

MOTIVATION

- Bulk BS-seq technology ignores epigenetic heterogeneity among individual cells.
- Single cell BS-seq¹ measures methylation at single-cell level (binary states).
- Only 20% to 40% CpG coverage: limits statistical analysis to a semi-quantitative level.
- Bayesian hierarchical model: jointly learn methylation profiles (i.e. predict uncovered methylation states) and cluster cells based on genome-wide methylation patterns.

SCBPR MODEL

Bayesian generalised linear model (GLM) of basis function regression coupled with a Bernoulli observation model², where y_i is the methylation state, h_i the genomic location and w the model parameters.

$$y_i = \begin{cases} 1 & \text{if } z_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$z_i \sim \mathcal{N}(h_i w, 1)$$

$$w = p(w)$$

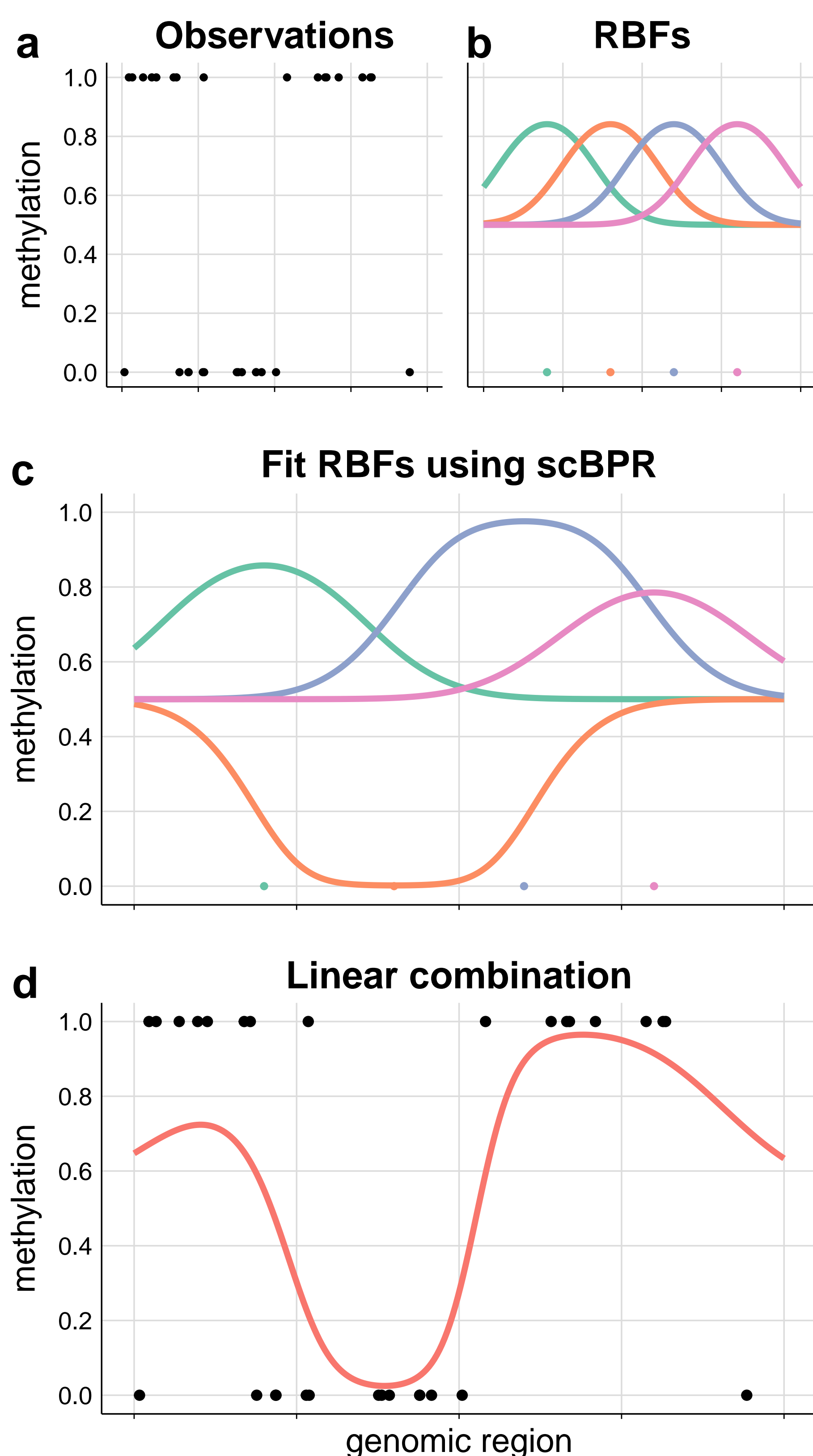


Figure 1: Process of learning methylation profiles using the single-cell Bernoulli Probit Regression (scBPR) model for a specific genomic region using 4 Radial Basis Functions (RBFs).

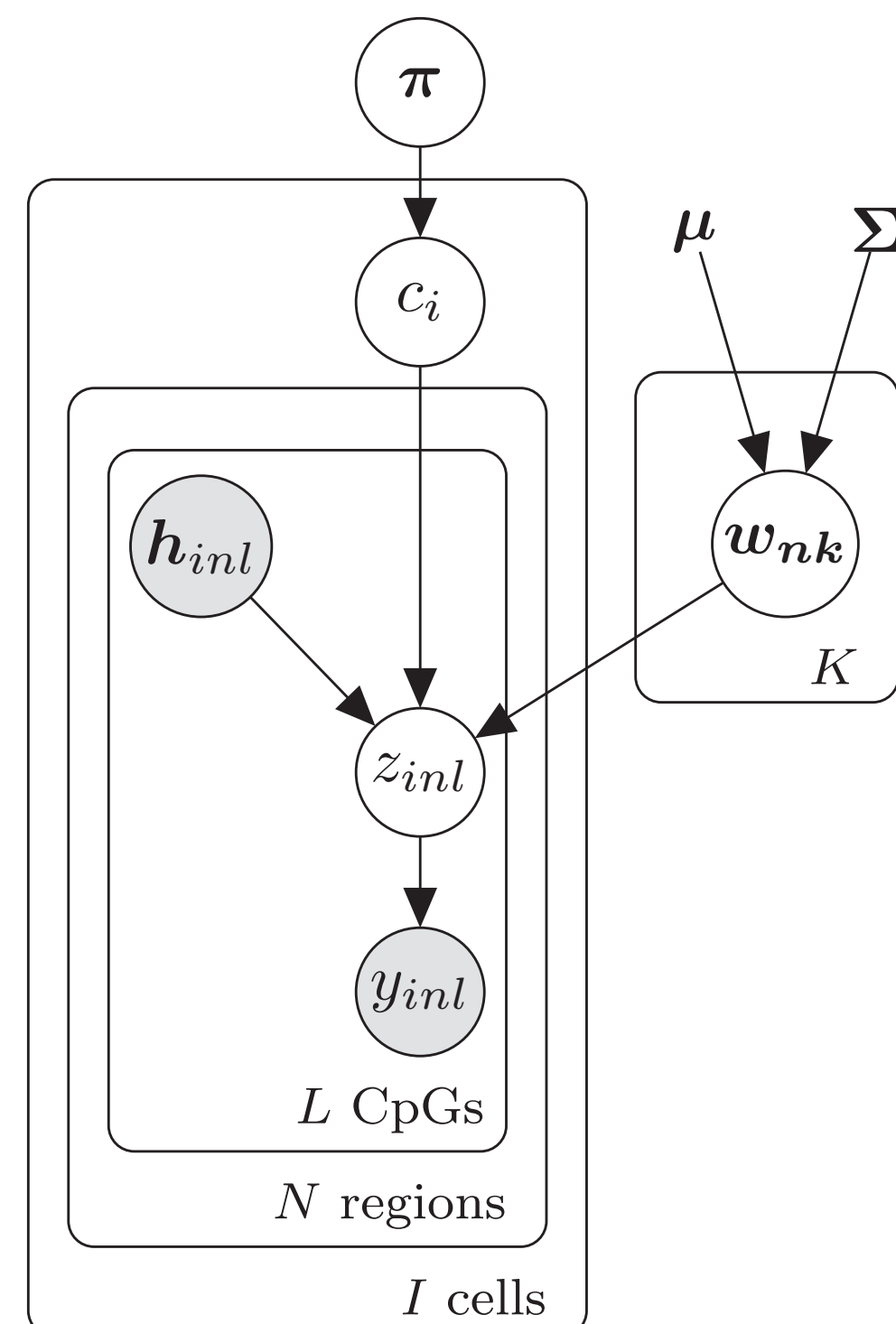
REFERENCES

1. Smallwood, S.A. *et al.*, 2014. Single-cell genome-wide bisulfite sequencing for assessing epigenetic heterogeneity. *Nature methods*, **11**(8), pp.817-820.
2. CAK & GS, 2016. Higher order methylation features for clustering and prediction in epigenomic studies. *Bioinformatics*, **32**(17), pp.i405-i412.
3. Angermueller, C. *et al.*, 2016. Parallel single-cell sequencing links transcriptional and epigenetic heterogeneity. *Nature methods.*, **13**(3), pp.229-232.

SCBPR FINITE DIRICHLET MIXTURE MODEL

Joint posterior distribution

$$p(\mathbf{W}, \mathbf{Z}, \mathbf{c}, \boldsymbol{\pi} | \mathbf{H}, \mathbf{Y}) \propto \left\{ \prod_i \prod_n p(y_{in} | z_{in}) p(z_{in} | \mathbf{w}_{n,c_i}, \mathbf{H}_{in}, c_i) p(c_i | \boldsymbol{\pi}) \right\} p(\boldsymbol{\pi}) \prod_n \prod_k p(\mathbf{w}_{nk})$$



Algorithm 1 Gibbs sampling for scBPR FDMM model

- 1: **initialize** Set $t \leftarrow 1$; set clusters K ; set parameters $\boldsymbol{\pi}^{(0)}, \mathbf{W}^{(0)}$
- 2: **while** $t \leq T$ **do**
- 3: Compute $\gamma(c_{ik})$, probability that cell i belongs to cluster k
- 4: Generate $c_i^{(t)} \sim \text{Discrete}(\gamma(c_i))$
- 5: Generate $\boldsymbol{\pi}^{(t)} \sim \text{Dir}(\{\alpha_k + \sum_i \mathbf{1}(c_i^{(t)} = k)\}_{k=1}^K)$
- 6: Generate $\mathbf{z}_{in}^{(t)} \sim \begin{cases} \mathcal{TN}(\mathbf{h}_{inl} \mathbf{w}_{nk}^{(t-1)}, 1, 0, \infty) & \text{if } y_{inl} = 1 \\ \mathcal{TN}(\mathbf{h}_{inl} \mathbf{w}_{nk}^{(t-1)}, 1, -\infty, 0) & \text{if } y_{inl} = 0 \end{cases}$
- 7: Generate $\mathbf{w}_{nk}^{(t)} \sim \mathcal{N}(\boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$
- 8: **end while**

IDENTIFYING CELL SUB-POPULATIONS

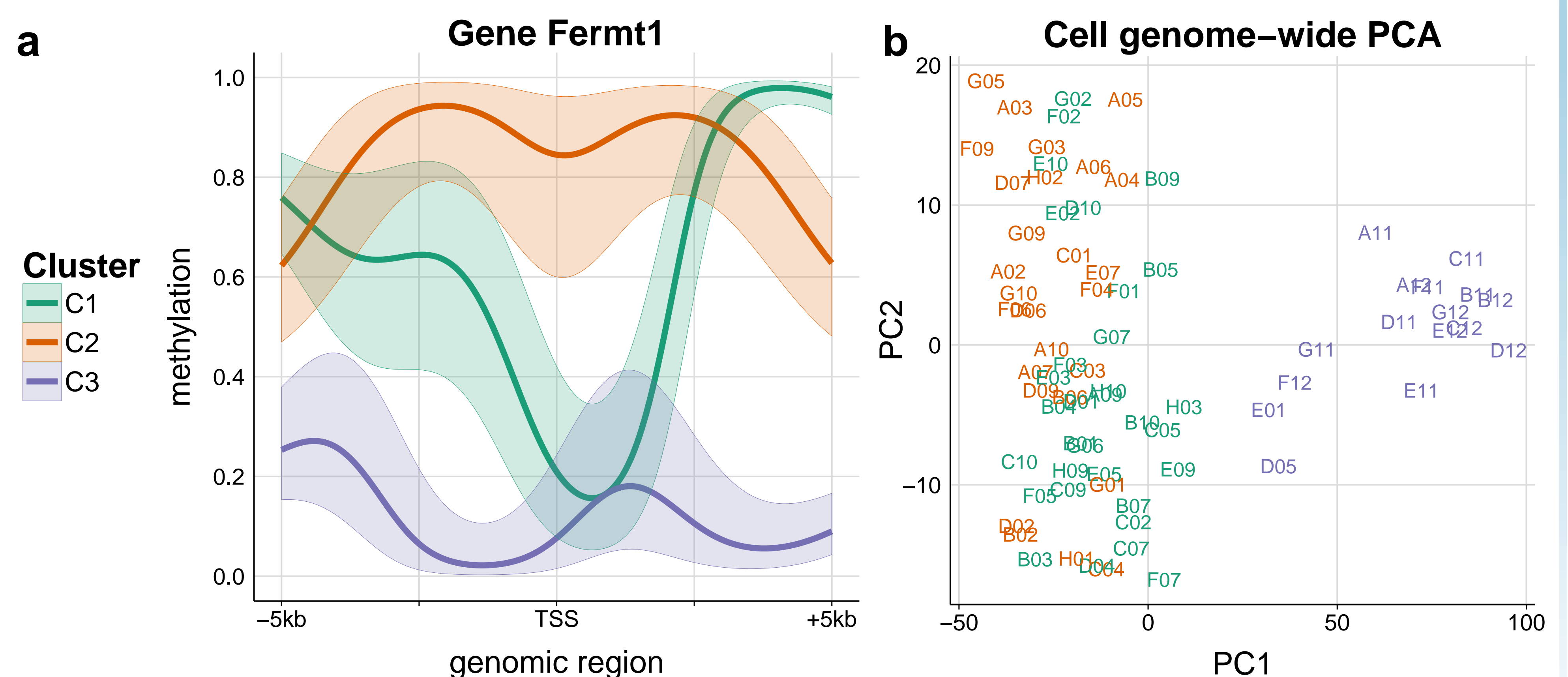


Figure 2: Clustering cells based on promoter methylation profiles on a real dataset³. (a) Methylation profiles for each sub-population of cells for gene *Fermt1*. (b) PCA of cells based on mean genome-wide methylation patterns; blue coloured cells are 2i cells which distinguish well from serum-cultured cells even based on mean methylation.

IMPUTING METHYLATION STATES

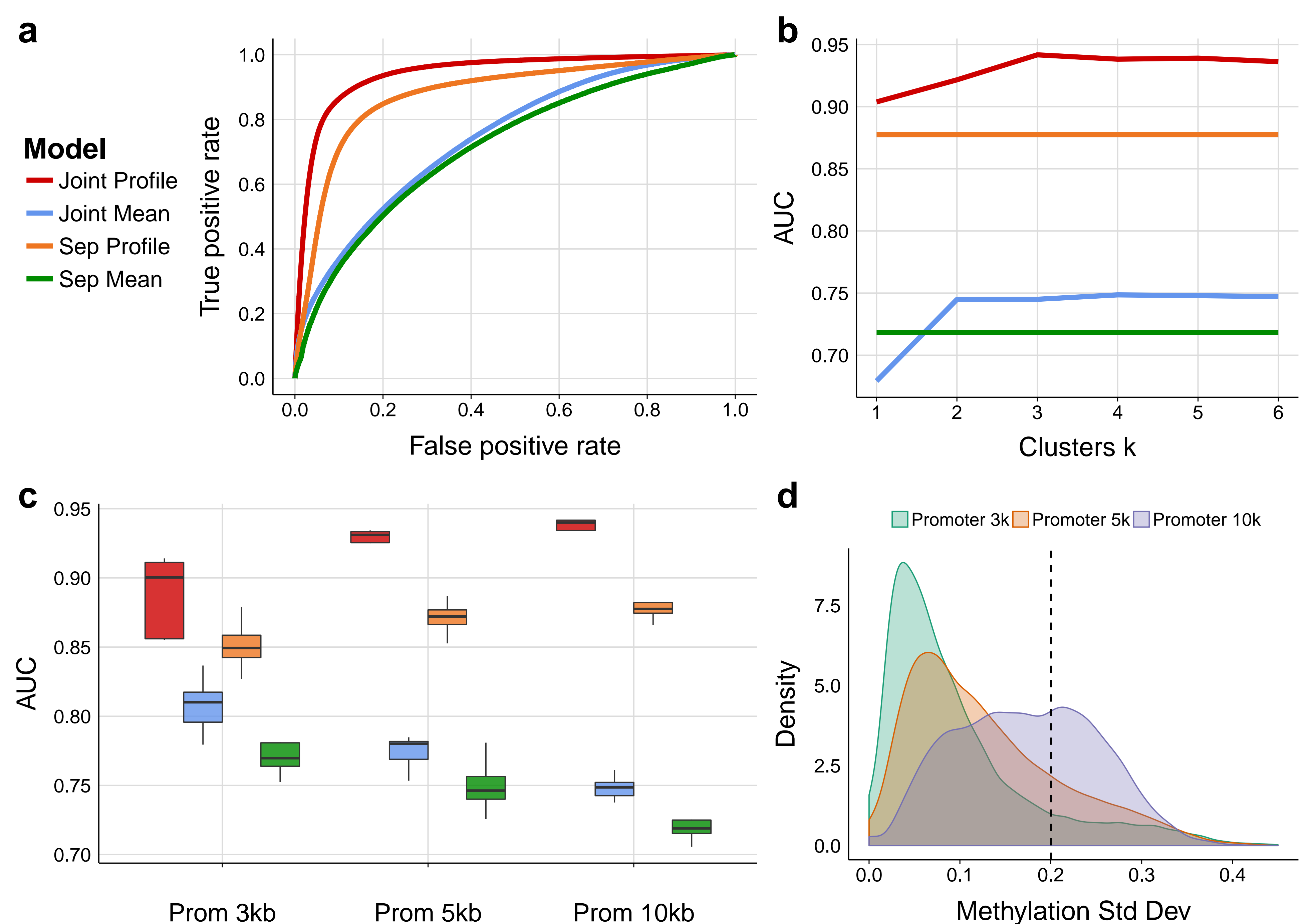


Figure 3: Imputation performance on real data³. (a) ROC curves for 10kb promoter windows and $K = 3$ clusters. (b) AUC while increasing the clusters for 10kb window. (c) Boxplot of AUCs for varying promoter windows, each dot represents a different experiment. (d) Mean methylation variability across cells on different promoter windows.