

# Pointwise HSIC: A Linear-Time Kernelized Co-occurrence Norm for Sparse Linguistic Expressions [EMNLP2018]

Sho Yokoi (Tohoku U./RIKEN AIP)\*, Sosuke Kobayashi (PFN), Kenji Fukumizu (ISM), Jun Suzuki\*, Kentaro Inui\*

✉ yokoi@ecei.tohoku.ac.jp 📄 1809.00800 🌐 https://bit.ly/2SfpZmv 🏠 github.com/cl-tohoku/phsic 📦 pip install phsic-cli

## Computing Co-occurrence Strength in NLP

### General Settings

**Observed:** paired data  $\mathcal{D} = \{(x_i, y_i)\}$

$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_n$	$y_n$

**Task:** compute co-occurrence strength of a pair  $(x, y)$

$x$	$y$	?
-----	-----	---

### Collocation Extraction

[Manning&Schütze,'99]

word bigrams from corpora

Have	you
you	ever
$\vdots$	$\vdots$
York	?

which pairs are collocation?

New	York	✓
in	the	■

### Dialogue Response Selection

[Lowe+, '15]

input-response message pairs from SNS

I'm hungry!	Let's have lunch
Will it rain today?	It's about to rain
$\vdots$	$\vdots$
I love this manga	I don't know

which response is the best?

I've lost my wallet	I saw it at the ...	✓
I've lost my wallet	I'm so sleepy	■
I've lost my wallet	I don't know	■

input by users    response candidates

## De facto Measure: Pointwise Mutual Information

### Mutual Information

$$MI(X, Y) = KL[\mathbf{P}_{XY} \parallel \mathbf{P}_X \mathbf{P}_Y]$$

$$= \mathbf{E}_{(x,y)} \left[ \log \frac{\mathbf{P}_{XY}(x,y)}{\mathbf{P}_X(x)\mathbf{P}_Y(y)} \right]$$

### Pointwise Mutual Information

$$PMI(x, y; X, Y) = \log \frac{\mathbf{P}_{XY}(x,y)}{\mathbf{P}_X(x)\mathbf{P}_Y(y)}$$

co-occurrence actual

co-occurrence by chance

### $\widehat{PMI}_{MLE}(x, y)$

$$= \log \frac{n \cdot \#(x, y)}{\#(x, \cdot) \#(\cdot, y)}$$

just counting

- ✓ easy to learn
- ✗ inapplicable to sparse expressions

### $\widehat{PMI}_{RNN}(x, y)$

$$= \log \frac{\widehat{\mathbf{P}}_{RNN}(y|x)}{\widehat{\mathbf{P}}_{RNN}(y)}$$

using RNNs [Li+, '16]

- ✗ tough to learn
- ✓ applicable to sparse expressions

## Proposed Measure: Pointwise HSIC

### HSIC [Gretton+, '05]

$$HSIC(X, Y; k, \ell) = MMD_{k, \ell}^2[\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y]$$

$$= \mathbf{E}_{(x,y)} [(\phi(x) - m_X)^\top C_{XY} (\psi(y) - m_Y)]$$

$$= \mathbf{E}_{(x,y)} \left[ \mathbf{E}_{(x',y')} [\tilde{k}(x, x') \tilde{\ell}(y, y')] \right]$$

$\phi(x) := k(x, \cdot), \psi(y) := \ell(y, \cdot)$   
 $m_X := \mathbf{E}_x[\phi(x)], m_Y := \mathbf{E}_y[\psi(y)]$   
 $C_{XY} := \mathbf{E}_{(x,y)} [(\phi(x) - m_X)(\psi(y) - m_Y)^\top]$

### Pointwise HSIC

$$PHSIC(x, y; X, Y, k, \ell) = (\phi(x) - m_X)^\top C_{XY} (\psi(y) - m_Y)$$

$$= \mathbf{E}_{(x',y')} [\tilde{k}(x, x') \tilde{\ell}(y, y')]$$

Kernels centered on data  
 $\tilde{k}(x, x') := k(x, x') - \mathbf{E}_x[k(x, x')] - \mathbf{E}_{x'}[k(x, x')] + \mathbf{E}_{x, x'}[k(x, x')]$

### Requires Very Short Learning Time

estimators are reduced to a simple matrix calculation

when using cosine similarity → can estimate in RKHS directly

$$PHSIC_{RKHS}(x, y; k, \ell) = (\phi(x) - \overline{\phi(x)})^\top \widehat{C}_{XY} (\psi(y) - \overline{\psi(y)})$$

when using non-linear kernels → use incomplete Cholesky decomposition

$$PHSIC_{ICD}(x, y; k, \ell) = (a - \bar{a})^\top \widehat{C}_{ICD} (b - \bar{b})$$

### Dependence of $(X, Y)$

the difference between  $\mathbf{P}_{XY}$  and  $\mathbf{P}_X \mathbf{P}_Y$

### Co-occurrence of $(x, y)$

the contribution of  $(x, y)$  to the dependence of  $(X, Y)$

## Applicable to Sparse Expressions

PHSIC smooths the matching using kernels

	add scores	deduct scores
$\widehat{PMI}(x, y; \mathcal{D}) = \log \frac{n \cdot \sum_i \mathbb{I}[x=x_i \wedge y=y_i]}{\sum_i \mathbb{I}[x=x_i] \sum_i \mathbb{I}[y=y_i]}$	$\mathcal{D} = \{ \dots, (x_i, y_i), \dots \}$	$\{ \dots, (x_i, y_i), \dots, (x_i, y_i), \dots \}$
$\widehat{PHSIC}(x, y; \mathcal{D}, k, \ell) = \frac{1}{n} \sum_i \widehat{k}(x, x_i) \widehat{\ell}(y, y_i)$	$\{ \dots, (x_i, y_i), \dots, (x_i, y_i), \dots \}$	$\{ \dots, (x_i, y_i), \dots, (x_i, y_i), \dots \}$

exact match

smooths the matching by kernels

## Experiments

### 1. Dialogue Response Selection [Lowe+, '15]

$n = 5 \times 10^5$

Method	Learning Time [s]	Predictive Performance [R@2]
PMI (w/RNN)	13054.9	0.52
PHSIC (cos of USE)	1.8	0.57

✓ thousands of times faster!    ✓ while outperforming PMI

$n = 10^3$

Method	Learning Time [s]	Predictive Performance [R@2]
PMI (w/RNN)	49.0	0.21
PHSIC (cos of USE)	0.0	0.56

✓ much more robust to data sparsity

- ### ✓ Kernels are Available
- **Cosine Similarity** between **Sentence Vectors**
    - Sentence vectors [Kiros+, '15; Dai&Le, '15; Iyyer+, '15; Hill+, '16; Cer+, '18]
    - Sum of word vectors [Mikolov+, '13; Pennington+, '14; Bojanowski+, '17]
    - Many pre-trained models are off-the-shelf!
  - **Structure Kernels**
    - [Collins&Duffy, '02; Bunescu&Mooney, '06; Moschitti, '06]
  - **Combinations**

### Kernels in Experiments

- **Cosine Similarity** between Sum of **Word Vectors**  
 $k(x, x') = \cos(\sum_{w \in x} v(w), \sum_{w \in x'} v(w))$ 

fastText [Bojanowski+, '17; Grave+, '18]
- **Gaussian Kernel** between Sum of **Word Vectors**  
 $k(x, x') = \exp(-\|\sum_{w \in x} v(w) - \sum_{w \in x'} v(w)\|^2 / 2\sigma^2)$
- **Cosine Similarity** between **Sentence Vectors**  
 $k(x, x') = \cos(v(x), v(x'))$ 

Universal Sentence Encoder [Cer+'18] w/ Deep Averaging Network [Iyyer+, '15]

### 2. Parallel Corpus Filtering

BLEU score after filtering (3M  $\rightarrow$  1M)

Method	BLEU score
random	39.82 (-1.20)
fast_align	40.56 (-0.46)
PHSIC (RBF kernel of fastText)	40.95 (-0.07)

using all (3M) training data  $\blacktriangle 41.02$

✓ PHSIC reduces # of training data to 1/3, almost without sacrificing BLEU