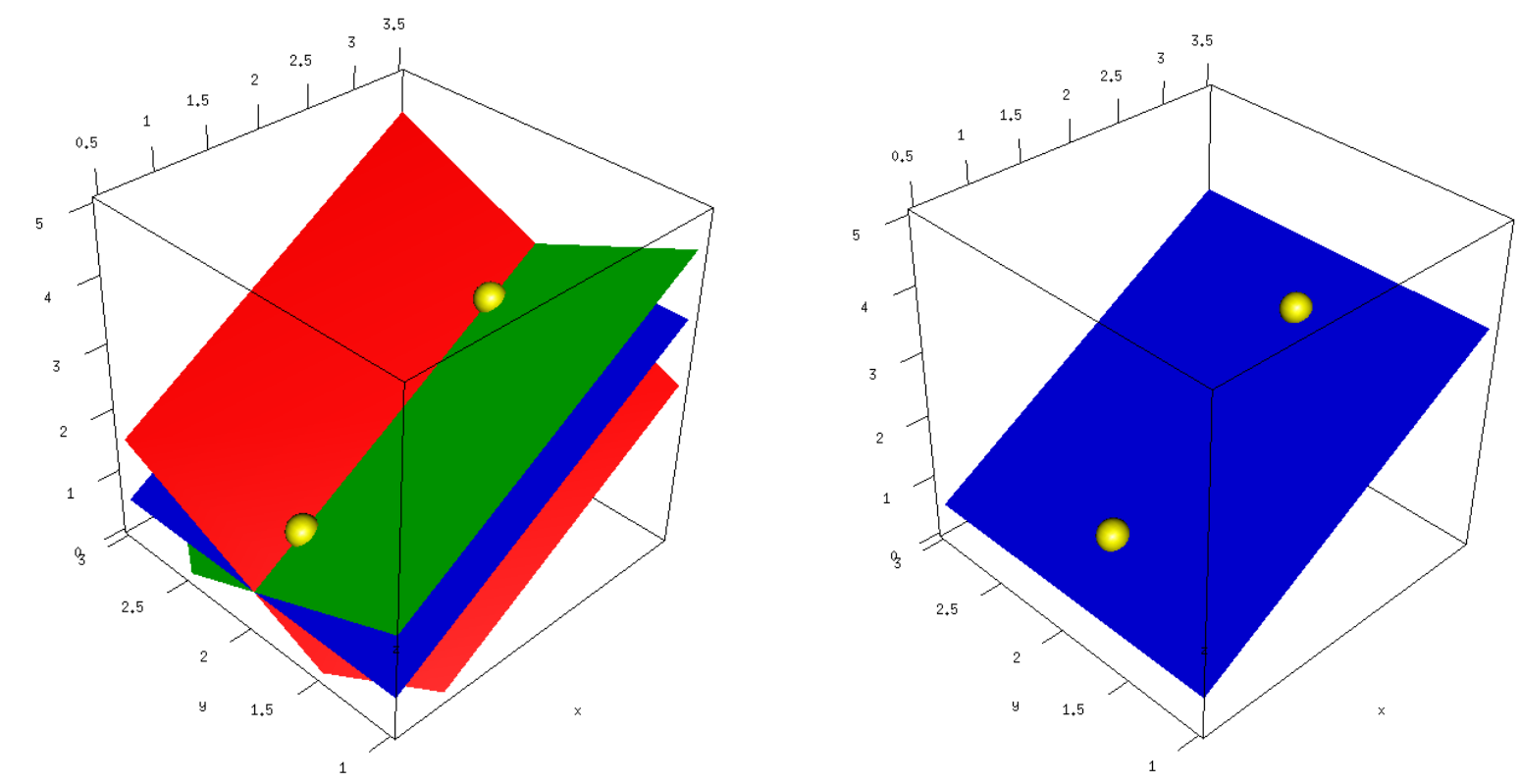


1. Introduction

We focus on linear regression problems that are **under-determined**, *i.e.*, we have the same or more attributes than observations ($n \leq d$).

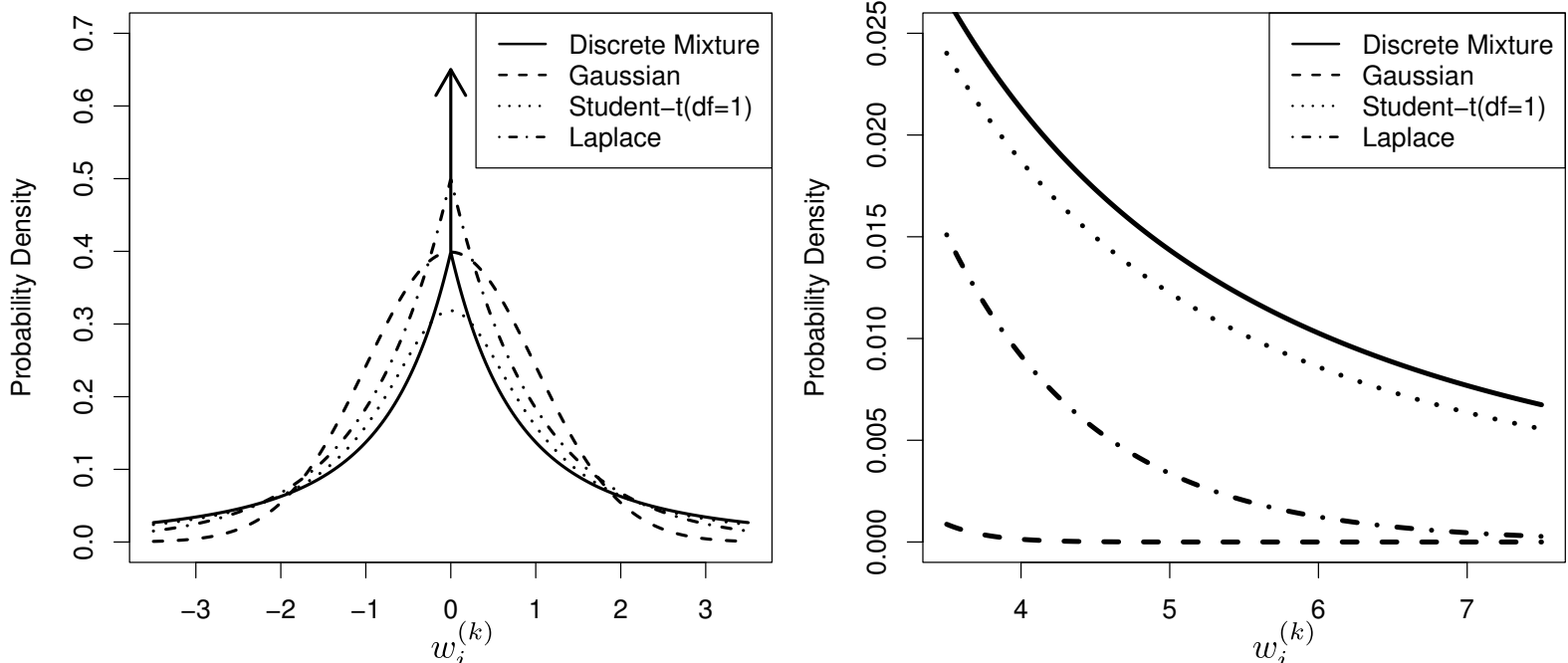
$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}\sigma^2).$$



A typical regularization assumes **sparsity** in \mathbf{w} .

2. Sparsity Assumption

Introduced by setting a **sparse enforcing prior** for \mathbf{w} , *e.g.*, Laplace, Student's T, Horseshoe or spike-and-slab.



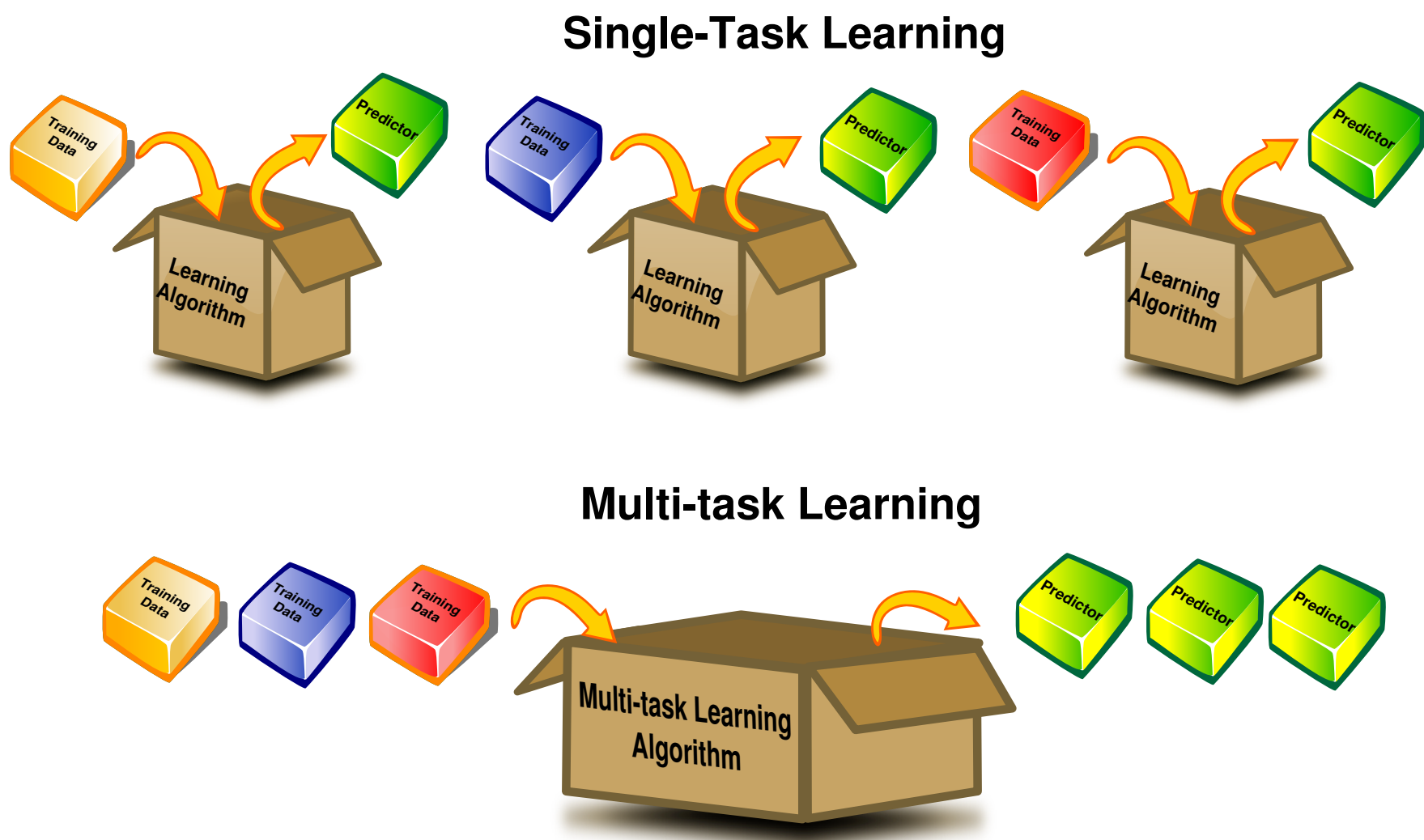
Discrete-mixture prior: $p(\mathbf{w}_i) = \rho\delta_0 + (1 - \rho)\pi(\mathbf{w}_i)$.

$$\pi(\mathbf{w}_i) = \int \mathcal{N}(\mathbf{w}_i | \mathbf{0}, \lambda_i^2) \frac{\lambda_i}{(\lambda_i^2 + 1)^{3/2}} d\lambda_i = \frac{1}{\sqrt{2\pi}} \left(1 - |\mathbf{w}_i| \frac{\Phi(-|\mathbf{w}_i|)}{\mathcal{N}(\mathbf{w}_i | \mathbf{0}, 1)} \right)$$

is the Strawderman-Bergen prior which has a **closed form convolution** with the Gaussian distribution.

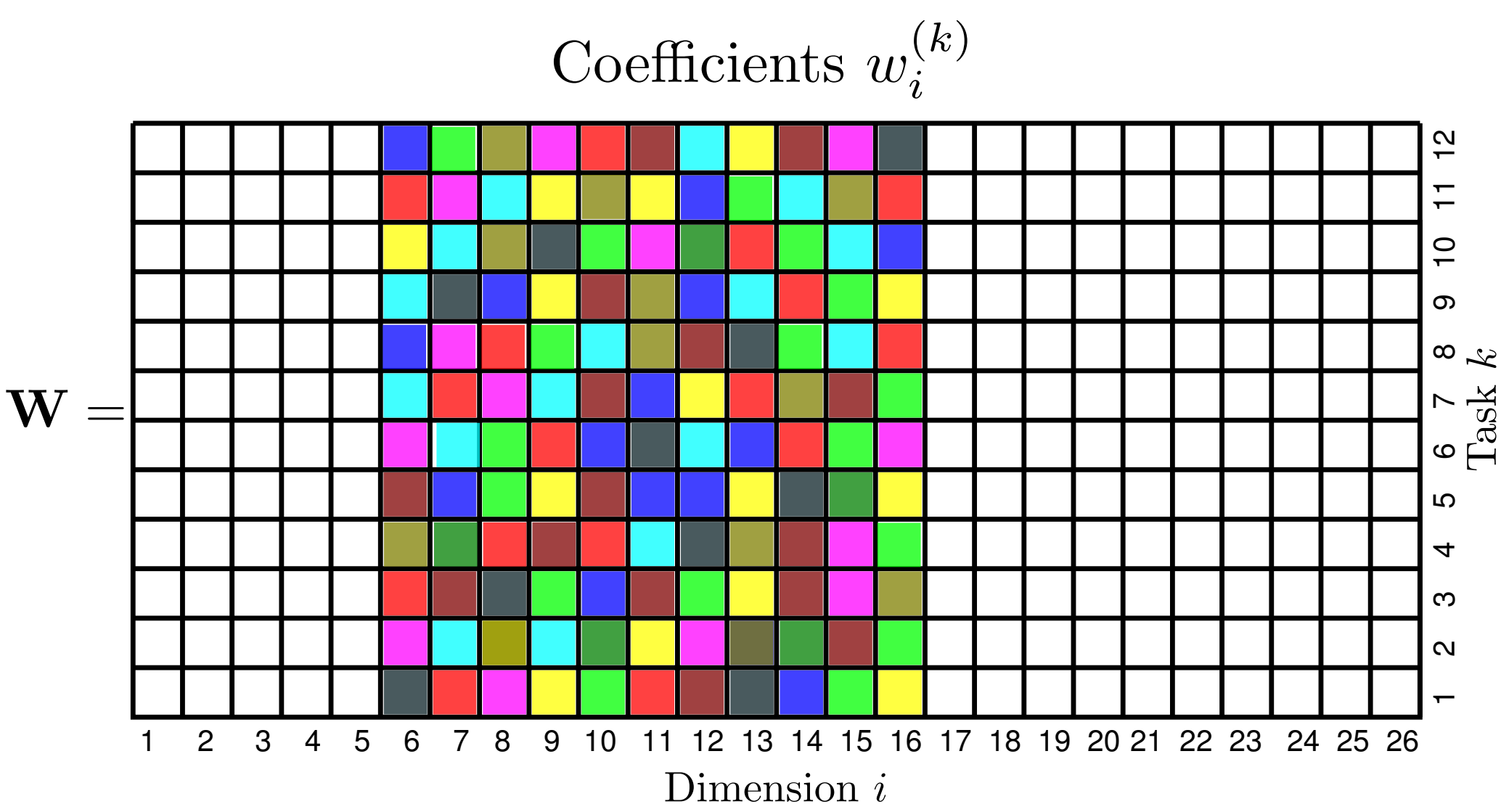
3. Multi-task Learning

There may be several learning tasks available for induction. Multi-task methods try to exploit **similarities** among tasks to improve the induction process.



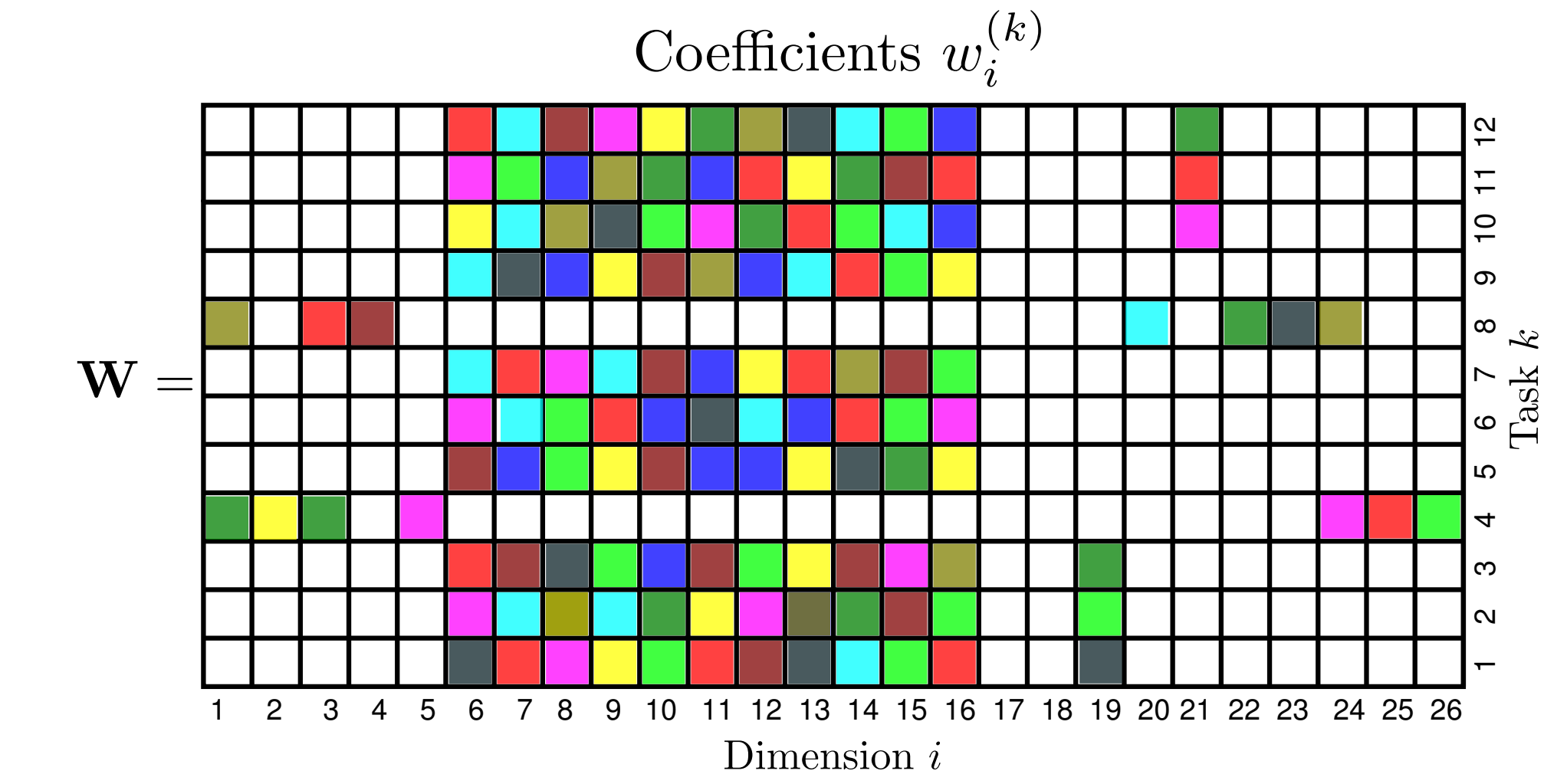
4. Typical Hypothesis

Tasks **share relevant and irrelevant** features.



5. Something More Reasonable

Outlier tasks and **outlier features**.



6. Robust Prior Distribution

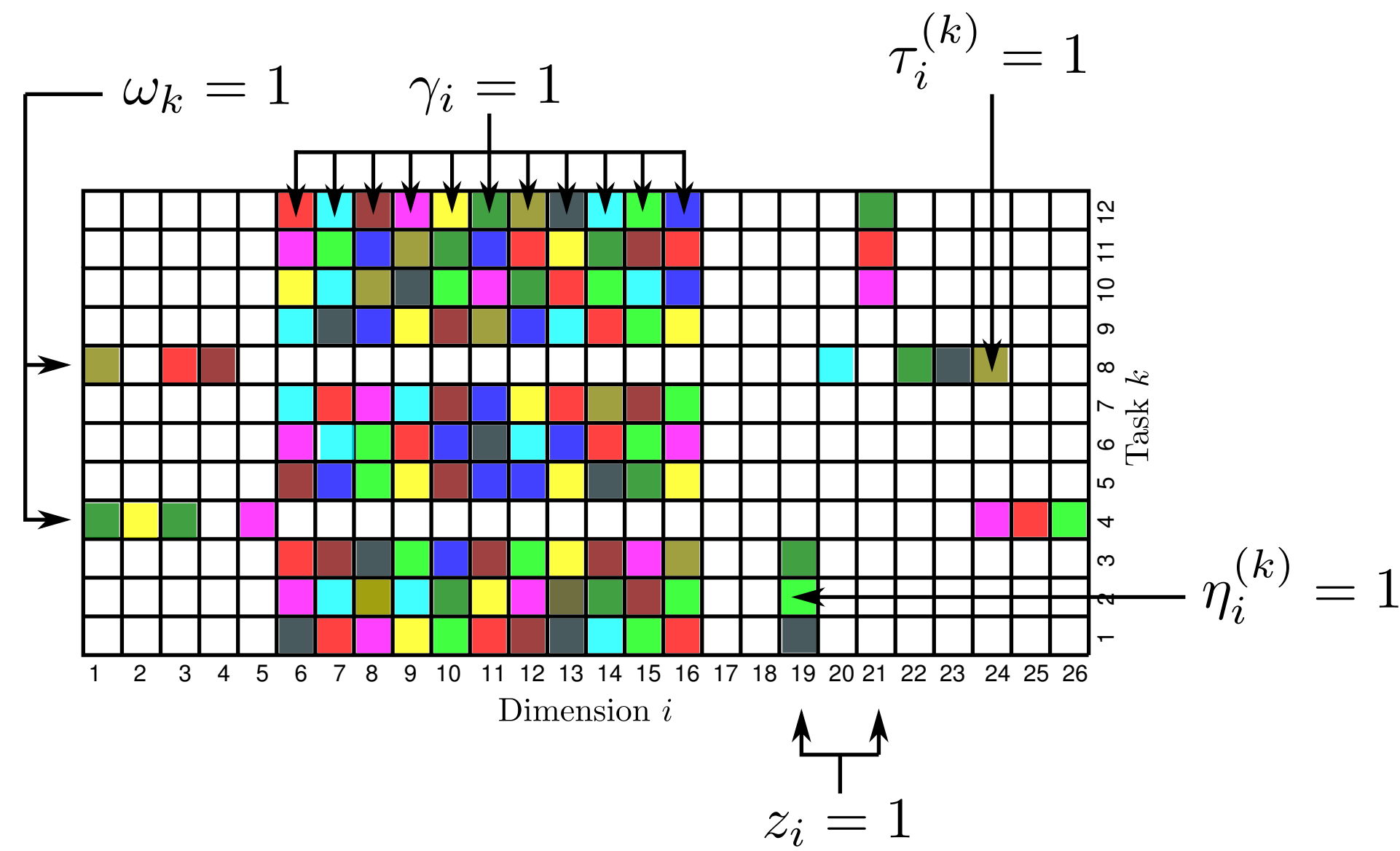
Define $\Omega = \{\mathbf{z}, \omega, \gamma, \{\tau^{(k)}\}_{k=1}^K, \{\eta^{(k)}\}_{k=1}^K\}$. The prior for \mathbf{W} is $p(\mathbf{W} | \Omega) = \prod_{i=1}^d \prod_{k=1}^K p(\mathbf{w}_i^{(k)} | \Omega)$, with $p(\mathbf{w}_i^{(k)} | \Omega) =$

$$\left\{ \left[\pi(\mathbf{w}_i^{(k)})^{\tau_i^{(k)}} \delta_0^{1-\tau_i^{(k)}} \right]^{\omega_k} \left[\pi(\mathbf{w}_i^{(k)})^{\gamma_i} \delta_0^{1-\gamma_i} \right]^{1-\omega_k} \right\}^{1-z_i} \times \left\{ \pi(\mathbf{w}_i^{(k)})^{\eta_i^{(k)}} \delta_0^{1-\eta_i^{(k)}} \right\}^{z_i},$$

Not outlier feature
Not outlier task

Outlier task
Outlier feature

where δ_0 is a point of probability mass at the origin.



5. Dirty Multi-task Feature Selection

Gaussian Likelihood

$$\mathbf{y}^{(k)} \sim \mathcal{N}(\mathbf{X}^{(k)} \mathbf{w}^{(k)}, \mathbf{0}, \mathbf{I}\sigma^2_{(k)}), \quad \forall k.$$

$$\mathbf{w}_i^{(k)} \sim \text{RobustPrior}(\mathbf{z}_i, \omega_k, \gamma_i, \tau_i^{(k)}, \eta_i^{(k)}), \quad \forall i, k.$$

$$\sigma^2_{(k)} \sim \text{InvGam}(5, 5), \quad \forall k.$$

Independence assumption

$$\mathbf{z}_i \sim \text{Bernoulli}(\rho_z), \quad \forall i,$$

$$\omega_k \sim \text{Bernoulli}(\rho_\omega), \quad \forall k,$$

$$\gamma_i \sim \text{Bernoulli}(\rho_\gamma), \quad \forall i,$$

$$\tau_i^{(k)} \sim \text{Bernoulli}(\rho_\tau), \quad \forall i, k,$$

$$\eta_i^{(k)} \sim \text{Bernoulli}(\rho_\eta), \quad \forall i, k,$$

Noise hyper-prior

$$\rho_z \sim \text{Beta}(1, 1),$$

$$\rho_\omega \sim \text{Beta}(1, 1),$$

$$\rho_\gamma \sim \text{Beta}(1, 1),$$

$$\rho_\tau \sim \text{Beta}(1, 1),$$

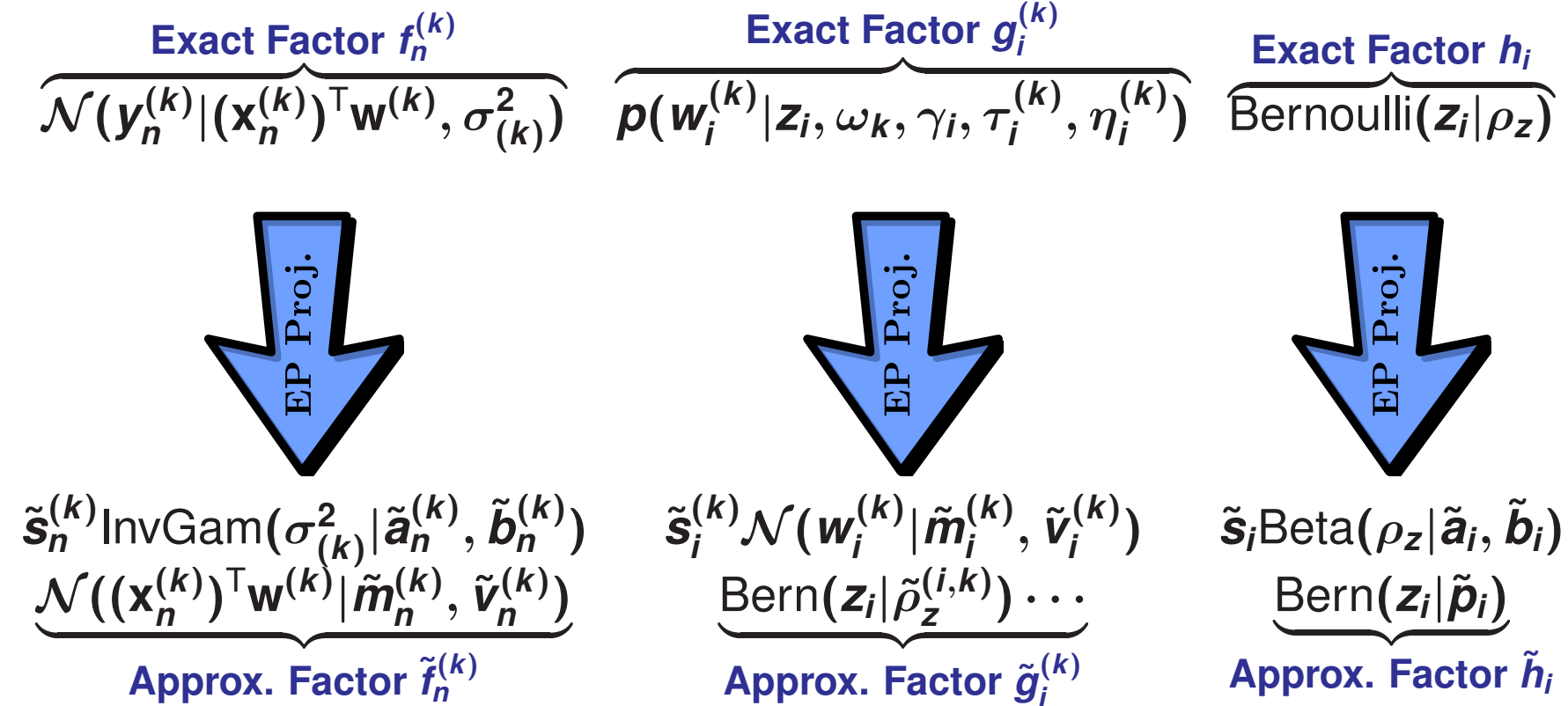
$$\rho_\eta \sim \text{Beta}(1, 1),$$

Prob. hyper-prior

7. Expectation Propagation

Approximates each factor in $p(\mathbf{Y}, \mathbf{W}, \Omega, \rho, \sigma^2 | \mathcal{X})$ with an unnormalized distribution inside an **exponential family** \mathcal{F} . $\mathcal{F} \rightarrow$ Gaussians on \mathbf{W} , Bernoullis on Ω , I. Gammas on σ^2 and Betas on ρ .

Example of factors that need approximation:



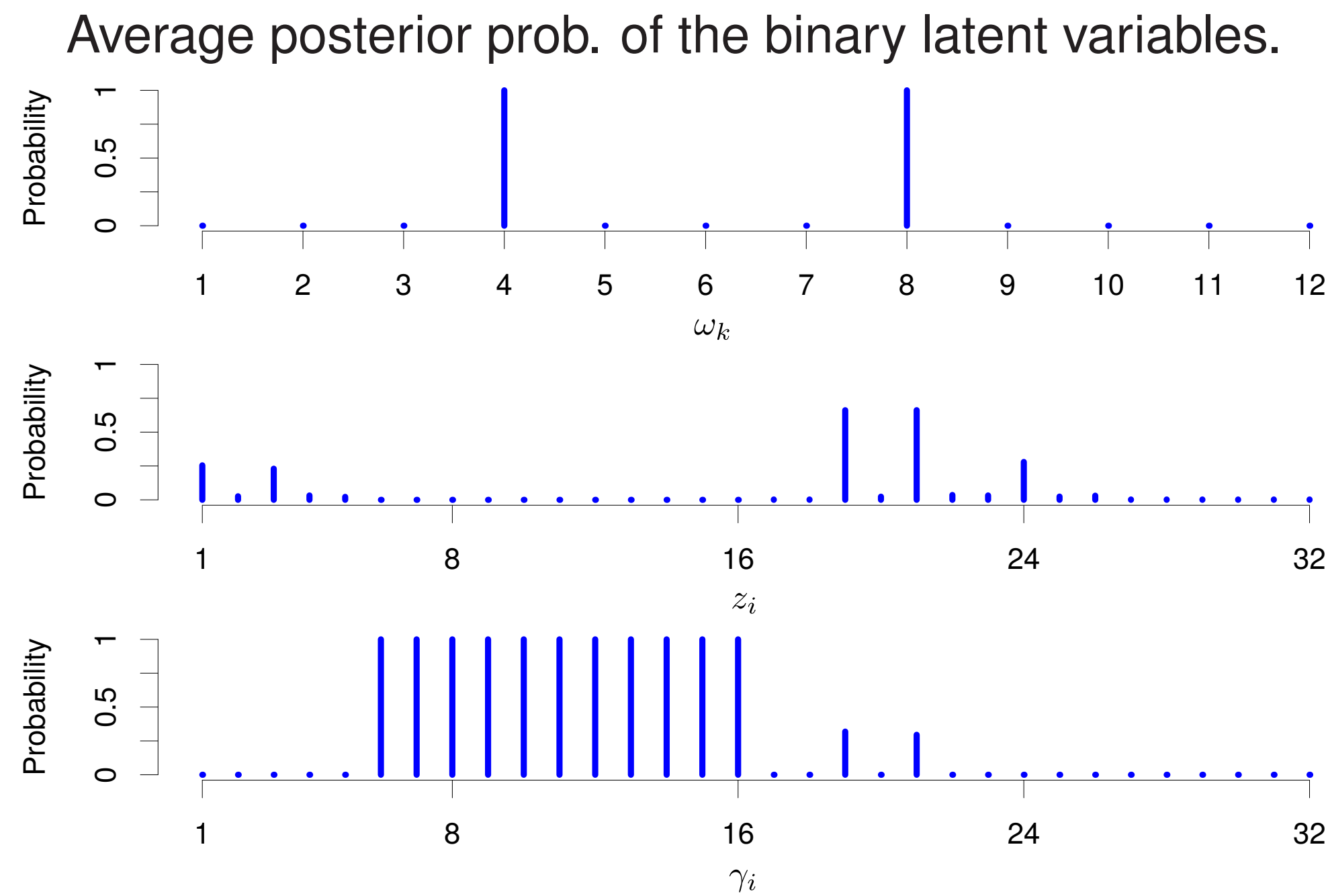
8. Experiments with Synthetic Data

$N = 200$, $d = 2,000$ and we use for \mathbf{W} the pattern above. $K = 12$, $\sigma^2_{(k)} = 0.5, \forall k$ and $\mathbf{w}_i^{(k)} \sim \text{Student}(\text{df} = 5)$.

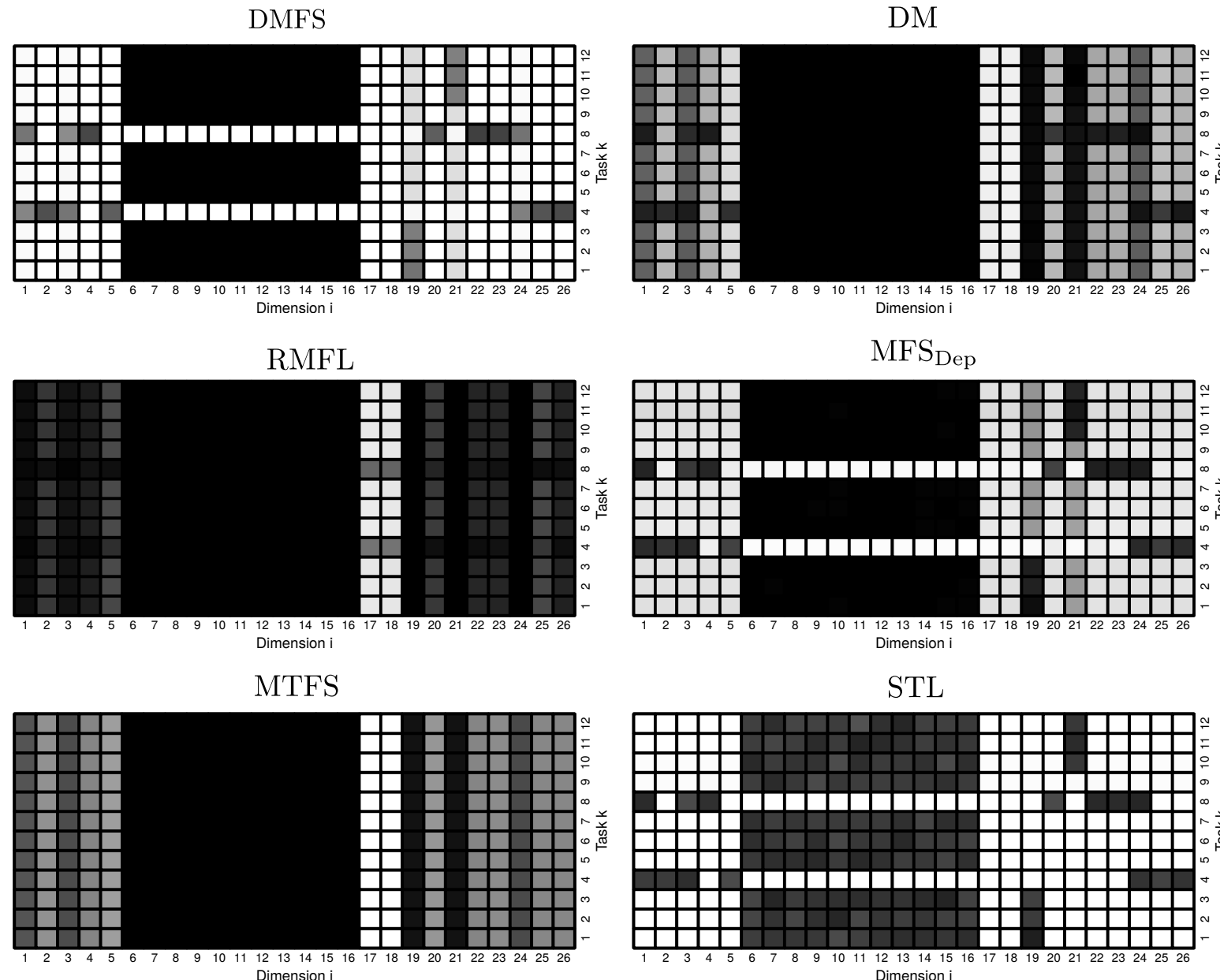
Method	Test RMSE	Rec. Error	Training Time
DMFS	0.73 ± 0.04	0.22 ± 0.02	21.29 ± 0.2
DM	0.86 ± 0.05	0.50 ± 0.03	150.35 ± 10.0
RMFL	0.90 ± 0.05	0.56 ± 0.03	95.42 ± 5.0
MFS _{Dep}	0.77 ± 0.06	0.32 ± 0.04	2 · 10³ ± 4 · 10²
MFS	0.81 ± 0.06	0.37 ± 0.04	6.7 ± 1.7
STL	0.78 ± 0.07	0.33 ± 0.06	4.76 ± 0.4

MFS and STL are **particular cases** of DMFS.

Several other experiments in the paper!



Average prob. that each coefficient is different from zero.



9. Summary and Conclusions of the Research Work

(1) - Most methods for multi-task feature selection assume jointly relevant and irrelevant features, which may be **too restrictive**. **(2)** - A robust prior allows tasks with **specific** relevant and irrelevant coefficients, and features to be **arbitrarily** relevant or irrelevant. **(3)** - Exact inference is infeasible under such a prior. However, a **quadrature-free** expectation propagation algorithm is possible. **(4)** - Several experiments show **gains** in the prediction performance and in the identification of relevant features for prediction. **(5)** - Our new prior is useful to **better understand** the data because it allows to identify outlier tasks and outlier features.