttamente la matrice Σ^{-1} del butterfly model ottengo:

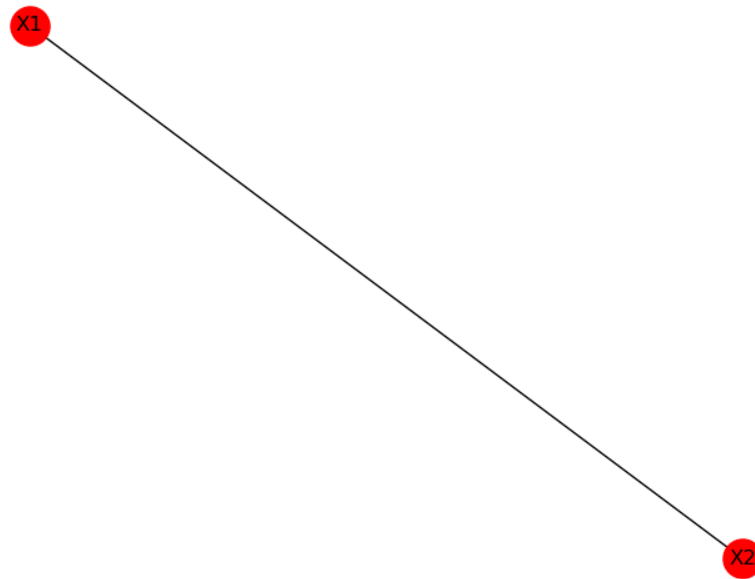


Figure 3: Butterfly model, $\alpha = 0.95$

Il che è senza dubbio sbagliato, solo con α molto elevato ottengo un grafo non vuoto che comunque dovrebbe essere il modello a farfalla, e non un segmento. Le volevo chiedere appunto un aiuto a capire dove fosse il, o più probabilmente gli, errori. Allego anche il codice completo.

pcalgorithm.py

```
1 import itertools
2 import copy
3 import numpy
4 import graphtools
5
6 def complete(n):
7     A = list()
8     for i in range(0,n):
9         row = list()
10        for j in range(0,n):
11            if i == j:
12                row.append(0)
13            else:
14                row.append(1)
15        A.append(row)
16    return A
17
18 def setdiff(A,B):
19     ret_set = copy.copy(A)
20     for x in B:
21         if x in A:
22             ret_set.remove(x)
23     return ret_set
24
25 def adj(i,G):
26     adjacents = list()
27     for j in range(0,len(G[i])):
28         if i!=j and G[i][j] == 1:
29             adjacents.append(j)
30     return adjacents
31
32 def findsubsets(S,l):
33     return list(itertools.combinations(S,l))
34
35 def erfinv(x):
```

```

36  sgn = 1
37  a = 0.147
38  PI = numpy.pi
39  if x<0:
40      sgn = -1
41  temp = 2/(PI*a) + numpy.log(1-x**2)/2
42  add_1 = temp**2
43  add_2 = numpy.log(1-x**2)/a
44  add_3 = temp
45  rt1 = (add_1-add_2)**0.5
46  rtarg = rt1 - add_3
47  return sgn*(rtarg**0.5)
48
49  def phi(p):
50      return (2**0.5)*erfinv(2*p-1)
51
52  def test(n,z,i,j,a,K):
53      root = (n-len(K)-3)**0.5
54      return root*abs(z(i,j)) <= phi(1-a/2)
55
56  def meanof(dataset):
57      n = len(dataset[0])
58      m = []
59      for i in range(0,n):
60          m.append(0.0)
61          datasize = len(dataset)
62          for i in range(0,datasize):
63              for j in range(0,n):
64                  m[j] = m[j] + float(dataset[i][j])
65          for i in range(0,n):
66              m[i] = m[i] / datasize
67      return m
68
69  def zeros(n,m):
70      zer = []
71      for i in range(0,n):
72          row = []
73          for j in range(0,m):
74              row.append(0)
75          zer.append(row)
76      return zer
77
78  def getcol(i,matrix):
79      col = []
80      for row in matrix:
81          col.append(row[i])
82      return col
83
84  def sigma(dataset,means):
85      n = len(means)
86      sigma = zeros(n,n)
87      for i in range(0,n):
88          for j in range(0,n):
89              dset_i = getcol(i,dataset)
90              dset_j = getcol(j,dataset)
91              means_i = means[i]
92              means_j = means[j]
93              sigma[i][j] = covar(dset_i,dset_j,means_i,means_j)
94      return sigma
95
96  def covar(X,Y,ux,uy):
97      n = len(X)
98      s = 0
99      for i in range(0,n):
100          s = s +(X[i] - ux)*(Y[i] - uy);
101      return float(s)/n
102
103  def getSigma(dataset):
104      means = meanof(dataset)
105      return sigma(dataset,means)
106
107  def getInverse(A):
108      return numpy.linalg.inv(A)
109

```

```

110 def pc_algorithm(a, sigma_inverse):
111     l = - 1
112     n = len(sigma_inverse)
113     z = lambda i, j : -sigma_inverse[i][j]/((sigma_inverse[i][i]*sigma_inverse[j][j])**0.5)
114     act_g = complete(n)
115     act = 0
116     while l<n-1:
117         l = l + 1
118         for i in range(0,n):
119             for j in range(i+1,n):
120                 if (act_g[i][j]!=1):
121                     continue
122                 adjacents = adj(i, act_g)
123                 if len(adjacents)==0:
124                     continue
125                 act_set = setdiff(adjacents,[j])
126                 all_k = findsubsets(act_set, 1)
127                 counter = 0
128                 while (act_g[i][j]!=0):
129                     if (counter >= len(all_k)):
130                         break
131                 K = all_k[counter];
132                 counter = counter+1
133                 if test(n,z,i,j,a,K):
134                     act_g[i][j] = 0
135         for i in range(0,n):
136             for j in range(0,n):
137                 if i>j:
138                     act_g[i][j] = act_g[j][i]
139     return act_g;
140
141 def rand_set(X,Var):
142     n = len(X)
143     rset = []
144     for i in range(0,n):
145         rset.append(numpy.random.normal(X[i],Var[i],n))
146     return numpy.transpose(rset);
147
148 def make_graph_from_dataset(dataset, alpha):
149     sigma = getSigma(dataset)
150     sigma_inverse = getInverse(sigma)
151     return pc_algorithm(alpha, sigma_inverse)
152
153 def plot(adj_matrix, varnames):
154     graphtools.plot_graph(adj_matrix, varnames)
155
156 def get_rand_dataset(dim):
157     return rand_set(range(0,dim),range(0,dim))
158
159 def butterfly():
160     return [[0.8, 0.5, 0, 0.6],[0.5, 1.4, -0.6, 0.4],[0, -0.6, 1.2, -0.3],[0.6, 0.4, -0.3, 1]]

```

main - butterfly

```

1 if __name__ == '__main__':
2     varnames = ['X1', 'X2', 'X3', 'X4']
3     alpha = 0.95
4     G = pc_algorithm(alpha, butterfly())
5     plot(G, varnames)

```

main - random variables

```

1 if __name__ == '__main__':
2     dataset = get_rand_dataset(7)
3     varnames = ['X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X7']
4     alpha = 0.50
5     G = make_graph_from_dataset(dataset, alpha)
6     plot(G, varnames)

```