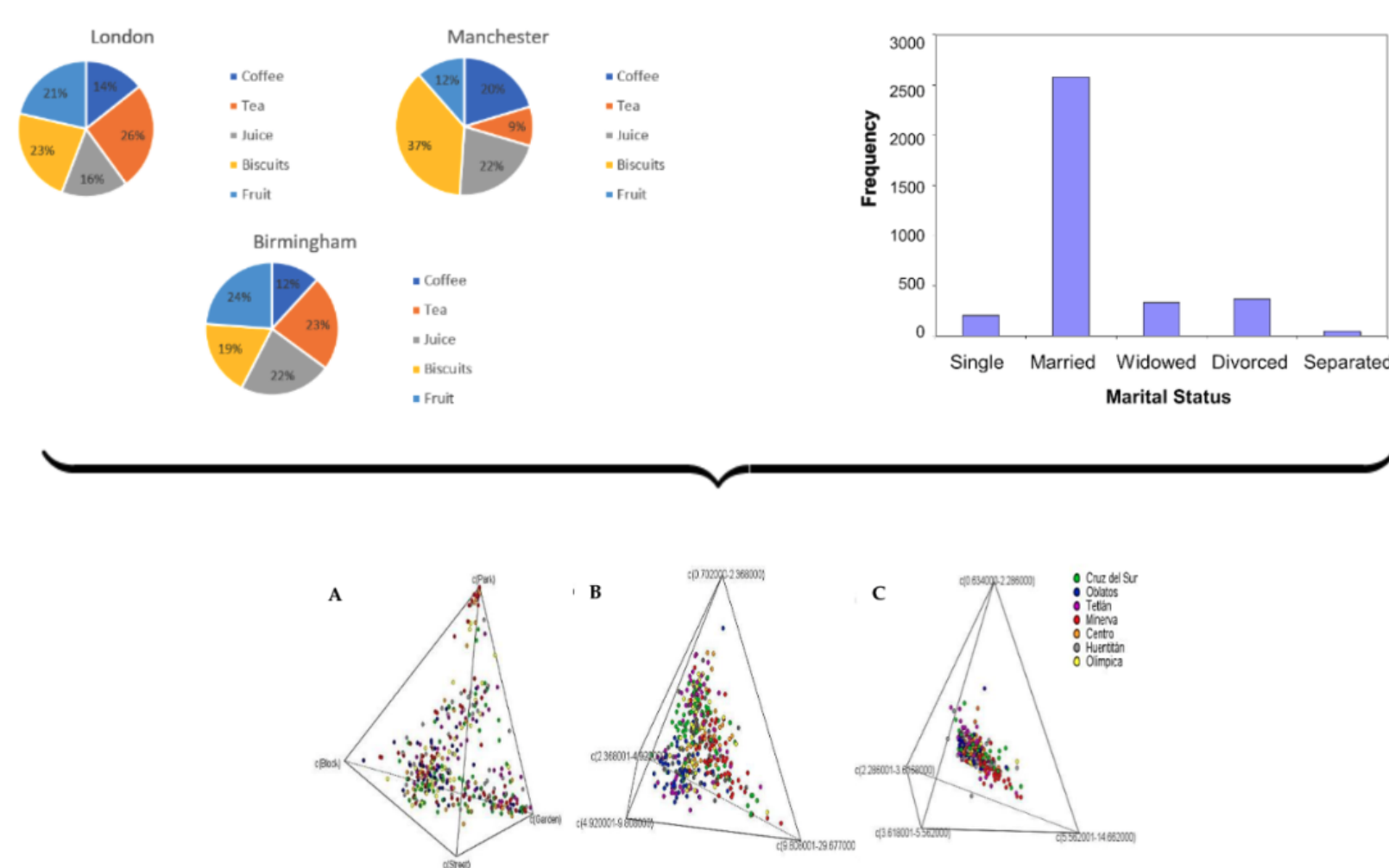


Motivation: compositional data



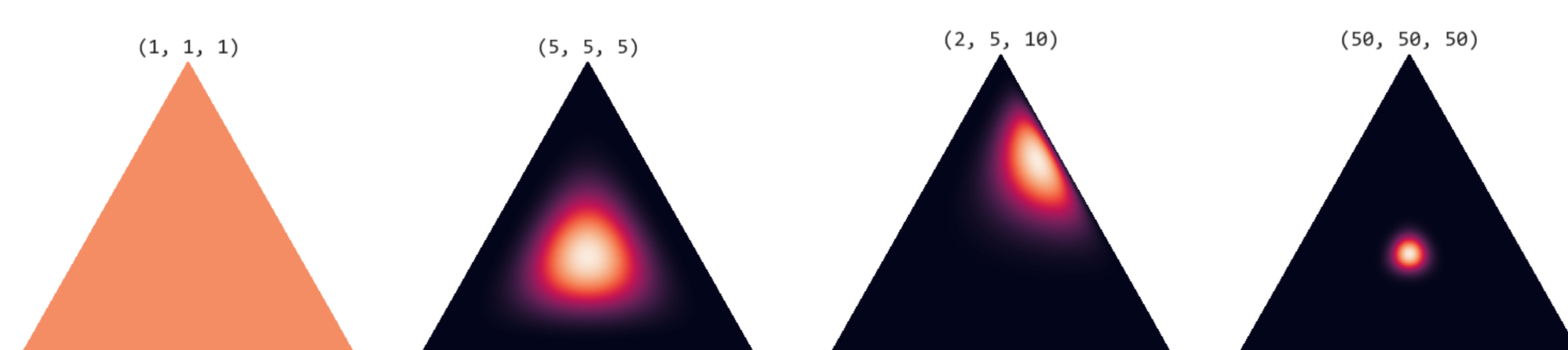
Simplex: $\mathbb{S}^K := \{x \in \mathbb{R}_+^K : \sum_{i=1}^K x_i = 1\}$

Background: shortcomings of the Dirichlet

Definition: $x \sim \text{Dirichlet}(\alpha)$ if $x \in \mathbb{S}^K$ with density:

$$p(x; \alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1}. \quad (1)$$

- ▶ **Extrema.** $\log p(x; \alpha) \rightarrow \pm\infty$ as $x_j \rightarrow 0$.
∴ log-likelihood is undefined in the presence of zeros.
- ▶ **Bias.** Re-write the density in canonical form $p(x; \alpha) = h(x) \exp\left(\sum_{i=1}^K \alpha_i \log x_i - A(\alpha)\right)$.
By theory of exponential families, MLE is unbiased for $\mathbb{E} \log x_j$.
∴ MLE is biased for the mean $\mu_j = \mathbb{E} x_j$.
- ▶ **Flexibility.** If $x_0 \in \mathbb{S}^K$ is a single datapoint, then $\log p(x_0; \alpha) \rightarrow \infty$ as $\alpha \rightarrow \infty$ along $\alpha = kx_0$.
∴ the Dirichlet log-likelihood is ill-behaved under flexible predictive models (e.g. GLMs, neural networks).

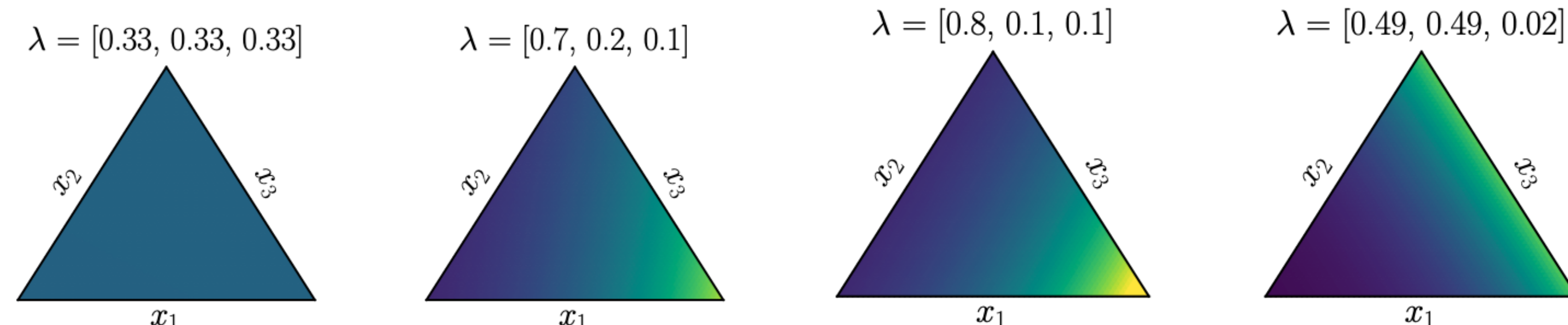


Solution: the continuous categorical

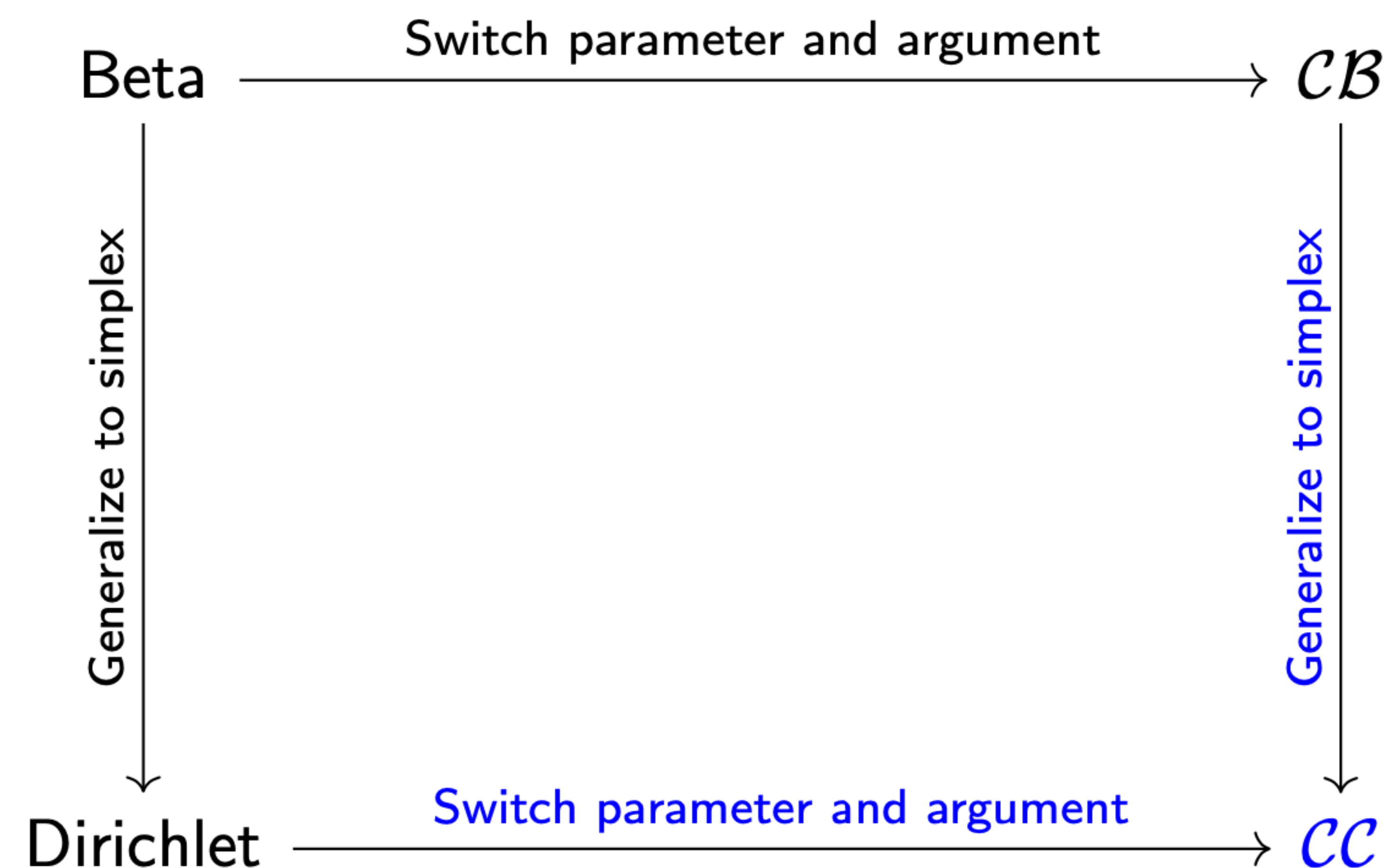
Definition: $x \in \mathbb{S}^K$ follows a *continuous categorical (CC)* distribution with parameter $\lambda \in \mathbb{S}^K$ if:

$$x \sim \mathcal{CC}(\lambda) \iff p(x; \lambda) \propto \prod_{i=1}^K \lambda_i^{x_i}$$

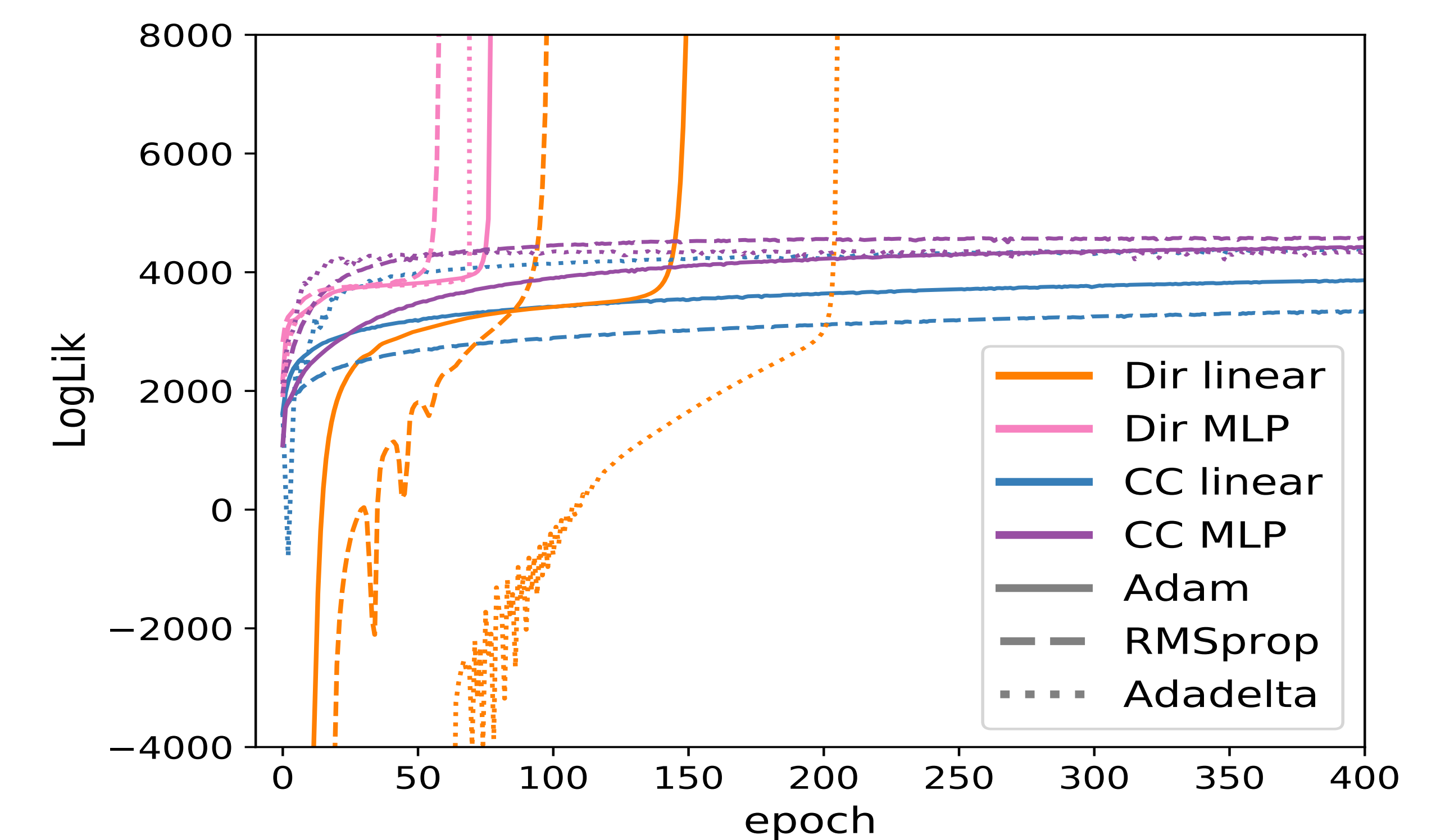
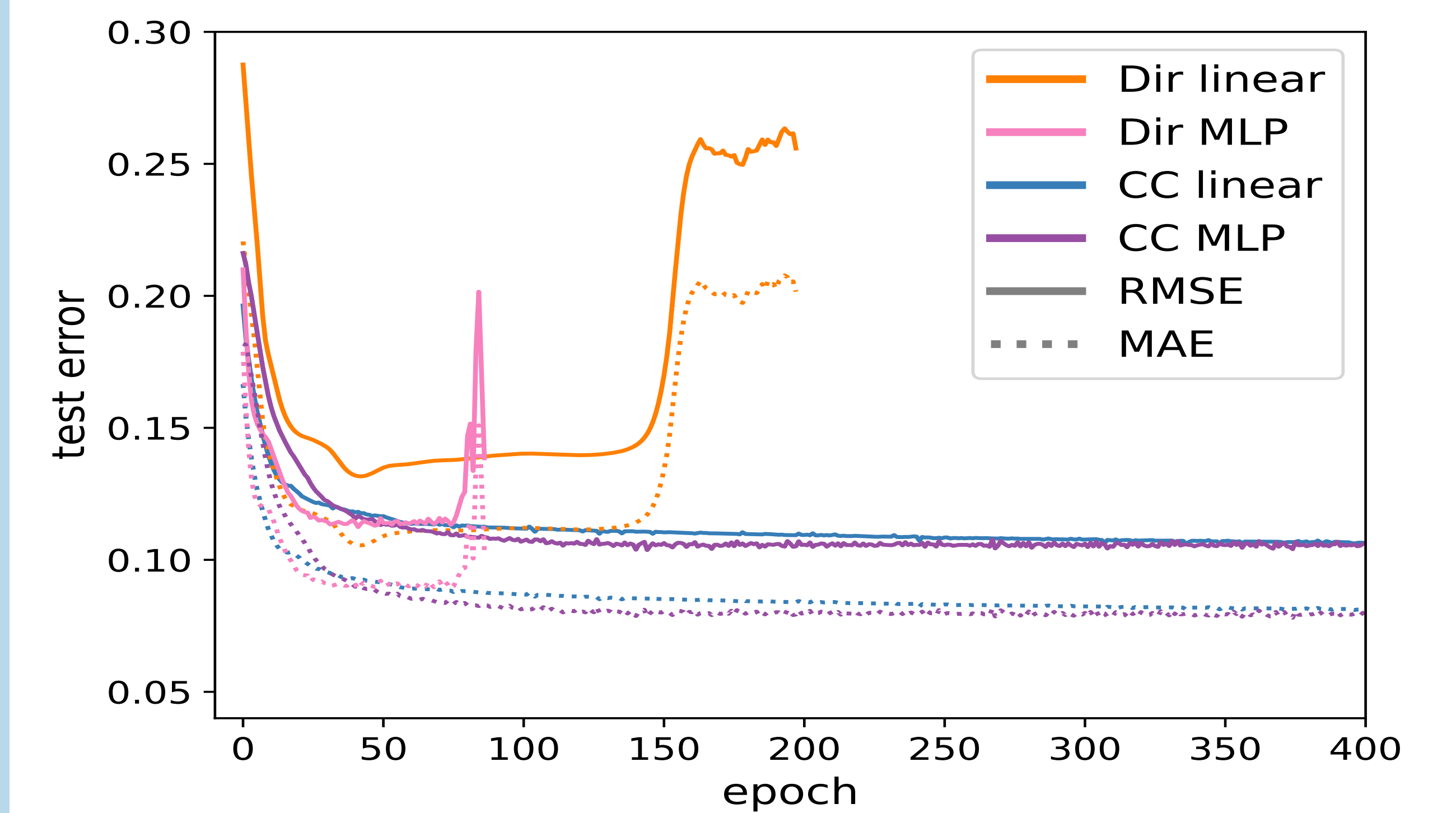
- ▶ **Extrema.** $\log p(x; \lambda)$ is finite at the extrema of the simplex.
∴ log-likelihood is well-defined in the presence of zeros.
- ▶ **Bias.** Re-write the \mathcal{CC} density in canonical form $p(x; \lambda) \propto \exp\left(\sum_{i=1}^K \log(\lambda_i) \cdot x_i\right)$.
∴ by theory of exponential families, MLE is unbiased for the mean $\mu_j = \mathbb{E} x_j$.
- ▶ **Flexibility.** The \mathcal{CC} density is convex in x .
∴ cannot represent interior modes, cannot concentrate mass on interior points and log-likelihood does not diverge.



Related distributions



Experiments: election data



Experiments: knowledge distillation

