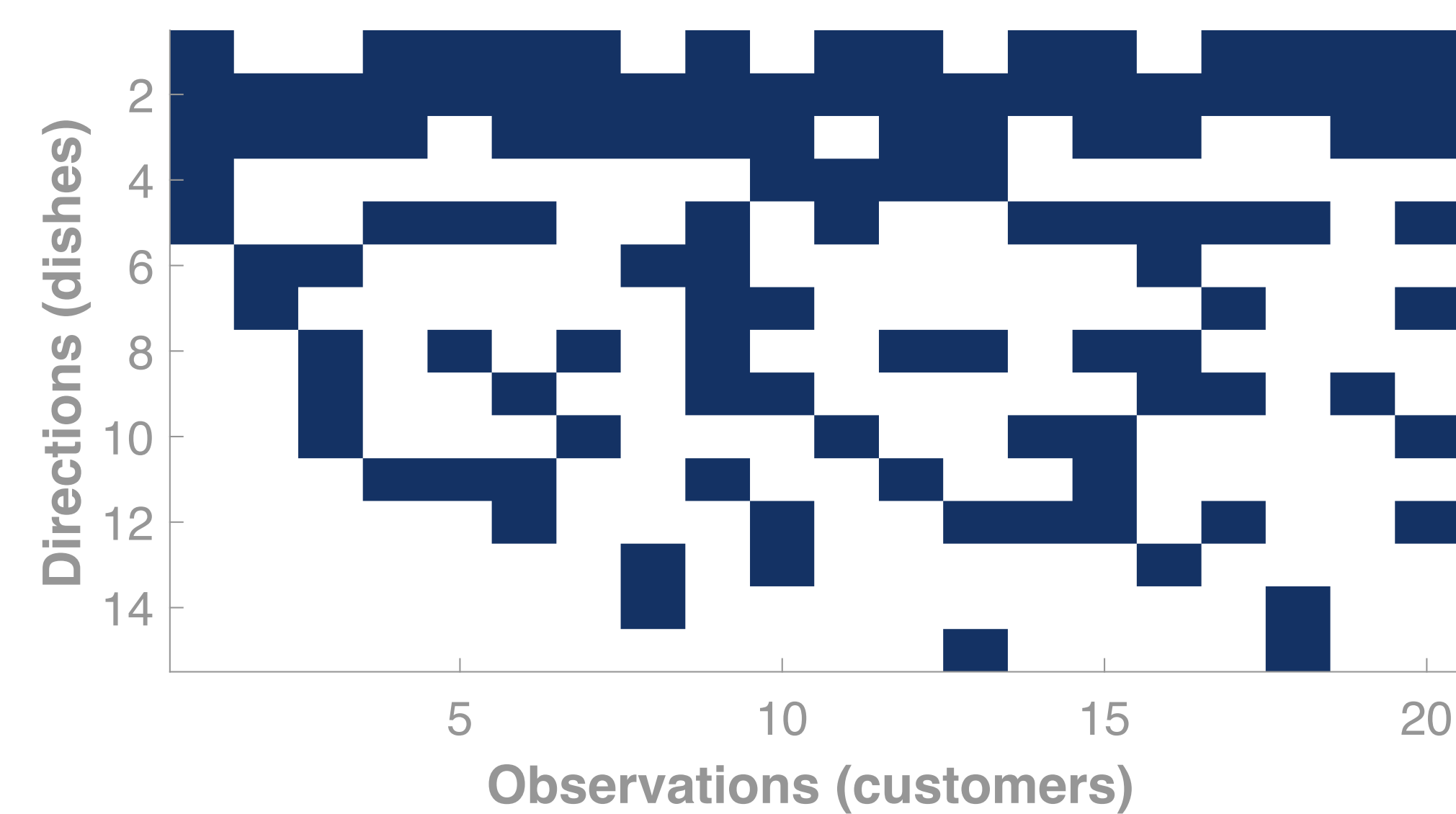


Context and contributions

Reduction of dimension is an ubiquitous pre-processing in numerous signal processing and machine learning tasks. The choice of the number of **principal components** K retained has a significant impact on performances. Existing methods able to infer K generally rely on RJMCMC methods, variational Bayes approximations or information criteria. **We propose a nonparametric Bayesian modeling of subspace estimation whose size is not known in advance.** We propose a Gibbs sampler as inference scheme. We build estimators w.r.t. to all possible subspaces and **discuss their consistency.** The method is validated on synthetic dataset and two real tasks.

Indian Buffet Process = prior over potentially infinite sparse binary matrices with $\mathbb{E}[K] = \alpha \log(N)$ (**regularizing effect**)



Distributions on the Stiefel manifold S_L

$$\text{vMF}(\mathbf{P}|\mathbf{F}) \sim \exp \text{tr}[\mathbf{F}^t \mathbf{P}]$$

$$\text{Bingham} \sim \exp \text{tr}[\Lambda \mathbf{P}^t \mathbf{A} \mathbf{P}]$$

Trick: write $\mathbf{p} = \mathbf{N}\mathbf{v}$ such that $\|\mathbf{v}\|_2 = 1$,

\mathbf{N} an orthonormal basis of $\mathbf{P}_{\setminus k}^\perp$, $\mathbf{p} = \mathbf{N}\mathbf{v}$

References

- [1] C. Elvira, P. Chainais, and N. Dobigeon, "Bayesian antisparsity coding," *IEEE Trans. Signal Process.*, vol. 65, pp. 1660–1672, April 2017.
- [2] C. Elvira, P. Chainais, and N. Dobigeon, "Bayesian nonparametric modeling of latent subspace," in prep.



¹This work was supported in part by the BNPSI ANR Project ANR-13-BS-03-0006-01 and the Fondation Centrale Initiative.

Bayesian model

$$\mathbf{y}_n = \mathbf{P}(\mathbf{z}_n \odot \mathbf{x}_n) + \mathbf{e}_n$$

$\mathbf{y}_n \in \mathbb{R}^D$, the observation vector

$$\mathbf{P} = [\mathbf{p}_1 \dots \mathbf{p}_D, \mathbf{0} \dots], \mathbf{P}^t \mathbf{P} = \mathbf{I}_D, \\ [\mathbf{p}_1 \dots \mathbf{p}_D] \sim \mathcal{U}_{S_D}$$

$$\mathbf{x}_n = [x_{1,n} \dots] \text{ where } x_{k,n} \sim \mathcal{N}(0, \delta_k^2 \sigma^2)$$

$$\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_n] \sim \text{IBP}(\alpha) \text{ a binary matrix}$$

$$\mathbf{e}_n \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_D)$$

$$\boldsymbol{\theta} = \{\delta^2, \sigma^2, \alpha\} \text{ vague conjugate priors}$$

$$p(\mathbf{P}, \mathbf{Z}, \boldsymbol{\theta} | \mathbf{Y}) = \int_{\mathbb{R}^{DN}} p(\mathbf{P}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{X} | \mathbf{Y}) d\mathbf{X}$$

Sampling active directions

$$\mathbf{p}_k | \mathbf{Y}, \mathbf{P}_{\setminus k} \stackrel{d}{\sim} \text{Bingham} \left(\sum z_{n,k,f}(\delta, \sigma^2) \mathbf{N}^t \mathbf{y}, \mathbf{y}^t \mathbf{N} \right)$$

Gibbs Sampler

foreach Iteration t do

Sample the non-singleton feature of \mathbf{Z} ;

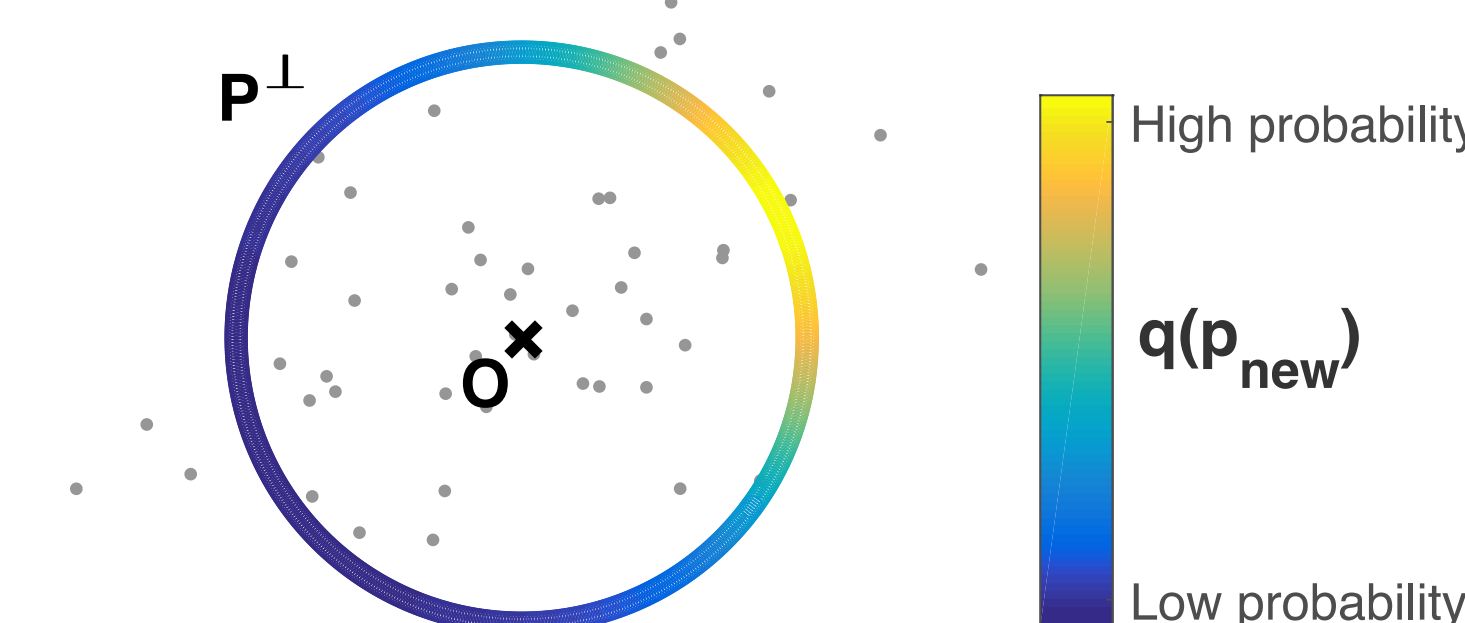
Add/suppress directions \sim von Mises Fisher;

$\mathbf{P} \sim$ Bingham;

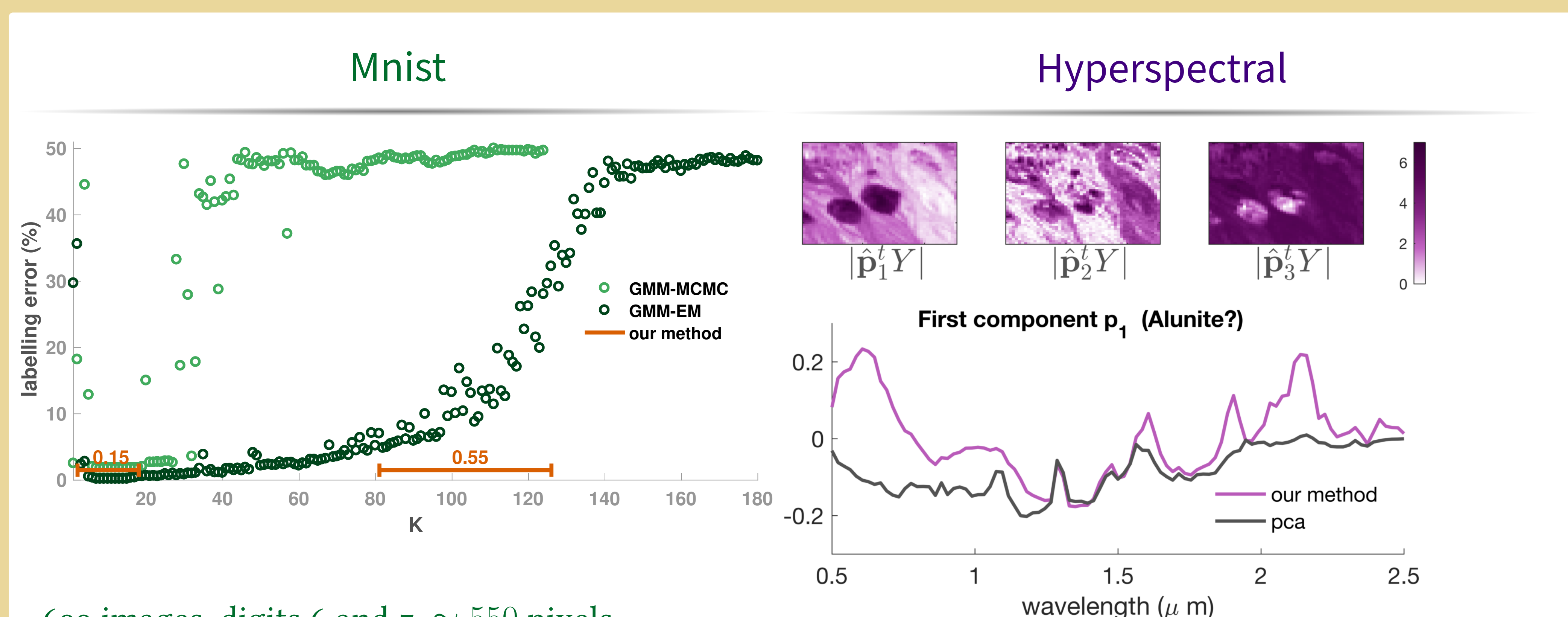
$\delta, \sigma^2, \alpha \sim$ conjugate distributions;

end

Explore new direction: Metropolis Hastings step with proposal $q \stackrel{d}{=} \text{vMF}(\mathbf{P}_{\text{new}} | [\mathbf{p}_1 \dots \mathbf{p}_K]^\perp, \mathbf{Y}, \sigma^2)$



Applications



- 600 images, digits 6 and 7, $\simeq 550$ pixels (after removing those with null variance)

- $\mathbf{x} \sim \pi_1 \mathcal{N}_1 + (1 - \pi_1) \mathcal{N}_2$ here

- compared with a MCMC- and EM-based inferences of a Gaussian mixture model for a varying number K of principal components

- our method reaches 1.5% of labeling error

- 2 areas are explored by the posterior $K | \mathbf{Y}$

- Cuprite-Hill dataset : $\simeq 100$ dimensions and 500 observations

- linear model hypothesis : voxel = $\mathbf{E}\mathbf{a}$ + noise, with $\|\mathbf{a}\|_1 = 1$ and $\mathbf{a} > 0$.

- $\hat{K} = 15$, ground truth ~ 10

- interpretability of the first directions

- next step: include unmixing!

Estimating the number of components

The natural estimator is inconsistent

$$\forall k, \limsup_{N \rightarrow +\infty} P(K_N = k | \mathbf{y}_1 \dots \mathbf{y}_N, \alpha) < 1 \text{ with probability 1}$$

If $\mathbf{y} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, a stronger result holds

$$P[K_N = 0 | \mathbf{y}_1 \dots \mathbf{y}_N, \alpha, \sigma^2] \xrightarrow[N \rightarrow +\infty]{a.s.} 0$$

Proposed methodology

All posteriors should be used to infer \hat{K}

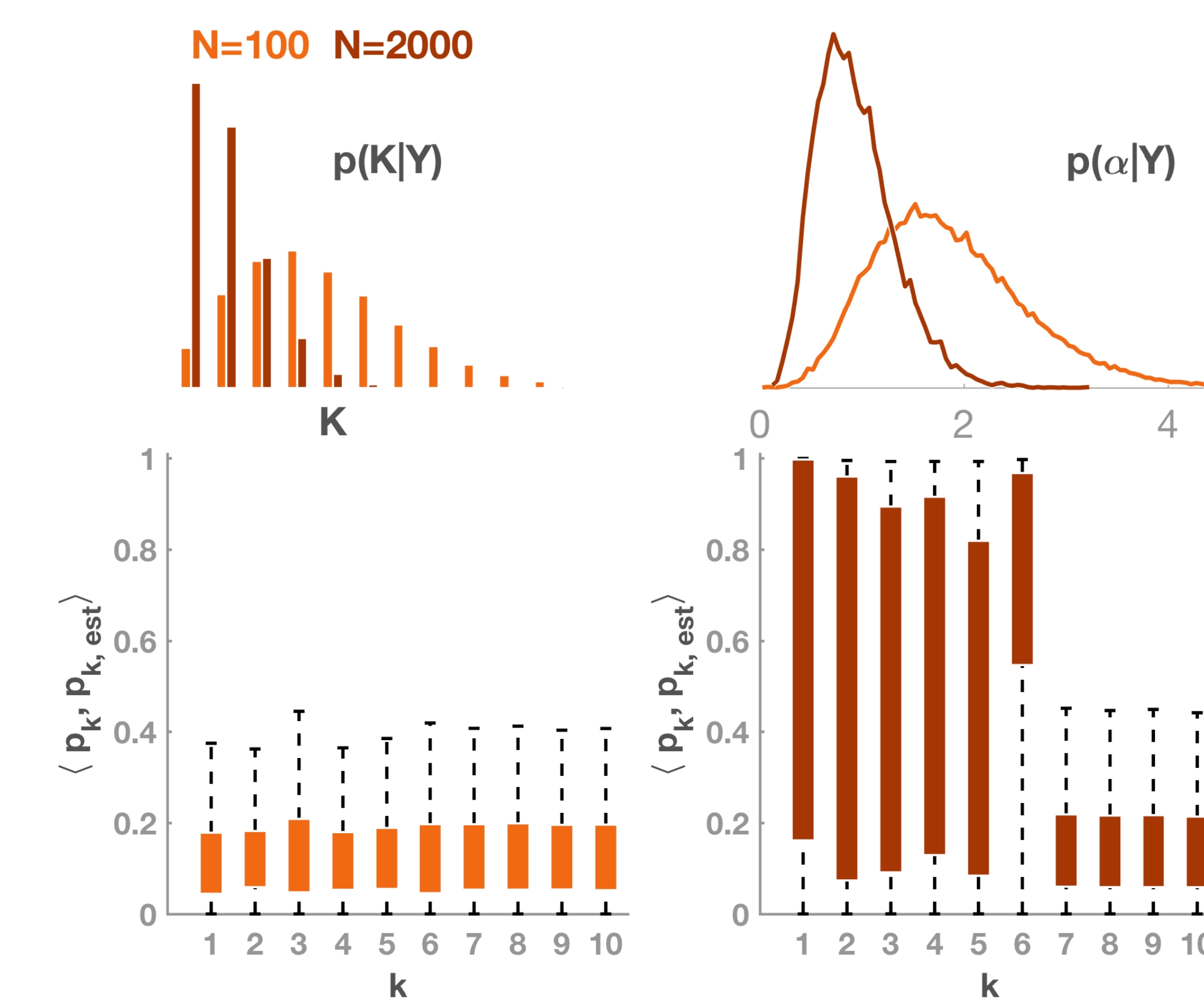
Let $\mathbf{u} \in \mathbb{R}^L, \|\mathbf{u}\|_2 = 1, \mathbf{V} \sim \mathcal{U}_{S_L}, W = |\langle \mathbf{u}, \mathbf{V} \rangle|$

Compare the marginal distributions $\langle \mathbf{u}, \mathbf{p}_k^{(t)} \rangle$ to $\langle \mathbf{u}, \mathbf{V} \rangle$ for various L using CDF

$$p_U(W \leq \lambda) = \frac{\text{vol}(\mathcal{S}_{L-2})}{\text{vol}(\mathcal{S}_{L-1})} 2 \int_0^\lambda (1 - w^2)^{(L-3)/2} dw$$

e.g., statistical tests.

Empirical observations



The estimator $\hat{K} | \mathbf{Y}$ empirically behaves correctly for large N , where $\mathbf{P}, \mathbf{Z}, \delta, \alpha$ and σ^2 have been marginalized out.

Several settings are tested : noise as input signal (✓), anisotropic noise (✓), scale factors δ^2 below the noise level (✗).