

SEQUENTIAL DEPENDENCE DECOMPOSITION & FEATURE LEARNING

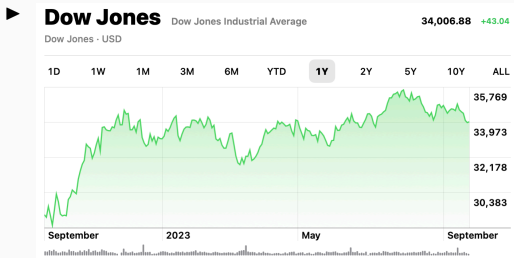
Xiangxiang Xu

Joint work with Prof. Lihong Zheng

Allerton Conference 2023



SEQUENCES



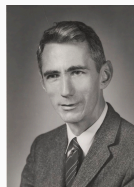
- ▶ Videos
- ▶ example of dependence in plain text



SHANNON'S EXPERIMENTS [1948]

Guessing Next Letter

1. **E** \implies solve **E** **?**



Claude E Shannon. A Mathematical Theory of Communication. 1948.

Guessing Next Letter

1. **E** \implies solve **E**?

- ▶ count the transition: **DEPENDENCE**
 - ▷ "E" followed by N: 0.5, P: 0.25, 'SPC': 0.25



Guessing Next Letter

1. **E** \implies solve **E** **?**

2. **E N** \implies solve **N** **?**

▶ count the transition: **DEPENDENCE**

▷ "E" followed by N: 0.5, P: 0.25, 'SPC': 0.25

▶ predict based on the last letter and digram frequency



Guessing Next Letter

1. **E** \implies solve **E** **?**

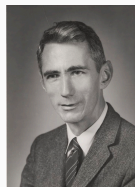
2. **E N** \implies solve **N** **?**

▶ count the transition: **DEPENDENCE**

▷ "E" followed by N: 0.5, P: 0.25, 'SPC': 0.25

▶ predict based on the last letter and digram frequency

▶ generate a Markov chain of letters



Guessing Next Letter

1. **E** \implies solve **E** **?**

2. **E N** \implies solve **N** **?**

▶ count the transition: **DEPENDENCE**

▷ "E" followed by N: 0.5, P: 0.25, 'SPC': 0.25

▶ predict based on the last letter and digram frequency

▶ generate a Markov chain of letters

▷ word-level construction

▷ trigram, n -gram (higher order Markov chains)



Guessing Next Letter

Second-order approximation (digram structure as in English).

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CTISBE

Third-order approximation (trigram structure as in English).

IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID
PONDENOME OF DEMONSTURES OF THE REPTAGIN IS
REGOACTIONA OF CRE



SHANNON'S EXPERIMENTS [1948]

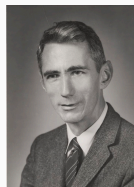
Guessing Next Letter

Second-order approximation (digram structure as in English).

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CTISBE

Third-order approximation (trigram structure as in English).

IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID
PONDENOME OF DEMONSTURES OF THE REPTAGIN IS
REGOACTIONA OF CRE



Operations

- frequency tables

264

SECRET AND URGENT

TABLE XII

ENGLISH TRIGRAMS

Figures represent approximate frequencies for 20,000 words, but are more important in their relation to one another.

A	agn.....	2	ame.....	53	aro.....	7	
Aam.....	1	ago.....	6	amh.....	1	arp.....	1
	agr.....	8	ami.....	6	arr.....	19	
	agu.....	2	aml.....	1	ars.....	46	
Abs.....	6	amn.....	1	art.....	48		
abd.....	1	Ahe.....	2	arv.....	3		
abe.....	1	aho.....	1	amp.....	10	ary.....	34
abl.....	3	ahu.....	1	ams.....	6		
abl.....	39	ahy.....	1	amy.....	1		
abo.....	28					Asa.....	1
abr.....	1	Aid.....	24	Ana.....	6	asc.....	3
abs.....	1	aig.....	3	anc.....	39	ase.....	20

Fletcher Pratt. Secret and Urgent. 1939.

Claude E Shannon. A Mathematical Theory of Communication. 1948.

SHANNON'S EXPERIMENTS [1948]

Guessing Next Letter

Second-order approximation (digram structure as in English).

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CTISBE

Third-order approximation (trigram structure as in English).

IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID
PONDENOME OF DEMONSTURES OF THE REPTAGIN IS
REGOACTIONA OF CRE



Operations

- ▶ frequency tables
- ▶ *opens a book...*
select a letter at random

264

SECRET AND URGENT

TABLE XII

ENGLISH TRIGRAMS

Figures represent approximate frequencies for 20,000 words, but are more important in their relation to one another.

A	agn.....	2	ame.....	53	aro.....	7	
Aam.....	1	ago.....	6	amh.....	1	arp.....	1
	agr.....	8	ami.....	6	arr.....	19	
Abs.....	-6	agu.....	2	aml.....	1	ars.....	46
abd.....	1	Ahe.....	2	amn.....	1	art.....	48
abe.....	1	Aho.....	1	amo.....	8	arv.....	3
abl.....	3	Ahu.....	1	amp.....	10	ary.....	34
abl.....	39	Ahy.....	1	ams.....	6		
abo.....	28			amy.....	1	Asa.....	1
abr.....	1	Aid.....	24	Ana.....	6	asc.....	3
abs.....	1	aig.....	3	anc.....	39	ase.....	20

Fletcher Pratt. Secret and Urgent. 1939.

Claude E Shannon. A Mathematical Theory of Communication. 1948.

- ▶ n -gram table on $\mathcal{X} \implies |\mathcal{X}|^n$ entries
- ▶ “book”: may not contain all combinations

- ▶ n -gram table on $\mathcal{X} \implies |\mathcal{X}|^n$ entries
- ▶ “book”: may not contain all combinations

Deep Learning Solution

- ▶ entries are structured

- ▶ n -gram table on $\mathcal{X} \implies |\mathcal{X}|^n$ entries
- ▶ “book”: may not contain all combinations

Deep Learning Solution

- ▶ entries are structured
- ▶ have “table” parameterized
 - ▷ deep neural nets: LSTM, Transformer
- ▶ learn parameters from “books”



- ▶ n -gram table on $\mathcal{X} \implies |\mathcal{X}|^n$ entries
- ▶ “book”: may not contain all combinations

Deep Learning Solution

- ▶ entries are structured
- ▶ have “table” parameterized
 - ▷ deep neural nets: LSTM, Transformer
- ▶ learn parameters from “books”

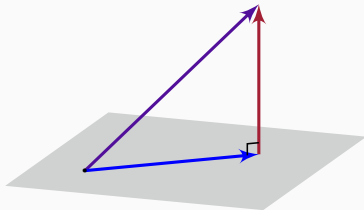


Hidden Parts?

- ▶ black-box features
- ▶ dependence structure, e.g., the “order”



FEATURE GEOMETRY



$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{\mathcal{X} \rightarrow \mathbb{R}\}$

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{\mathcal{X} \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.
 - ▷ P_X : Metric (distribution)

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.
 - ▷ P_X : Metric (distribution)
- ▶ Inner product spaces $\mathcal{F}_X[P_X], \mathcal{F}_Y[P_Y]$

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.
 - ▷ P_X : Metric (distribution)
- ▶ Inner product spaces $\mathcal{F}_X[P_X], \mathcal{F}_Y[P_Y]$

Joint functions $\mathcal{F}_{X \times Y}[P_X P_Y]$

- ▶ Metric $P_X P_Y$ helps decouple the dependence

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.
 - ▷ P_X : Metric (distribution)
- ▶ Inner product spaces $\mathcal{F}_X[P_X], \mathcal{F}_Y[P_Y]$

Joint functions $\mathcal{F}_{X \times Y}[P_X P_Y]$

- ▶ Metric $P_X P_Y$ helps decouple the dependence
- ▶ Product of $f \in \mathcal{F}_X, g \in \mathcal{F}_Y$: $f \otimes g: (x, y) \mapsto f(x)g(y)$

$(X, Y) \sim P_{X,Y}$ variables of interest, e.g., input and output pair

Feature Space

- ▶ Features of X : $\mathcal{F}_X \triangleq \{X \rightarrow \mathbb{R}\}$
- ▶ Inner Product $\langle f_1, f_2 \rangle \triangleq \mathbb{E}_{P_X} [f_1(X)f_2(X)]$
 - ▷ Induced geometry: norm, orthogonality, projection, etc.
 - ▷ P_X : Metric (distribution)
- ▶ Inner product spaces $\mathcal{F}_X[P_X], \mathcal{F}_Y[P_Y]$

Joint functions $\mathcal{F}_{X \times Y}[P_X P_Y]$

- ▶ Metric $P_X P_Y$ helps decouple the dependence
- ▶ Product of $f \in \mathcal{F}_X, g \in \mathcal{F}_Y$: $f \otimes g: (x, y) \mapsto f(x)g(y)$

- ▶ $f \otimes g \triangleq \sum_{i=1}^k f_i \otimes g_i$ for k -dim features

$$\begin{aligned} f &= (f_1, \dots, f_k) \\ g &= (g_1, \dots, g_k) \end{aligned}$$

Canonical Dependence Kernel (CDK)

$$i_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

- ▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$
- ▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp Y$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

- ▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$
- ▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp\!\!\!\perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

- ▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$
- ▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp\!\!\!\perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ learning $\mathbf{i}_{X;Y}$ by maximizing the H-score

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{\mathcal{X} \times \mathcal{Y}}$

▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp\!\!\!\perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ learning $\mathbf{i}_{X;Y}$ by maximizing the H-score
- ▶ efficiently computable from data samples $\{(x_i, y_i)\}_{i=1}^n$

$$\mathcal{H} \stackrel{k=1}{=} \text{cov}(f(X), g(Y)) - \frac{1}{2} \cdot \mathbb{E}[f^2(X)] \cdot \mathbb{E}[g^2(Y)]$$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

$$\triangleright \mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$$

$$\triangleright \|\mathbf{i}_{X;Y}\| = 0 \text{ iff } X \perp\!\!\!\perp Y$$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ learning $\mathbf{i}_{X;Y}$ by maximizing the H-score
- ▶ efficiently computable from data samples $\{(x_i, y_i)\}_{i=1}^n$

$$\mathcal{H} \stackrel{k=1}{=} \underbrace{\text{cov}(f(X), g(Y))}_{\frac{1}{n} \sum_i f(x_i) g(y_i)} - \frac{1}{2} \cdot \underbrace{\mathbb{E}[f^2(X)]}_{\frac{1}{n} \sum_i f^2(x_i)} \cdot \mathbb{E}[g^2(Y)]$$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

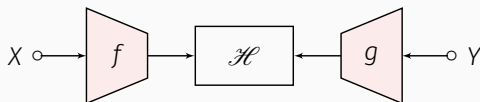
▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$

▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ apply neural feature extractors to process data



Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

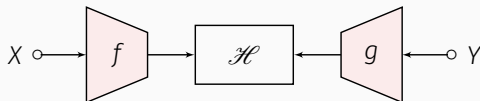
▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$

▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp\!\!\!\perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ apply neural feature extractors to process data



- ▶ maximize $\mathcal{H}(f, g) \implies f \otimes g = \mathbf{i}_{X;Y}$

Canonical Dependence Kernel (CDK)

$$\mathbf{i}_{X;Y} = \frac{P_{X,Y}}{P_X P_Y} - 1$$

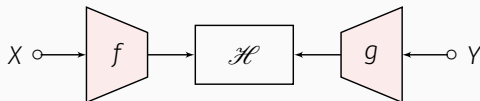
▷ $\mathbf{i}_{X;Y} \in \mathcal{F}_{X \times Y}$

▷ $\|\mathbf{i}_{X;Y}\| = 0$ iff $X \perp\!\!\!\perp Y$

Learning CDK from Features

$$\mathcal{H}(f, g) \triangleq \frac{1}{2} \left(\|\mathbf{i}_{X;Y}\|^2 - \|\mathbf{i}_{X;Y} - f \otimes g\|^2 \right)$$

- ▶ apply neural feature extractors to process data



- ▶ maximize $\mathcal{H}(f, g) \implies f \otimes g = \mathbf{i}_{X;Y}$

▷ compute $\|\mathbf{i}_{X;Y}\|$ from features

▷ prediction/estimation

$$P_{X|Y}(x|y) = P_X(x) \cdot [1 + f(x) \cdot g(y)]$$

ONE MORE VARIABLE

Learn X based on Y : $i_{X;Y}$

ONE MORE VARIABLE

Learn X based on

- ▶ Y : $\mathbf{i}_{X;Y}$
- ▶ (Y, Z) : $\mathbf{i}_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z}[P_X P_{Y,Z}]$

Learn X based on

- ▶ Y : $\mathbf{i}_{X;Y}$
- ▶ (Y, Z) : $\mathbf{i}_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z}[P_X P_{Y,Z}]$

Contribution of Z

$$\mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$$

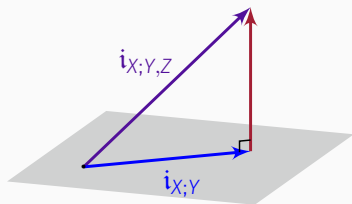
ONE MORE VARIABLE

- Learn X based on
- ▶ Y : $i_{X;Y}$
 - ▶ (Y, Z) : $i_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z} [P_X P_{Y,Z}]$

Contribution of Z

$$i_{X;Y,Z} - i_{X;Y}$$



Markov Plane

- ▶ i : $X - Y - Z$

ONE MORE VARIABLE

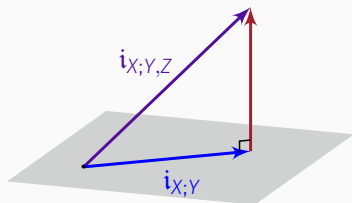
- Learn X based on
- ▶ Y : $i_{X;Y}$
 - ▶ (Y, Z) : $i_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z} [P_X P_{Y,Z}]$

Contribution of Z

$$i_{X;Y,Z} - i_{X;Y}$$

$i_{X;Y}$ Markov Component



Markov Plane

▶ $i: X - Y - Z$

ONE MORE VARIABLE

- Learn X based on
- ▶ Y : $\mathbf{i}_{X;Y}$
 - ▶ (Y, Z) : $\mathbf{i}_{X;Y,Z}$

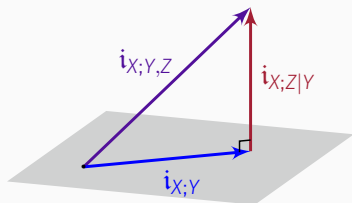
Space $\mathcal{F}_{x \times y \times z} [P_X P_{Y,Z}]$

Contribution of Z

$$\mathbf{i}_{X;Z|Y} \triangleq \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$$

$\mathbf{i}_{X;Y}$ Markov Component

$\mathbf{i}_{X;Z|Y}$ Conditional Dependence



Markov Plane

▶ $\mathbf{i}: X - Y - Z$

ONE MORE VARIABLE

- Learn X based on
- ▶ Y : $\mathbf{i}_{X;Y}$
 - ▶ (Y, Z) : $\mathbf{i}_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z} [P_X P_{Y,Z}]$

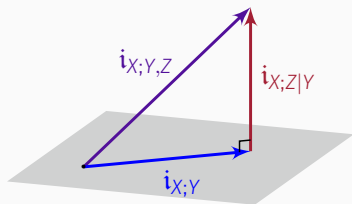
Contribution of Z

$$\mathbf{i}_{X;Z|Y} \triangleq \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$$

$\mathbf{i}_{X;Y}$ Markov Component

$\mathbf{i}_{X;Z|Y}$ Conditional Dependence

- ▶ $\|\mathbf{i}_{X;Z|Y}\| = 0$ iff $X \perp\!\!\!\perp Z|Y$



Markov Plane

- ▶ $\mathbf{i} : X - Y - Z$

ONE MORE VARIABLE

- Learn X based on
- ▶ Y : $\mathbf{i}_{X;Y}$
 - ▶ (Y, Z) : $\mathbf{i}_{X;Y,Z}$

Space $\mathcal{F}_{X \times Y \times Z} [P_X P_{Y,Z}]$

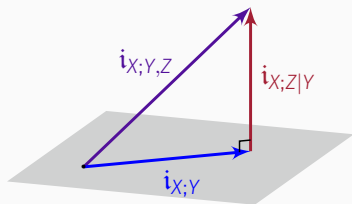
Contribution of Z

$$\mathbf{i}_{X;Z|Y} \triangleq \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$$

$\mathbf{i}_{X;Y}$ Markov Component

$\mathbf{i}_{X;Z|Y}$ Conditional Dependence

- ▶ $\|\mathbf{i}_{X;Z|Y}\| = 0$ iff $X \perp\!\!\!\perp Z|Y$



Markov Plane

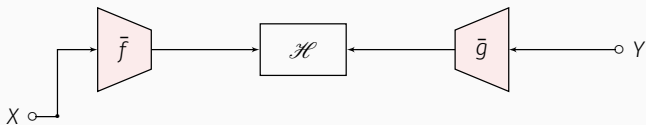
▷ $\mathbf{i} : X - Y - Z$

$$\|\mathbf{i}_{X;Y,Z}\|^2 = \|\mathbf{i}_{X;Y}\|^2 + \|\mathbf{i}_{X;Z|Y}\|^2$$

$\mathbf{i}_{X;Z|Y} = \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$: dependence Y cannot capture

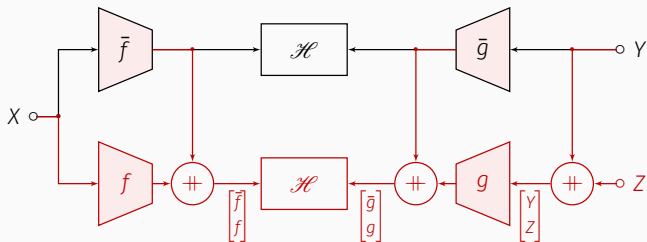
LEARNING CONDITIONAL DEPENDENCE

$\mathbf{i}_{X;Z|Y} = \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$: dependence Y cannot capture



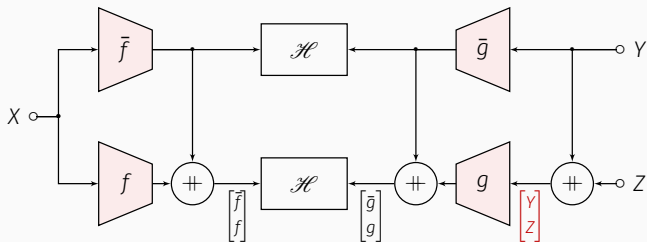
LEARNING CONDITIONAL DEPENDENCE

$\mathbf{i}_{X;Z|Y} = \mathbf{i}_{X;Y,Z} - \mathbf{i}_{X;Y}$: dependence Y cannot capture



LEARNING CONDITIONAL DEPENDENCE

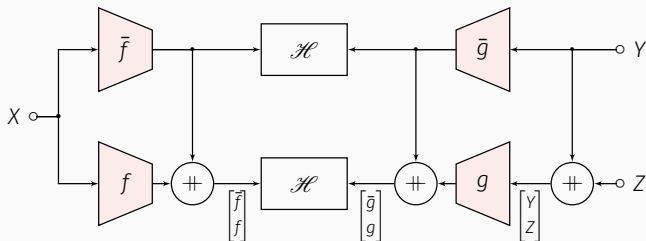
$i_{X;Z|Y} = i_{X;Y,Z} - i_{X;Y}$: dependence Y cannot capture



- nesting: separate conditional dependence from the joint

LEARNING CONDITIONAL DEPENDENCE

$i_{X;Z|Y} = i_{X;Y,Z} - i_{X;Y}$: dependence Y cannot capture



- ▶ nesting: separate conditional dependence from the joint
- ▶ training: maximize the sum of two H-scores
 - ▷ optimal solution: $\bar{f} \otimes \bar{g} = i_{X;Y}$, $f \otimes g = i_{X;Z|Y}$
 - ▷ measure the strength of conditional dependence

... X_{-3} X_{-2} X_{-1} X_0 ...

X_{-1} X_0 \dots

“Looking Back”

- ▶ previous state: $i_{X_0;X_{-1}}$

Learn X_0 based on

X_{-2} X_{-1} X_0 \dots

“Looking Back”

Learn X_0 based on

- ▶ previous state: $\mathbf{i}_{X_0;X_{-1}}$
- ▶ past 2 states: $\mathbf{i}_{X_0;(X_{-1},X_{-2})}$

... X_{-3} X_{-2} X_{-1} X_0 ...

“Looking Back”

Learn X_0 based on

- ▶ previous state: $i_{X_0; X_{-1}}$
- ▶ past 2 states: $i_{X_0; (X_{-1}, X_{-2})}$
- ▶ \vdots
- ▶ past n states: $i_{X_0; (X_{-1}, \dots, X_{-n})}$

... X_{-3} X_{-2} X_{-1} X_0 ...

“Looking Back”

Learn X_0 based on

- ▶ previous state: $\mathbf{i}_{X_0; X_{-1}}$
- ▶ past 2 states: $\mathbf{i}_{X_0; (X_{-1}, X_{-2})}$
- ⋮
- ▶ past n states: $\mathbf{i}_{X_0; (X_{-1}, \dots, X_{-n})}$

- ▶ Gain from ℓ -th layer $\mathbf{i}_\ell \triangleq \mathbf{i}_{X_0; X_{-\ell} | (X_{-1}, \dots, X_{-\ell+1})}$
 - conditional dependence at lag ℓ

... X_{-3} X_{-2} X_{-1} X_0 ...

“Looking Back”

Learn X_0 based on

- ▶ previous state: $\mathbf{i}_{X_0; X_{-1}}$
- ▶ past 2 states: $\mathbf{i}_{X_0; (X_{-1}, X_{-2})}$
- ⋮
- ▶ past n states: $\mathbf{i}_{X_0; (X_{-1}, \dots, X_{-n})}$

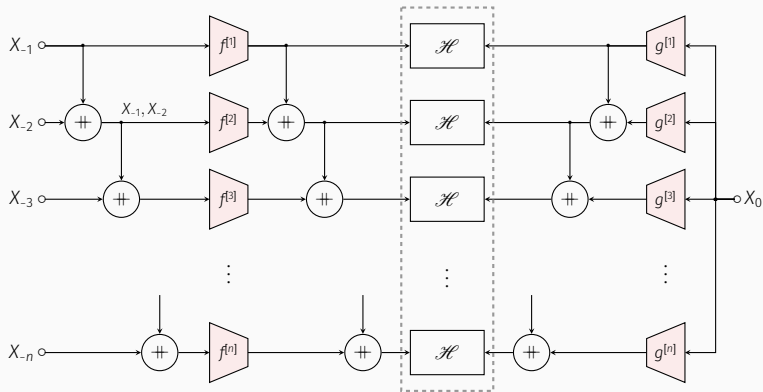
- ▶ Gain from ℓ -th layer $\mathbf{i}_\ell \triangleq \mathbf{i}_{X_0; X_{-\ell} | (X_{-1}, \dots, X_{-l+1})}$

- ▷ conditional dependence at lag ℓ

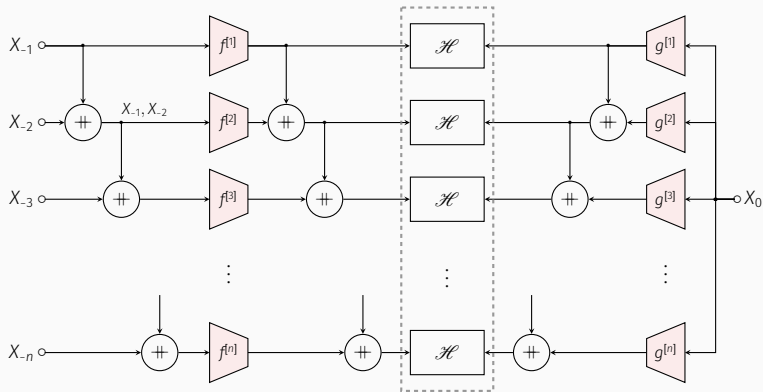
- ▶ Orthogonal decomposition

- ▷ Dependence between X_0 and past n states = $\sum_{\ell=1}^n \mathbf{i}_\ell$

LEARNING THE DECOMPOSITION

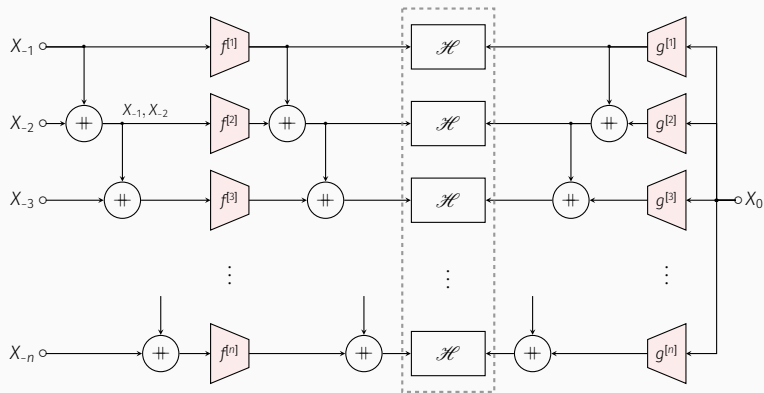


LEARNING THE DECOMPOSITION



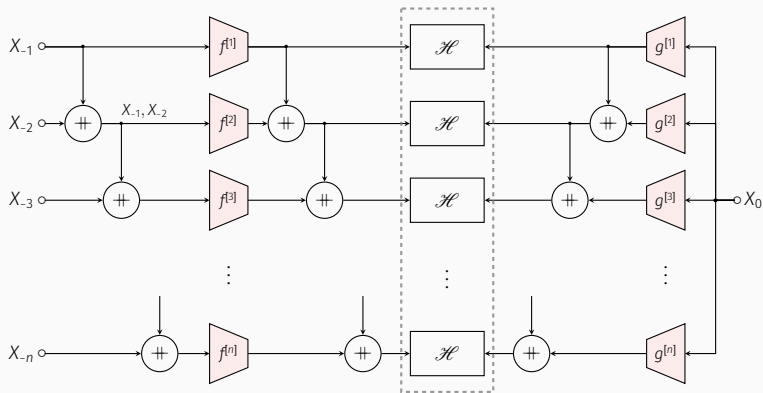
- l -th branch learns i_ℓ : conditional dependence at lag l

LEARNING THE DECOMPOSITION



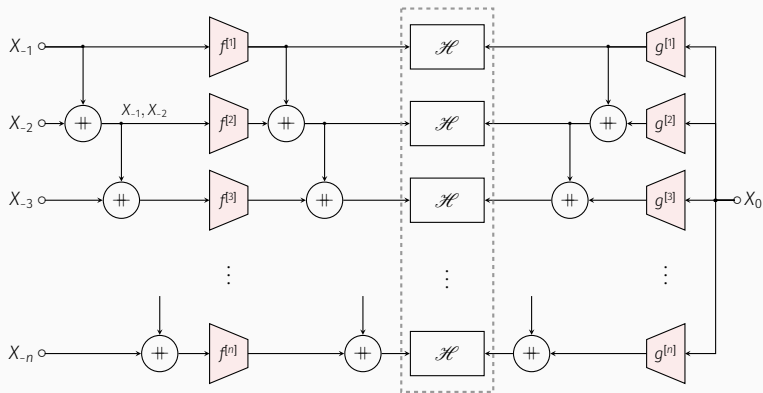
- ▶ l -th branch learns i_ℓ : conditional dependence at lag l
 - ▷ top l branches: dependence between X_0 and past l states

LEARNING THE DECOMPOSITION



- ▶ l -th branch learns \mathbf{i}_l : conditional dependence at lag l
 - ▷ top l branches: dependence between X_0 and past l states
- ▶ dependence “spectrum” over lags: $\{\|\mathbf{i}_l\|^2, l \geq 1\}$

LEARNING THE DECOMPOSITION



- ▶ l -th branch learns \mathbf{i}_l : conditional dependence at lag l
 - ▷ top l branches: dependence between X_0 and past l states
- ▶ dependence “spectrum” over lags: $\{\|\mathbf{i}_l\|^2, l \geq 1\}$
 - ▷ Markov Chain of Order $M \implies$ cutoff at $l = M$

Sequence Observations

- ▶ Dependence on the History?



Sequence Observations

- ▶ Dependence on the History?
- ▶ First/Second/Third-Order?



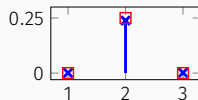
Sequence Observations

- ▶ Dependence on the History?
- ▶ First/Second/Third-Order?
- ★ Plot Dependence Spectrum $\|\mathbf{i}_\ell\|^2, \ell \geq 1$



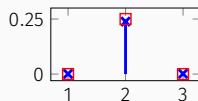
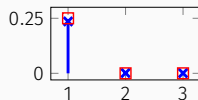
Sequence Observations

- ▶ Dependence on the History?
- ▶ First/Second/Third-Order?
- ★ Plot Dependence Spectrum $\|i_\ell\|^2, \ell \geq 1$



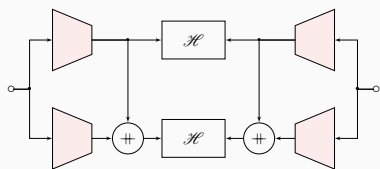
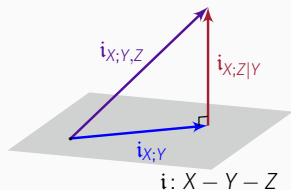
Sequence Observations

- ▶ Dependence on the History?
- ▶ First/Second/Third-Order?
- ★ Plot Dependence Spectrum $\|i_\ell\|^2, \ell \geq 1$



SUMMARY

SUMMARY



- ▶ Feature Geometry
 - ▷ Feature Learning \leftrightarrow Geometric Operations
 - ▷ Nesting Technique
- ▶ Case Study: Learning Random Processes
 - ▷ Decompose Sequential Dependence

LEARN MORE

- ▶ arXiv: 2309.10140

