

# Cancer prognosis on the symmetric group

Jean-Philippe Vert

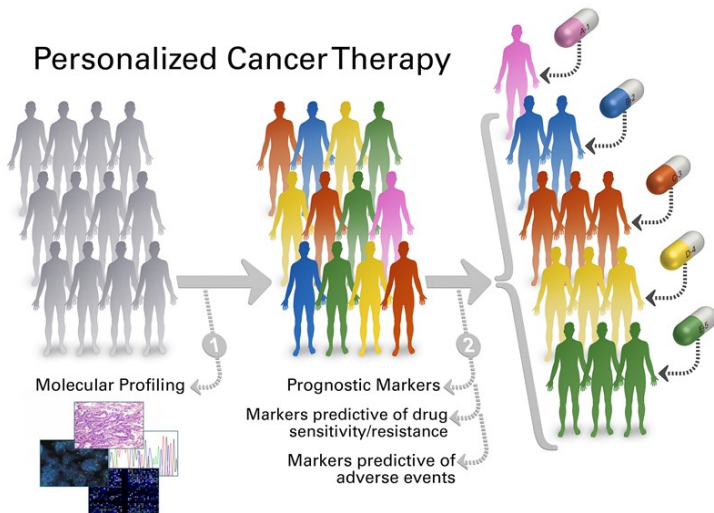


ENS Paris, December 15, 2016

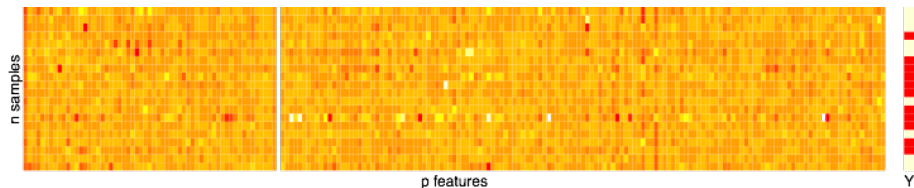


# Opportunities

## Personalized Cancer Therapy

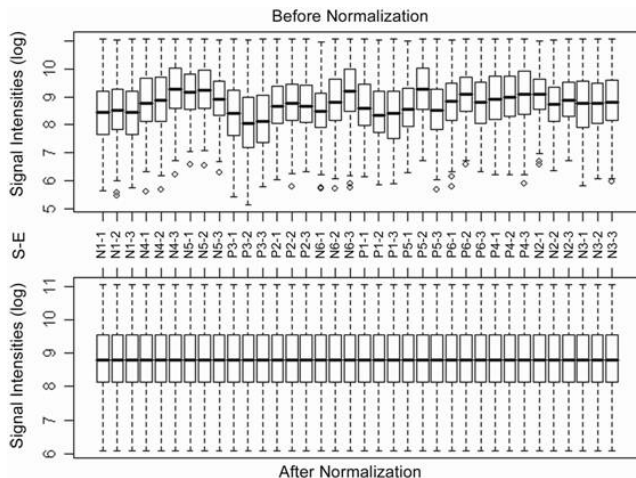


# Example: cancer prognosis from gene expression data



- $X$  gene expression profile of each patient
- $Y$  survival information of each patient
- $n = 10^2 \sim 10^4$
- $p = 2 \times 10^4$
- Goal: learn to predict  $Y$  from  $X$
- But... where does  $X$  come from?

# From raw data to $X$

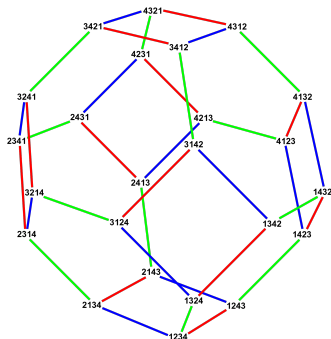


Quantile normalization (per sample) to remove various technical artefacts

# Working on the symmetric group

After QN, each sample  $X_i$  is:

- a **target distribution**  $d \in \mathbb{R}^p$ ,
- permuted by a **samples-specific permutation**  $\sigma_i \in \mathcal{S}_p$ , the symmetric group over the set of features



Can we directly estimate a model  $Y = f(\sigma)$  ?

# Outline

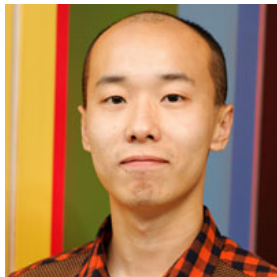
- 1 The Kendall and Mallows kernels
- 2 Supervised quantile normalization
- 3 Conclusion

# Outline

- 1 The Kendall and Mallows kernels
- 2 Supervised quantile normalization
- 3 Conclusion



## Joint work with

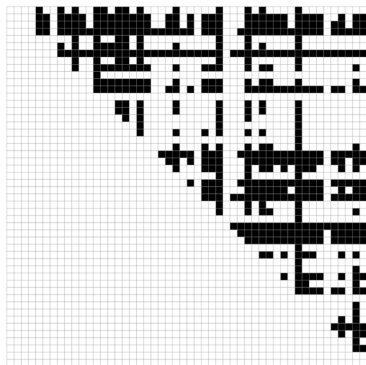


Yunlong Jiao

# An idea: all pairwise comparisons

Replace  $x \in \mathbb{R}^p$  by  $\Phi(x) \in \{0, 1\}^{p(p-1)/2}$ :

$$\Phi_{i,j}(x) = \begin{cases} 1 & \text{if } x_i \leq x_j, \\ 0 & \text{otherwise.} \end{cases}$$



**One sample  $x$   
 $p$  features**

**Mapping  $f(x)$   
 $p(p-1)/2$  bits**

# Related work: Top scoring pairs (TSP)



(Geman et al., 2004; Tan et al., 2005; Leek, 2009)

# Practical challenge



- Need to store  $O(p^2)$  bits per sample
- Need to train a model in  $O(p^2)$  dimensions

## Theorem (Wahba, Schölkopf, ...)

Training a linear model over a representation  $\Phi(x) \in \mathbb{R}^Q$  of the form:

$$\min_{w \in \mathbb{R}^Q} \frac{1}{n} \sum_{i=1}^n \ell(w^\top \Phi(x_i), y_i) + \lambda \|w\|^2$$

can be done efficiently, independently of  $Q$ , if the kernel

$$K(x, x') = \Phi(x)^\top \Phi(x')$$

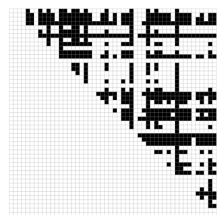
can be computed efficiently.

Ex: ridge regression,  $O(Q^3 + nQ^2)$  becomes  $O(n^3 + n^2 T)$

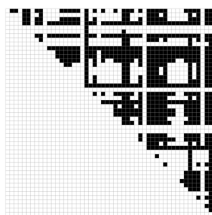
Other: SVM, logistic regression, Cox model, survival SVM, ...

# Kernel trick for us: Kendall's $\tau$

$$\Phi(x)^\top \Phi(x') = \tau(x, x') \quad (\text{up to a scaling})$$



$\times$



$$= \tau \left( \begin{array}{c} \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \end{array} , \begin{array}{c} \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \\ \text{red} \\ \text{green} \end{array} \right)$$

$O(p^2)$

$O(p \log(p))$

Good news for SVM and kernel methods!

## More formally

- For two permutations  $\sigma, \sigma'$  let  $n_c(\sigma, \sigma')$  (resp.  $n_d(\sigma, \sigma')$ ) the number of **concordant** (resp. **discordant**) pairs.
- The **Kendall kernel** (a.k.a. **Kendall tau coefficient**) is defined as

$$K_\tau(\sigma, \sigma') = \frac{n_c(\sigma, \sigma') - n_d(\sigma, \sigma')}{\binom{p}{2}}.$$

- The **Mallows kernel** is defined for any  $\lambda \geq 0$  by

$$K_M^\lambda(\sigma, \sigma') = e^{-\lambda n_d(\sigma, \sigma')}.$$

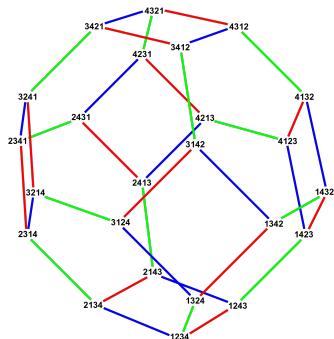
### Theorem (Jiao and V., 2015)

*The Kendall and Mallows kernels are **positive definite**.*

### Theorem (Knight, 1966)

*These two kernels for permutations can be evaluated in  $O(p \log p)$  time.*

# Related work



Cayley graph of  $S_4$

- Kondor and Barbarosa (2010) proposed the **diffusion kernel** on the Cayley graph of the symmetric group generated by adjacent transpositions.
- Computationally intensive ( $O(p^p)$ )
- Mallows kernel is written as

$$K_M^\lambda(\sigma, \sigma') = e^{-\lambda n_d(\sigma, \sigma')},$$

where  $n_d(\sigma, \sigma')$  is the **shortest path distance** on the Cayley graph.

- It can be computed in  $O(p \log p)$



# Application: supervised classification

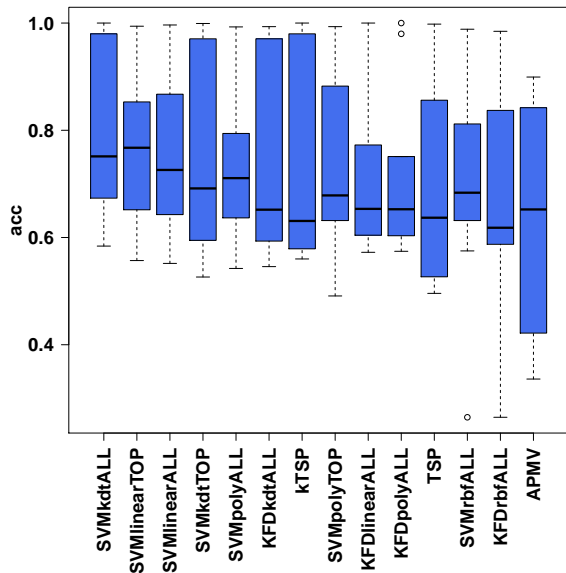
## Datasets

Dataset	No. of features	No. of samples (training/test)	
		$C_1$	$C_2$
Breast Cancer 1	23624	44/7 (Non-relapse)	32/12 (Relapse)
Breast Cancer 2	22283	142 (Non-relapse)	56 (Relapse)
Breast Cancer 3	22283	71 (Poor Prognosis)	138 (Good Prognosis)
Colon Tumor	2000	40 (Tumor)	22 (Normal)
Lung Cancer 1	7129	24 (Poor Prognosis)	62 (Good Prognosis)
Lung Cancer 2	12533	16/134 (ADCA)	16/15 (MPM)
Medulloblastoma	7129	39 (Failure)	21 (Survivor)
Ovarian Cancer	15154	162 (Cancer)	91 (Normal)
Prostate Cancer 1	12600	50/9 (Normal)	52/25 (Tumor)
Prostate Cancer 2	12600	13 (Non-relapse)	8 (Relapse)

## Methods

- Kernel machines Support Vector Machines (SVM) and Kernel Fisher Discriminant (KFD) with Kendall kernel, linear kernel, Gaussian RBF kernel, polynomial kernel.
- Top Scoring Pairs (TSP) classifiers [?].
- Hybrid scheme of SVM + TSP feature selection algorithm.

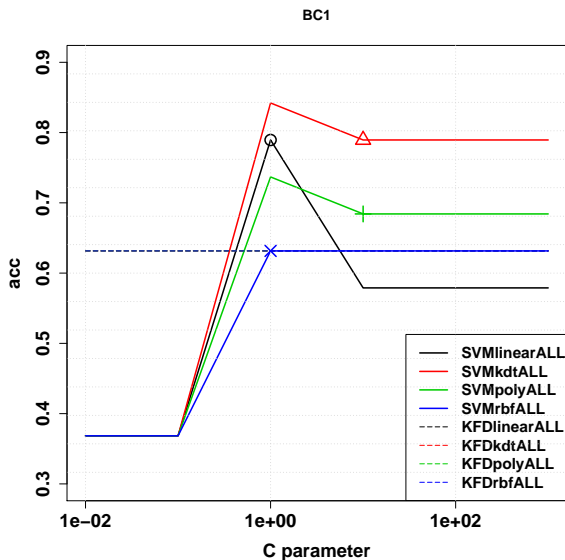
# Results



Kendall kernel SVM

- **Competitive accuracy!**
- Less sensitive to regularization parameter!
- No need for feature selection!

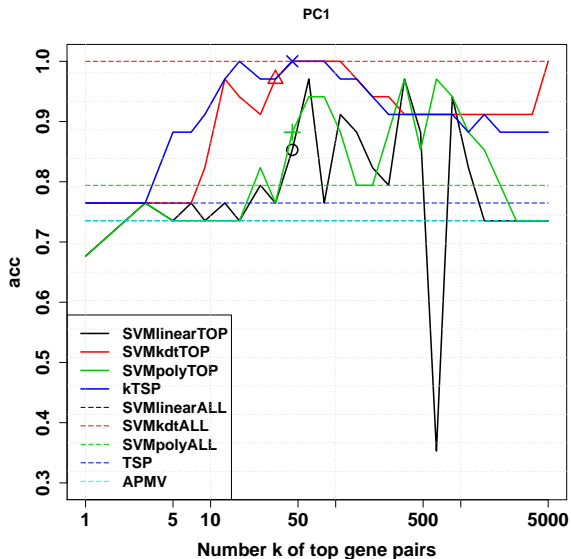
# Results



Kendall kernel SVM

- Competitive accuracy!
- **Less sensitive to regularization parameter!**
- No need for feature selection!

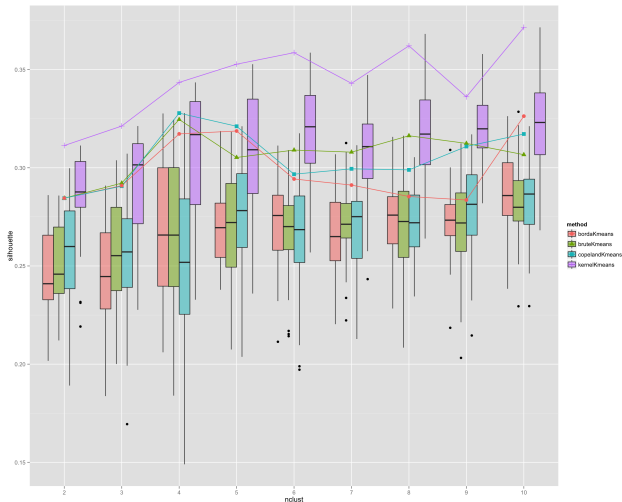
# Results



Kendall kernel SVM

- Competitive accuracy!
- Less sensitive to regularization parameter!
- **No need for feature selection!**

# Application: clustering



- APA data (full rankings)
- $n = 5738$ ,  $p = 5$
- (new) Kernel k-means vs (standard) k-means in  $\mathbb{S}_5$
- Show silhouette as a function of number of clusters (higher better)

## Extension to partial rankings

- Two interesting types of partial rankings are **interleaving partial ranking**

$$x_{i_1} \succ x_{i_2} \succ \dots \succ x_{i_k}, \quad k \leq n.$$

and **top-k partial ranking**

$$x_{i_1} \succ x_{i_2} \succ \dots \succ x_{i_k} \succ X_{\text{rest}}, \quad k \leq n.$$

- Partial rankings can be **uniquely represented** by a set of permutations compatible with all the observed partial orders.

### Theorem

*For these two particular types of partial rankings, the convolution kernel (Haussler, 1999) induced by Kendall kernel*

$$K_{\tau}^*(R, R') = \frac{1}{|R||R'|} \sum_{\sigma \in R} \sum_{\sigma' \in R'} K_{\tau}(\sigma, \sigma')$$

*can be evaluated in  $O(k \log k)$  time.*

## Extension to partial rankings

- Two interesting types of partial rankings are **interleaving partial ranking**

$$x_{i_1} \succ x_{i_2} \succ \dots \succ x_{i_k}, \quad k \leq n.$$

and **top-k partial ranking**

$$x_{i_1} \succ x_{i_2} \succ \dots \succ x_{i_k} \succ X_{\text{rest}}, \quad k \leq n.$$

- Partial rankings can be **uniquely represented** by a set of permutations compatible with all the observed partial orders.

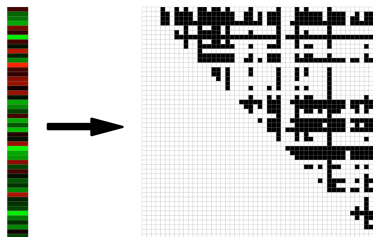
### Theorem

*For these two particular types of partial rankings, the convolution kernel (Haussler, 1999) induced by Kendall kernel*

$$K_{\tau}^*(R, R') = \frac{1}{|R||R'|} \sum_{\sigma \in R} \sum_{\sigma' \in R'} K_{\tau}(\sigma, \sigma')$$

*can be evaluated in  $O(k \log k)$  time.*

# Extension to smoother, continuous representations



One sample  $x$   
 $p$  features

Mapping  $f(x)$   
 $p(p-1)/2$  bits

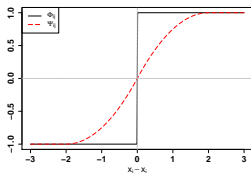
- Instead of  $\Phi : \mathbb{R}^p \rightarrow \{0, 1\}^{p(p-1)/2}$ , consider the continuous mapping  $\Psi_a : \mathbb{R}^p \rightarrow \mathbb{R}^{p(p-1)/2}$ :

$$\Psi_a(x) = \mathbb{E}\Phi(x + \epsilon) \quad \text{with} \quad \epsilon \sim (\mathcal{U}[-\frac{a}{2}, \frac{a}{2}])^n$$

- Corresponding kernel  $G_a(x, x') = \Psi_a(x)^\top \Psi_a(x')$



# Computation of $G(x, x')$



- $G_a(x, x')$  can be computed **exactly** in  $O(p^2)$  by explicit computation of  $\Psi_a(x)$  in  $\mathbb{R}^{p(p-1)/2}$

- $G_a(x, x')$  can be computed **approximately** in  $O(D^2 p \log p)$  by Monte-Carlo approximation:

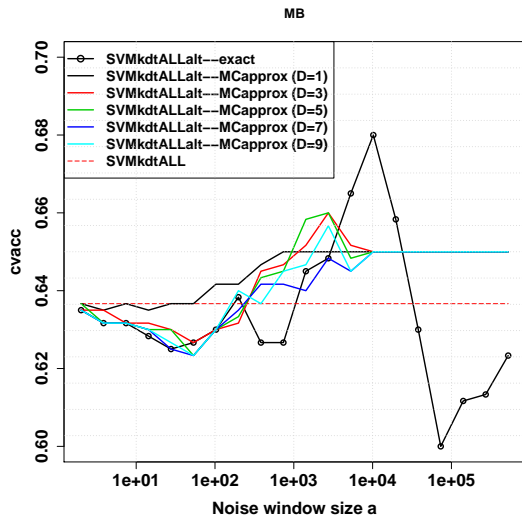
$$\tilde{G}_a(x, x') = \frac{1}{D^2} \sum_{i,j=1}^D K(x + \epsilon_i, x' + \epsilon'_j)$$

- Theorem: for supervised learning, Monte-Carlo approximation is better<sup>1</sup> than exact computation when  $n = o(p^{1/3})$

---

<sup>1</sup> faster for the same accuracy

# Performance of $G_a(x, x)$



# Outline

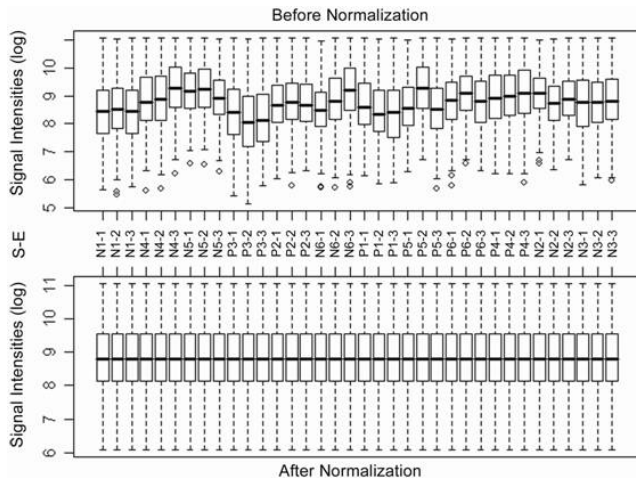
- 1 The Kendall and Mallows kernels
- 2 Supervised quantile normalization**
- 3 Conclusion

## Joint work with



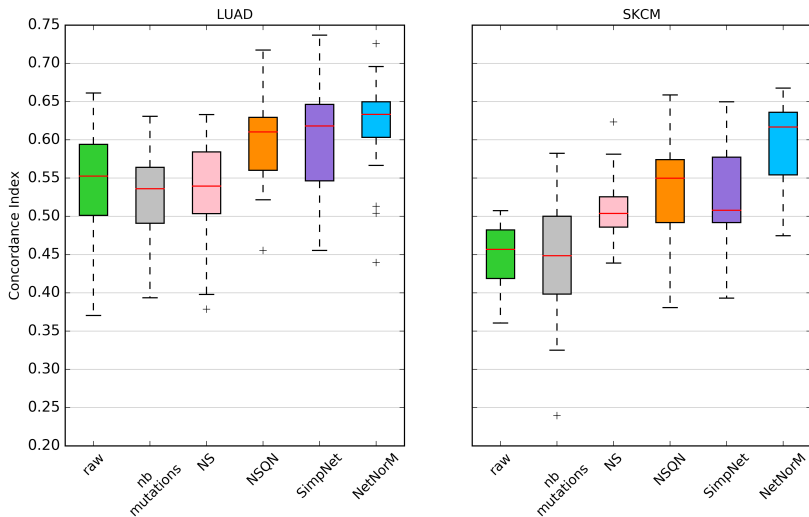
Marine Le Morvan

# Standard full quantile normalization



Typically followed by a predictive model  $f(X)$  on the normalized data

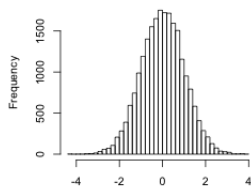
# Choosing a "good" target distributions is important



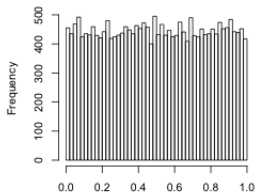
Cancer prognosis from somatic mutations

# How to choose a "good" target distribution?

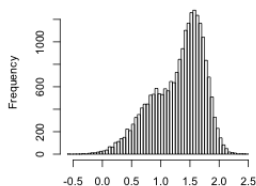
gaussian distribution (mean=0, sd=1)



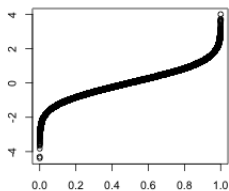
uniform distribution



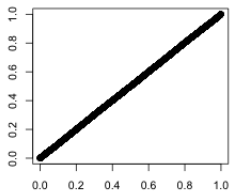
bigaussian distribution



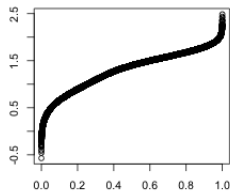
quantile function (-> gaussian)



quantile function (-> uniform)



quantile function (-> bigaussian)



# Learning the target distribution

- $x_1, \dots, x_n$  a set of  $p$ -dimensional samples
- $f \in \mathbb{R}^p$  a non-decreasing target distribution (CDF)
- For  $x \in \mathbb{R}^p$ , let  $\Phi_f(x) \in \mathbb{R}^p$  be the data after QN with target distribution  $f$
- **Standard approaches** (NSQN, NetNorM, ...)
  - 1 Fix  $f$  arbitrarily
  - 2 QN all samples to get  $\Phi_f(x_1), \dots, \Phi_f(x_n)$
  - 3 Learn a generalized linear model  $(w, b)$  on normalized data:

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^n \ell_i(w^\top \Phi_f(x_i) + b) + \lambda \Omega(w)$$

- **SUQUAN: jointly learn  $f$  and  $(w, b)$ :**

$$\min_{w, b, f} \frac{1}{n} \sum_{i=1}^n \ell_i(w^\top \Phi_f(x_i) + b) + \lambda \Omega(w) + \gamma \Omega_2(f)$$



# SUQAN: supervised quantile normalization

- For  $x \in \mathbb{R}^p$ , let  $\Pi_x \in \mathbb{R}^{p \times p}$  the permutation matrix of  $x$ 's entries

$$x = \begin{pmatrix} 4.5 \\ 1.2 \\ 10.1 \\ 8.9 \end{pmatrix} \quad \Pi_x = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix} \quad f = \begin{pmatrix} 0 \\ 1 \\ 3 \\ 4 \end{pmatrix}$$

- Quantile normalized  $x$  with target distribution  $f$  is:

$$\Phi_f(x) = \Pi_x f$$

- SUQUAN solves

$$\begin{aligned} \min_{w,b,f} \frac{1}{n} \sum_{i=1}^n \ell \left( w^\top \Pi_{x_i} f + b \right) + \lambda \Omega(w) + \gamma \Omega_2(f) \\ = \min_{w,b,f} \frac{1}{n} \sum_{i=1}^n \ell \left( \langle w f^\top, \Pi_{x_i} \rangle + b \right) + \lambda \Omega(w) + \gamma \Omega_2(f) \end{aligned} \tag{1}$$

- A particular **rank-1 matrix optimization**,  $x$  is **replaced by  $\Pi_x$**
- Solved by alternatively optimizing  $f$  and  $w$

# Experiments

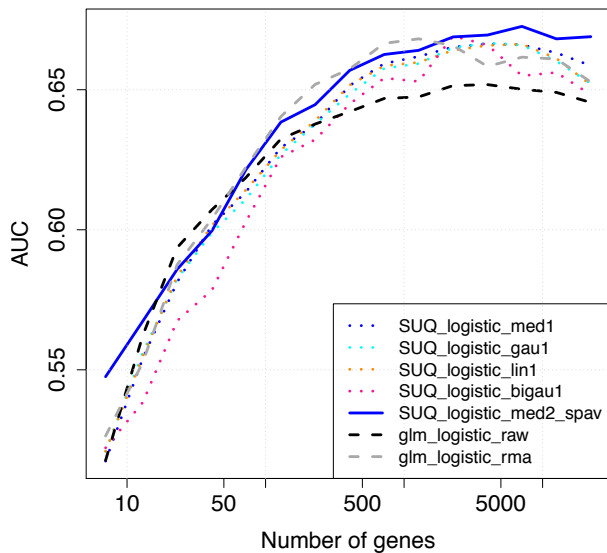
$$\min_{w,b,f} \frac{1}{n} \sum_{i=1}^n \ell_i \left( w^\top \Phi_f(x_i) + b \right) + \frac{\lambda}{2} \|w\|_2^2 + \frac{\gamma}{2} \sum_{j=1}^{p-1} (f_{j+1} - f_j)^2$$

- Breast cancer prognosis from gene expression data.
- Two classes of patients: those who relapsed within 6 years of diagnosis and those who did not.

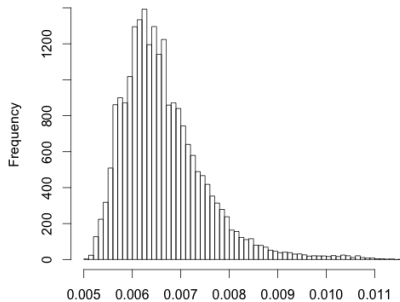
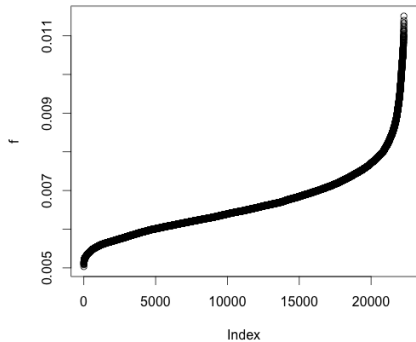
Dataset name	# genes	# patients	# positives	% positives
GSE7390	22283	189	58	0.31
GSE4922	22283	225	73	0.32
GSE2990	22283	106	32	0.30
GSE2034	22283	271	104	0.38
GSE1456	22283	141	37	0.26

# Performance

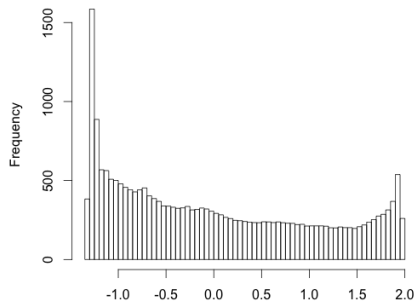
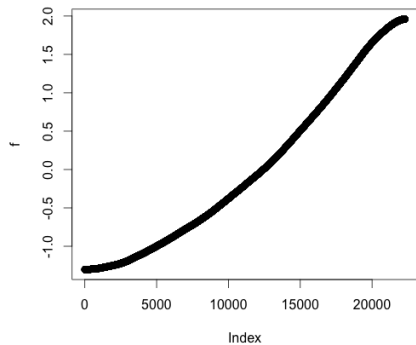
average over all datasets



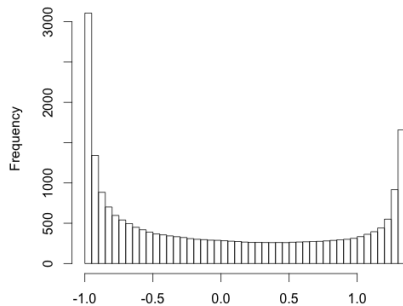
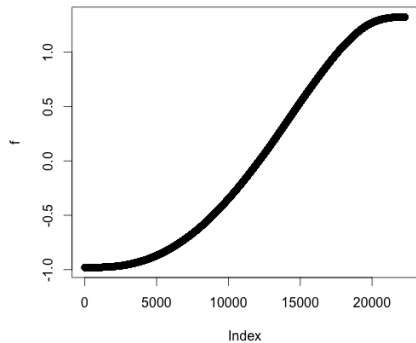
# Estimated distribution: iteration=0



# Estimated distribution: iteration=1



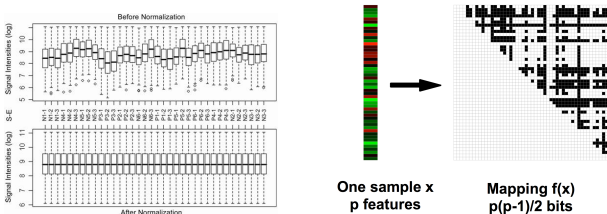
# Estimated distribution: iteration=2



# Outline

- 1 The Kendall and Mallows kernels
- 2 Supervised quantile normalization
- 3 Conclusion**

# Conclusion



- Representing omics data as **permutations** has some potential
  - **Kendall and Mallows** kernel in  $O(p \ln(p))$
  - **SUQUAN** supervised quantile normalization as matrix regression
- Understanding the **benefits and cost** of different representations remains very heuristic and sometimes counterintuitive
- **Learning representation** may help



# Thanks



The Adolph C. and Mary Sprague  
Miller Institute for Basic  
Research in Science  
*University of California, Berkeley*



SIMONS  
INSTITUTE  
for the Theory of Computing