

Some kernels for structured data

- Goal: construct a similarity score for object such as
 - sequences
 - with variable length
 - by their interpretations
 - labeled graphs (or trees)
 - different size
 - different structure
 - other objects
 - by their interpretations
- The similarity must be a PD kernel

Rational kernels

- Compare sequences $\mathbf{x}, \mathbf{y} \in \Sigma^*$

- $\mathbf{x} = (0, 1, 1, 0, 1, 1, 0, 0)$

- $\mathbf{y} = (1, 1, 0, 1)$

- Transducer

- maps a seq. \mathbf{x} sto seq. \mathbf{z} with a weight

- defines a “weighed relation” $T(\mathbf{x}, \mathbf{z}) \rightarrow \mathbb{R}$

- is implemented by a *finite state automaton*

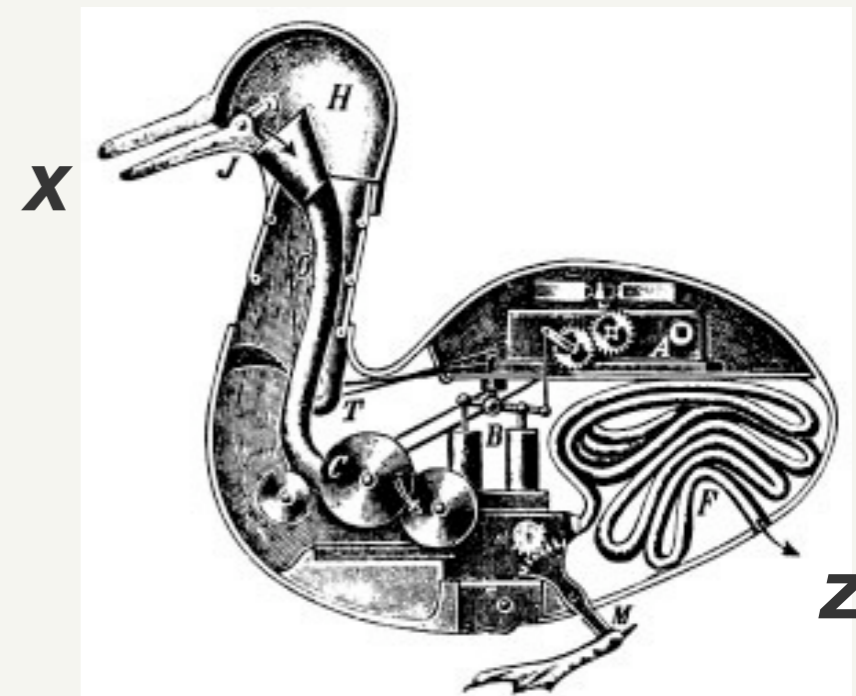
- Kernel

\mathbf{x}, \mathbf{y} are similar if they are transduced often to the same \mathbf{z}

- $K(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{z}} T(\mathbf{x}, \mathbf{z}) T(\mathbf{y}, \mathbf{z})$

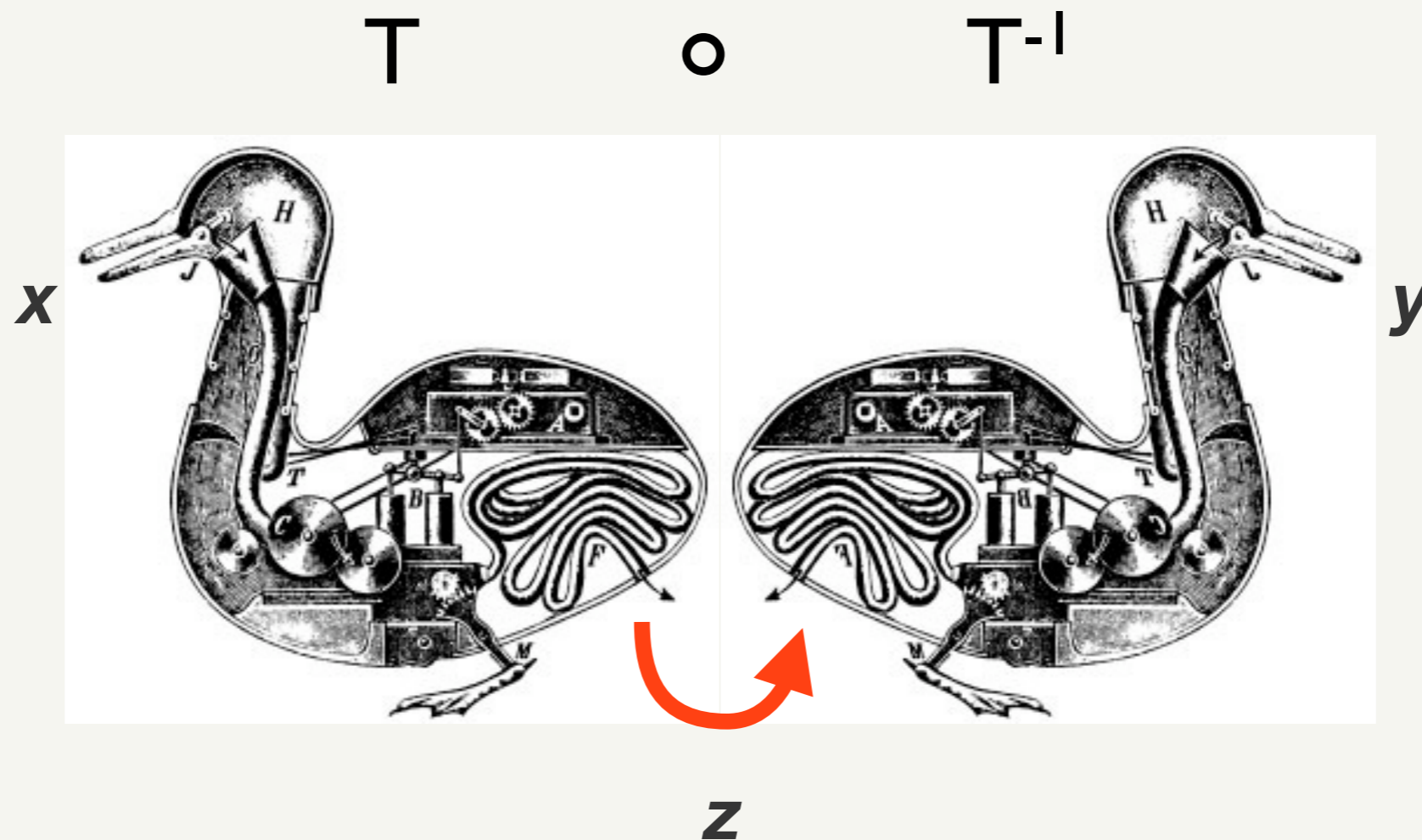
- Advantage

Given an automaton for T , can construct an automaton for K



Rational kernels: Implementation

- Automaton for $K(x, y)$
 - invert T
 - compose T and T^{-1}



Rational kernels: Examples

- Bag-of-subsequences
 - \mathbf{x} binary sequence
 - \mathbf{z} binary sequence of 4 characters
 - $T(\mathbf{x}, \mathbf{z}) = \#$ occurrences of \mathbf{z} in \mathbf{x}
 - $K(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{z}} T(\mathbf{x}, \mathbf{z}) T(\mathbf{y}, \mathbf{z})$ is large iff \mathbf{x}, \mathbf{y} contain similar subsequences

$$\begin{array}{r} \mathbf{x} = (0, 1, 1, 0, 1, 1, 0) \\ \quad \mathbf{T} \quad \quad \mathbf{z} \\ \mathbf{2} \quad 0, 1, 1, 0 \\ \mathbf{1} \quad 1, 1, 0, 1 \\ \mathbf{1} \quad 1, 0, 1, 1 \end{array}$$

$$\begin{array}{r} \mathbf{y} = (1, 1, 0, 1) \\ \quad \mathbf{T} \quad \quad \mathbf{z} \\ \mathbf{1} \quad 1, 1, 0, 1 \end{array}$$

$$K(\mathbf{x}, \mathbf{y}) = 1$$

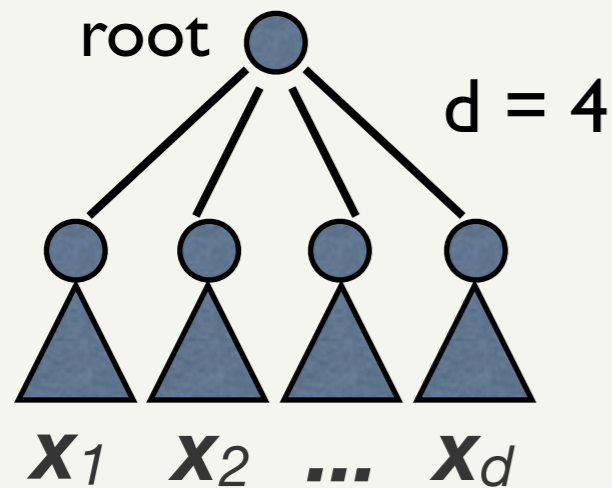
- Normalization
 $K(\mathbf{x}, \mathbf{y}) / (K(\mathbf{x}, \mathbf{x}) K(\mathbf{y}, \mathbf{y}))^{1/2}$

- Other examples
 - HMM-like models

Convolution kernels

- To compare objects \mathbf{x}, \mathbf{y}
 - decompose each object in d components
 - compare components and combine results
 - (repeat recursively until atomic components)
- Example: tree

\mathbf{x} is a d -degree tree



subtrees are
 d -degree trees or leaves

Subpart relation

$$R(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d, \mathbf{x})$$

$$K(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^d K_i(\mathbf{x}_i, \mathbf{y}_i)$$

Convolution kernels

- Example: string

- \mathbf{x} is a string
- Subpart relation



$R(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x})$ iif

$\mathbf{x}_1, \mathbf{x}_2$ are (non-empty) strings such that $\mathbf{x} = \text{concat}(\mathbf{x}_1, \mathbf{x}_2)$

- Multiple decompositions are possible

- $R^{-1}(\mathbf{x}) = \{ (\mathbf{x}_1, \mathbf{x}_2) : R(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}) \}$

- Convolution kernel

$$k(x, y) = \sum_{x' \in R^{-1}(x)} \sum_{y' \in R^{-1}(y)} \prod_{i=1}^r k_i(x'_i, y'_i)$$

Convolution kernels

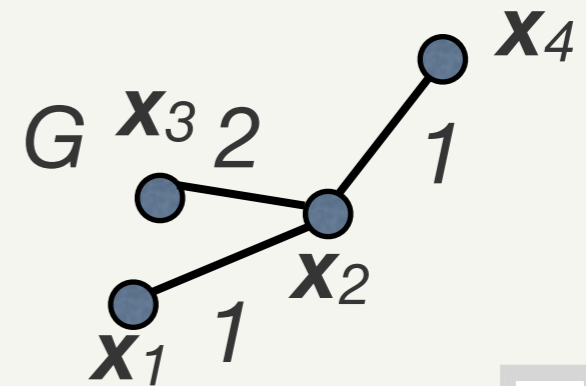
We can represent the relation " x_1, \dots, x_d are the parts of x " by a relation R on the set $X_1 \times \dots \times X_D \times X$, where $R(x_1, \dots, x_D, x)$ is true iff x_1, \dots, x_D are the parts of x . For brevity, let $\vec{x} = x_1, \dots, x_D$, and denote $R(x_1, \dots, x_D, x)$ by $R(\vec{x}, x)$. Let $R^{-1}(x) = \{\vec{x} : R(\vec{x}, x)\}$. We say R is *finite* if $R^{-1}(x)$ is finite for all $x \in X$. Here are some examples:

1. If x is a D -tuple in $X = X_1 \times \dots \times X_D$, and each component of $x \in X$ is a part of x , then $R(\vec{x}, x)$ iff $\vec{x} = x$.
2. If $X_1 = X_2 = X$, where X is the set of all finite strings over a finite alphabet \mathcal{A} , then we can define $R(x_1, x_2, x)$ iff $x_1 \circ x_2 = x$, where $x_1 \circ x_2$ denotes the concatenation of strings x_1 and x_2 .
3. Continuing the previous example, if the alphabet \mathcal{A} has only one letter, then a finite string can be represented by the nonnegative integer n that is its length, so $X_1 = X_2 = X = \{0, 1, \dots\}$ and $R(n_1, n_2, n)$ iff $n_1 + n_2 = n$.
4. If $X_1 = \dots = X_D = X$, where X is the set of all D -degree ordered and rooted trees, then we can define $R(\vec{x}, x)$ iff x_1, \dots, x_D are the D subtrees of the root of the tree $x \in X$.

Kernels based on local info

- Given

- $\{ \mathbf{x}_1, \dots, \mathbf{x}_n \}$ collection of objects
- “local” distances
formally: G undirected weighed DAG



D	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3
\mathbf{x}_1	0	1	3
\mathbf{x}_2		0	2
\mathbf{x}_3			0

- Get geodesic distances D

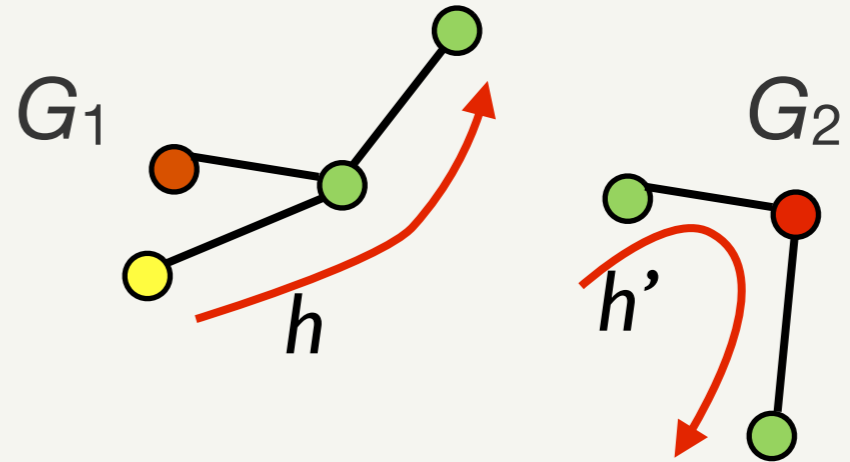
- all shortest-paths D
- regularize by finding low-dimensional embedding (ISOMAP)

- Get a kernel

- Use identity $D(\mathbf{x}_1, \mathbf{x}_2) = K(\mathbf{x}_1, \mathbf{x}_1) + K(\mathbf{x}_2, \mathbf{x}_2) - 2K(\mathbf{x}_1, \mathbf{x}_2)$
- Make *positive definite* by incrementing the diagonal
 $K \leftarrow K + \lambda I$

Graph kernels

- Compare labeled graphs $\mathbf{x}, \mathbf{y} \in \Sigma^*$
 - given a kernel on *paths*
 $k_{\text{path}}(h, h')$
 - extend to kernel on graphs
 - try to capture “topology”
- Compare all paths $W(G_1), W(G_2)$



$$k_G(G_1, G_2) = \sum_{h \in W(G_1)} \sum_{h' \in W(G_2)} k_{\text{path}}(h, h').$$

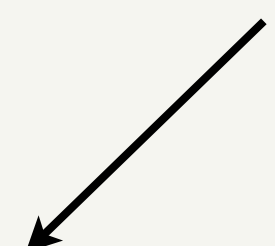
- walks (any path)
 - proper paths (no self intersection)
 - shortest paths
 - random walks
- Suggested reading :-)
Vedaldi Soatto, “Relaxed Matching Kernels”, 2008

Fisher kernels

- Compare objects \mathbf{x} , \mathbf{y} by a generative model
 - given $p(\mathbf{x} | \theta)$
 - map points \mathbf{x} to maximum-likelihood parameters $\theta_{\mathbf{x}}$
 - compare $K(\theta_{\mathbf{x}}, \theta_{\mathbf{y}})$
- Local analysis
 - log-likelihood function $L(\mathbf{x}, \theta) = \log p(\mathbf{x} | \theta)$
 - assume $\mathbf{x} \sim p(\mathbf{x} | \theta)$
 - maximum likelihood is consistent $\forall \hat{\theta} : E[L(x, \hat{\theta})] \leq E[L(x, \theta)]$
- Fisher score

$$U(\mathbf{x}, \theta) = \nabla_{\theta} L(\mathbf{x}, \theta) \quad E[U(\mathbf{x}, \theta)] = \frac{\partial}{\partial \theta} E[L(\mathbf{x}, \theta)] = 0$$

- Fisher information

$$I(\theta) = E[U(\mathbf{x}, \theta)^2] = \text{var } U(\mathbf{x}, \theta)$$


Fisher kernels

- Fisher information matrix as approx. second derivative

$$\begin{aligned} E \left[\frac{\partial^2}{\partial \theta^2} L(\mathbf{x}, \theta) \right] &= E \left[\frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2}{\partial \theta^2} p(\mathbf{x}|\theta) \right] - E \left[\left(\frac{1}{p(\mathbf{x}|\theta)} \frac{\partial}{\partial \theta} p(\mathbf{x}|\theta) \right)^2 \right] \\ &\approx -E \left[\left(\frac{\partial}{\partial \theta} \log p(\mathbf{x}|\theta) \right)^2 \right] \\ &= -E[U(\mathbf{x}, \theta)^2] = -I(\theta) \end{aligned}$$

- Approx. ML estimate

$$\begin{aligned} L(\mathbf{x}, \theta + \delta\theta) &\approx L(\mathbf{x}, \theta) + U(\mathbf{x}, \theta)\delta\theta - \frac{1}{2}I(\theta)(\delta\theta)^2 \\ \delta\theta_{\mathbf{x}} &\approx I(\theta)^{-1}U(\mathbf{x}, \theta) \end{aligned}$$

- Fisher kernel

$$K(\mathbf{x}, \mathbf{y}) = \delta\theta_{\mathbf{x}} I(\theta) \delta\theta_{\mathbf{y}} = U(\mathbf{x}, \theta) I(\theta)^{-1} U(\mathbf{y}, \theta)$$

Invariance

- Why weighting by I ?

$$K(\mathbf{x}, \mathbf{y}) = \delta\theta_{\mathbf{x}} I(\theta) \delta\theta_{\mathbf{y}} = U(\mathbf{x}, \theta) I(\theta)^{-1} U(\mathbf{y}, \theta)$$

- Reparametrization $\theta = \phi(\lambda)$

$$L'(\mathbf{x}, \lambda) = L(\mathbf{x}, \phi(\lambda)) \quad U'(\mathbf{x}, \lambda) = U(\mathbf{x}, \phi(\lambda)) \dot{\phi}(\lambda)$$

$$I'(\lambda) = \dot{\phi}(\lambda) I(\phi(\lambda)) \dot{\phi}(\lambda)$$

- Fisher kernel is invariant to reparametrization

$$K(\mathbf{x}, \mathbf{y}) = U' (I')^{-1} U' = U \phi \phi^{-1} I \phi^{-1} \phi U = U I^{-1} U$$

Tutorial

- MediaLandscape Player