

Transfer Learning of Genome Wide Transcription Dynamics during Malaria Infection

Venelin Mitov
ETH Zürich

February 5, 2014
Thesis presentation

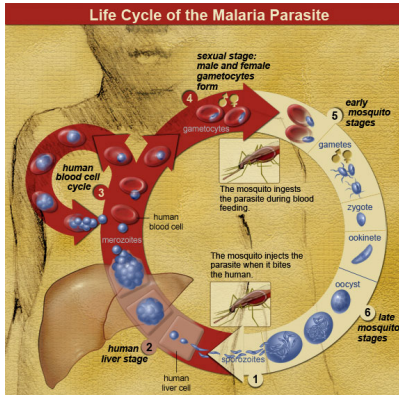
Outline

Malaria Host Transcription Dynamics

Post-Infection Time Inference in Mice

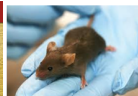
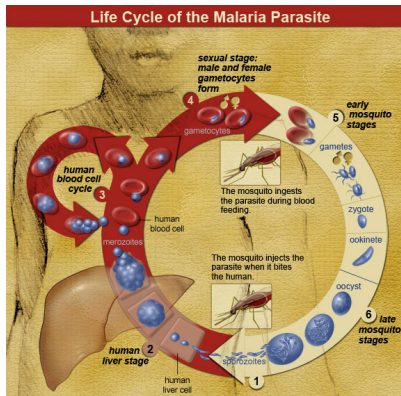
Transfer Learning To Human Data

Discussion



Courtesy: National Institute of Allergy and
Infectious Diseases

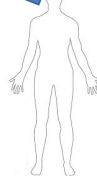
A Transfer Learning Approach



1. Find genes in infected mice, which have informative time-course dynamics for the inference of post-infection time in mice



Transfer learning



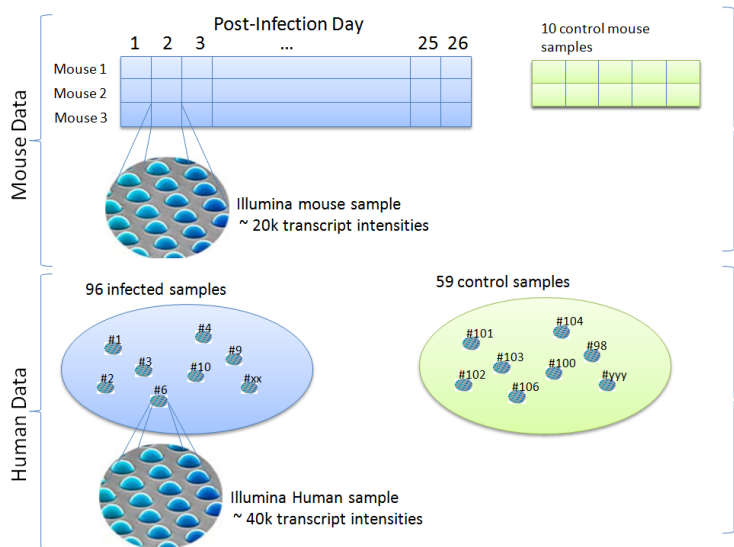
2. Map these genes to their human homologs and narrow down this set to genes that are relevant for malaria progression in the human context

M. musculus → H. sapiens

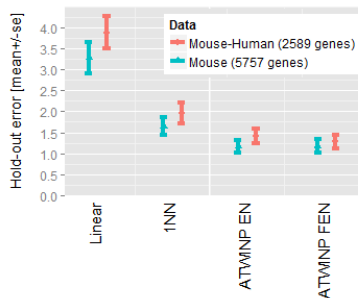
EIF6
PDCD2
ACHE
...

Courtesy: National Institute of Allergy and

Infectious Diseases

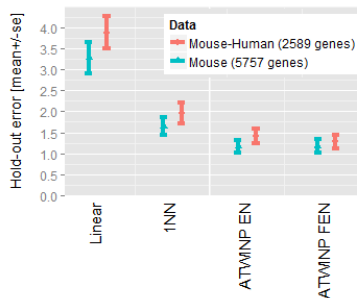
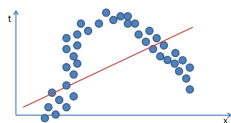


Results from Leave-One-Mouse-Out Cross Validation



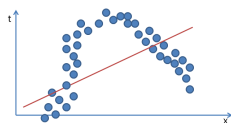
Results from Leave-One-Mouse-Out Cross Validation

Linear regression

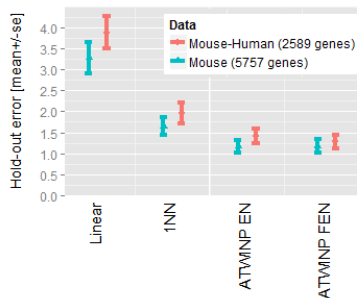
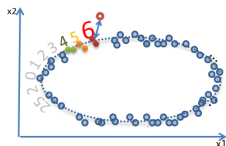


Results from Leave-One-Mouse-Out Cross Validation

Linear regression

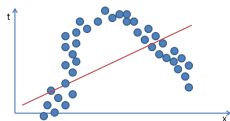


Classification (1st NN)

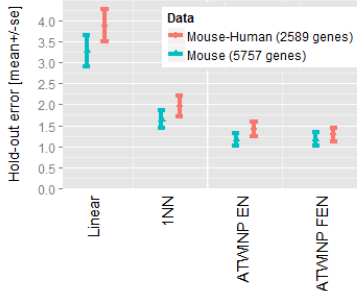
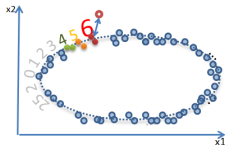


Results from Leave-One-Mouse-Out Cross Validation

Linear regression



Classification (1st NN)

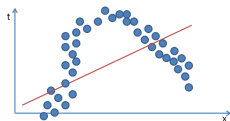


Aggregated Time Windows

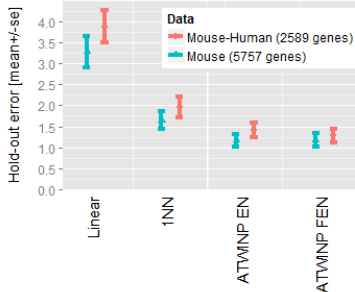
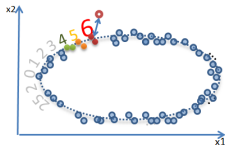


Results from Leave-One-Mouse-Out Cross Validation

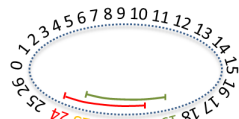
Linear regression



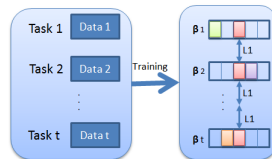
Classification (1st NN)



Aggregated Time Windows



+ fused logistic regression





One-Against-All Linear Logistic Regression



Model the logit function, $\text{logit}(\pi) := \log(\pi/(1 - \pi))$, as a linear function of \mathbf{x} :

$$\text{logit}(\pi^{(j)}(\mathbf{x})) \approx \mathbf{x}^T \boldsymbol{\beta}^{(j)}.$$

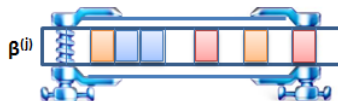
The negative log-likelihood is defined as:

$$-\ell^{(j)}(\boldsymbol{\beta}^{(j)}; [X|\mathbf{y}_j]) = \sum \log \left(\mathbf{1} + \exp(-\mathbf{y}_j \odot X\boldsymbol{\beta}^{(j)}) \right), \quad j = 1, \dots, t.$$

Maximum likelihood fit for $\boldsymbol{\beta}^{(j)}$:

$$\boldsymbol{\beta}^{(j)*} := \arg \min_{\boldsymbol{\beta}^{(j)} \in \mathbb{R}^{(1+d)}} \left\{ -\ell^{(j)}(\boldsymbol{\beta}^{(j)}; [X|\mathbf{y}_j]) \right\}$$

Regularization and Automatic Variable Selection



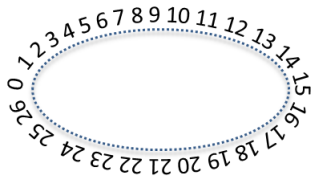
- ▶ L2-penalty (Ridge): $\frac{1}{2} \lambda_2 ||\beta^{(j)}||_2^2 = \frac{1}{2} \lambda_2 \sum_{k=1}^d \beta_k^{(j)2}$
- ▶ L1-penalty (Lasso): $\lambda_1 ||\beta^{(j)}||_1 = \lambda_1 \sum_{k=1}^d |\beta_k|$
- ▶ Elastic Net penalty (Lasso+Ridge): $\lambda_1 ||\beta^{(j)}||_1 + \frac{1}{2} \lambda_2 ||\beta^{(j)}||_2^2$

Maximum A-Posteriori fit for $\beta^{(j)}$:

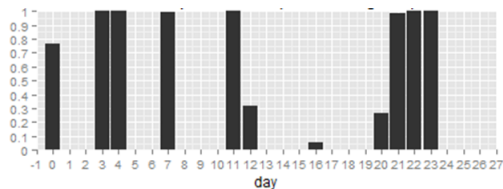
$$\beta^{(j)*} := \arg \min_{\beta^{(j)} \in \mathbb{R}^{(1+d)}} \left\{ -\ell^{(j)}(\beta^{(j)}; [X|\mathbf{y}_j]) + \lambda_1 ||\beta||_1 + \frac{1}{2} \lambda_2 ||\beta^{(j)}||_2^2 \right\}$$

Single Day versus Time Window Prediction

Single day

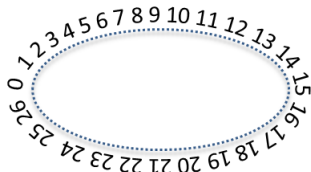


Predicted probabilities

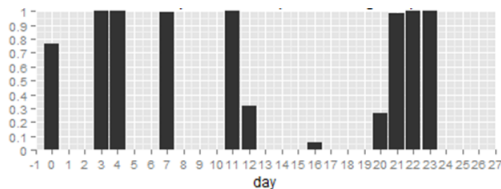


Single Day versus Time Window Prediction

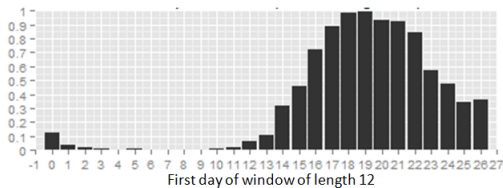
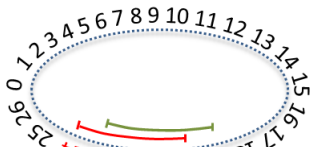
Single day



Predicted probabilities

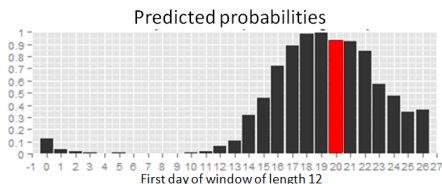
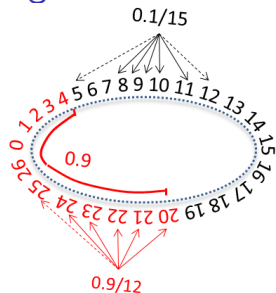


Time Window



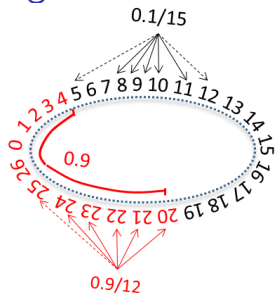


Aggregated Time Window Predictor (ATWINP)

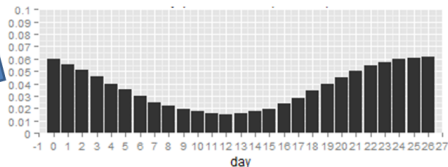
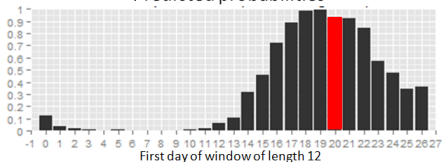




Aggregated Time Window Predictor (ATWINP)

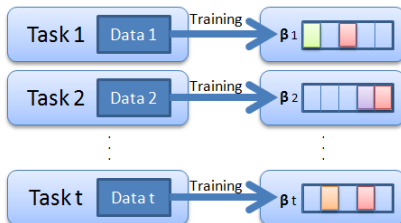


Predicted probabilities



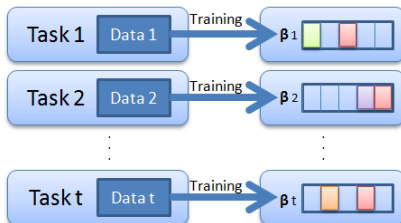
The Idea of Multi-Task Learning

Single Task Learning

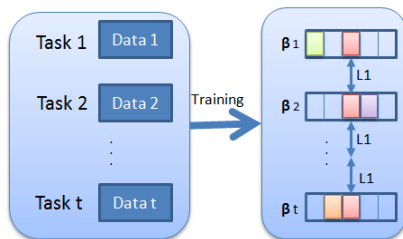


The Idea of Multi-Task Learning

Single Task Learning



Multi-Task Learning





Fused Elastic Net Logistic Regression (FLR)

Let $B := [\beta^{(1)}, \dots, \beta^{(t)}] \in \mathbb{R}^{(1+d) \times t}$ be the coefficient matrix for all tasks and let $R \in \mathbb{R}^{t \times t}$ be a matrix defined in the following way:

$$R_{ij} := \begin{cases} 1 & \text{if } j = i - 1 \text{ or } (i, j) = (1, t) \\ 0 & \text{otherwise} \end{cases}, \quad i, j = 1, \dots, t.$$

The multi-task fused elastic net negative log-likelihood is defined as:

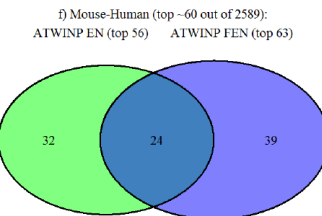
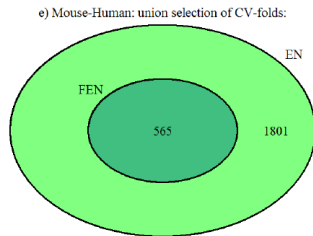
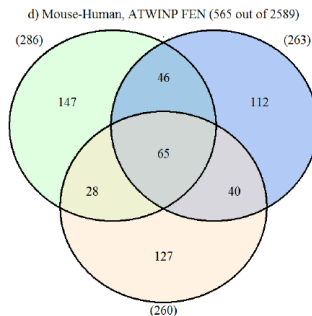
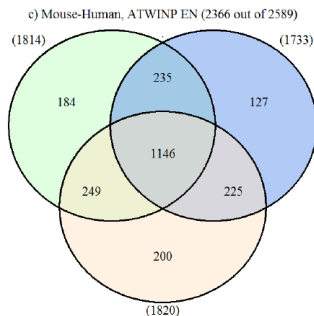
$$\begin{aligned} -\ell^{MT}(B; [X|Y]) &:= \sum \log([1] + \exp(-Y \odot XB)) \\ &\quad + ||[\lambda_1] \odot B||_1 + \frac{1}{2} ||[\lambda_2] \odot B||_2^2 \\ &\quad + ||[\nu] \odot B(I - R)||_1 \end{aligned}$$

The Fused Elastic Net Logistic Regression (FENLR) fit for B is obtained by solving

$$B^* = \arg \min_{B \in \mathbb{R}^{(1+d) \times t}} -\ell^{MT}(B; [X|Y]).$$

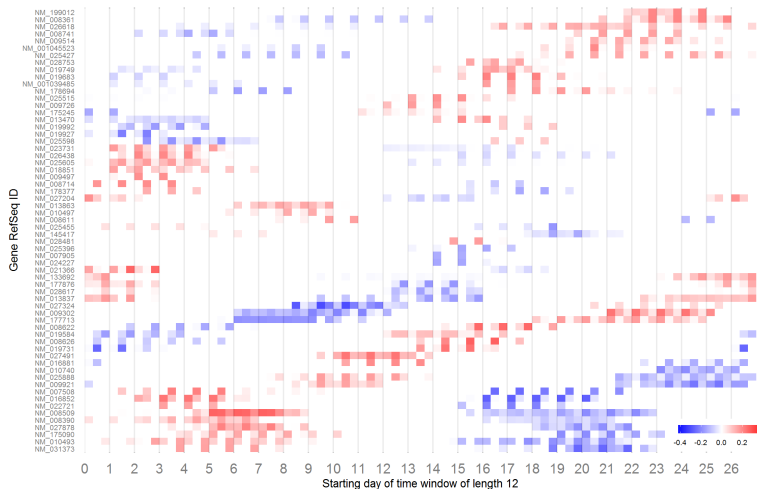


Selected Genes

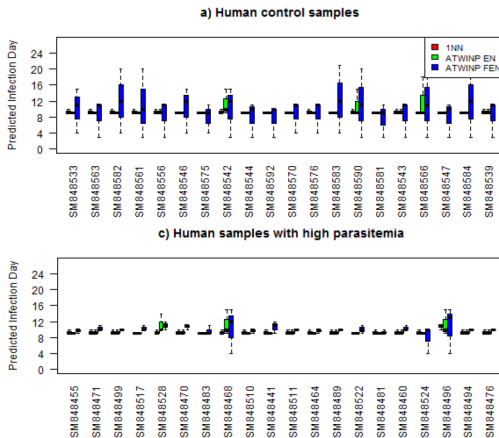




Top 60 Genes, ATWINP FEN



Post-Infection Time Prediction in Human Patients



Discussion

- ▶ Our model can predict the post-infection time of an unlabeled infected mouse-sample with expected deviation of 1.28 days from the true post-infection time.
- ▶ The gene-expression profile of an infected host-organism preserves information with respect to the beginning of the infection, and can be used to characterize the disease progression on a fine time-scale.
- ▶ We were able to identify a set of genes that are informative for the disease progression in mice and we could quantify the effect of each selected gene at all points in the time-course of the infection.
- ▶ At the current time knowledge transfer from mouse to human patients cannot provide a valuable estimation of the post-infection time in humans.

Acknowledgements

- ▶ prof. Manfred Claassen - Advisor of the master thesis project
- ▶ Stefan, Eirini, Anita, Ana - colleagues in the Claassen's Group
- ▶ David and Brenda - Stanford Microbiology and Immunology Lab