

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

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



1809.00800



<https://bit.ly/2SfpZmv>



[github.com/cl-tohoku/phsic](https://github.com/cl-tohoku/phsic)



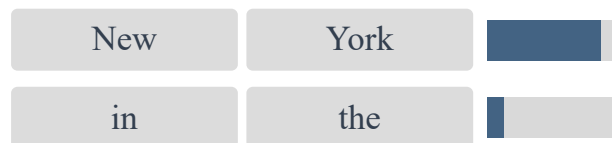
`pip install phsic-cli`

## Computing Co-occurrence Strength in NLP

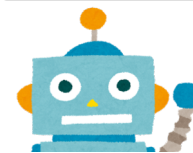
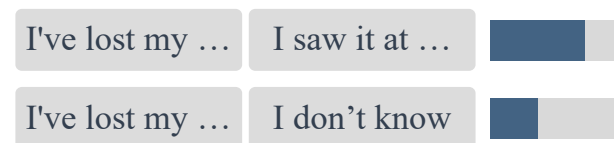
observed: paired data  $\mathcal{D} = \{(x_i, y_i)\} \sim \mathbf{P}_{XY}$

task: compute “co-occurrence/association/relation strength” of  $(x, y) \in \mathcal{X} \times \mathcal{Y}$

### Collocation Extraction



### Dialogue Response Selection



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

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



1809.00800



<https://bit.ly/2SfpZmv>



[github.com/cl-tohoku/phsic](https://github.com/cl-tohoku/phsic)



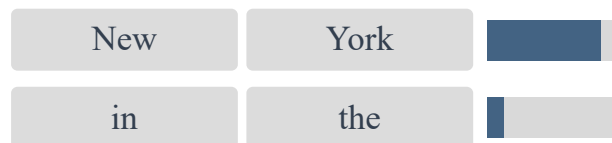
`pip install phsic-cli`

## Computing Co-occurrence Strength in NLP

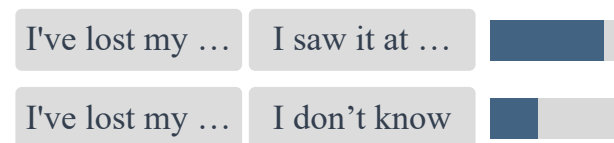
observed: paired data  $\mathcal{D} = \{(x_i, y_i)\} \sim \mathbf{P}_{XY}$

task: compute “co-occurrence/association/relation strength” of  $(x, y) \in \mathcal{X} \times \mathcal{Y}$

### Collocation Extraction



### Dialogue Response Selection



## De facto Measure: Pointwise Mutual Information

co-occurrence  
actual

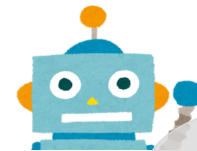
$$\text{PMI}(x, y) = \log \frac{\mathbf{P}_{XY}(x, y)}{\mathbf{P}_X(x)\mathbf{P}_Y(y)}$$

co-occurrence  
by chance

just  
counting

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

- ✓ easy to learn
- ✗ inappropriate to sparse data



using  
RNNs

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

- ✗ tough to learn
- ✓ applicable to sparse data

# Proposed Measure for Computing Co-occurrence Strength: Pointwise HSIC

Dependence of

$X$

and

$Y$

Co-occurrence of  $x$  and  $y$

## Mutual Information

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

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

contribute

## Pointwise Mutual Information

$$PMI(x, y; X, Y)$$

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

## HSIC [Gretton+'05]

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

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

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

contribute

## 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

$$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$



Requires very short learning time

1000 times faster than RNN-based PMI



Applicable to sparse objects

PHSIC allows various and available similarity metrics to be plugged in as kernels

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

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



1809.00800



<https://bit.ly/2SfpZmv>



[github.com/cl-tohoku/phsic](https://github.com/cl-tohoku/phsic)



`pip install phsic-cli`

**input:** paired data  $\mathcal{D} = \{(x_i, y_i)\} \sim \mathbf{P}_{XY}$  **goal:** compute “co-occurrence strength” of  $(x, y) \in \mathcal{X} \times \mathcal{Y}$

## Collocation Extraction

New

York

in

the

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

✓ easy to learn

✗ inappropriate to sparse data

## Dialogue Response Selection

I've lost my ...

I saw it at ...

I've lost my ...

I'm so sleepy

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

✗ tough to learn

✓ applicable to sparse data

## Proposed Measure (PHSIC)

$$\begin{aligned} \widehat{\text{PHSIC}}(x, y; k, \ell) \\ = \frac{1}{n} \sum_i \widehat{k}(x, x_i) \widehat{\ell}(y, y_i) \end{aligned}$$

✓ easy to learn

✓ applicable to sparse data