

# Supplementary materials for “Drug susceptibility prediction against a panel of drugs using kernelized Bayesian multitask learning”

Mehmet Gönen and Adam A. Margolin

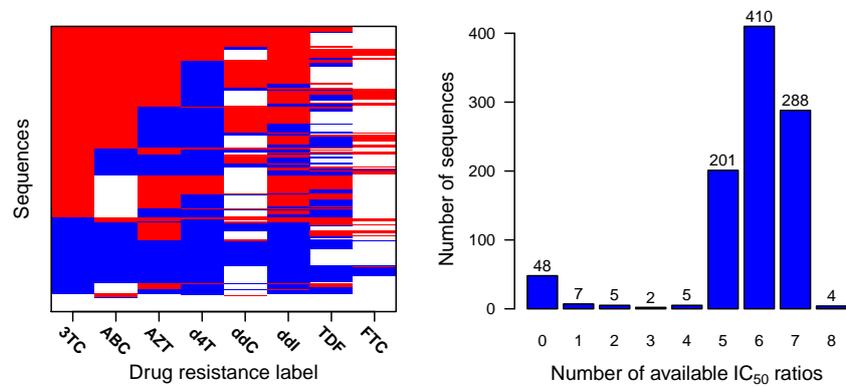


Figure 1: Drug resistance labels and histogram of available IC<sub>50</sub> ratios for 970 reverse transcriptase sequences. For drug resistance labels, red: resistant; blue: susceptible; white: missing

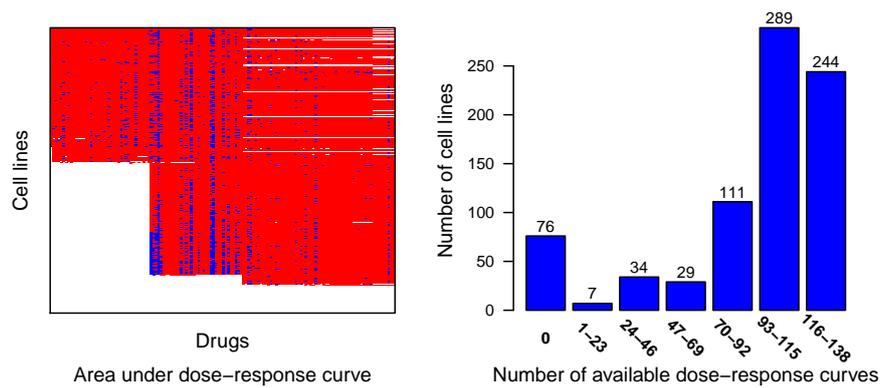


Figure 2: AUC values and histogram of available dose-response curves for 790 cell lines. For AUC values, red: higher; blue: lower; white: missing

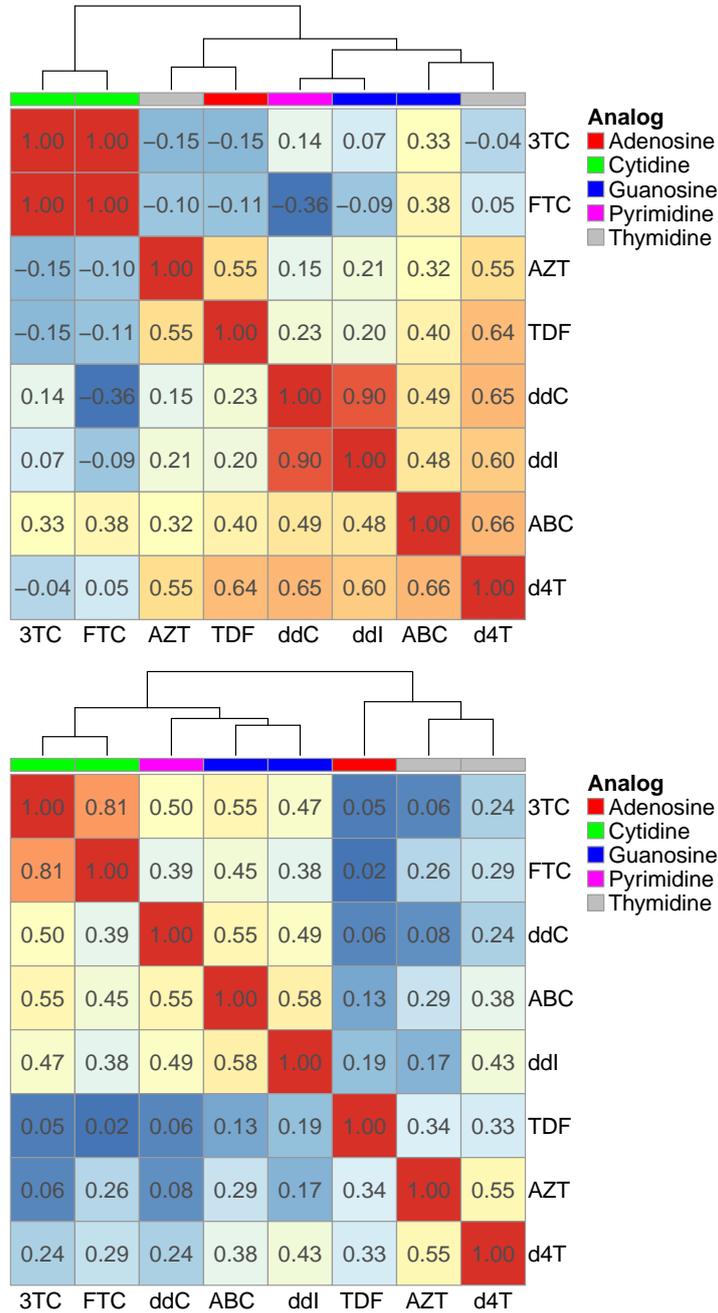


Figure 3: Hierarchical clustering of HIV-1 drugs based on correlation matrices. Top: Using the correlation matrix calculated over the measured  $IC_{50}$  ratios. Bottom: Using the correlation matrix calculated over the task-specific classification parameters found by kernelized Bayesian multitask learning

# 1 Inference details of kernelized Bayesian multitask learning

The approximate posterior distributions of the precision priors can be updated as

$$\begin{aligned}\alpha(\lambda_s^i) &= \alpha_\lambda + 1/2 \\ \beta(\lambda_s^i) &= (1/\beta_\lambda + \langle (a_s^i)^2 \rangle / 2)^{-1},\end{aligned}$$

where  $\langle g(\cdot) \rangle$  denotes the posterior expectation as usual, i.e.  $E_{q(\cdot)}[g(\cdot)]$ .

The approximate posterior distribution of the projection matrix can be updated as

$$\begin{aligned}\Sigma(\mathbf{a}_s) &= (\text{diag}(\langle \boldsymbol{\lambda}_s \rangle) + \mathbf{K}\mathbf{K}^\top / \sigma_h^2)^{-1} \\ \mu(\mathbf{a}_s) &= \Sigma(\mathbf{a}_s)(\mathbf{K}\langle \mathbf{h}^s \rangle^\top / \sigma_h^2),\end{aligned}$$

and the approximate posterior distribution of the hidden representation for each data point can be updated as

$$\begin{aligned}\Sigma(\mathbf{h}_i) &= \left( \mathbf{I} / \sigma_h^2 + \sum_{t \in \mathcal{T}_i} \langle \mathbf{w}_t \mathbf{w}_t^\top \rangle \right)^{-1} \\ \mu(\mathbf{h}_i) &= \Sigma(\mathbf{h}_i) \left( \langle \mathbf{A}^\top \rangle \mathbf{k}_i / \sigma_h^2 + \sum_{t \in \mathcal{T}_i} [\langle f_{t,i} \rangle \langle \mathbf{w}_t \rangle] \right),\end{aligned}$$

where  $\mathcal{T}_i$  gives the indices of tasks with given class labels for data point  $i$ . Note that, for each data point, we use the tasks with given class labels only to update the shared hidden representation.

The approximate posterior distribution of the classification parameters for each task can be updated as

$$\begin{aligned}\Sigma(\mathbf{w}_t) &= \left( \mathbf{I} / \sigma_w^2 + \sum_{i \in \mathcal{I}_t} \langle \mathbf{h}_i \mathbf{h}_i^\top \rangle \right)^{-1} \\ \mu(\mathbf{w}_t) &= \Sigma(\mathbf{w}_t) \sum_{i \in \mathcal{I}_t} [\langle f_{t,i} \rangle \langle \mathbf{h}_i \rangle],\end{aligned}$$

where it can be seen that the inference mechanism transfers information between the tasks because the task-specific classification parameters are updated using the shared hidden representations.

The approximate posterior distribution of the predicted outputs can be updated as

$$\begin{aligned}\Sigma(f_{t,i}) &= 1 \\ \mu(f_{t,i}) &= \langle \mathbf{w}_t^\top \rangle \langle \mathbf{h}_i \rangle \\ \rho(f_{t,i}) &\triangleq f_{t,i} y_{t,i} > \nu,\end{aligned}$$

where we can fortunately calculate the expectation of the truncated normal distribution in closed-form.

The inference procedure summarized in Algorithm ?? sequentially updates the approximate posterior distributions of the priors, latent variables and model parameters until convergence, which can be checked by monitoring the lower bound. The first term of the lower bound corresponds to the sum of exponential form expectations of the distributions in the joint likelihood. The second term is the sum of negative entropies of the approximate posteriors in the ensemble. The only nonstandard distribution in these terms is the truncated normal distribution used for the predicted outputs; nevertheless, the truncated normal distribution has a closed-form formula also for its entropy.

---

**Algorithm 1** Kernelized Bayesian Multitask Learning (KBMTL)

---

**Require:**  $\mathbf{K}$ ,  $\{\mathbf{y}_t\}_{t=1}^T$ ,  $R$ ,  $\alpha_\lambda$ ,  $\beta_\lambda$ ,  $\sigma_h$ ,  $\sigma_w$  and  $\nu$   
1: Initialize  $q(\mathbf{A})$ ,  $q(\mathbf{H})$  and  $\{q(\mathbf{w}_t), q(\mathbf{f}_t)\}_{t=1}^T$  randomly  
2: **repeat**  
3:   Update  $q(\mathbf{\Lambda})$ ,  $q(\mathbf{A})$  and  $q(\mathbf{H})$   
4:   Update  $\{q(\mathbf{w}_t), q(\mathbf{f}_t)\}_{t=1}^T$   
5: **until** convergence  
6: **return**  $q(\mathbf{A})$  and  $\{q(\mathbf{w}_t)\}_{t=1}^T$

---

## 2 Extension to regression problems

In regression problems, for each task, we are given an output vector  $\mathbf{y}_t = \{y_{t,i} \in \mathbb{R}\}_{i \in \mathcal{I}_t}$  instead of a label vector. The binary classification part of our method is replaced with the following distributional assumptions:

$$\begin{aligned} w_{t,s} &\sim \mathcal{N}(w_{t,s}; 0, \sigma_w^2) && \forall (t, s) \\ \epsilon_t &\sim \mathcal{G}(\epsilon_t; \alpha_\epsilon, \beta_\epsilon) && \forall t \\ y_{t,i} | \mathbf{h}_i, \mathbf{w}_t, \epsilon_t &\sim \mathcal{N}(y_{t,i}; \mathbf{w}_t^\top \mathbf{h}_i, \epsilon_t^{-1}) && \forall (t, i \in \mathcal{I}_t), \end{aligned}$$

where we do not need to use a bias term by centering the output vector to zero mean. We again used a variational approximation for inference, and the factorable ensemble approximation of the required posterior becomes

$$p(\Theta | \mathbf{K}, \{\mathbf{y}_t\}_{t=1}^T) \approx q(\Theta) = q(\mathbf{\Lambda})q(\mathbf{A})q(\mathbf{H}) \prod_{t=1}^T [q(\mathbf{w}_t)q(\epsilon_t)],$$

where  $q(\mathbf{\Lambda})$ ,  $q(\mathbf{A})$ ,  $q(\mathbf{H})$  and  $q(\mathbf{w}_t)$  remain intact. One additional factor in the ensemble is defined as

$$q(\epsilon_t) = \mathcal{G}(\epsilon_t; \alpha(\epsilon_t), \beta(\epsilon_t)).$$

The approximate posterior distribution of the regression parameters for each task can be updated as

$$\begin{aligned} \Sigma(\mathbf{w}_t) &= \left( \mathbf{I} / \sigma_w^2 + \sum_{i \in \mathcal{I}_t} \langle \epsilon_t \rangle \langle \mathbf{h}_i \mathbf{h}_i^\top \rangle \right)^{-1} \\ \mu(\mathbf{w}_t) &= \Sigma(\mathbf{w}_t) \sum_{i \in \mathcal{I}_t} \left[ \langle \epsilon_t \rangle \langle \mathbf{f}_{t,i} \rangle \langle \mathbf{h}_i \rangle \right], \end{aligned}$$

and the approximate posterior distribution of the additional precision priors can be updated as

$$\alpha(\epsilon_t) = \alpha_\epsilon + N_t/2$$

$$\beta(\epsilon_t) = \left( 1/\beta_\epsilon + \sum_{i \in \mathcal{I}_t} \langle (y_{t,i} - \mathbf{w}_t \mathbf{h}_i)^2 \rangle / 2 \right)^{-1},$$

where other update equations for  $q(\boldsymbol{\Lambda})$ ,  $q(\mathbf{A})$  and  $q(\mathbf{H})$  are not changed.

### 3 Baseline algorithms

#### 3.1 Bayesian probit classifier

Its distributional assumptions are defined as

$$\begin{aligned} \gamma &\sim \mathcal{G}(\gamma; \alpha_\gamma, \beta_\gamma) \\ b|\gamma &\sim \mathcal{N}(b; 0, \gamma^{-1}) \\ \eta_f &\sim \mathcal{G}(\eta_f; \alpha_\eta, \beta_\eta) && \forall f \\ w_f|\eta_f &\sim \mathcal{N}(w_f; 0, \eta_f^{-1}) && \forall f \\ f_i|b, \mathbf{w}, \mathbf{x}_i &\sim \mathcal{N}(f_i; \mathbf{w}^\top \mathbf{x}_i + b, 1) && \forall i \\ y_i|f_i &\sim \delta(f_i y_i > \nu) && \forall i, \end{aligned}$$

where the predicted outputs of data points are modeled as a linear function of input features (i.e.  $\mathbf{w}^\top \mathbf{x}_i + b$ ). We learn the priors  $\{\gamma, \boldsymbol{\eta}\}$ , latent variables  $\mathbf{f}$  and model parameters  $\{b, \mathbf{w}\}$  using a deterministic variational approximation as we do for our method.

#### 3.2 Bayesian relevance vector machine

Its distributional assumptions are defined as

$$\begin{aligned} \gamma &\sim \mathcal{G}(\gamma; \alpha_\gamma, \beta_\gamma) \\ b|\gamma &\sim \mathcal{N}(b; 0, \gamma^{-1}) \\ \lambda_i &\sim \mathcal{G}(\lambda_i; \alpha_\lambda, \beta_\lambda) && \forall i \\ a_i|\lambda_i &\sim \mathcal{N}(a_i; 0, \lambda_i^{-1}) && \forall i \\ f_i|\mathbf{a}, b, \mathbf{k}_i &\sim \mathcal{N}(f_i; \mathbf{a}^\top \mathbf{k}_i + b, 1) && \forall i \\ y_i|f_i &\sim \delta(f_i y_i > \nu) && \forall i, \end{aligned}$$

where the predicted outputs of data points are modeled as a linear function of their kernel representations (i.e.  $\mathbf{a}^\top \mathbf{k}_i + b$ ). We again learn the priors  $\{\gamma, \boldsymbol{\lambda}\}$ , latent variables  $\mathbf{f}$  and model parameters  $\{\mathbf{a}, b\}$  using a deterministic variational approximation as we do for our method.