

# Additive Noise Models: Identifiability, Learning Algorithms, Hidden Variables and Time Series

Behrad Moniri

Sharif University of Technology

*bemoniri@ee.sharif.edu*

January 22, 2019

# Model Definition

We call the SCM,  $C$ , an additive noise model if each observed variable  $X_j$  is associated with a node  $j$  in a directed acyclic graph  $G$ , and the value of  $X_j$  is obtained as a function of its parents in  $G$ , plus independent additive noise  $N_j$ , i.e.

$$X_j = f_j(\mathbf{PA}_j) + N_j, \quad j = 1, \dots, p \quad (1)$$

with jointly independent variables  $N_j$ . We will assume that the noise variables have strictly positive density.

## Causal Minimality

For those models with strictly positive density, causal minimality reduces to the condition that each function  $f_j$  is not constant in any of its arguments.

# Identifiability: The Bivariate Case

Hoyer et al. (2009) proves the following theorem about the identifiability of bivariate additive noise models.

# Identifiability: The Bivariate Case

## Theorem

An additive noise model with two variables, i.e.,  $X_1 = N_1$  and  $X_2 = f_j(X_1) + N_2$ , with  $N_1 \perp\!\!\!\perp N_2$ , is identifiable if it does not solve the following differential equation for all  $x_i, x_j$  with  $\nu''(x_j - f(x_i))f'(x_i) \neq 0$ :

$$\xi''' = \xi'' \left( -\frac{\nu''' f'}{\nu''} + \frac{f''}{f'} \right) - 2\nu'' f'' f' + \nu' f''' + \frac{\nu' \nu''' f'' f'}{\nu''} - \frac{\nu' (f'')^2}{f'},$$

Here  $\xi := \log p_{X_1}$  and  $\nu := \log p_{N_2}$  and we have skipped the arguments  $x_2 - f(x_1)$ ,  $x_1$ , and  $x_1$  for  $\nu$ ,  $\xi$ , and  $f$  and their derivatives, respectively.

## Corollary

*Gaussian Noise: Assume that  $\nu''' = \xi''' = 0$  everywhere. If a backward model exists, then  $f$  is linear.*

## Corollary

*Assume that  $f(x) = x$  and  $p_x(x) = e^{-x-e^{-x}}$  and  $p_n(n) = e^{-n-e^{-n}}$ . With  $p_y(y) = e^{-y-2\log(1+e^y)}$ ,  $\tilde{p}_n(\tilde{n}) = e^{-2\tilde{n}-e^{-\tilde{n}}}$  and  $g(y) = \log(1 + e^{-y})$ , one obtains:*

$$p(x, y) = p_n(y - f(x))p_x(x) = \tilde{p}_n(x - g(y))p_y(y)$$

*so the model is not identifiable.*

# Identifiability: From Bivariate to Multivariate

## Definition

Consider an ANM with  $p$  variables. We call this SCM a *restricted additive noise model* if for all  $j \in \mathbf{V}$ ,  $i \in \mathbf{PA}_j$  and all sets  $\mathbf{S} \subseteq \mathbf{V}$  with  $\mathbf{PA}_j \setminus \{i\} \subseteq \mathbf{S} \subseteq \mathbf{ND}_j \setminus \{i, j\}$ , there is an  $x_{\mathbf{S}}$  with  $p_{\mathbf{S}}(x_{\mathbf{S}}) > 0$ , such that

$$\left( f_j(x_{\mathbf{PA}_j \setminus \{i\}}, X_i), \mathcal{L}(X_i | X_{\mathbf{S}} = x_{\mathbf{S}}), \mathcal{L}(N_j) \right)$$

satisfies the bivariate identifiability conditions.

We assume that the noise variables to have non-vanishing densities and the functions  $f_j$  are three times differentiable.

# Identifiability: The Multivariate Case

Peters et al. (2014) proves a very interesting theorem. This theorem states how we can generalize a bivariate identifiability to the multivariate case, in this case ANM identifiability.

## Theorem

*Let  $X_1, \dots, X_p$  be generated by a restricted additive noise model with graph  $G_0$  and assume that  $P_{\mathbf{X}}$  satisfies causal minimality with respect to  $G_0$ , i.e., the functions  $f_j$  are not constant. Then,  $G_0$  is identifiable from the joint distribution.*

# Identifiability: Post Non Linear (PNL) Models

## Definition

PNL Models are introduced in Zhang and Hyvärinen (2009). A PNL is an SCM where each expresses each variable  $X_i$  as

$$X_i = g_i(f_i(\mathbf{PA}_i) + N_i), \quad i = 1, \dots, n$$



## Theorem (Bivariate Identifiability)

Assume that  $x_2 = f_2(f_1(x_1) + e_2)$  and  $x_1 = g_2(g_1(x_2) + e_1)$ . Densities and nonlinear functions are three times differentiable. We then have the following equation for every  $(x_1, x_2)$  satisfying  $\eta'' h' \neq 0$ :

$$t_1 = g_2^{-1}(x_1), \quad z_2 = f_2^{-1}(x_2), \quad h = f_1 \circ g_2, \quad h_1 = g_1 \circ f_2$$

$$\eta_1(t_1) = \log p_{t_1}(t_1) \quad \eta_2(e_2) = \log p_{e_2}(e_2)$$

$$\eta_1''' - \frac{\eta_1'' h''}{h'} = \left( \frac{\eta_2' \eta_2'''}{\eta_2''} - 2\eta_2'' \right) \cdot h' h'' - \frac{\eta_2'''}{\eta_2''} \cdot h' \eta_1'' + \eta_2' \cdot \left( h''' - \frac{h''^2}{h'} \right)$$

and  $h_1$  depends on  $\eta_1$ ,  $\eta_2$ , and  $h$  in the following way:

$$\frac{1}{h_1'} = \frac{\eta_1'' + \eta_2'' h'^2 - \eta_2' h''}{\eta_2'' h'}$$

	$p_{e_2}$	$p_{t_1} (t_1 = g_2^{-1}(x_1))$	$h = f_1 \circ g_2$	Remark
I	Gaussian	Gaussian	linear	$h_1$ also linear
II	log-mix-lin-exp	log-mix-lin-exp	linear	$h_1$ strictly monotonic, and $h'_1 \rightarrow 0$ , as $z_2 \rightarrow +\infty$ or as $z_2 \rightarrow -\infty$
III	log-mix-lin-exp	one-sided asymptotically exponential (but not log-mix-lin-exp)	$h$ strictly monotonic, and $h' \rightarrow 0$ , as $t_1 \rightarrow +\infty$ or as $t_1 \rightarrow -\infty$	—
IV	log-mix-lin-exp	generalized mixture of two exponentials	Same as above	—
V	generalized mixture of two exponentials	two-sided asymptotically exponential	Same as above	—

Figure: All unidentifiable cases with the assumptions made above

# Score Based Method

# RESIT Algorithm

- First proposed in Peters et al. (2014)
- Assumption : Multivariate ANM + Causal Sufficiency
- Idea :  $X_i$  is sink  $\iff N_i \perp\!\!\!\perp \mathbf{X} \setminus \{X_i\}$
- There are two stages in the algorithm:
  - Stage 1 : Finding a causal order
  - Stage 2 : Estimating DAG by removing edges
- Number of Tests (Less than PC)
  - Stage 1 :  $O(n^2)$
  - Stage 2 :  $O(n)$

# RESIT Algorithm

---

**Algorithm 1** Regression with subsequent independence test (RESIT)
 

---

```

1: Input: I.i.d. samples of a  $p$ -dimensional distribution on  $(X_1, \dots, X_p)$ 
2:  $S := \{1, \dots, p\}, \pi := []$ 
3: PHASE 1: Determine topological order.
4: repeat
5:   for  $k \in S$  do
6:     Regress  $X_k$  on  $\{X_i\}_{i \in S \setminus \{k\}}$ .
7:     Measure dependence between residuals and  $\{X_i\}_{i \in S \setminus \{k\}}$ .
8:   end for
9:   Let  $k^*$  be the  $k$  with the weakest dependence.
10:   $S := S \setminus \{k^*\}$ 
11:   $\text{pa}(k^*) := S$ 
12:   $\pi := [k^*, \pi]$  ( $\pi$  will be the topological order, its last component being a sink)
13: until  $\#S = 0$ 
14: PHASE 2: Remove superfluous edges.
15: for  $k \in \{2, \dots, p\}$  do
16:   for  $\ell \in \text{pa}(\pi(k))$  do
17:     Regress  $X_{\pi(k)}$  on  $\{X_i\}_{i \in \text{pa}(\pi(k)) \setminus \{\ell\}}$ .
18:     if residuals are independent of  $\{X_i\}_{i \in \{\pi(1), \dots, \pi(k-1)\}}$  then
19:        $\text{pa}(\pi(k)) := \text{pa}(\pi(k)) \setminus \{\ell\}$ 
20:     end if
21:   end for
22: end for
23: Output:  $(\text{pa}(1), \dots, \text{pa}(p))$ 

```

---

# RESIT Algorithm : Performance (Linear Setting)

$$\beta_{jk} \sim [-2, -0.1] \cup [0.1, 2] \quad N_j \sim K_j \cdot \text{sign}(M_j) \cdot |M_j|^{\alpha_j} \text{ such that } M_j \sim N(0, 1), \\ K_j \sim U(0.1, 0.5) \text{ and } \alpha_j \sim U([2, 4]).$$

	GDS	BF	RESIT	LINGAM	PC	CPC	GES	RAND
$p = 4, n = 100$								
DAG	$0.7 \pm 0.9$	$0.6 \pm 0.8$	$1.2 \pm 1.3$	$1.9 \pm 1.2$	$3.5 \pm 1.5$	$3.6 \pm 1.4$	$3.1 \pm 1.7$	$4.4 \pm 1.0$
CPDAG	$1.1 \pm 1.5$	$0.9 \pm 1.4$	$1.5 \pm 1.7$	$2.4 \pm 1.5$	$2.4 \pm 1.7$	$2.3 \pm 1.6$	$2.0 \pm 2.0$	$4.3 \pm 1.4$
$p = 4, n = 500$								
DAG	$0.2 \pm 0.6$	$0.1 \pm 0.3$	$0.6 \pm 0.8$	$0.5 \pm 0.8$	$3.1 \pm 1.4$	$3.2 \pm 1.4$	$2.9 \pm 1.6$	$4.1 \pm 1.2$
CPDAG	$0.3 \pm 0.9$	$0.2 \pm 0.5$	$0.9 \pm 1.3$	$0.8 \pm 1.2$	$1.9 \pm 1.8$	$1.6 \pm 1.7$	$1.6 \pm 1.9$	$3.9 \pm 1.4$
$p = 15, n = 100$								
DAG	$12.2 \pm 5.3$	—	$25.2 \pm 8.3$	$11.1 \pm 3.7$	$13.0 \pm 3.6$	$13.7 \pm 3.7$	$12.7 \pm 4.2$	$57.4 \pm 26.4$
CPDAG	$13.2 \pm 5.4$	—	$27.0 \pm 8.5$	$12.4 \pm 3.9$	$10.7 \pm 3.5$	$10.8 \pm 3.8$	$12.4 \pm 4.9$	$58.5 \pm 27.1$
$p = 15, n = 500$								
DAG	$6.1 \pm 6.4$	—	$51.2 \pm 17.8$	$3.4 \pm 2.8$	$10.2 \pm 3.8$	$10.8 \pm 4.2$	$8.7 \pm 4.6$	$57.6 \pm 24.2$
CPDAG	$6.8 \pm 6.9$	—	$54.5 \pm 18.5$	$4.5 \pm 3.8$	$8.2 \pm 4.6$	$7.5 \pm 4.4$	$7.1 \pm 5.6$	$58.9 \pm 25.0$

Figure: Structural Hamming Distance of Estimated Graph

# RESIT Algorithm : Performance (Non Linear Setting)

Functions sampled from a Gaussian process with  $BW = 1$ . Gaussian Noise with random variance.

	GDS	BF	RESIT	LiNGAM	PC	CPC	GES	RAND
$p = 4, n = 100$								
DAG	$1.5 \pm 1.4$	$1.0 \pm 1.0$	$1.7 \pm 1.3$	$3.5 \pm 1.2$	$3.5 \pm 1.5$	$3.8 \pm 1.4$	$3.5 \pm 1.3$	$4.0 \pm 1.3$
CPDAG	$1.7 \pm 1.7$	$1.2 \pm 1.4$	$2.0 \pm 1.6$	$3.0 \pm 1.4$	$2.9 \pm 1.5$	$2.7 \pm 1.4$	$3.4 \pm 1.7$	$3.9 \pm 1.4$
$p = 4, n = 500$								
DAG	$0.5 \pm 0.9$	$0.3 \pm 0.5$	$0.8 \pm 0.9$	$3.7 \pm 1.2$	$3.5 \pm 1.5$	$3.8 \pm 1.5$	$3.3 \pm 1.5$	$4.1 \pm 1.2$
CPDAG	$0.6 \pm 1.1$	$0.6 \pm 1.0$	$1.0 \pm 1.3$	$3.0 \pm 1.7$	$3.1 \pm 1.9$	$2.8 \pm 1.8$	$3.4 \pm 1.9$	$3.8 \pm 1.6$
$p = 15, n = 100$								
DAG	$14.3 \pm 4.9$	—	$15.4 \pm 5.7$	$15.4 \pm 3.6$	$14.2 \pm 3.5$	$15.5 \pm 3.6$	$24.8 \pm 6.3$	$56.8 \pm 24.1$
CPDAG	$15.1 \pm 5.4$	—	$16.5 \pm 5.9$	$15.3 \pm 4.0$	$13.3 \pm 3.6$	$13.3 \pm 4.0$	$26.4 \pm 6.5$	$58.0 \pm 24.7$
$p = 15, n = 500$								
DAG	$13.0 \pm 8.4$	—	$10.1 \pm 5.7$	$21.4 \pm 6.9$	$13.9 \pm 4.5$	$15.1 \pm 4.8$	$26.8 \pm 8.5$	$56.1 \pm 26.8$
CPDAG	$14.2 \pm 9.2$	—	$11.3 \pm 6.3$	$21.1 \pm 7.3$	$13.7 \pm 4.9$	$13.4 \pm 5.1$	$28.6 \pm 8.8$	$57.0 \pm 27.3$

Figure: Structural Hamming Distance of Estimated Graph

# References

- Patrik O. Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *Advances in Neural Information Processing Systems 21*, pages 689–696. 2009.
- Jonas Peters, Joris M. Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with continuous additive noise models. *J. Mach. Learn. Res.*, 15(1):2009–2053, 2014.
- Kun Zhang and Aapo Hyvärinen. On the identifiability of the post-nonlinear causal model. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI '09*, pages 647–655, 2009.



# The End