

Kernels on Histograms through the Transportation Polytope

Marco Cuturi, cuturi@ism.ac.jp

The Institute of Statistical Mathematics, Tokyo

August 1st - MLSS Workshop

Outline

- 1 Kernels on Histograms**
 - Why histograms?...
 - Why non-standard kernels?...
 - What about standard distances?...
- 2 Kernels derived from permanents and the transportation polytope**
 - Some definitions
 - Transportation of clouds of points
 - The volume of the transportation polytope
 - Generating function (or weighted volume) of the transportation polytope

Outline

- 1 Kernels on Histograms**
 - Why histograms?...
 - Why non-standard kernels?...
 - What about standard distances?...
- 2 Kernels derived from permanents and the transportation polytope**
 - Some definitions
 - Transportation of clouds of points
 - The volume of the transportation polytope
 - Generating function (or weighted volume) of the transportation polytope

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Some solutions

- Use fingerprints (feature vectors defined by experts)
- Slice the objects if they are discrete and convolute kernels on subparts [Haussler, 1998]
- Use the power of statistical models [Jaakkola, 1998]
- In the most simple setting, use bags-of-components to describe your object, that is histograms.

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Some solutions

- Use fingerprints (feature vectors defined by experts)
- Slice the objects if they are discrete and convolute kernels on subparts [Haussler, 1998]
- Use the power of statistical models [Jaakkola, 1998]
- In the most simple setting, use bags-of-components to describe your object, that is histograms.

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Some solutions

- Use fingerprints (feature vectors defined by experts)
- Slice the objects if they are discrete and convolute kernels on subparts [Haussler, 1998]
- Use the power of statistical models [Jaakkola, 1998]
- In the most simple setting, use bags-of-components to describe your object, that is histograms.

Problem

- Kernel methods, SVM's are well understood for **vector inputs** with a challenge: which $\gamma, \Sigma, a, b, d \dots$?
- However, there is **no natural dot-product or norm** for strings, graphs, texts... so we cannot use Gaussian $e^{-\gamma \|\cdot - \cdot\|^2}$, polynomial $(\langle \cdot, \cdot \rangle + b)^d$ etc..

Some solutions

- Use fingerprints (feature vectors defined by experts)
- Slice the objects if they are discrete and convolute kernels on subparts [Haussler, 1998]
- Use the power of statistical models [Jaakkola, 1998]
- In the most simple setting, use bags-of-components to describe your object, that is histograms.

Bag-of-components, histogram representations

Example

A long sequence seen as a bag of n -grams:

AABHLK FHGH ... HAABGJYHLK ... AT

$\mapsto \{(AAB, 2), (HLK, 2), (FHG, 1) \dots\}$

Bag-of-components, histogram representations

Example

A text as a bag of words:

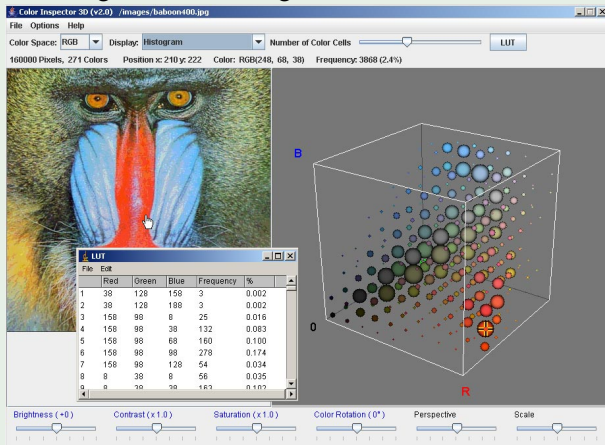
the cat eats the mouse $\mapsto \{(the, 2), (cat, 1), \dots\}$.

Why histograms?...

Bag-of-components, histogram representations

Example

An image as an histogram of RGB colors, or any other features:



Probability representation

A **composite** object z (a text) can be seen as an aggregation of basic **components** (words) taken from \mathcal{X} (a dictionary).

$$z \mapsto \mu \in M_1^+(\mathcal{X})$$

when \mathcal{X} is finite $\text{card}(\mathcal{X}) = d$, this is a **multinomial** distribution from Σ_d , that is a vector $\theta \in [0, 1]^d$ such that $\sum_{i=1}^d \theta_i = 1$.

kernels for structured objects \rightarrow kernels for histograms

In practice: complex objects \rightarrow histograms $\xrightarrow{k_{\Sigma_d} + \text{SVM}}$ results.

Probability representation

A **composite** object z (a text) can be seen as an aggregation of basic **components** (words) taken from \mathcal{X} (a dictionary).

$$z \mapsto \mu \in M_1^+(\mathcal{X})$$

when \mathcal{X} is finite $\text{card}(\mathcal{X}) = d$, this is a **multinomial** distribution from Σ_d , that is a vector $\theta \in [0, 1]^d$ such that $\sum_{i=1}^d \theta_i = 1$.

kernels for structured objects \rightarrow kernels for histograms

In practice: complex objects \rightarrow histograms $\xrightarrow{k_{\Sigma_d} + \text{SVM}}$ results.

$$\Sigma_d \subset \mathbb{R}_+^d \subset \mathbb{R}^d$$

- Sure, so Gaussian, polynomials etc.. can be used.
- In fact most early applications of SVM [Joachims '98, Chapelle '99] used this framework.
- However... counts are **positive** [Hein, Bousquet '05], which make multinomial vectors lie in the positive orthant,
- and **normalized** so their mass is 1, so we can use the information geometry of the manifold of probability measures [Lafferty & Lebanon '05, Lebanon '06, Kondor Jebara '03, C Fukumizu Vert 05 etc.]

$$\Sigma_d \subset \mathbb{R}_+^d \subset \mathbb{R}^d$$

- Sure, so Gaussian, polynomials etc.. can be used.
- In fact most early applications of SVM [Joachims '98, Chapelle '99] used this framework.
- However... counts are **positive** [Hein, Bousquet '05], which make multinomial vectors lie in the positive orthant,
- and **normalized** so their mass is 1, so we can use the information geometry of the manifold of probability measures [Lafferty & Lebanon '05, Lebanon '06, Kondor Jebara '03, C Fukumizu Vert 05 etc.]

$$\Sigma_d \subset \mathbb{R}_+^d \subset \mathbb{R}^d$$

- Sure, so Gaussian, polynomials etc.. can be used.
- In fact most early applications of SVM [Joachims '98, Chapelle '99] used this framework.
- However... counts are **positive** [Hein, Bousquet '05], which make multinomial vectors lie in the positive orthant,
- and **normalized** so their mass is 1, so we can use the information geometry of the manifold of probability measures [Lafferty & Lebanon '05, Lebanon '06, Kondor Jebara '03, C Fukumizu Vert 05 etc.]

$$\Sigma_d \subset \mathbb{R}_+^d \subset \mathbb{R}^d$$

- Sure, so Gaussian, polynomials etc.. can be used.
- In fact most early applications of SVM [Joachims '98, Chapelle '99] used this framework.
- However... counts are **positive** [Hein, Bousquet '05], which make multinomial vectors lie in the positive orthant,
- and **normalized** so their mass is 1, so we can use the information geometry of the manifold of probability measures [Lafferty & Lebanon '05, Lebanon '06, Kondor Jebara '03, C Fukumizu Vert 05 etc.]

Dozens of classic probability metrics

- Some well-known divergences and metrics can be used directly as kernels: Jensen-Divergence, Hellinger distance, χ^2 distance etc... (long list, but **not** KL)
- Some kernels aim at incorporating **a priori information (a kernel κ)** on \mathcal{X} ... as if they were actually computing these metrics in a feature space (the rkhs \mathcal{H}_κ).

The transport or Monge-Kantorovich distance

- Very old distance (goes back to Monge, late 18th), known as Earth Movers' in vision.

Dozens of classic probability metrics

- Some well-known divergences and metrics can be used directly as kernels: Jensen-Divergence, Hellinger distance, χ^2 distance etc... (long list, but **not** KL)
- Some kernels aim at incorporating **a prior information (a kernel κ)** on \mathcal{X} ... as if they were actually computing these metrics in a feature space (the rkhs \mathcal{H}_κ).

The transport or Monge-Kantorovich distance

- Very old distance (goes back to Monge, late 18th), known as 'Earth Movers' in vision.
- It is not directly a (nu) kernel (remember the edit distance?) but we can take a closer look at it to build a kernel.

Dozens of classic probability metrics

- Some well-known divergences and metrics can be used directly as kernels: Jensen-Divergence, Hellinger distance, χ_2 distance etc... (long list, but **not** KL)
- Some kernels aim at incorporating **a prior information (a kernel κ)** on \mathcal{X} ... as if they were actually computing these metrics in a feature space (the rkhs \mathcal{H}_κ).

The transport or Monge-Kantorovich distance

- Very old distance (goes back to Monge, late 18th), known as Earth Movers' in vision.
- It is not directly a (n.d) kernel (remember the edit distance?) but we can take a closer look at it to build a kernel.

Dozens of classic probability metrics

- Some well-known divergences and metrics can be used directly as kernels: Jensen-Divergence, Hellinger distance, χ_2 distance etc... (long list, but **not** KL)
- Some kernels aim at incorporating **a prior information (a kernel κ)** on \mathcal{X} ... as if they were actually computing these metrics in a feature space (the rkhs \mathcal{H}_κ).

The transport or Monge-Kantorovich distance

- Very old distance (goes back to Monge, late 18th), known as Earth Movers' in vision.
- It is not directly a (n.d) kernel (remember the edit distance?) but we can take a closer look at it to build a kernel.

Outline

- 1 **Kernels on Histograms**
 - Why histograms?...
 - Why non-standard kernels?...
 - What about standard distances?...
- 2 **Kernels derived from permanents and the transportation polytope**
 - Some definitions
 - Transportation of clouds of points
 - The volume of the transportation polytope
 - Generating function (or weighted volume) of the transportation polytope

MK distance between two histograms

- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^T\mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.
- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$

MK distance between two histograms

- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^T\mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.
- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$.

MK distance between two histograms

- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^T\mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.
- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$.

MK distance between two histograms

- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^T\mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.
- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$

MK distance between two histograms

- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^\top \mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.

- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$

MK distance between two histograms

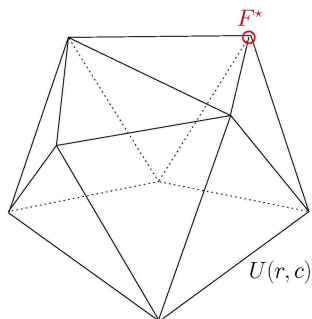
- Consider two discrete histograms $r = (r_1, \dots, r_d)$ and $c = (c_1, \dots, c_d)$ of equal sum N .
- $T \in \mathbb{R}_{d,d}$ is an arbitrary distance matrix between bins.
- $U(r, c) = \{F \in \mathbb{N}_{d \times d} \mid F\mathbf{1}_d = r, F^\top \mathbf{1}_d = c\}$, transportation matrices between r and c .

- $F = f_{ij}$ is such that

$$\left(\begin{array}{cccc|c} f_{11} & f_{12} & \cdots & f_{1d} & r_1 \\ f_{21} & f_{22} & \cdots & f_{2d} & r_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} & r_d \\ \hline c_1 & c_2 & \cdots & c_d & N \end{array} \right)$$

- Then, $d_{\text{MK}}(r, c) = \min_{F \in U(r, c)} \langle F, T \rangle$.
- We write $F^* = \operatorname{argmin}_{F \in U(r, c)} \langle F, T \rangle$

Some definitions



Within the **whole transportation polytope**, and for a given T , the MK distance only considers $\langle F^*, T \rangle$. How can we use the extra information of $U(r, c)$?

Permanent

For a $n \times n$ matrix $M = [m_{ij}]$, the permanent of M , per M is

$$\text{per } M = \sum_{\sigma \in \mathcal{S}_n} \prod_{i=1}^n m_{i\sigma(i)}$$

Permanent kernel between clouds of points

- Two finite clouds of points, $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_n\}$, x_i and $y_j \in \mathcal{X}$.
- Consider further that \mathcal{X} has a kernel κ .
- Then $\text{per } \kappa(x, y) = \text{per}[\kappa(x_i, y_j)]_{i,j=1}^n$.
- If \mathcal{X} is a distance metric and $\kappa = \exp(-\beta d)$

Permanent

For a $n \times n$ matrix $M = [m_{ij}]$, the permanent of M , $\text{per } M$ is

$$\text{per } M = \sum_{\sigma \in \mathcal{S}_n} \prod_{i=1}^n m_{i\sigma(i)}$$

Permanent kernel between clouds of points

- Two finite **clouds of points**, $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_n\}$, x_i and $y_j \in \mathcal{X}$.
- Consider further that \mathcal{X} has a kernel κ .
- Then $k_{\text{per}} : (x, y) \mapsto \text{per}([\kappa(x_i, y_j)]_{1 \leq i, j \leq n})$ is p.d.
- If T is a distance matrix and $\kappa = e^{-T}$,

$$k_{\text{per}}(x, y) = \sum_{\sigma \in \mathcal{S}_n} e^{-\sum_{i=1}^n t(x_{\sigma_i}, y_i)}.$$

Permanent

For a $n \times n$ matrix $M = [m_{ij}]$, the permanent of M , $\text{per } M$ is

$$\text{per } M = \sum_{\sigma \in \mathcal{S}_n} \prod_{i=1}^n m_{i\sigma(i)}$$

Permanent kernel between clouds of points

- Two finite **clouds of points**, $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_n\}$, x_i and $y_j \in \mathcal{X}$.
- Consider further that \mathcal{X} has a **kernel** κ .
- Then $k_{\text{per}} : (x, y) \mapsto \text{per}([\kappa(x_i, y_j)]_{1 \leq i, j \leq n})$ is p.d.
- If T is a distance matrix and $\kappa = e^{-T}$,

$$k_{\text{per}}(x, y) = \sum_{\sigma \in \mathcal{S}_n} e^{-\sum_{i=1}^n t(x_{\sigma_i}, y_i)}.$$

Permanent

For a $n \times n$ matrix $M = [m_{ij}]$, the permanent of M , $\text{per } M$ is

$$\text{per } M = \sum_{\sigma \in \mathcal{S}_n} \prod_{i=1}^n m_{i\sigma(i)}$$

Permanent kernel between clouds of points

- Two finite **clouds of points**, $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_n\}$, x_i and $y_i \in \mathcal{X}$.
- Consider further that \mathcal{X} has a **kernel** κ .
- Then $k_{\text{per}} : (x, y) \mapsto \text{per}([\kappa(x_i, y_j)]_{1 \leq i, j \leq n})$ is p.d.
- If T is a distance matrix and $\kappa = e^{-T}$,

$$k_{\text{per}}(x, y) = \sum_{\sigma \in \mathcal{S}_n} e^{-\sum_{i=1}^n t(x_{\sigma_i}, y_i)}.$$

Permanent

For a $n \times n$ matrix $M = [m_{ij}]$, the permanent of M , per M is

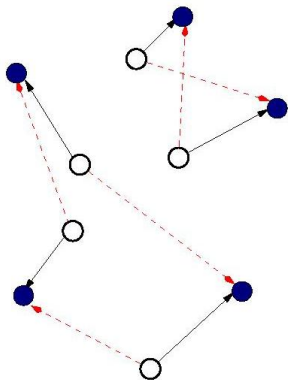
$$\text{per } M = \sum_{\sigma \in \mathcal{S}_n} \prod_{i=1}^n m_{i\sigma(i)}$$

Permanent kernel between clouds of points

- Two finite **clouds of points**, $x = \{x_1, \dots, x_n\}$ and $y = \{y_1, \dots, y_n\}$, x_i and $y_j \in \mathcal{X}$.
- Consider further that \mathcal{X} has a **kernel** κ .
- Then $k_{\text{per}} : (x, y) \mapsto \text{per}([\kappa(x_i, y_j)]_{1 \leq i, j \leq n})$ is p.d.
- If T is a distance matrix and $\kappa = e^{-T}$,

$$k_{\text{per}}(x, y) = \sum_{\sigma \in \mathcal{S}_n} e^{-\sum_{i=1}^n t(x_{\sigma_i}, y_i)}.$$

Transportation of clouds of points



—→ Optimal permutation σ^*

- - - → Other permutation $\sigma \in S_n$

$$k_{\text{per}}(X, Y) = e^{-d_{\sigma^*}^*} + e^{-d_{\sigma}} + \dots$$

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

- $k_{ij}(r, c) = |U(r, c)|$ is p.d.

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

- $k_{ij}(r, c) = |U(r, c)|$ is p.d.
- Proof: Birkhoff one-to-one correspondence between every $T \in U(r, c)$ and pairs of Gen. Young Tableaux (η, τ) , (η, σ) with identical shape η .

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

- $k_{\text{vol}}(r, c) = |U(r, c)|$ is p.d.
- Proof: RSK one-to-one correspondence between every $T \in U(r, c)$ and pairs of Gen. Young Tableaux $(\eta, r), (\eta, c)$ with identical shape η .
- Hence $|U(r, c)| = \sum_{\eta} K_{\eta, r} K_{\eta, c}$

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

- $k_{\text{vol}}(r, c) = |U(r, c)|$ is p.d.
- Proof: RSK one-to-one correspondence between every $T \in U(r, c)$ and pairs of Gen. Young Tableaux $(\eta, r), (\eta, c)$ with identical shape η .
- Hence $|U(r, c)| = \sum_{\eta} K_{\eta, r} K_{\eta, c}$

Volume of $U(r, c)$ for general histograms

- $|U(r, c)|$ quantifies the number of transportation plans
- Extremely difficult to compute exactly (active research topic) [Lorea & al. 2004, Barvinok 1993, Beck...]
- Some approximations exists [Diaconis & al. 1994]

Positive definiteness of the volume

- $k_{\text{vol}}(r, c) = |U(r, c)|$ is p.d.
- Proof: RSK one-to-one correspondence between every $T \in U(r, c)$ and pairs of Gen. Young Tableaux $(\eta, r), (\eta, c)$ with identical shape η .
- Hence $|U(r, c)| = \sum_{\eta} K_{\eta, r} K_{\eta, c}$

Weighted volume of $U(r, c)$

- $f(T) = \sum_{F \in U(r, c)} e^{-\langle F, T \rangle}$: generating function
- Also very difficult to compute although some SIS approximations have been proposed

Positive definiteness of weighted volumes

• define $k_{U, T} : (r, c) \mapsto \sum_{F \in U(r, c)} \mathcal{K}(F) e^{-\langle F, T \rangle}$

• $k_{U, T}$ is positive definite if $\langle k_{U, T}(r, c), k_{U, T}(r', c') \rangle \geq 0$

• $k_{U, T}$ is positive definite if $\langle k_{U, T}(r, c), k_{U, T}(r, c) \rangle \geq 0$

• $k_{U, T}$ is positive definite if $\langle k_{U, T}(r, c), k_{U, T}(r, c) \rangle \geq 0$

• $k_{U, T}$ is positive definite if $\langle k_{U, T}(r, c), k_{U, T}(r, c) \rangle \geq 0$

• $k_{U, T}$ is positive definite if $\langle k_{U, T}(r, c), k_{U, T}(r, c) \rangle \geq 0$

Weighted volume of $U(r, c)$

- $f(T) = \sum_{F \in U(r, c)} e^{-\langle F, T \rangle}$: generating function
- Also very difficult to compute although some SIS approximations have been proposed

Positive definiteness of weighted volumes

• defining $k_{a,r}(T, \Theta) = \sum_{F \in U(r, c)} \nu(F) e^{-\langle F, T \rangle}$

• Θ and ν as above \rightarrow K as

$$\nu(F) = \prod_{i=1}^r \frac{(2c_i)!}{(c_i!)^2} \prod_{j=1}^c \nu_j^{2a-1}$$

• for $0 \leq a \leq 2$ we have that $k_{a,r}$ is p.d.

Weighted volume of $U(r, c)$

- $f(T) = \sum_{F \in U(r, c)} e^{-\langle F, T \rangle}$: generating function
- Also very difficult to compute although some SIS approximations have been proposed

Positive definiteness of weighted volumes

- define $k_{\varphi, T} : (r, c) \mapsto \sum_{F \in U(r, c)} \varphi(F) e^{-\langle T, F \rangle}$
- and $\varphi_a : \mathbb{N}_{d \times d} \rightarrow \mathbb{R}$ as

$$\varphi_a(F) = \prod_i \frac{(2f_{ij})!^a}{f_{ij}!} \prod_{i \neq j} f_{ij}!^{2a-1}$$

for $0 \leq a \leq 2$ we have that $k_{\varphi, T}$ is p.d

- In particular: $a = 0$ Fisher-Yates, $a = \frac{1}{2}$ only diagonal weights.

Weighted volume of $U(r, c)$

- $f(T) = \sum_{F \in U(r, c)} e^{-\langle F, T \rangle}$: generating function
- Also very difficult to compute although some SIS approximations have been proposed

Positive definiteness of weighted volumes

- define $k_{\varphi, T} : (r, c) \mapsto \sum_{F \in U(r, c)} \varphi(F) e^{-\langle T, F \rangle}$
- and $\varphi_a : \mathbb{N}_{d \times d} \rightarrow \mathbb{R}$ as

$$\varphi_a(F) = \prod_i \frac{(2f_{ii})!^a}{f_{ii}!} \prod_{i \neq j} f_{ij}!^{2a-1}$$

for $0 \leq a \leq 2$ we have that $k_{\varphi, T}$ is p.d

- In particular: $a = 0$ Fisher-Yates, $a = \frac{1}{2}$ only diagonal weights.

Weighted volume of $U(r, c)$

- $f(T) = \sum_{F \in U(r, c)} e^{-\langle F, T \rangle}$: generating function
- Also very difficult to compute although some SIS approximations have been proposed

Positive definiteness of weighted volumes

- define $k_{\varphi, T} : (r, c) \mapsto \sum_{F \in U(r, c)} \varphi(F) e^{-\langle T, F \rangle}$
- and $\varphi_a : \mathbb{N}_{d \times d} \rightarrow \mathbb{R}$ as

$$\varphi_a(F) = \prod_i \frac{(2f_{ii})!^a}{f_{ii}!} \prod_{i \neq j} f_{ij}!^{2a-1}$$

for $0 \leq a \leq 2$ we have that $k_{\varphi, T}$ is p.d

- In particular: $a = 0$ Fisher-Yates, $a = \frac{1}{2}$ only diagonal weights.

2000 images, 10 classes, 20 pixels per image

σ	Gaussian	Permanent
.1	34.3 (± 1.4)	32.3 (± 1.2)
.2	33.45 (± 1.0)	31.3 (± 1.3)
.3	37.3 (± 1.0)	33.2 (± 1.2)

Table: Misclassification rate expressed in percents for the 2 considered kernels along with their standard errors averaged over 3 cross-validation folds.