

# Fast Discovery of Pairwise Interactions in High Dimensions using Gaussian Processes

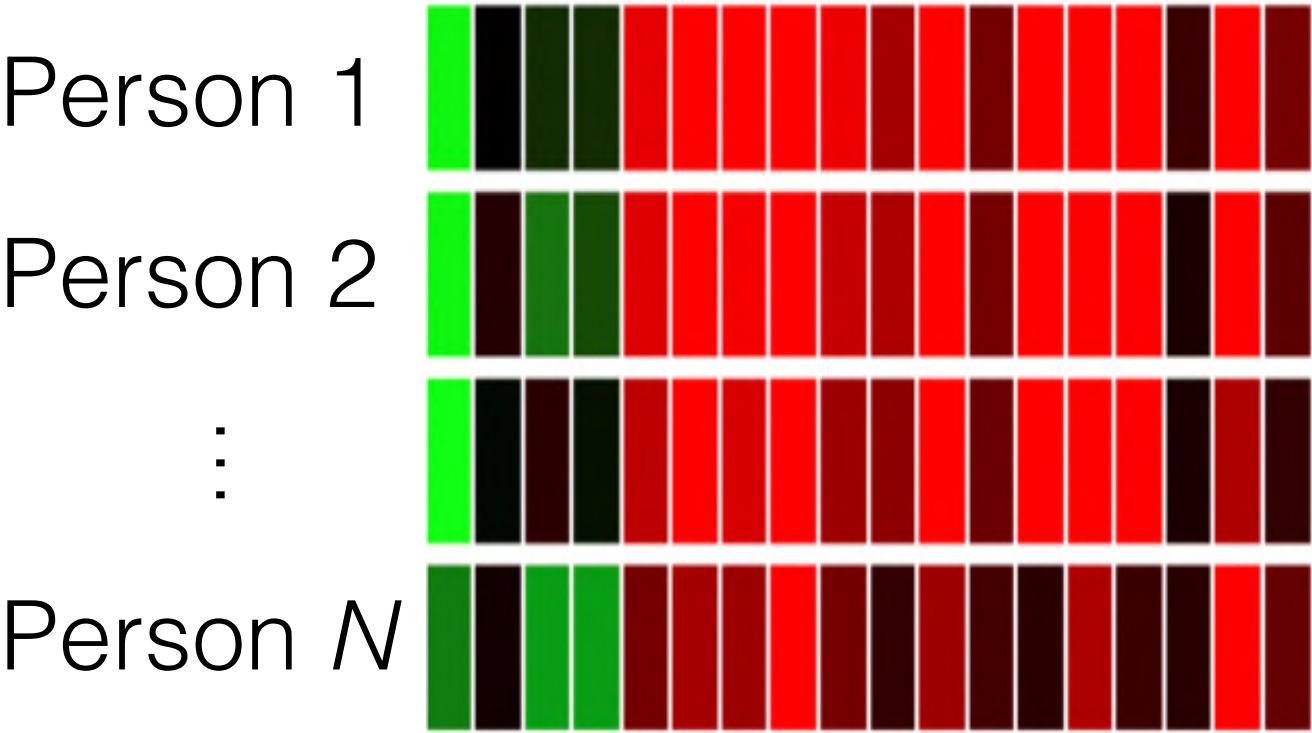
Tamara Broderick

Associate Professor  
EECS, MIT

Raj Agrawal, Jonathan H. Huggins, Brian L. Trippe



Gene expression levels



# Environmental factors

## Gene expression levels

Person 1



Person 2



⋮

Person  $N$



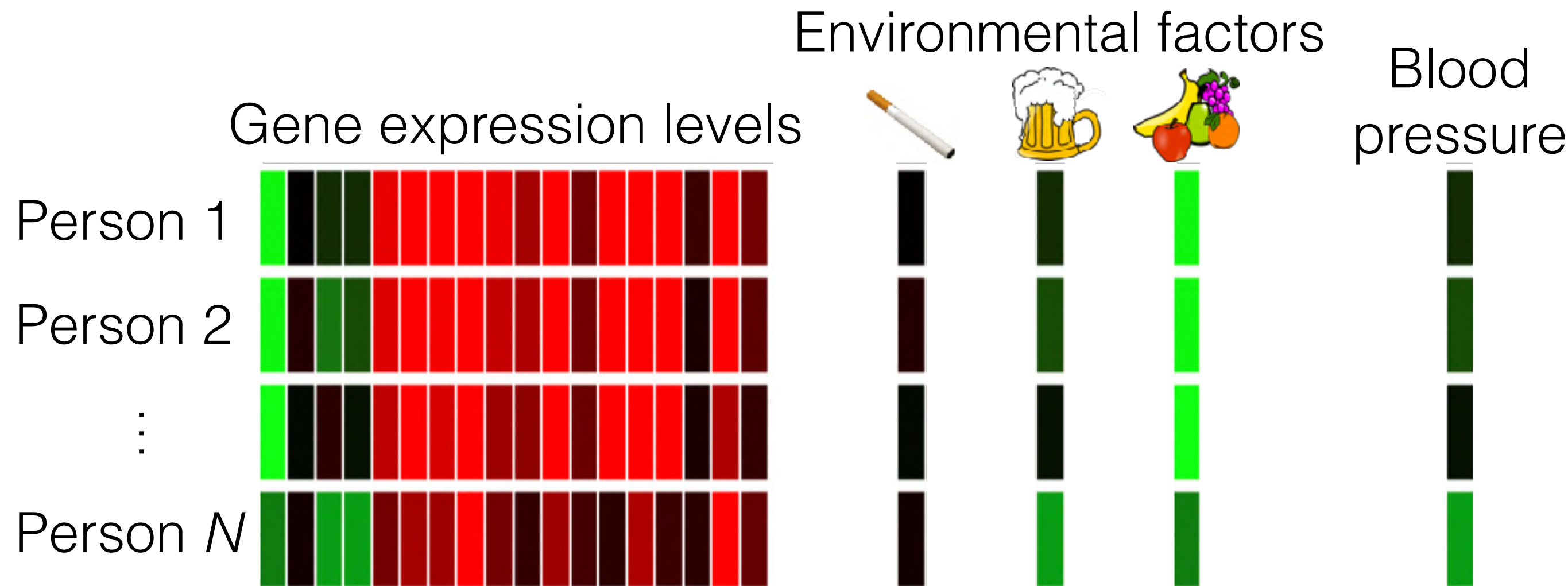


The figure displays a 4x20 grid of colored rectangles. The colors are categorized into three main groups: bright green (low frequency), dark green (medium frequency), and dark red (high frequency). The grid shows a clear pattern of high-frequency pairs (red) and low-frequency pairs (green) across the 20 columns and 4 rows.

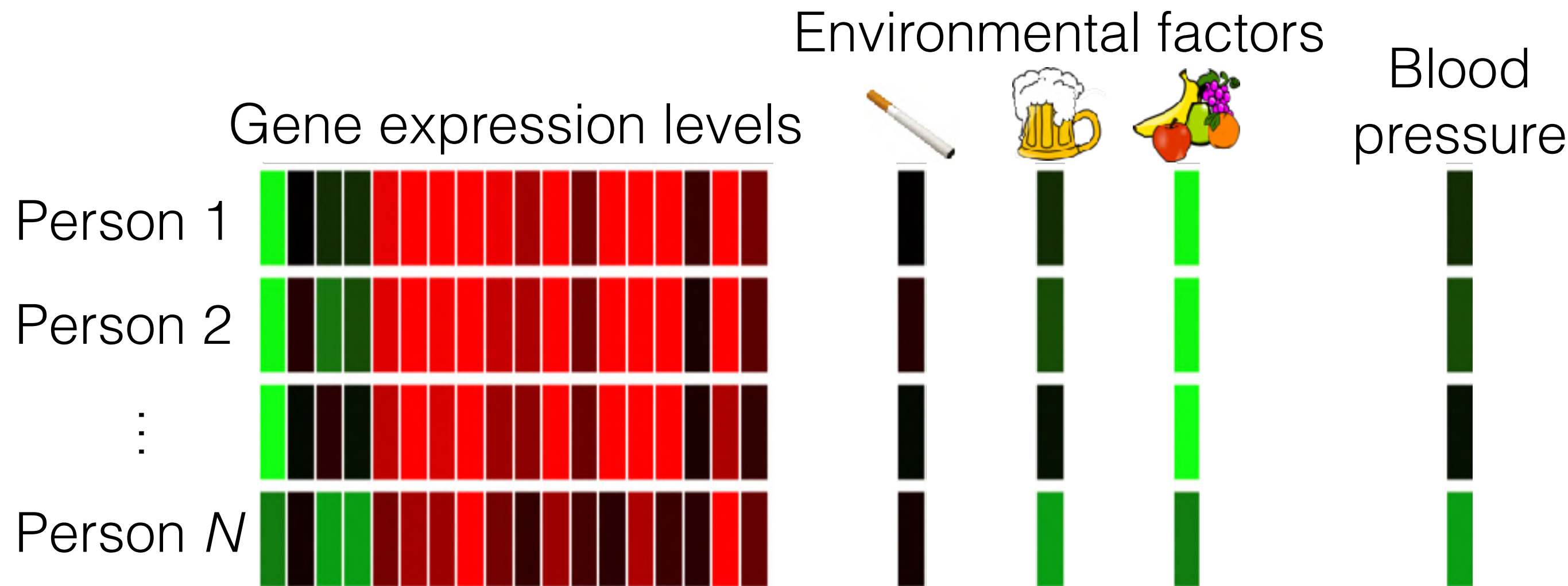
## Person 2

•  
•  
•

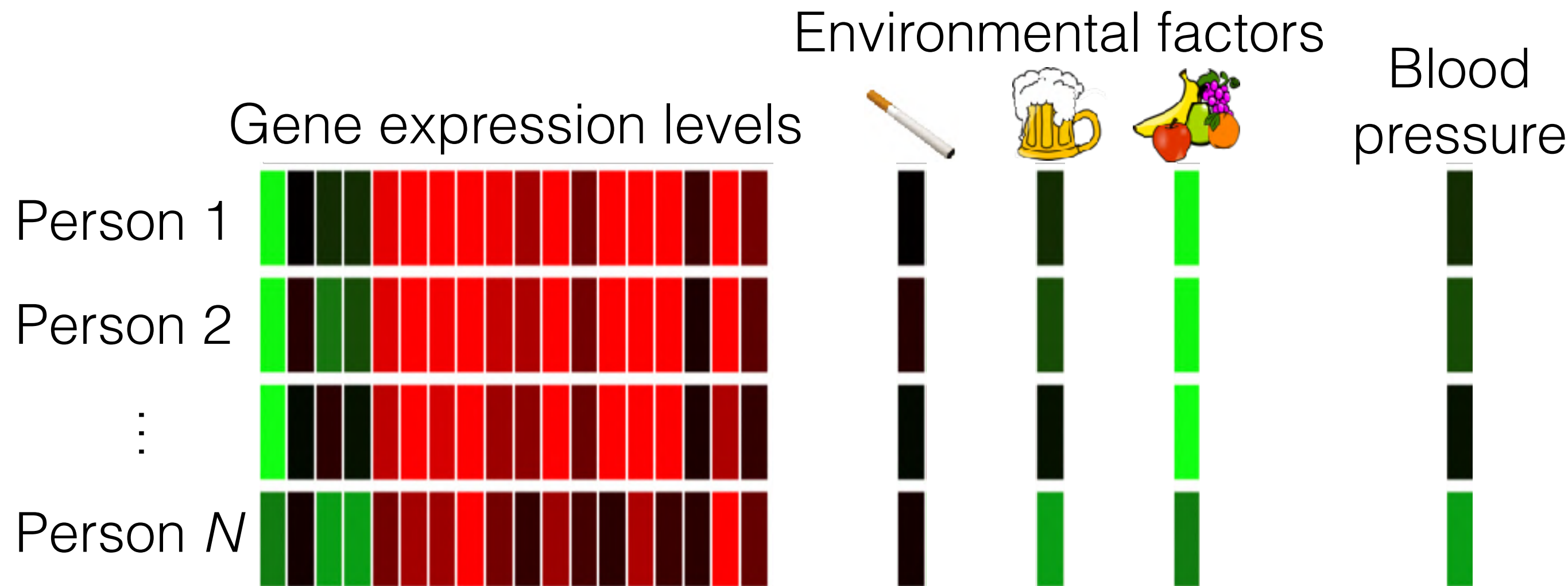
Person  $N$



- Which genes/factors are associated with a health issue?

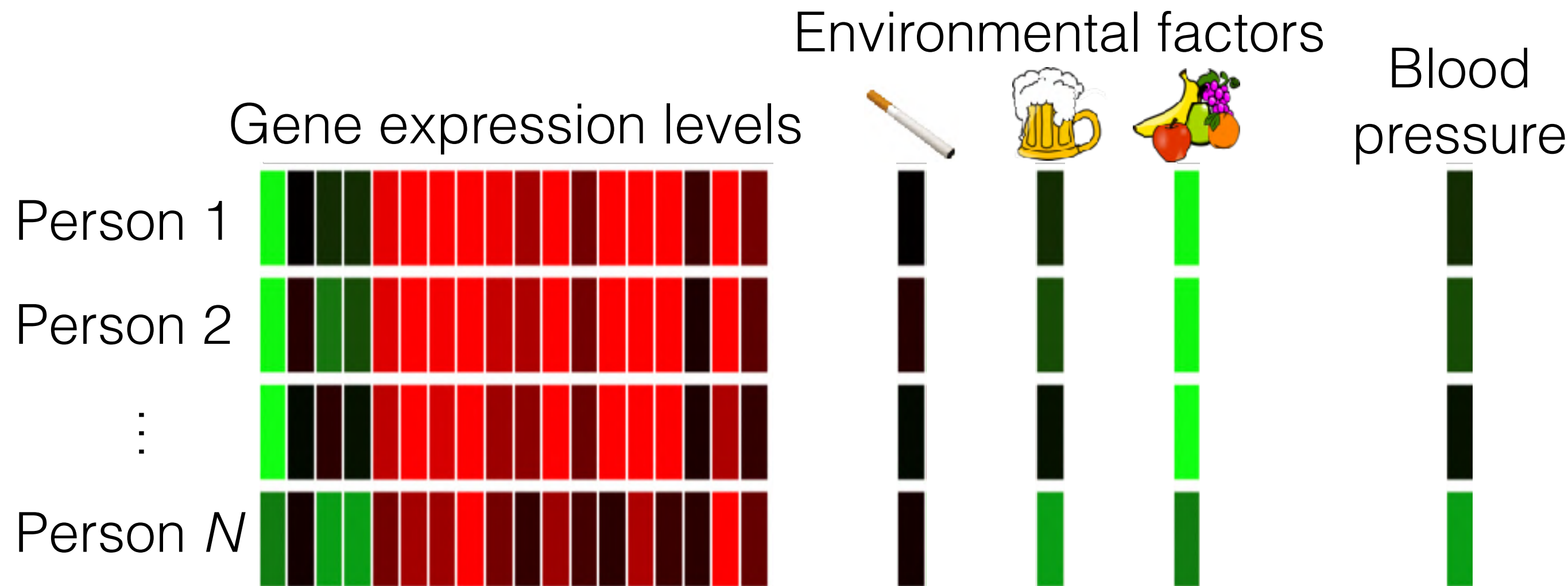


- Which genes/factors are associated with a health issue?
- Want small subset of  $p (> N)$  covariates



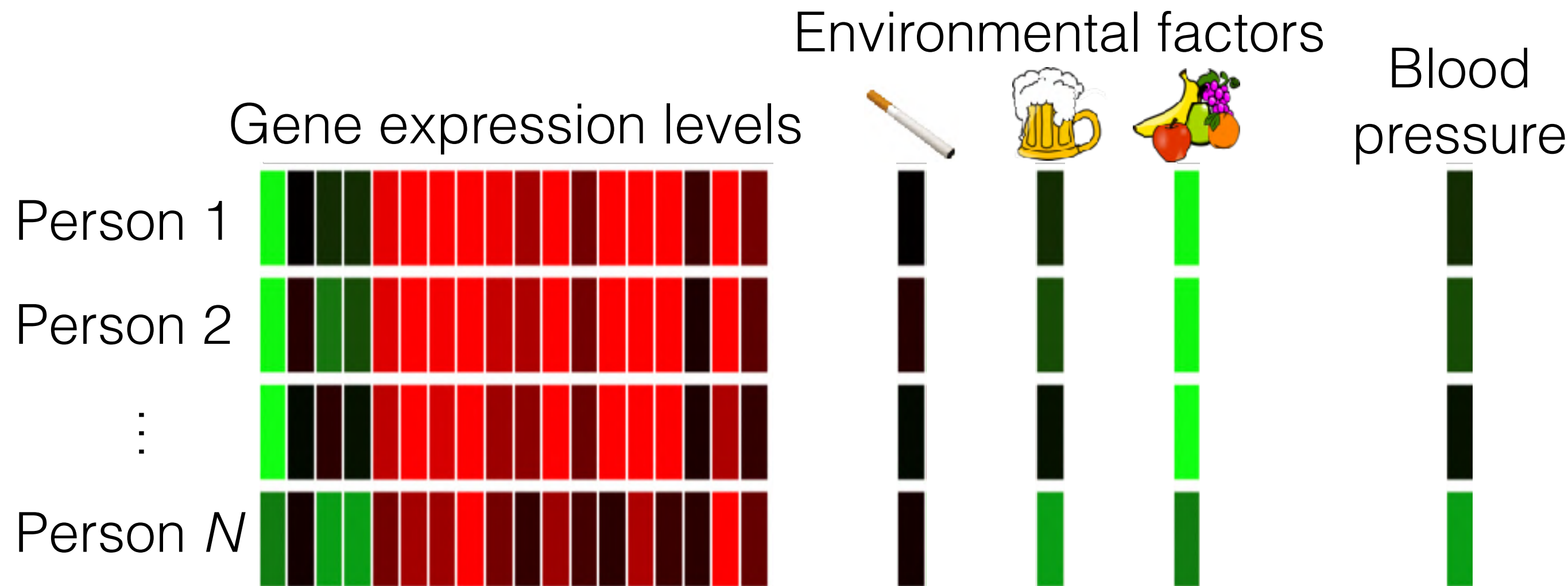
- Which genes/factors are associated with a health issue?
- Want small subset of  $p (> N)$  covariates (cf. LASSO)





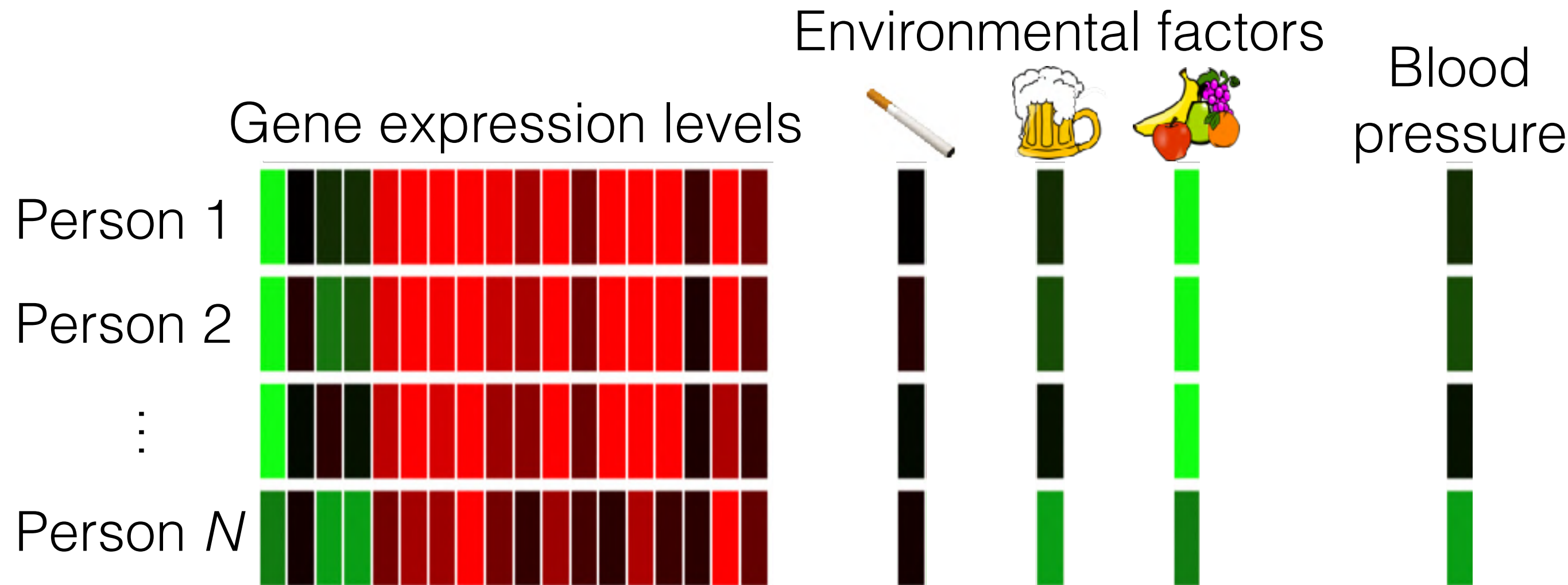
- Which genes/factors are associated with a health issue?
- Want small subset of  $p (> N)$  covariates (cf. LASSO)
- Additive model often not enough: need interactions





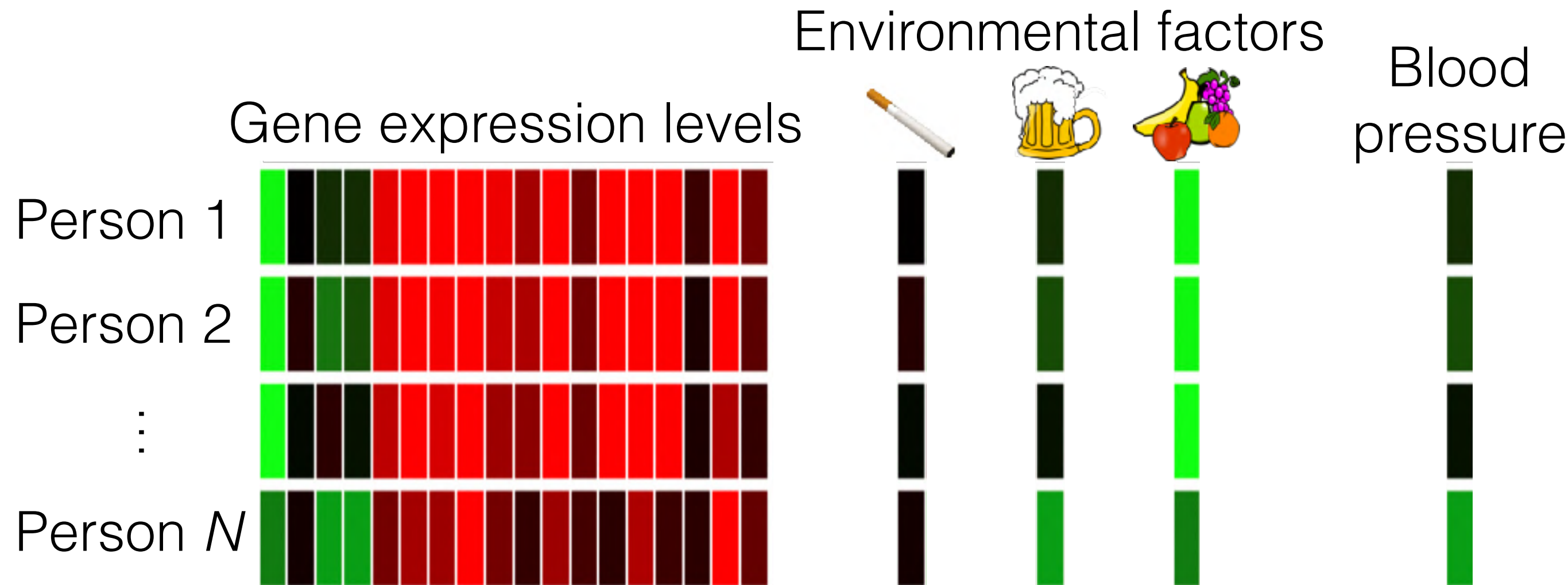
- Which genes/factors are associated with a health issue?
- Want small subset of  $p$  ( $> N$ ) covariates (cf. LASSO)
- Additive model often not enough: need interactions (now  $p^2$  dims!)

# Pairwise interactions in high dimensions



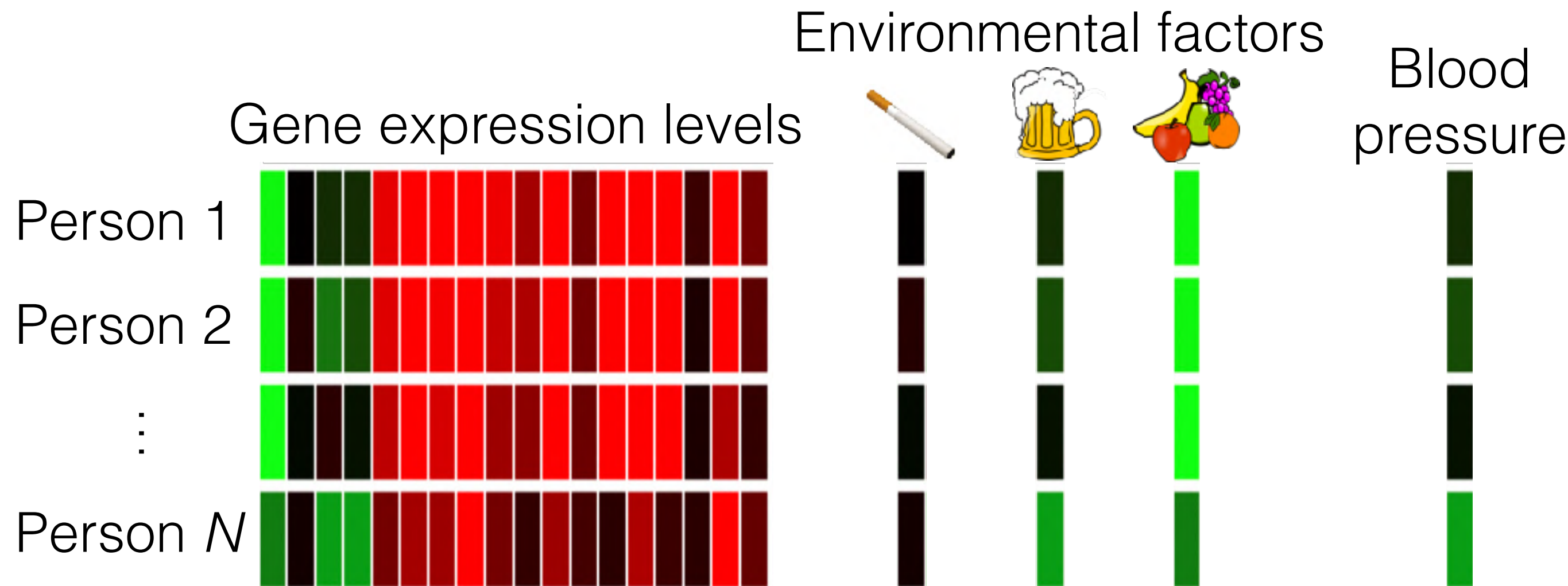
- Which genes/factors are associated with a health issue?
- Want small subset of  $p$  ( $> N$ ) covariates (cf. LASSO)
- Additive model often not enough: need interactions (now  $p^2$  dims!)

# Pairwise interactions in high dimensions



- Which genes/factors are associated with a health issue?
- Want small subset of  $p$  ( $> N$ ) covariates (cf. LASSO)
- Additive model often not enough: need interactions (now  $p^2$  dims!)
- **We provide:** Fast, accurate (Bayes) method for interaction discovery

# Pairwise interactions in high dimensions



- Which genes/factors are associated with a health issue?
- Want small subset of  $p$  ( $> N$ ) covariates (cf. LASSO)
- Additive model often not enough: need interactions (now  $p^2$  dims!)
- **We provide:** Fast, accurate (Bayes) method for interaction discovery
  - Better scaling in  $p$  & better accuracy than LASSO-based methods. Orders of magnitude faster than naive Bayesian inference

# Roadmap

# Roadmap

- Setup: Discovering main and interaction effects

# Roadmap

- Setup: Discovering main and interaction effects
- Our method



# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model

# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)

# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)
  - Fast reporting of results

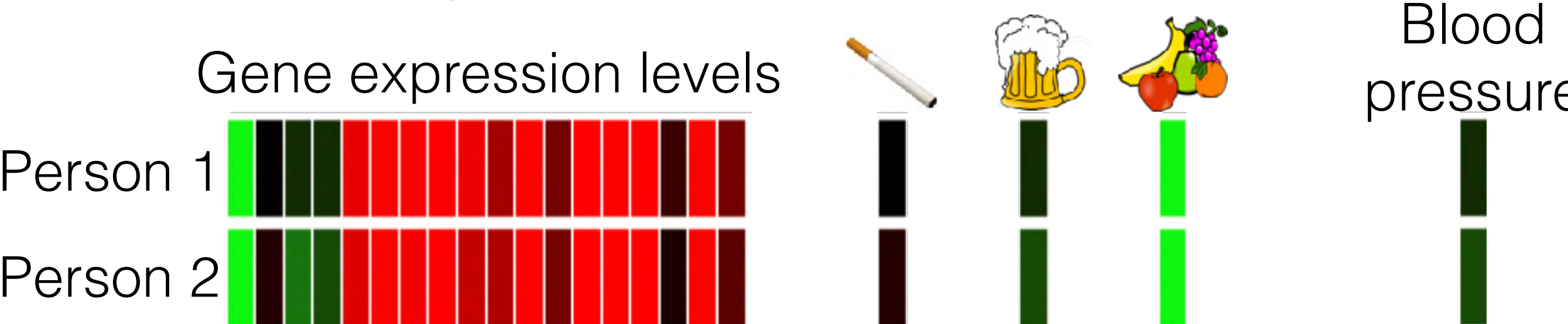
# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)
  - Fast reporting of results
- Experiments on simulated and real data

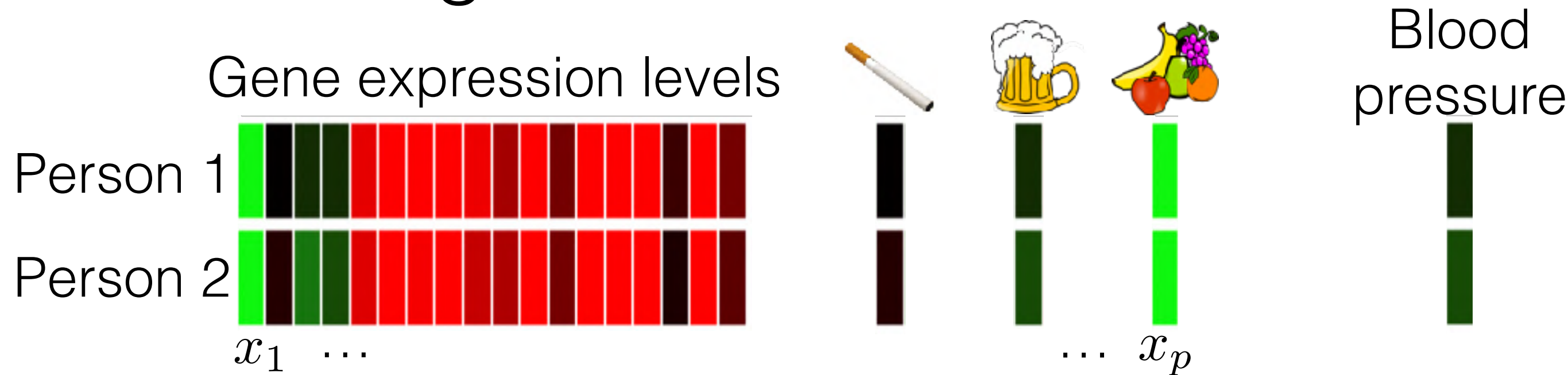
# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)
  - Fast reporting of results
- Experiments on simulated and real data

# Discovering main and interaction effects

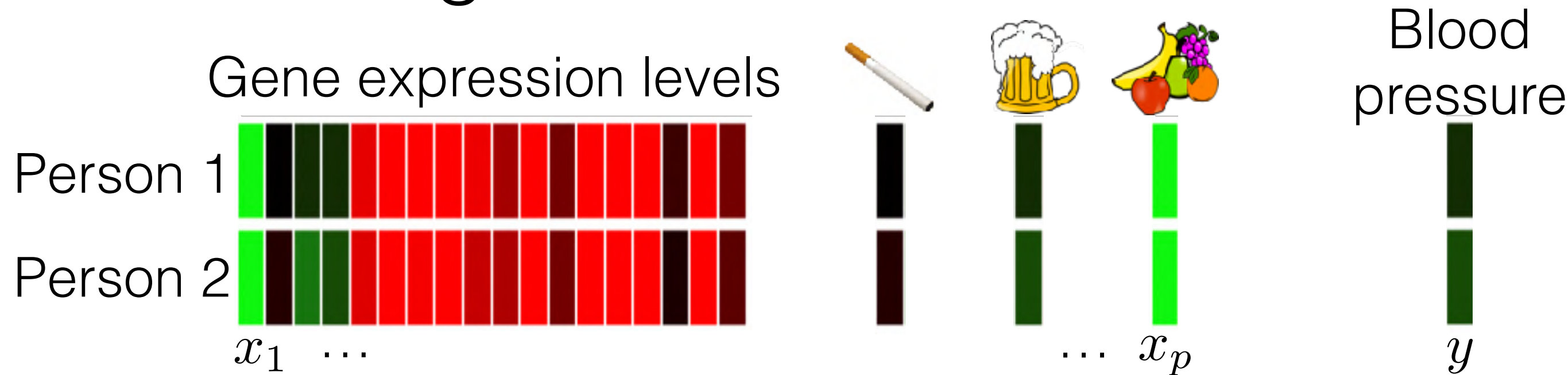


# Discovering main and interaction effects

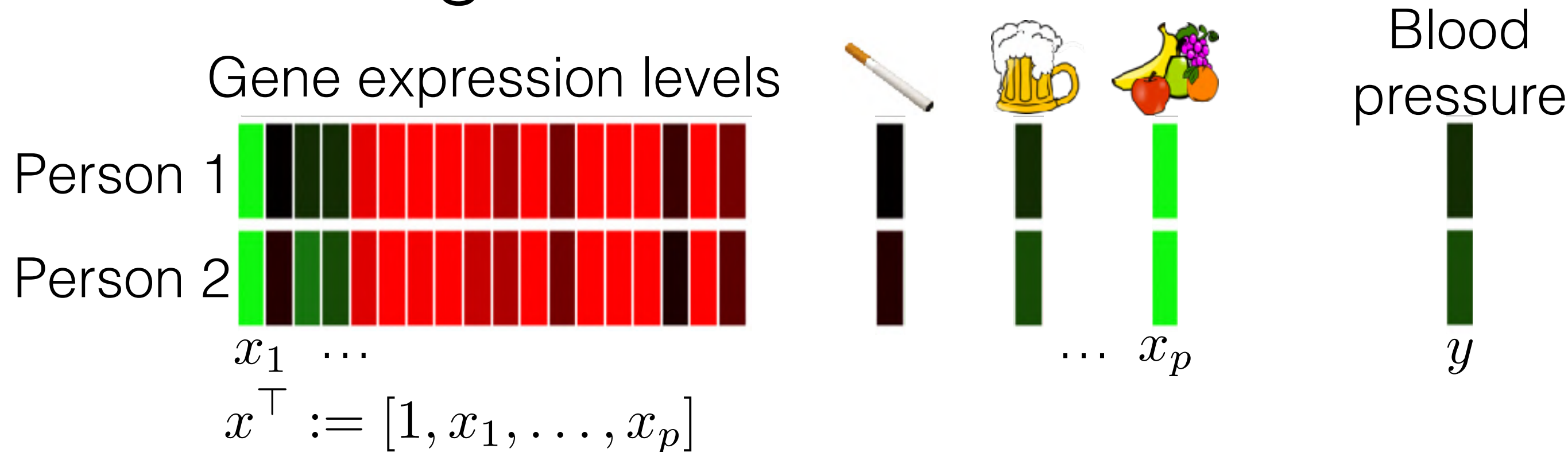




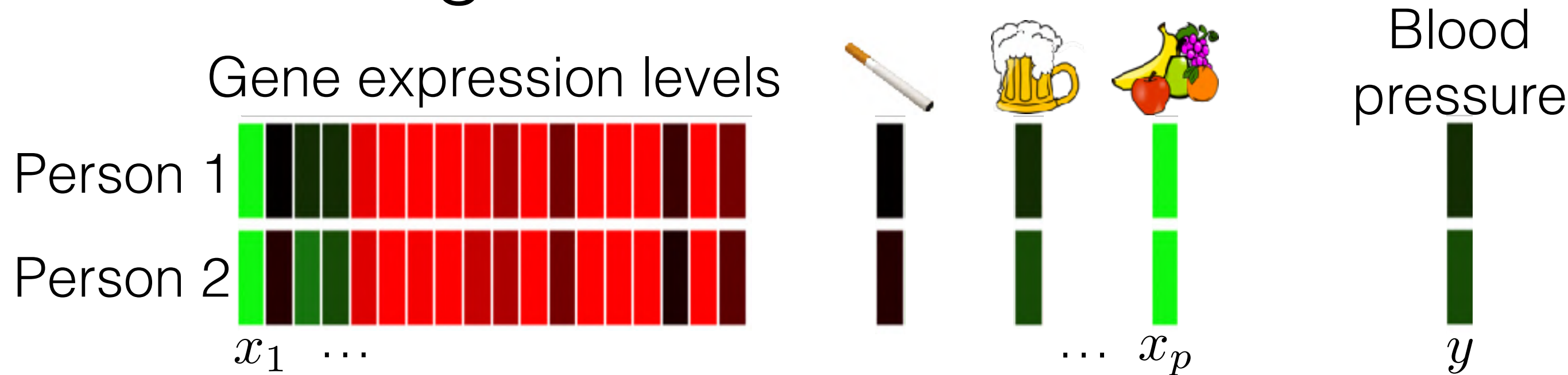
# Discovering main and interaction effects



# Discovering main and interaction effects



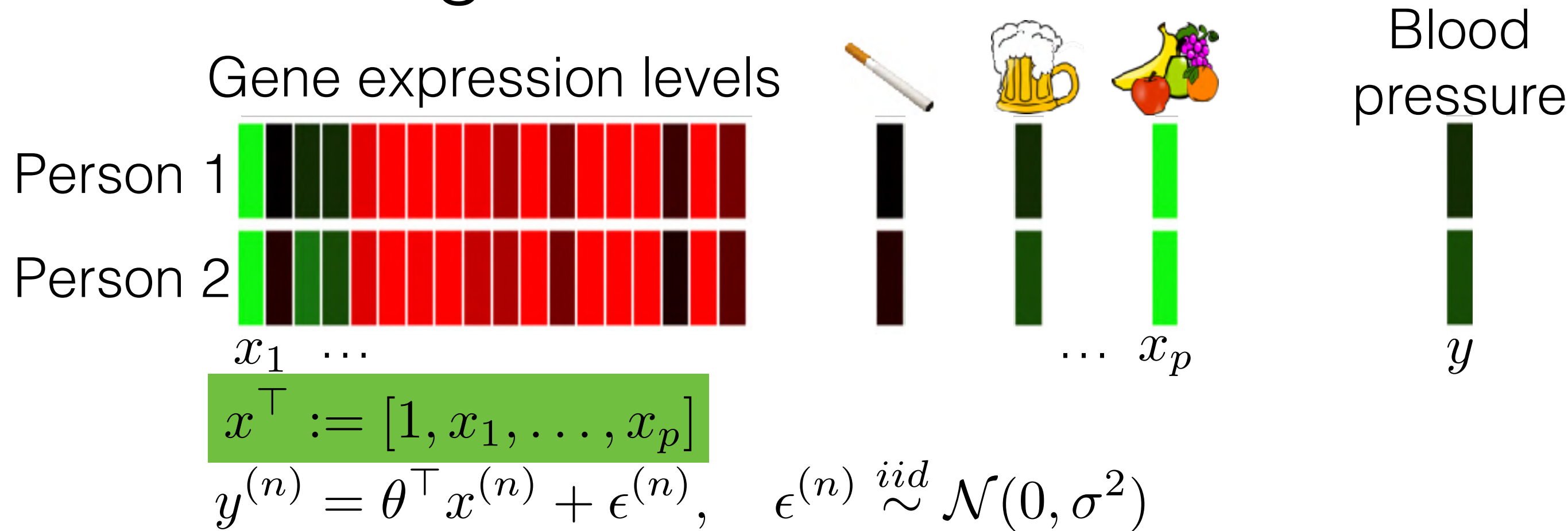
# Discovering main and interaction effects



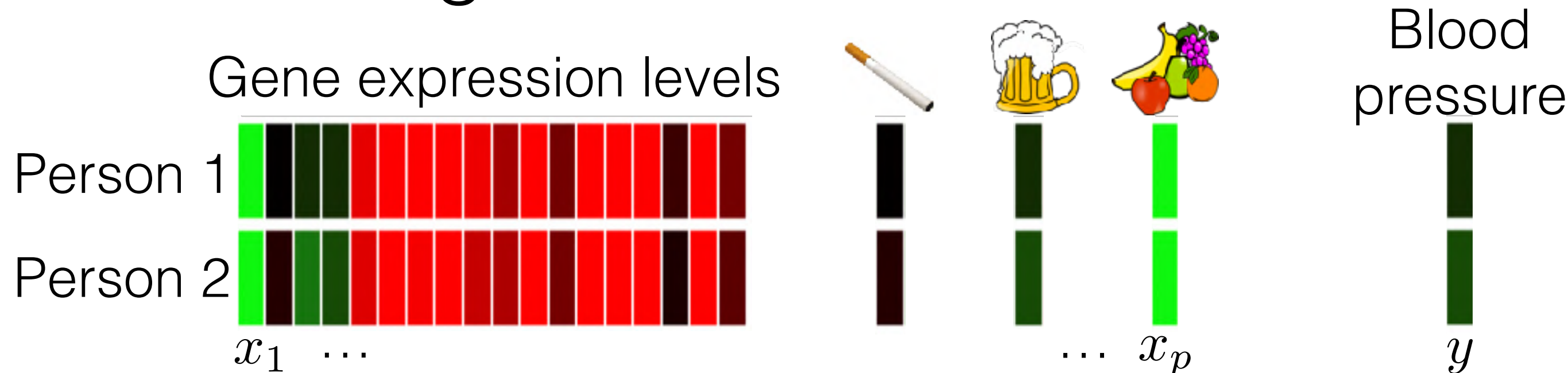
$$x^\top := [1, x_1, \dots, x_p]$$

$$y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

# Discovering main and interaction effects



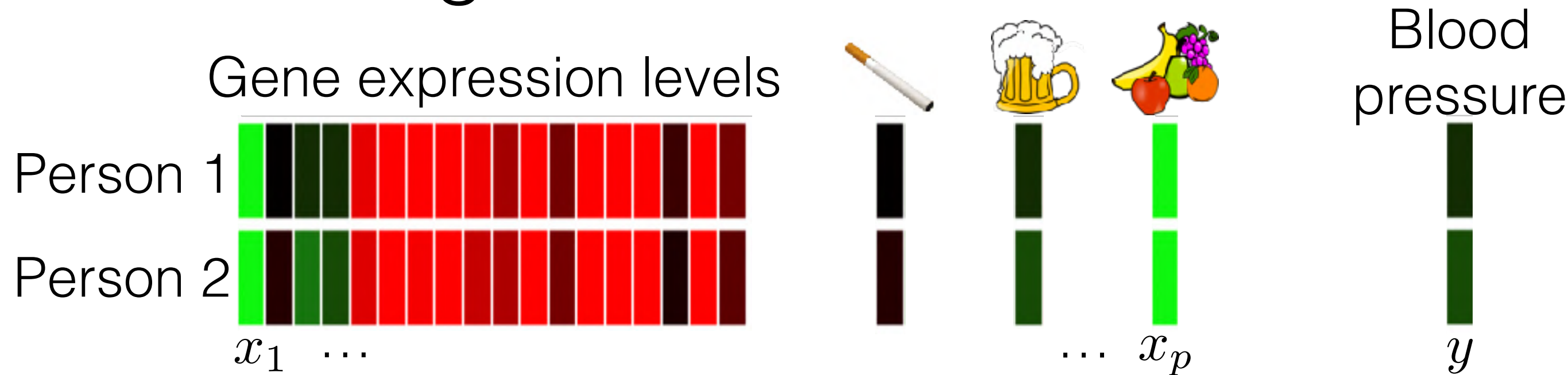
# Discovering main and interaction effects



$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

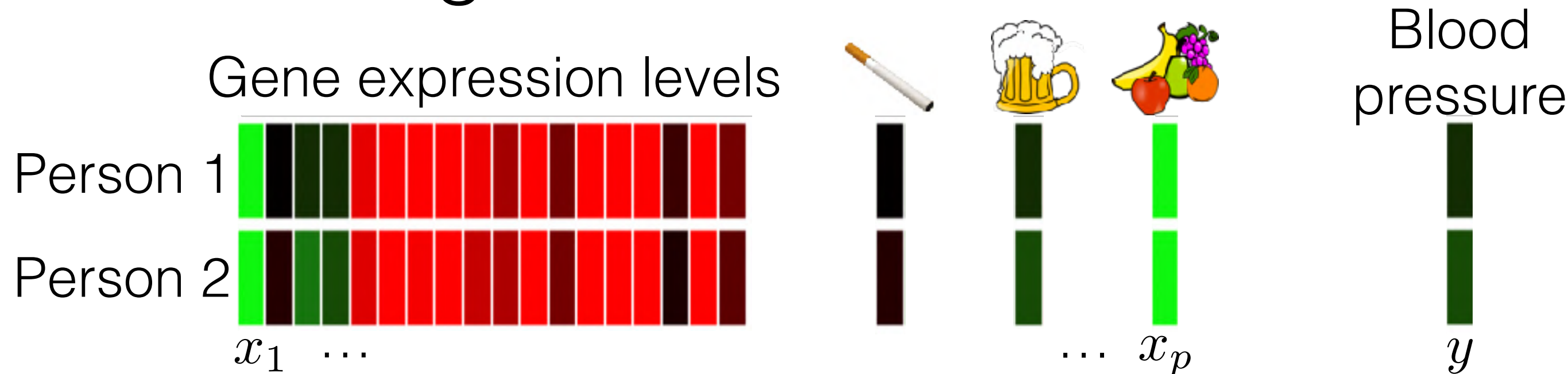
# Discovering main and interaction effects



$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

# Discovering main and interaction effects

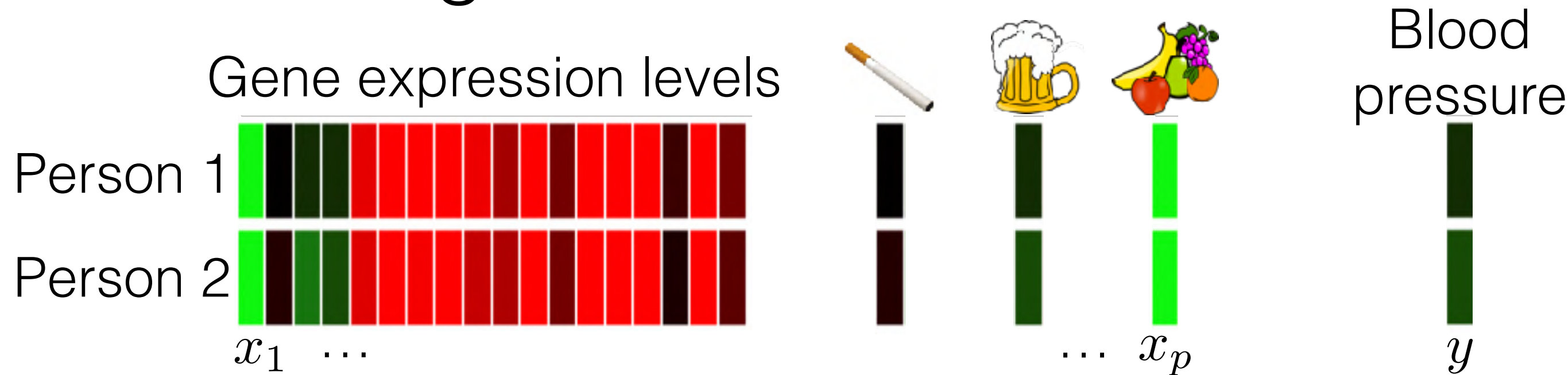


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$



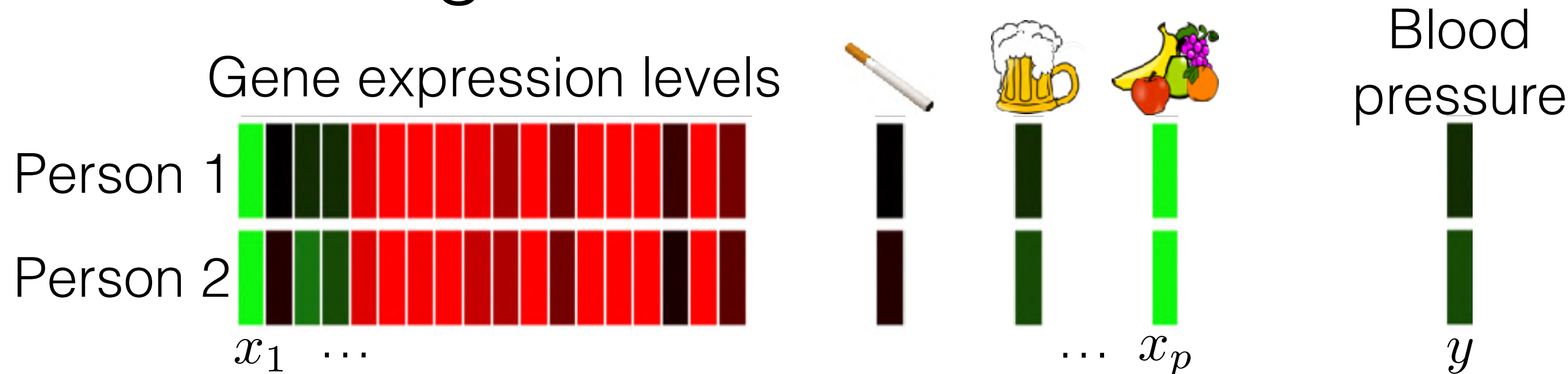
# Discovering main and interaction effects



$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

# Discovering main and interaction effects

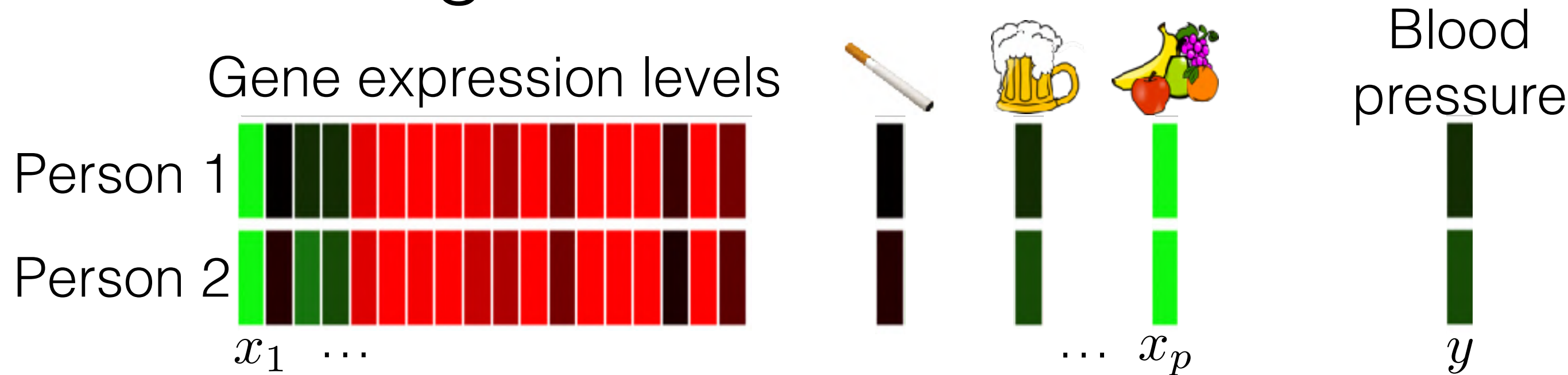


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation

# Discovering main and interaction effects

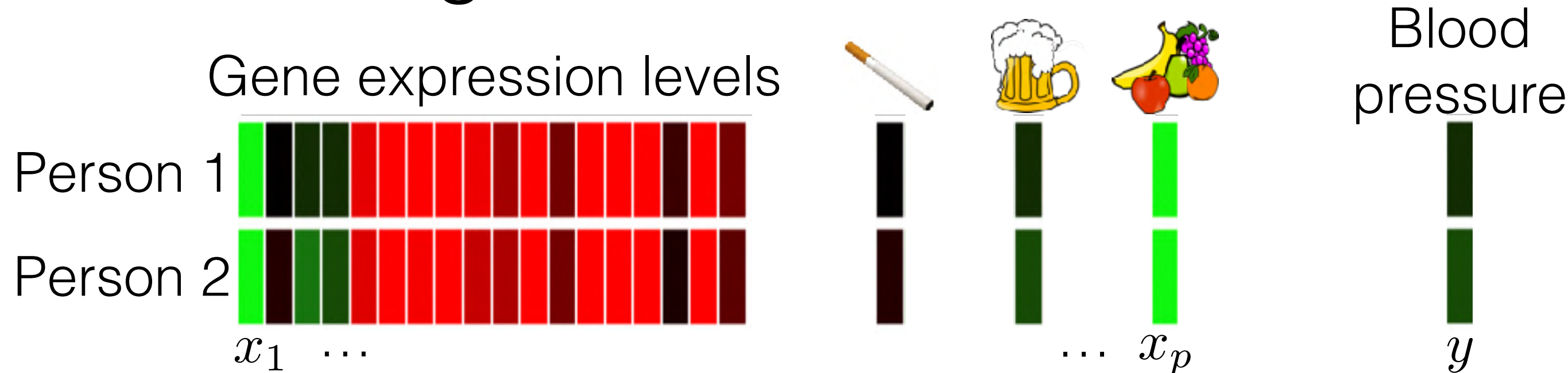


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:

# Discovering main and interaction effects

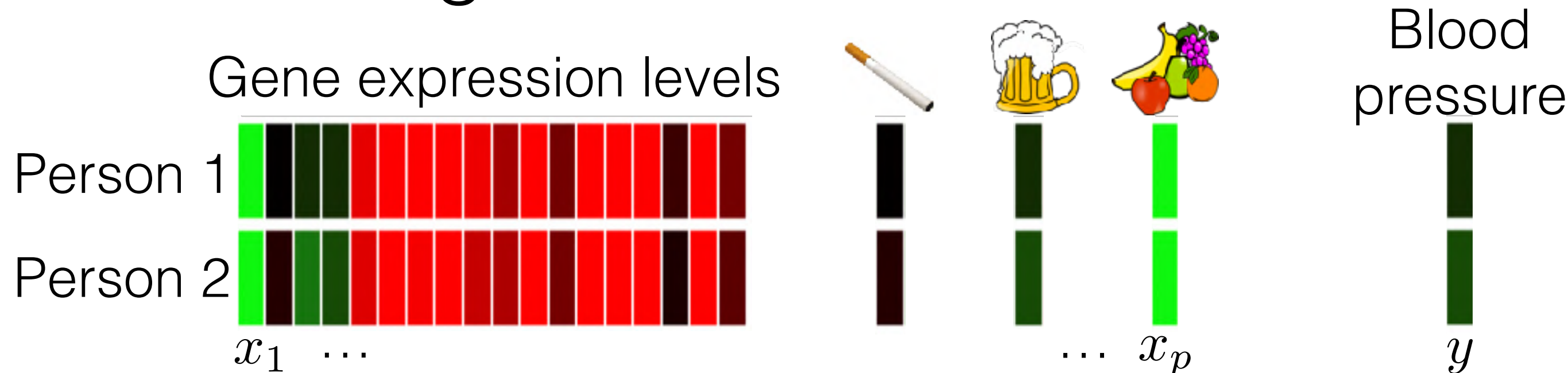


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1 x_2, \dots, x_{p-1} x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:
  - *Sparsity*: most main effects are negligible (interpretable)

# Discovering main and interaction effects

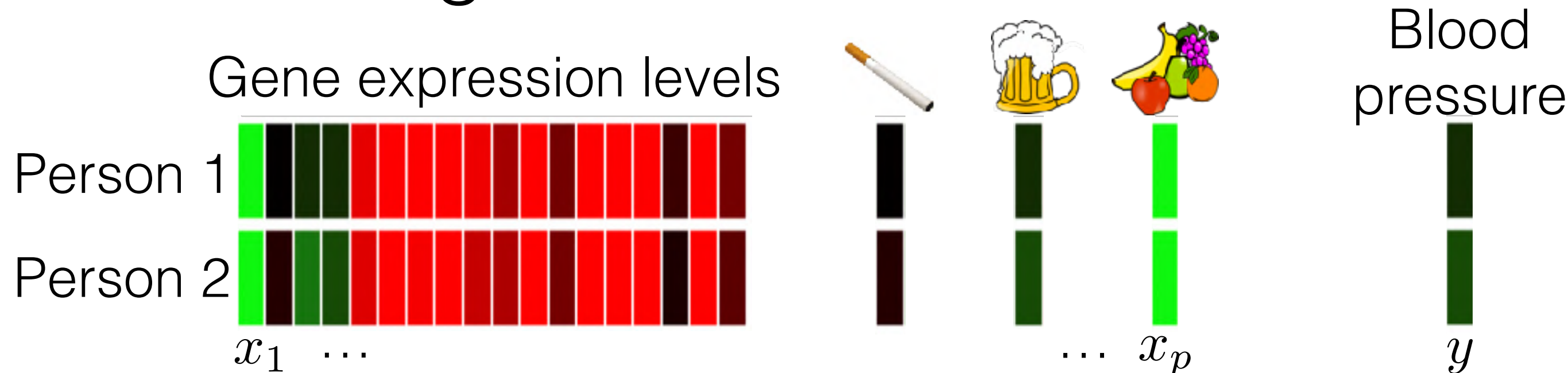


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:
  - *Sparsity:* most main effects are negligible (interpretable)
  - *Strong hierarchy:* Interaction only if main effects are present

# Discovering main and interaction effects

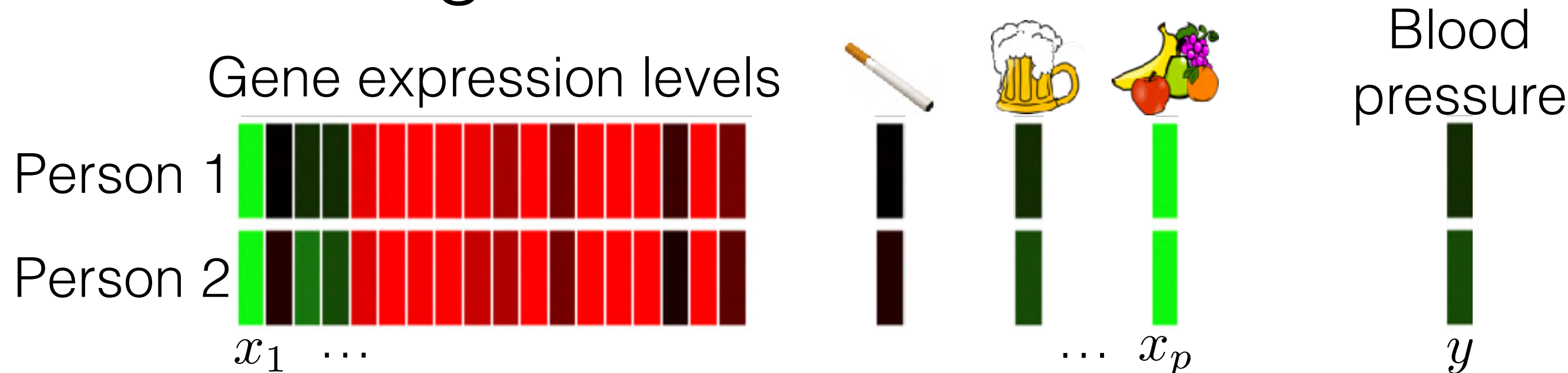


$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:
  - *Sparsity:* most main effects are negligible (interpretable)
  - *Strong hierarchy:* Interaction only if main effects are present  
[We are able to lose this assumption in new work; see discussion]

# Discovering main and interaction effects



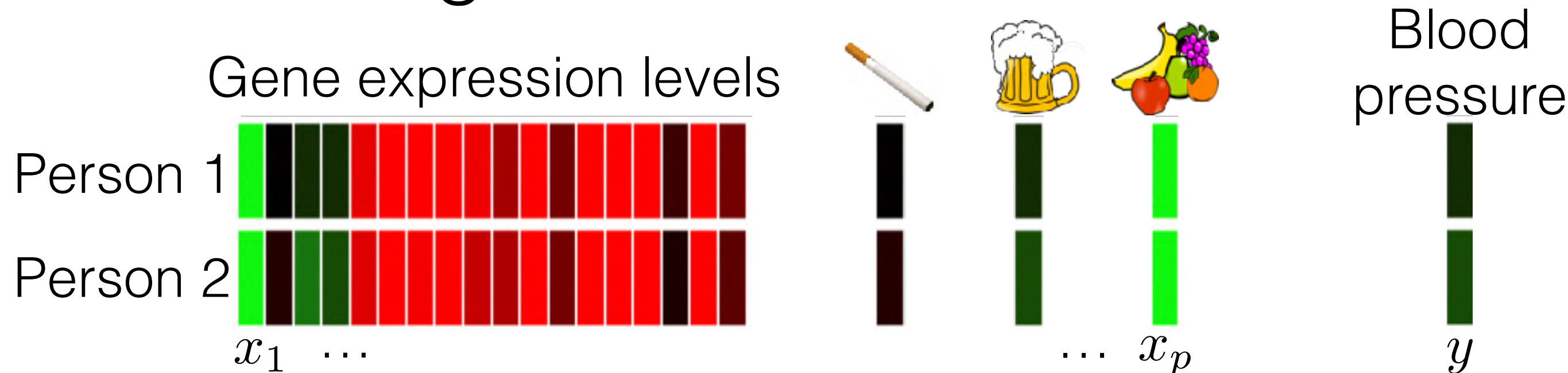
$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:
  - *Sparsity:* most main effects are negligible (interpretable)
  - *Strong hierarchy:* Interaction only if main effects are present  
[We are able to lose this assumption in new work; see discussion]
- $p^2$  covariates: large  $p \rightarrow$  statistical & computational challenge



# Discovering main and interaction effects



$$\Phi_2^\top(x) := [1, x_1, \dots, x_p, x_1x_2, \dots, x_{p-1}x_p]$$

$$y^{(n)} = \theta^\top \Phi_2(x^{(n)}) + \epsilon^{(n)}, \quad \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- **Goal:** Parameter selection/estimation under assumptions:
  - *Sparsity:* most main effects are negligible (interpretable)
  - *Strong hierarchy:* Interaction only if main effects are present  
[We are able to lose this assumption in new work; see discussion]
- $p^2$  covariates: large  $p \rightarrow$  statistical & computational challenge
- **Our solution:** using structure in covariates + sparsity assumptions to reduce to a problem *linear* in  $p$

# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference
  - Fast reporting of results
- Experiments on simulated and real data

# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference
  - Fast reporting of results
- Experiments on simulated and real data

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. Choose generative model

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. Choose generative model
2. Compute posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. Choose generative model
2. Compute posterior
3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. Choose generative model

2. Compute posterior

3. Report relevant summaries of the posterior



# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy
2. Compute posterior
3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. Compute posterior
3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**) to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]**

2. Compute posterior

3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. **Kernel Interaction Sampler (**KIS**)**: Use kernel trick to run MCMC in  $O(p)$  time per iteration
3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. **Kernel Interaction Sampler (**KIS**)**: Use kernel trick to run MCMC in  $O(p)$  time per iteration
3. Report relevant summaries of the posterior

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. **Kernel Interaction Sampler (**KIS**)**: Use kernel trick to run MCMC in  $O(p)$  time per iteration
3. **Kernel Interaction Trick (**KIT**)**: Use kernel trick to report *all* non-negligible main and interaction effects in  $O(p)$  time

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. **Kernel Interaction Sampler (**KIS**)**: Use kernel trick to run MCMC in  $O(p)$  time per iteration
3. **Kernel Interaction Trick (**KIT**)**: Use kernel trick to report *all* non-negligible main and interaction effects in  $O(p)$  time

# Our approach

A Bayesian method: expert information, uncertainty quantification, regularization, flexibility

1. New Bayesian generative model: **Sparse Kernel Interaction Model (**SKIM**)** to encode sparsity and strong hierarchy [Carvalho et al 2009; Pironen, Vehtari 2017; Chipman 1996, Griffin & Brown 2017]
2. **Kernel Interaction Sampler (**KIS**)**: Use kernel trick to run MCMC in  $O(p)$  time per iteration
3. **Kernel Interaction Trick (**KIT**)**: Use kernel trick to report *all* non-negligible main and interaction effects in  $O(p)$  time

Not just for SKIM



# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 1: sample  $\theta$

# Kernel Interaction Sampler vs. Naive MCMC

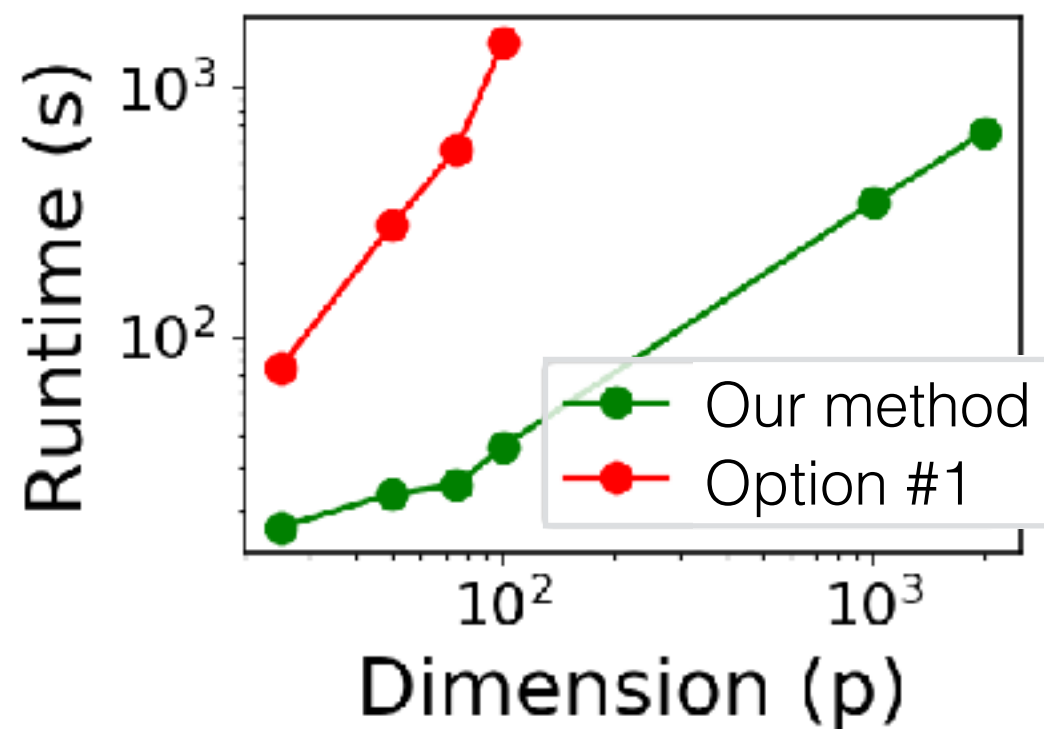
- MCMC option 1: sample  $\theta$  ( $p^2$  parameters)

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 1: sample  $\theta$  ( $p^2$  parameters)
  - Time cost:  $O(p^2N)$

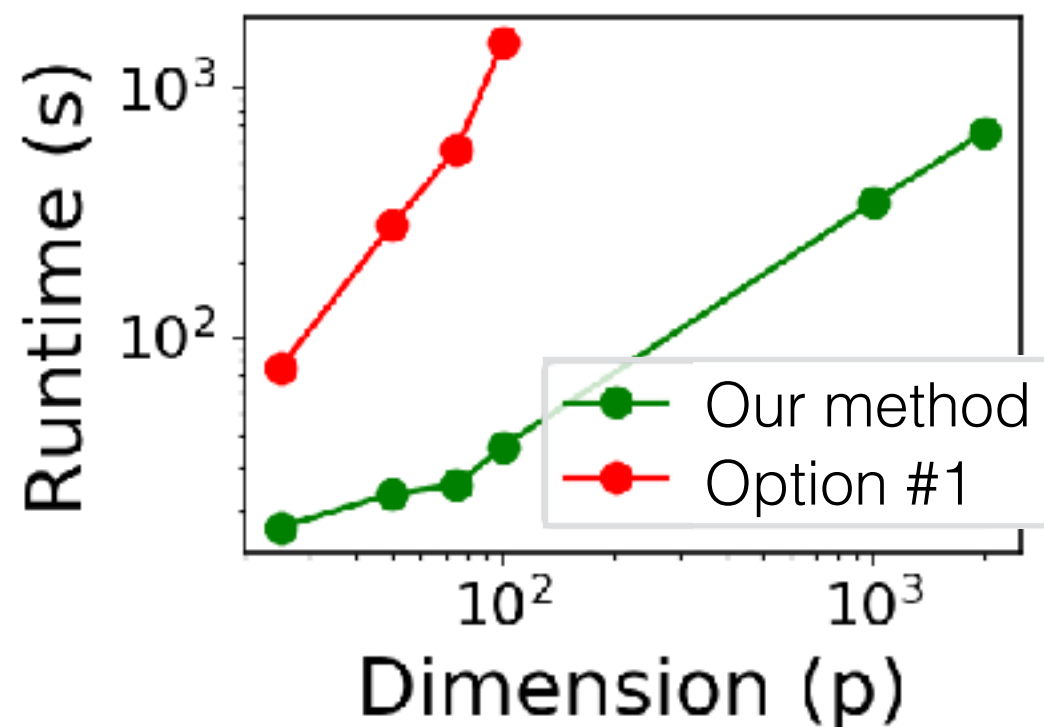
# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 1: sample  $\theta$  ( $p^2$  parameters)
  - Time cost:  $O(p^2N)$



# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 1: sample  $\theta$  ( $p^2$  parameters)
  - Time cost:  $O(p^2N)$



- Mixing (1000 iters Stan):
  - Option #1: all  $\hat{R} > 1.05$
  - Our method: all  $\hat{R} < 1.05$

# Kernel Interaction Sampler vs. Naive MCMC

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$X^{\top} X$$



# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$X^{\top} X \quad + \quad \text{prior precision matrix}$$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$X^{\top} X$$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$ 
  - Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X$ :  $N \times p$

$\Phi_2$ :  $N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$

$N$    $p^2$

# Kernel Interaction Sampler vs. Naive MCMC

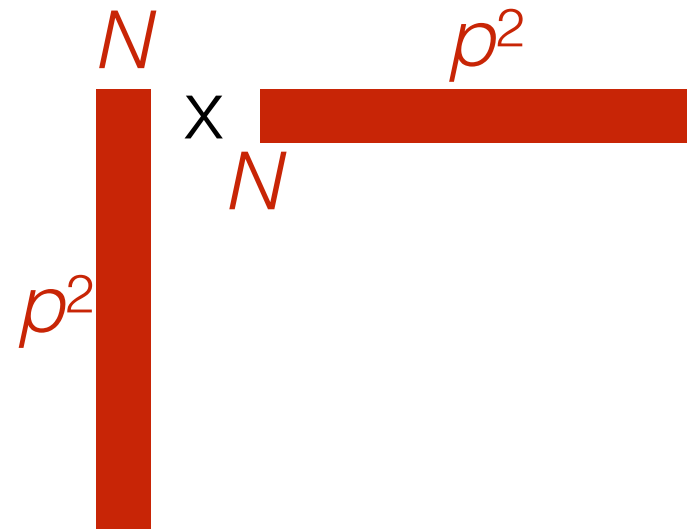
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$

$N \times N =$



# Kernel Interaction Sampler vs. Naive MCMC

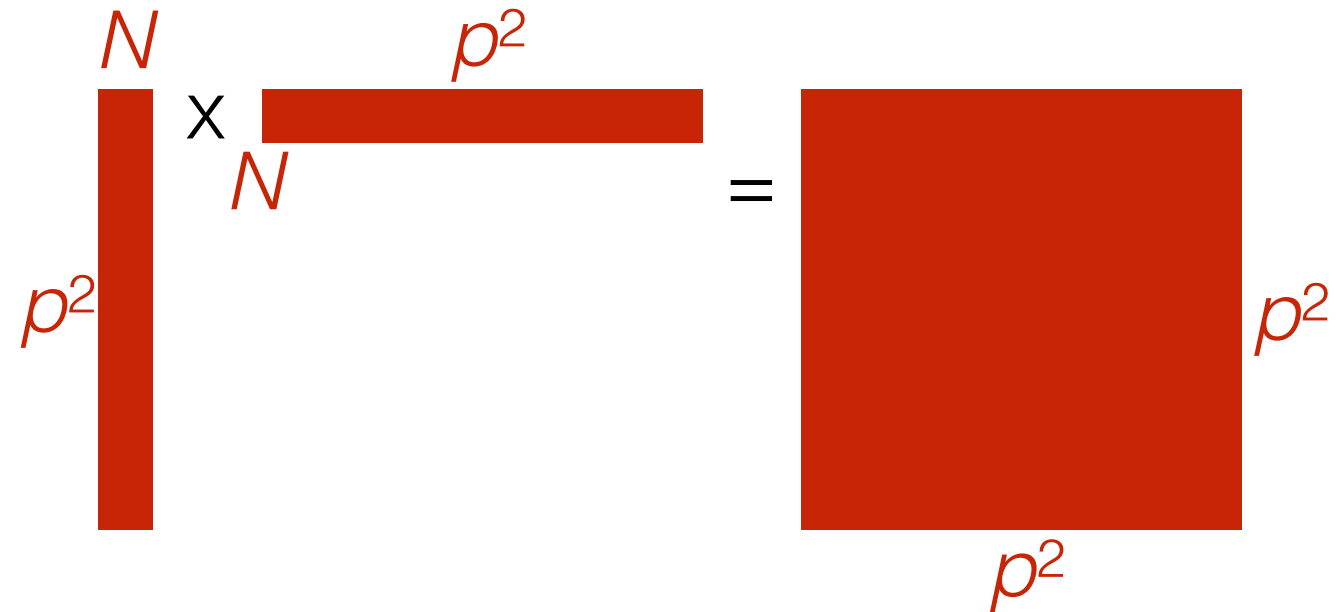
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

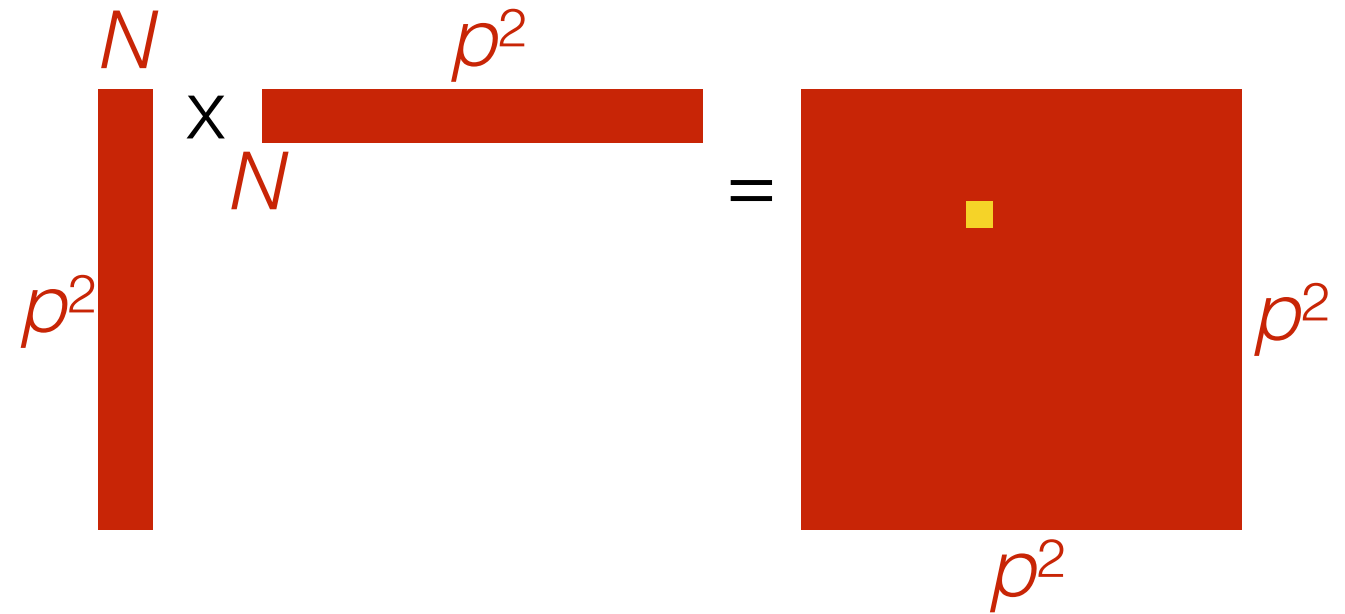
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

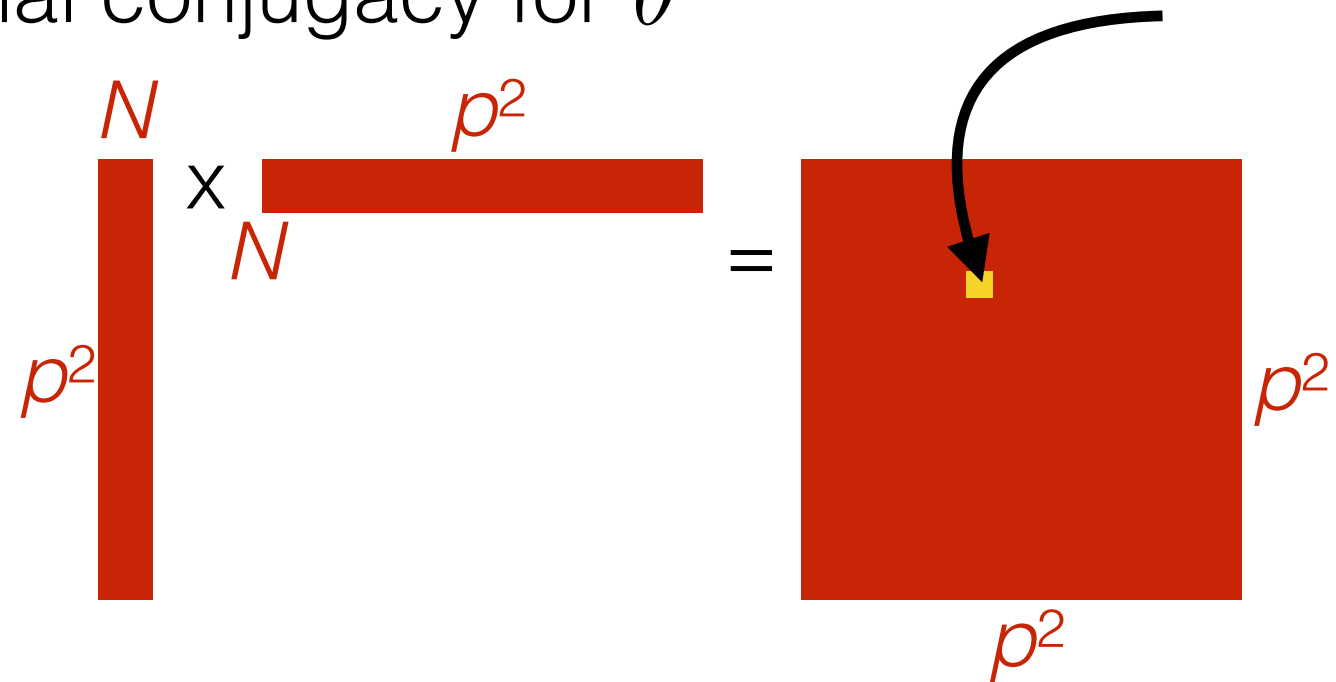
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

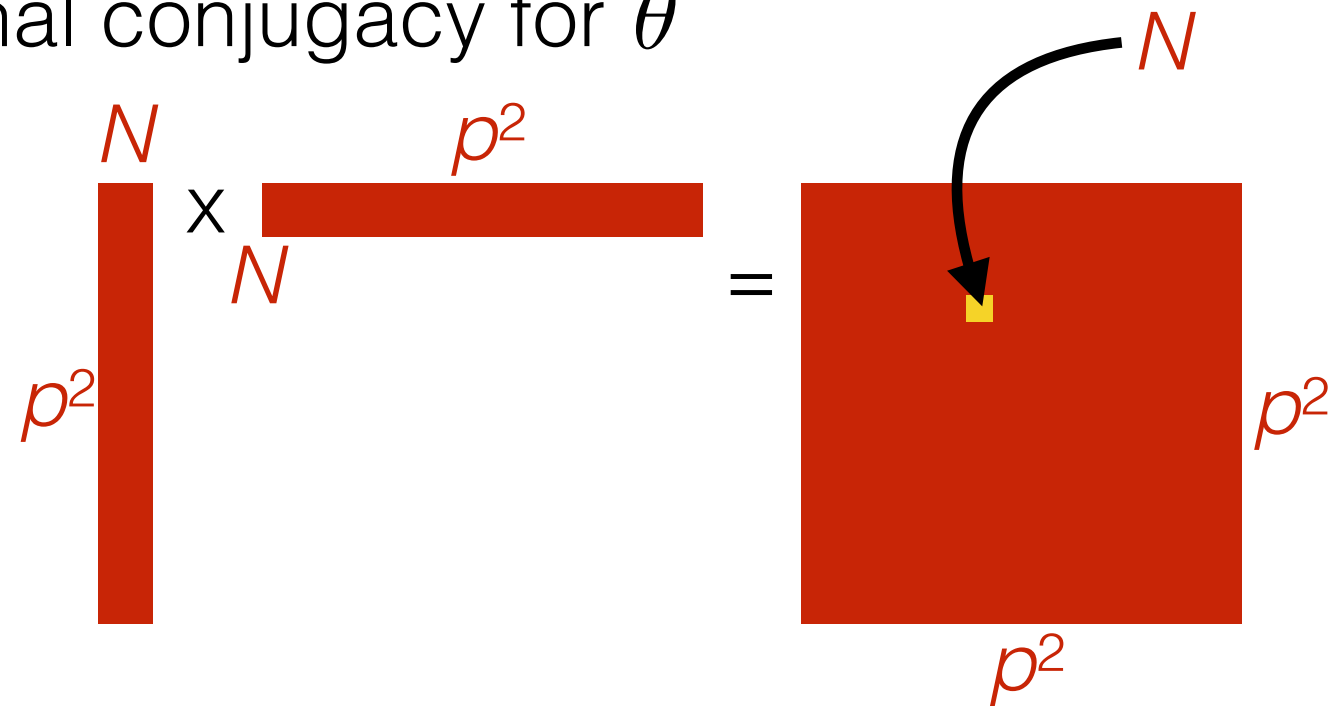
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

- MCMC option 2: use conditional conjugacy for  $\theta$

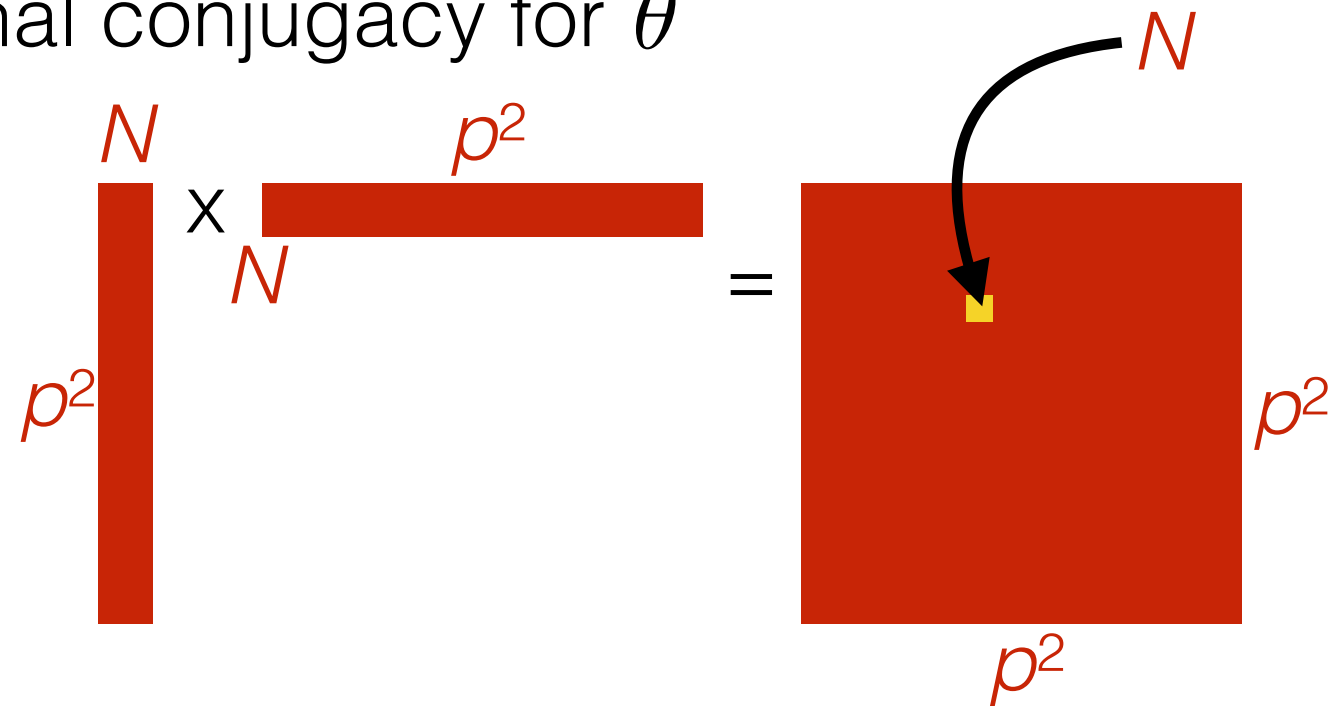
- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$

- Naive time cost:  $O(p^4 N + p^6)$



# Kernel Interaction Sampler vs. Naive MCMC

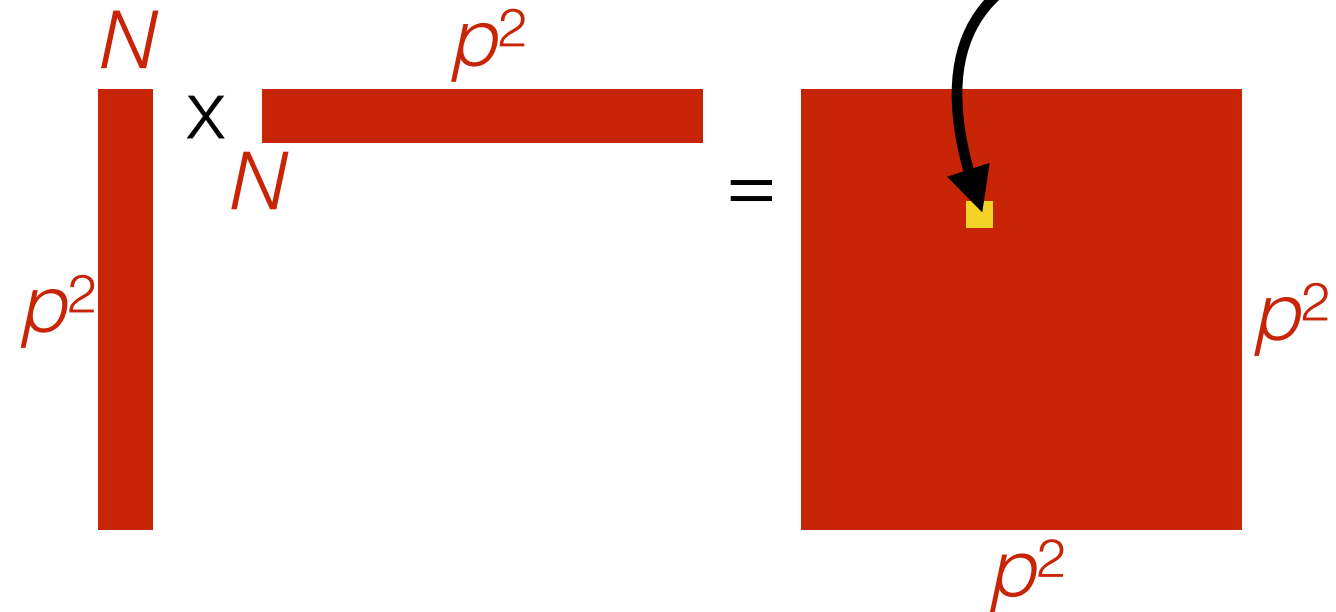
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



- Naive time cost:  $O(p^4 N + p^6)$
- Woodbury time cost:  $O(p^2 N^2 + N^3)$

# Kernel Interaction Sampler vs. Naive MCMC

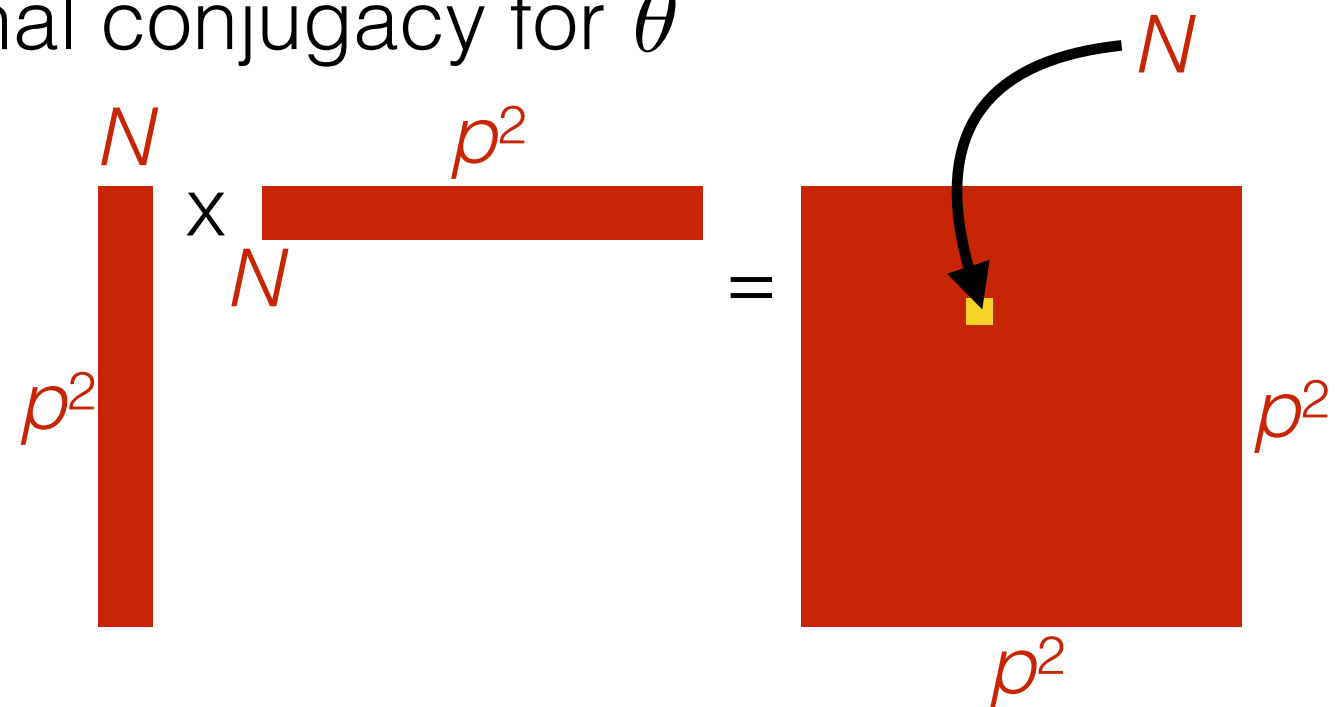
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

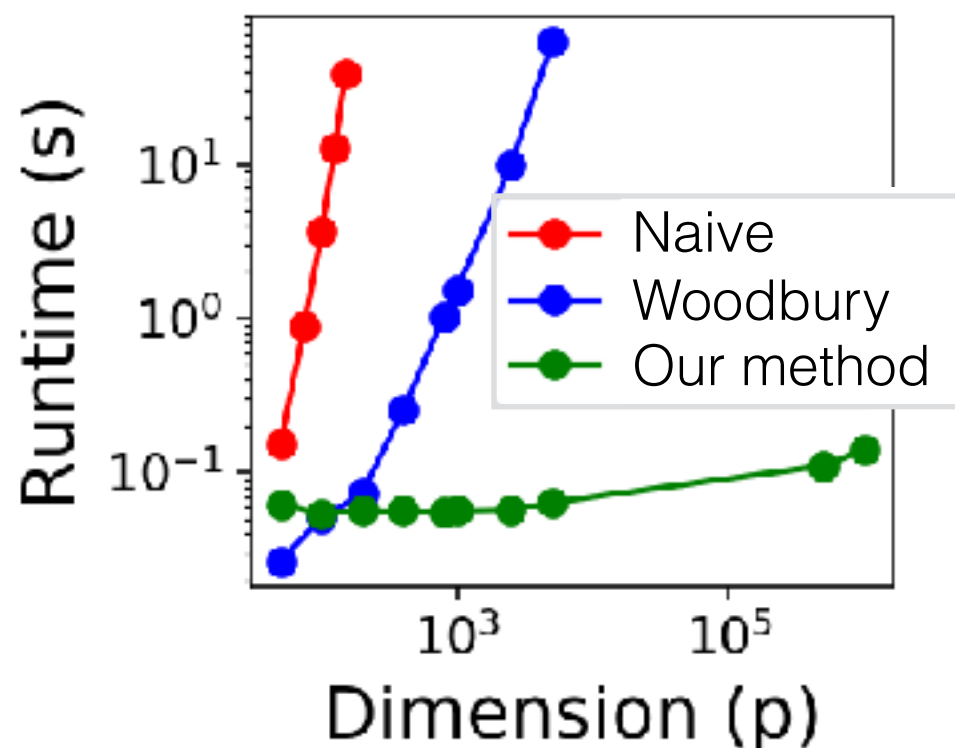
$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



- Naive time cost:  $O(p^4 N + p^6)$
- Woodbury time cost:  $O(p^2 N^2 + N^3)$



# Kernel Interaction Sampler vs. Naive MCMC

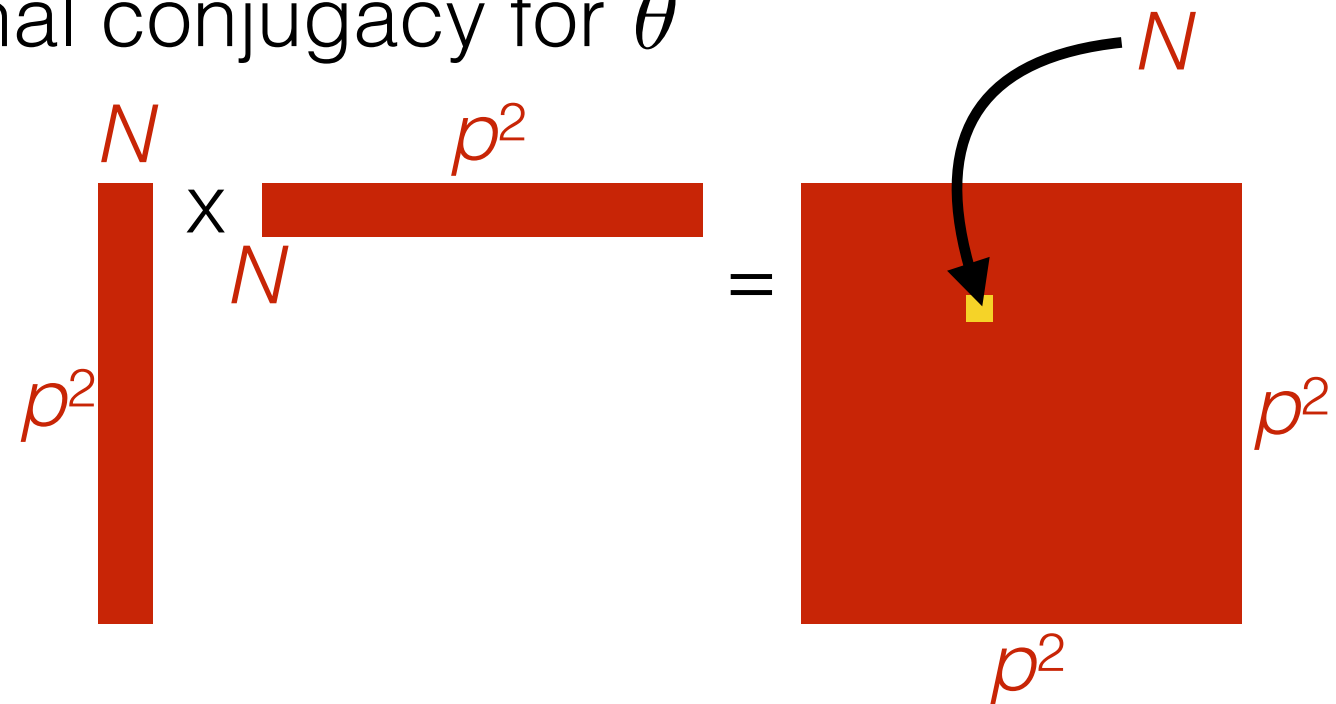
- MCMC option 2: use conditional conjugacy for  $\theta$

- Compute and invert

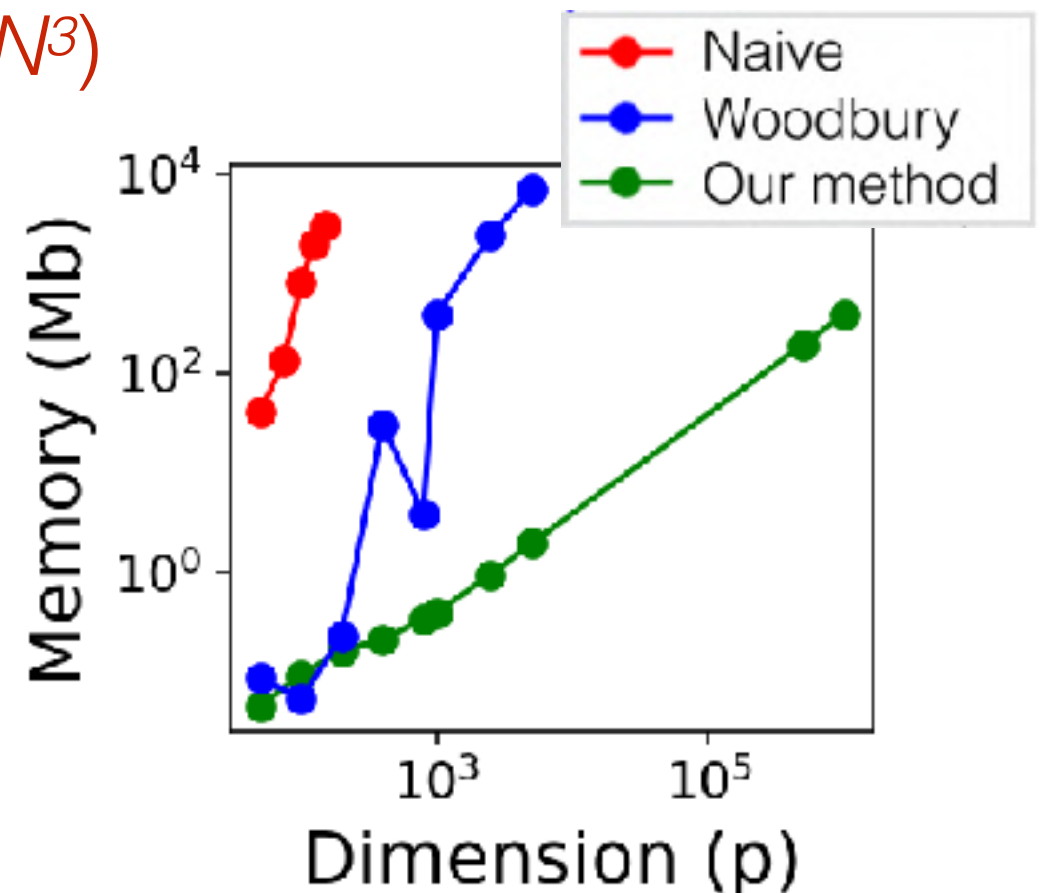
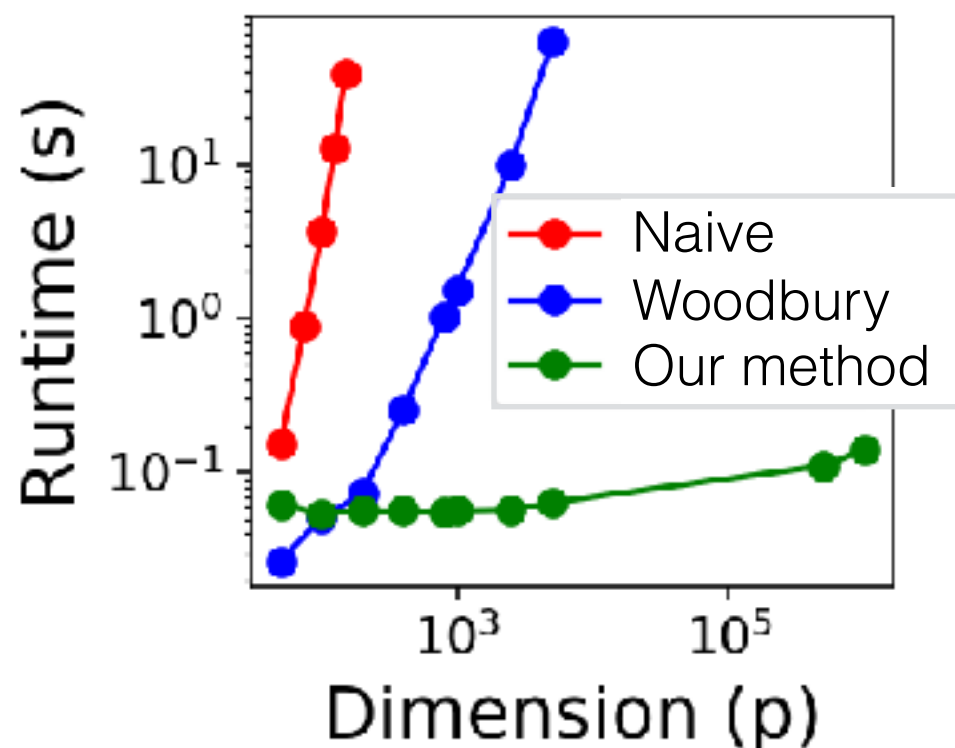
$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$



- Naive time cost:  $O(p^4 N + p^6)$
- Woodbury time cost:  $O(p^2 N^2 + N^3)$





# Kernel Interaction Sampler vs. Naive MCMC


- Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$$X: N \times p$$


$$\Phi_2: N \times p^2$$

# Kernel Interaction Sampler vs. Naive MCMC

- Compute and invert  
  $\Phi_2(X)^\top \Phi_2(X)$   
 $X: N \times p$   
 $\Phi_2: N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC


use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert  
  $\Phi_2(X)^\top \Phi_2(X)$   
 $X: N \times p$   
 $\Phi_2: N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert


$$\Phi_2(X)^\top \Phi_2(X)$$

$X: N \times p$

$\Phi_2: N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X)^\top \Phi_2(X)$$

$X$ :  $N \times p$

$\Phi_2$ :  $N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top \Phi_2(X) \Phi_2(X)^\top$$

$X$ :  $N \times p$

$\Phi_2$ :  $N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top \Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top \Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X) \Phi_2(X)^T$$

$X: N \times p$

$\Phi_2: N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X$ :  $N \times p$

$\Phi_2$ :  $N \times p^2$

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X$ :  $N \times p$

$\Phi_2$ :  $N \times p^2$

$$N \times p^2$$

# Kernel Interaction Sampler vs. Naive MCMC

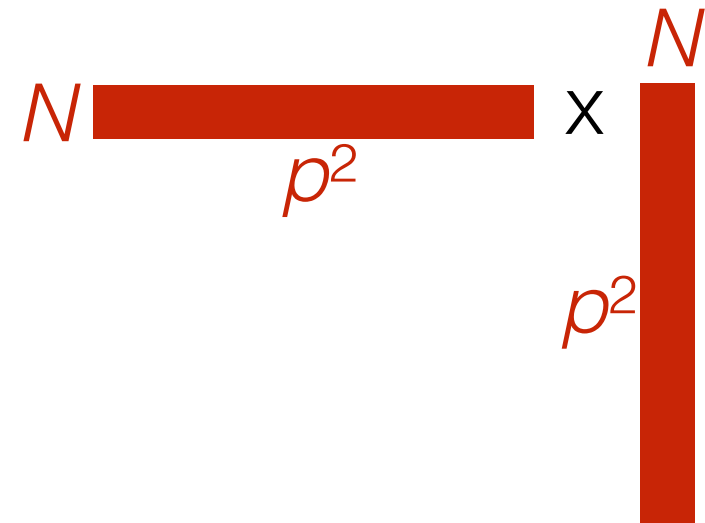
- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^{\top}$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

A diagram illustrating the dimensions of the matrix multiplication  $\Phi_2(X) \Phi_2(X)^\top$ . It consists of a horizontal red bar representing a matrix of size  $N \times p^2$ , with  $N$  above it and  $p^2$  below it. To its right is a vertical red bar representing a matrix of size  $N \times p^2$ , with  $N$  to its left and  $p^2$  to its right. The two bars are separated by a multiplication symbol 'x', and the entire expression is followed by an equals sign '='.

# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

$N \times p^2 \times p^2 \times N = N \times N$

# Kernel Interaction Sampler vs. Naive MCMC

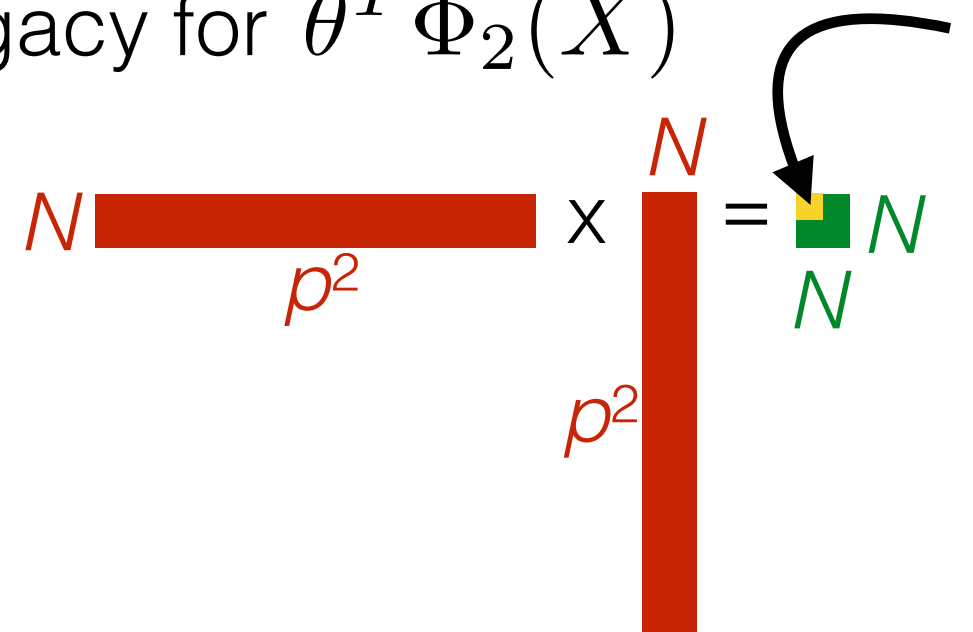
- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$



# Kernel Interaction Sampler vs. Naive MCMC

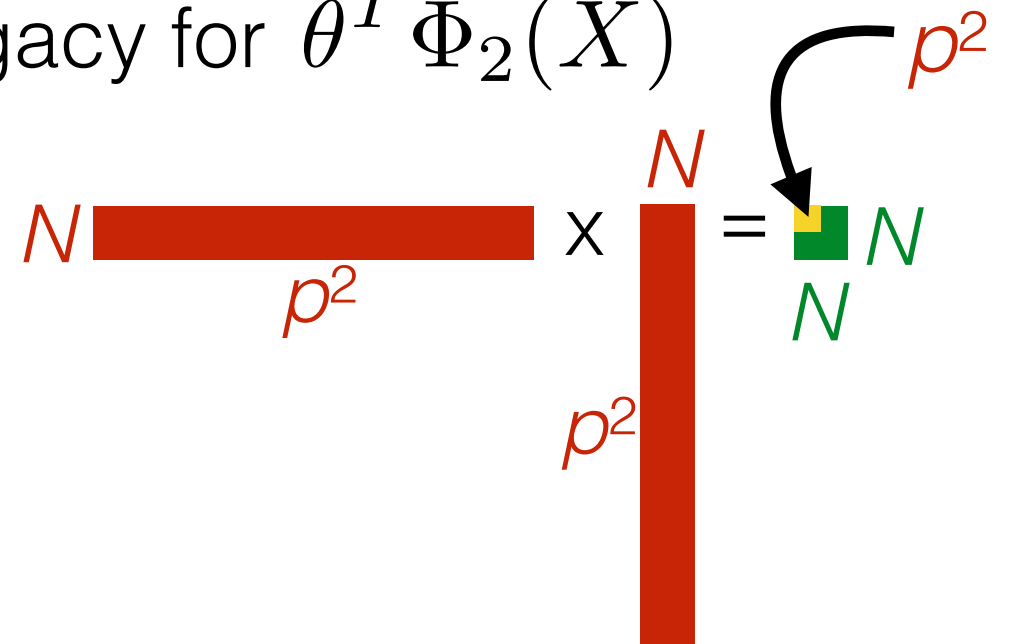
- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$





# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

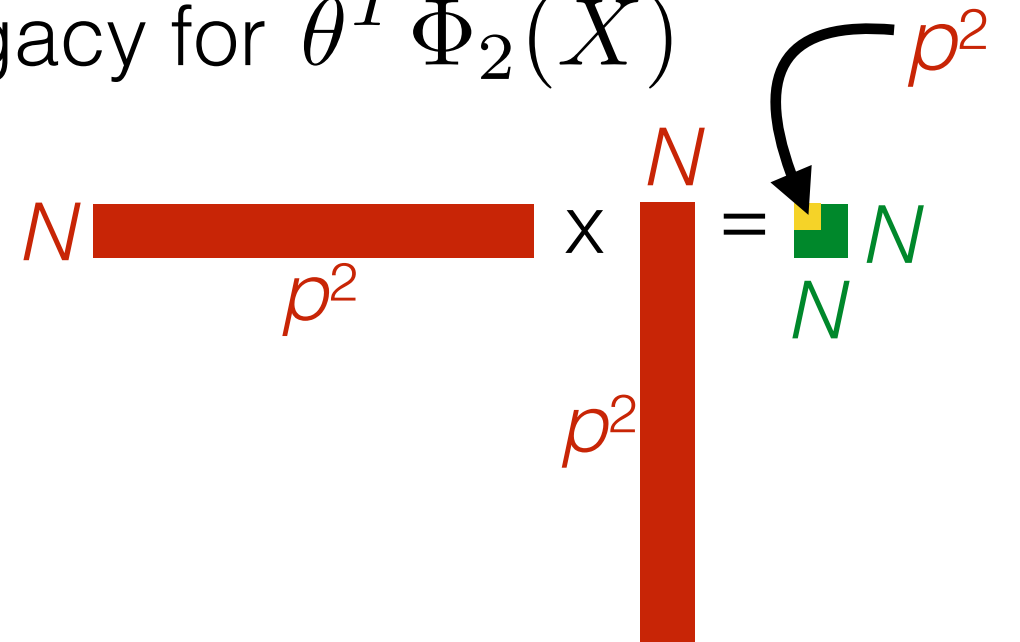
- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

- Kernel trick:  $O(p)$  cost



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

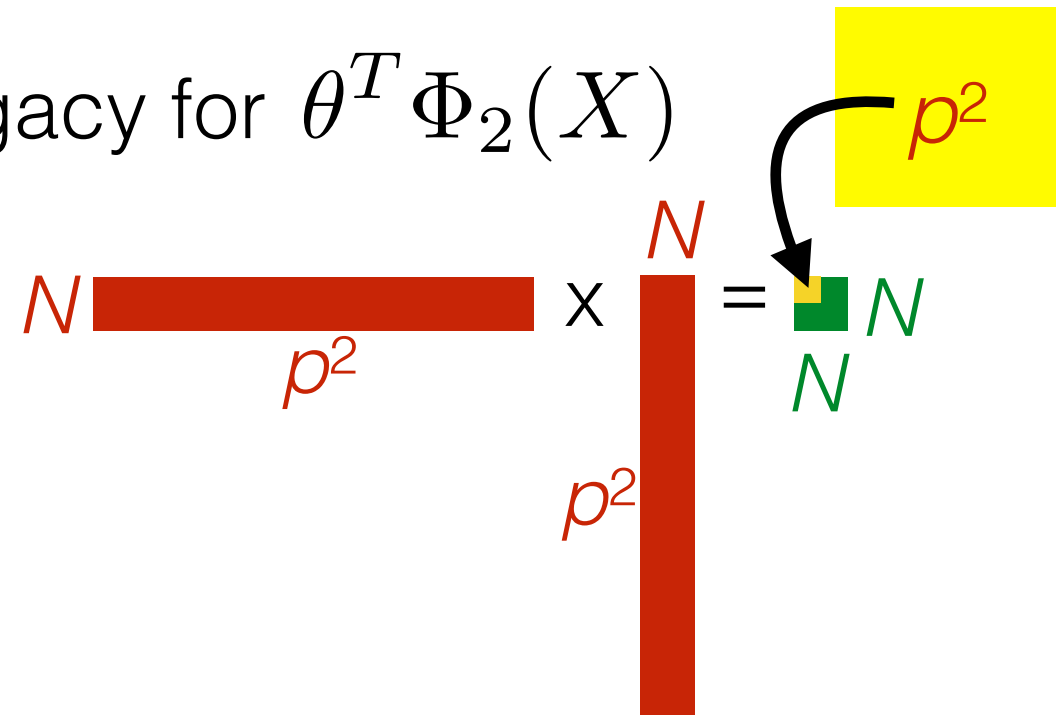
- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

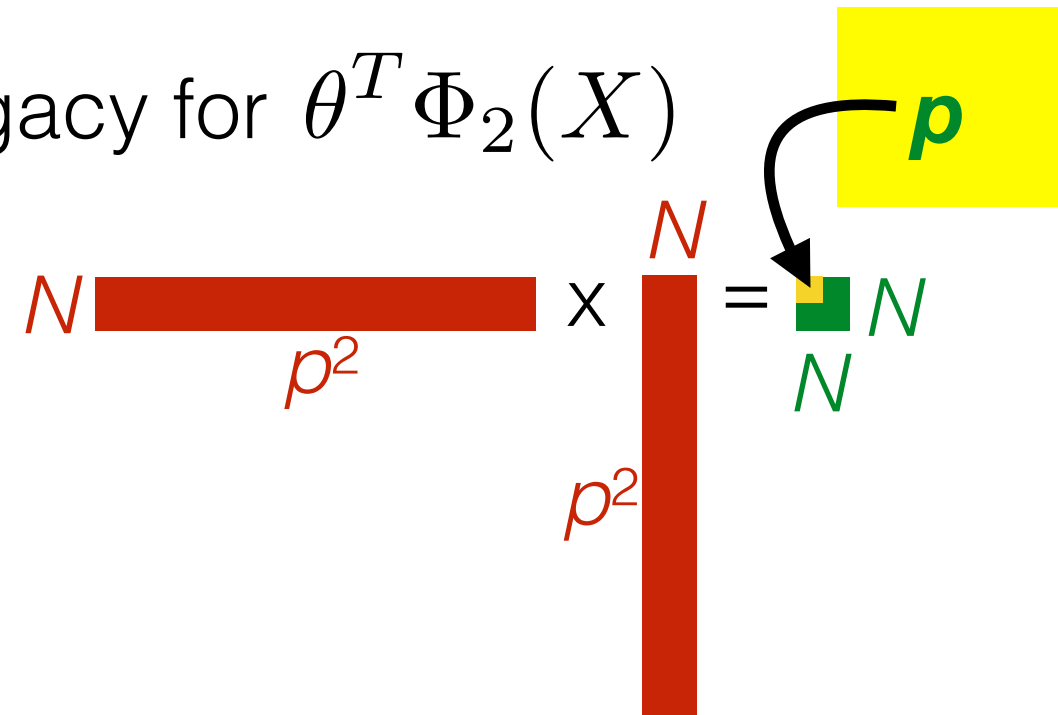
$\Phi_2: N \times p^2$

- Kernel trick:  $O(p)$  cost



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$ 
  - Compute and invert  $\Phi_2(X) \Phi_2(X)^\top$   
 $X: N \times p$   
 $\Phi_2: N \times p^2$
  - Kernel trick:  $O(p)$  cost



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

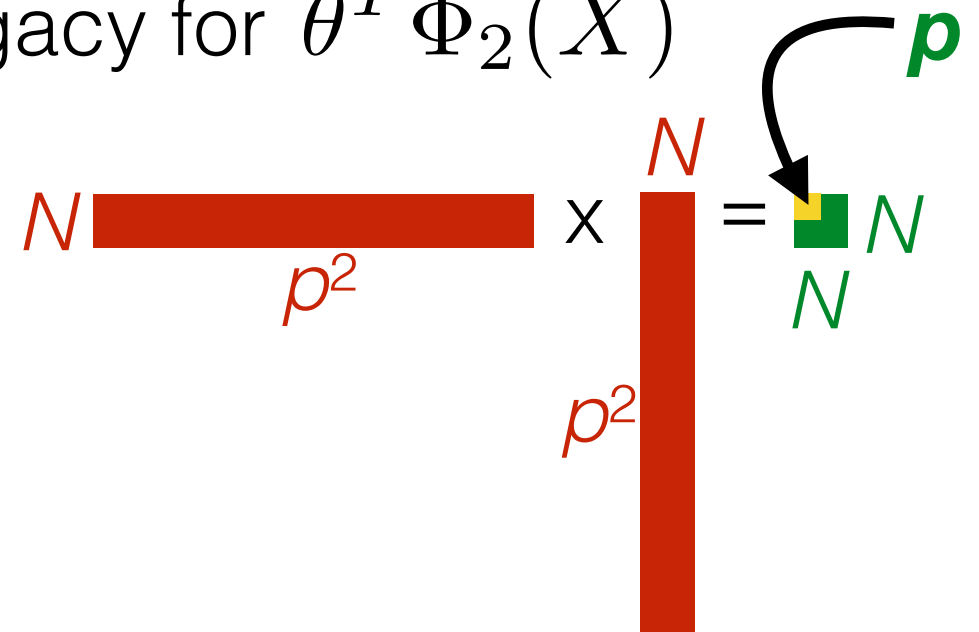
- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

- Kernel trick:  $O(p)$  cost



# Kernel Interaction Sampler vs. Naive MCMC

- Our approach: use conditional conjugacy for  $\theta^T \Phi_2(X)$

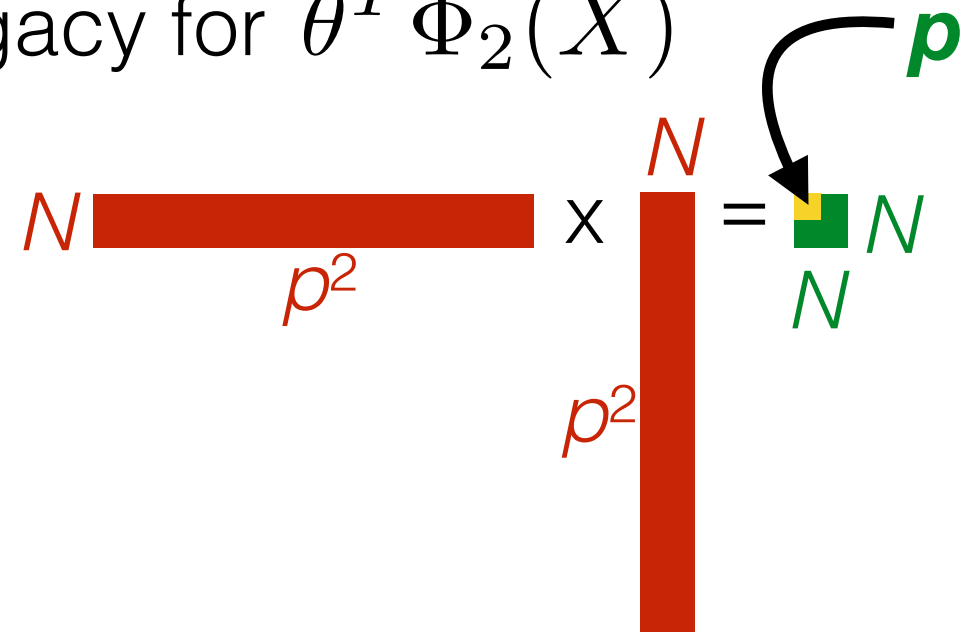
- Compute and invert

$$\Phi_2(X) \Phi_2(X)^\top$$

$X: N \times p$

$\Phi_2: N \times p^2$

- Kernel trick:  $O(p)$  cost
  - Our time cost:  $O(\mathbf{p}N^2 + N^3)$



# Reporting: Kernel Interaction Trick

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects




# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$

  
 $i^{\text{th}}$  position

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$


  
 $i^{\text{th}}$  position

$$g(e_i) = \theta_{x_i} + \theta_{x_i^2}$$

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$

  
 $i^{\text{th}}$  position

$$g(e_i) = \theta_{x_i} + \theta_{x_i^2}$$

$$g(-e_i) = -\theta_{x_i} + \theta_{x_i^2}$$

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$

  
 $i^{\text{th}}$  position

$$g(e_i) = \theta_{x_i} + \theta_{x_i^2}$$

$$g(-e_i) = -\theta_{x_i} + \theta_{x_i^2}$$

$$\frac{g(e_i) - g(-e_i)}{2} = \theta_{x_i}$$

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$

  
 $i^{\text{th}}$  position

$$g(e_i) = \theta_{x_i} + \theta_{x_i^2}$$

$$g(-e_i) = -\theta_{x_i} + \theta_{x_i^2}$$

$$\frac{g(e_i) - g(-e_i)}{2} = \theta_{x_i}$$

- Step B: Find  $k \ll p$  sparse main effects: takes  $O(p)$  time

# Reporting: Kernel Interaction Trick

- Can access posterior of  $g(x) = \theta^T \Phi_2(x)$  in  $O(p)$  time per iteration
- *But* our goal is to find main and interaction effects
- Step A: Report posterior of  $\theta_{x_i}$  or  $\theta_{x_i x_j}$  in  $O(1)$  time

$$e_i = [0, 0, \dots, 0, 1, 0, \dots, 0]$$

  
 $i^{\text{th}}$  position

$$g(e_i) = \theta_{x_i} + \theta_{x_i^2}$$

$$g(-e_i) = -\theta_{x_i} + \theta_{x_i^2}$$

$$\frac{g(e_i) - g(-e_i)}{2} = \theta_{x_i}$$

- Step B: Find  $k \ll p$  sparse main effects: takes  $O(p)$  time
- Step C: Report just the  $k^2$  strong-hierarchy interaction effects: takes  $O(k^2)$  time

# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)
  - Fast reporting of results
- Experiments on simulated and real data



# Roadmap

- Setup: Discovering main and interaction effects
- Our method
  - A Bayesian generative model
  - Fast inference (using Gaussian processes)
  - Fast reporting of results
- Experiments on simulated and real data

# Timing vs. LASSO-based methods

# Timing vs. LASSO-based methods

- LASSO (pairs, hierarchical):  $O(p^2)$  per iteration

# Timing vs. LASSO-based methods

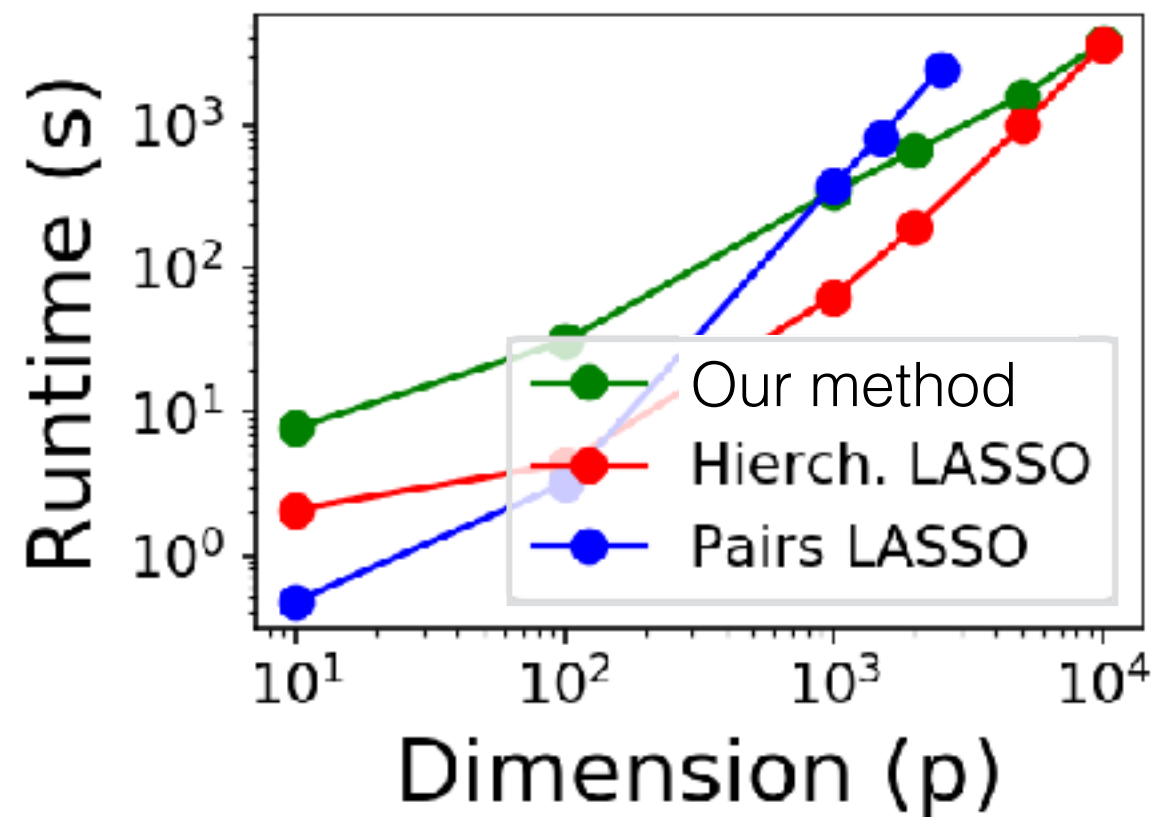
- LASSO (pairs, hierarchical):  $O(p^2)$  per iteration [Lim, Hastie 2015]

# Timing vs. LASSO-based methods

- LASSO (pairs, hierarchical):  $O(p^2)$  per iteration [Lim, Hastie 2015]
- Our method:  $O(p)$  per iteration

# Timing vs. LASSO-based methods

- LASSO (pairs, hierarchical):  $O(p^2)$  per iteration [Lim, Hastie 2015]
- Our method:  $O(p)$  per iteration
- Competitive empirically for moderate  $p$ :



# Experiments: Simulated

# Experiments: Simulated

- 36 different simulated data sets (so know true effects)



# Experiments: Simulated

- 36 different simulated data sets (so know true effects)
  - Up to  $p = 500 \rightarrow \approx 125,000$  total parameters

# Experiments: Simulated, Selection

- 36 different simulated data sets (so know true effects)
- Up to  $p = 500 \rightarrow \approx 125,000$  total parameters

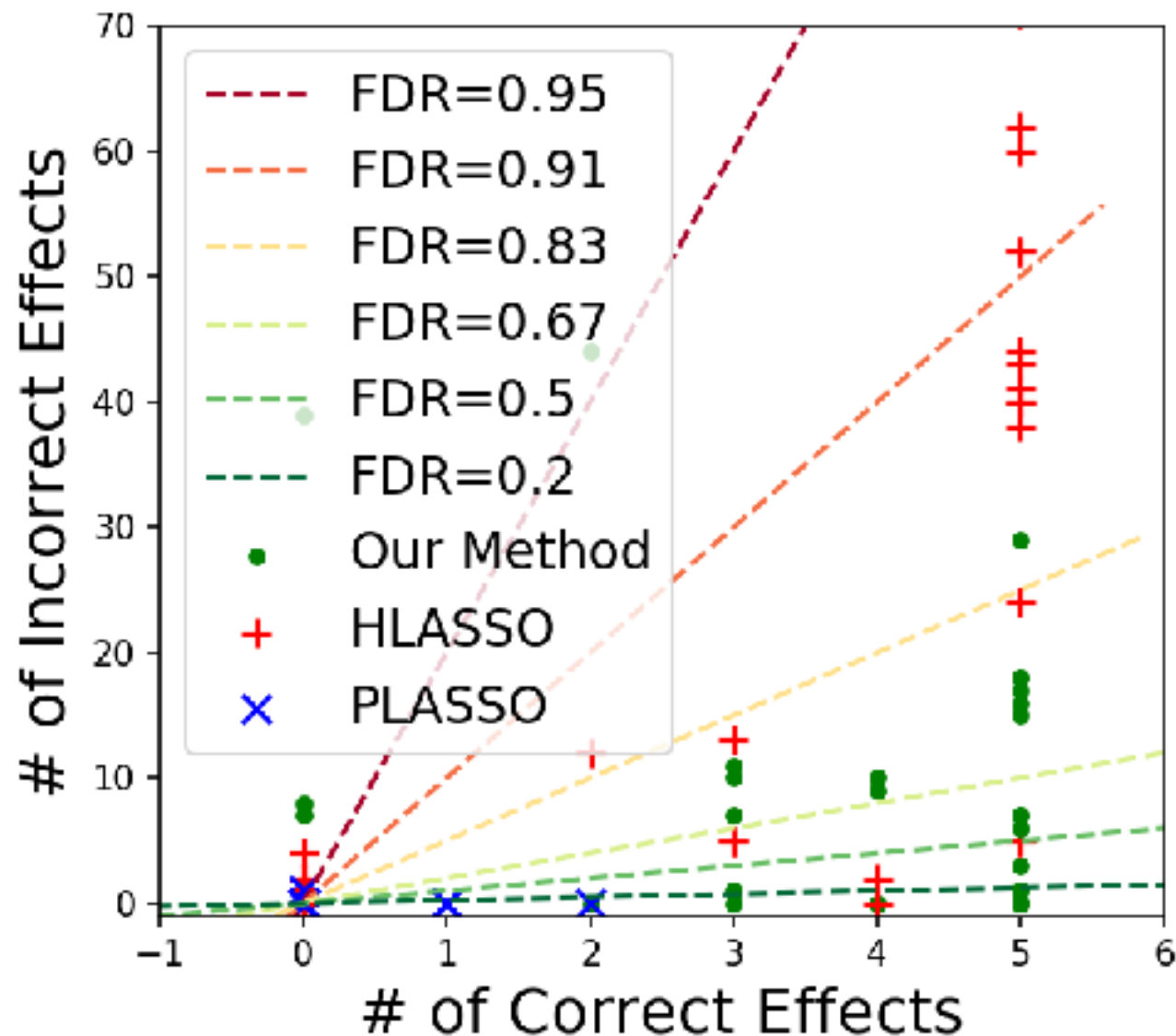
# Experiments: Simulated, Selection

- 36 different simulated data sets (so know true effects)
  - Up to  $p = 500 \rightarrow \approx 125,000$  total parameters
- False discovery rate (FDR): proportion incorrect

# Experiments: Simulated, Selection

- 36 different simulated data sets (so know true effects)
  - Up to  $p = 500 \rightarrow \approx 125,000$  total parameters
- False discovery rate (FDR): proportion incorrect

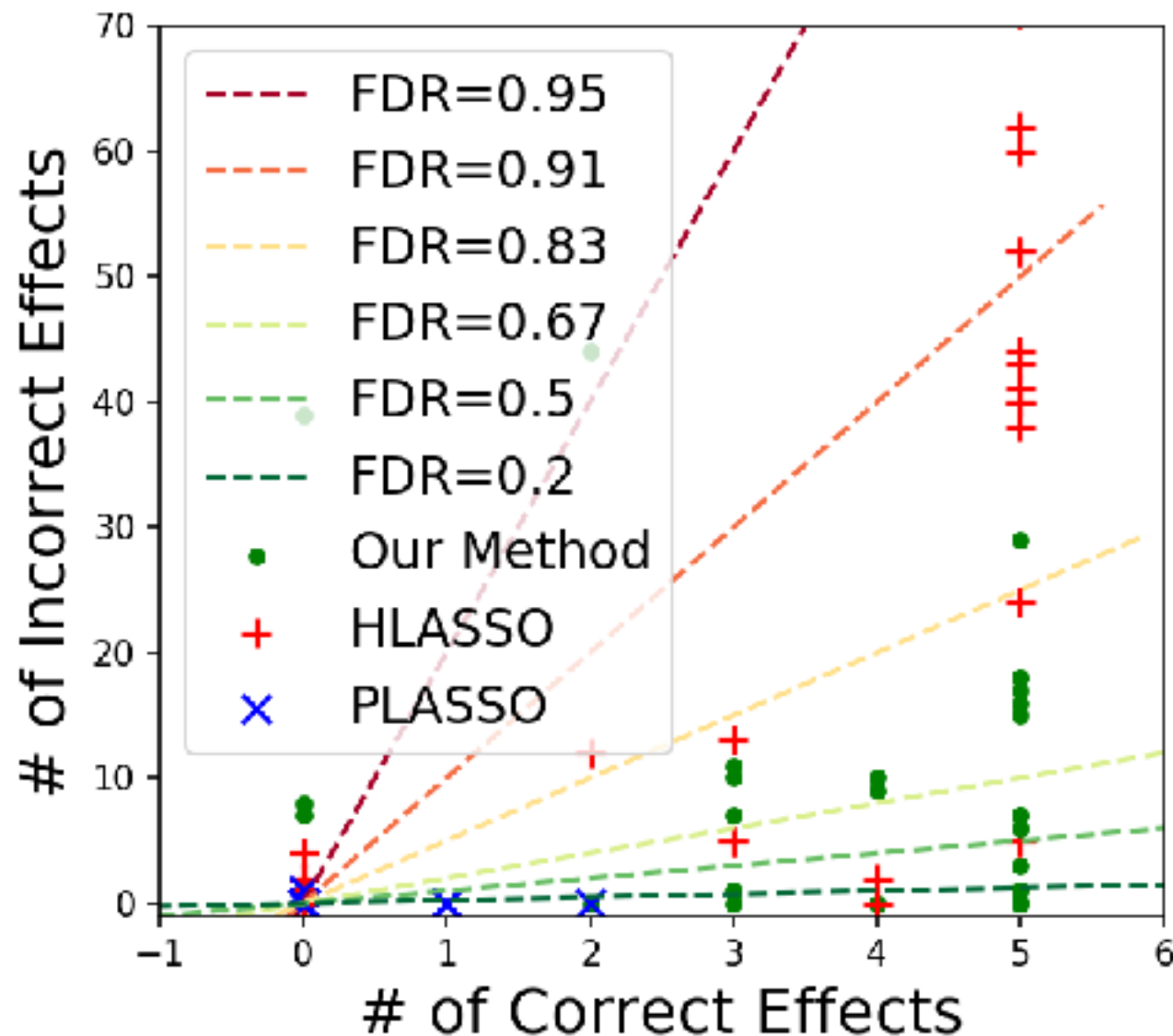
Main effects



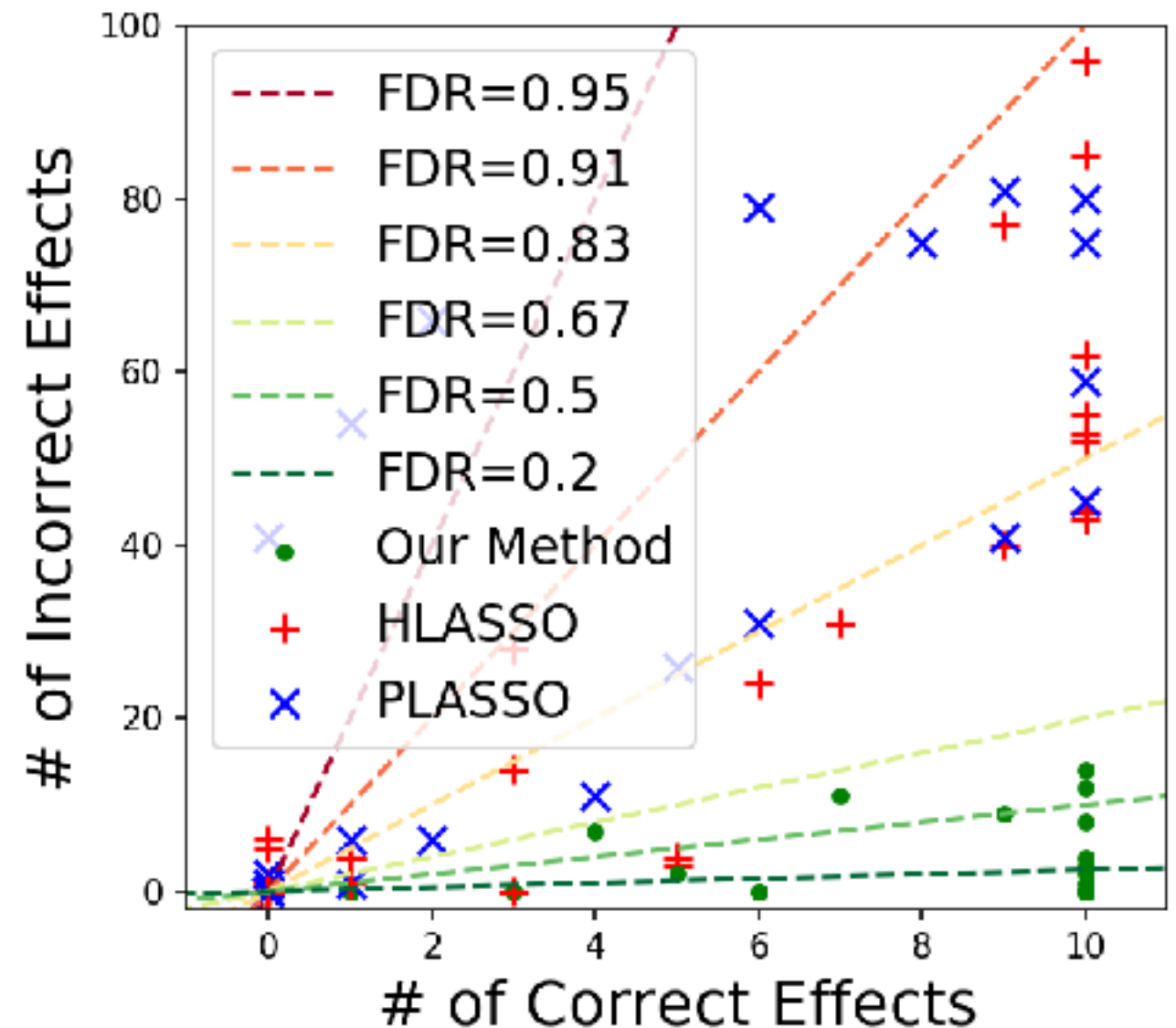
# Experiments: Simulated, Selection

- 36 different simulated data sets (so know true effects)
  - Up to  $p = 500 \rightarrow \approx 125,000$  total parameters
- False discovery rate (FDR): proportion incorrect

Main effects



Pairwise effects



# Experiments: Real covariates

# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction

# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction
- Covariates: Residential Building Data Set
  - Highly correlated: 20 of 105 capture 99% of variance



# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction
- Covariates: Residential Building Data Set
  - Highly correlated: 20 of 105 capture 99% of variance
- Key: (# correct effects): (# of incorrect effects)
  - **Higher** green is better: **lower** red is better

# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction
- Covariates: Residential Building Data Set
  - Highly correlated: 20 of 105 capture 99% of variance
- Key: (# correct effects): (# of incorrect effects)
  - **Higher** green is better: **lower** red is better

METHOD	#MAIN	#PAIR
PLASSO	2 : 5	3 : 21

# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction
- Covariates: Residential Building Data Set
  - Highly correlated: 20 of 105 capture 99% of variance
- Key: (# correct effects): (# of incorrect effects)
  - **Higher** green is better: **lower** red is better

METHOD	#MAIN	#PAIR
PLASSO	2 : 5	3 : 21
HLASSO	3 : 19	3 : 18

# Experiments: Real covariates

- Simulated effects: 5 main, 10 interaction
- Covariates: Residential Building Data Set
  - Highly correlated: 20 of 105 capture 99% of variance
- Key: (# correct effects): (# of incorrect effects)
  - **Higher** green is better: **lower** red is better

METHOD	#MAIN	#PAIR
Our method	3 : 0	3 : 0
PLASSO	2 : 5	3 : 21
HLASSO	3 : 19	3 : 18

# Experiments: Real data

# Experiments: Real data

- Covariates and response: Auto MPG

# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...

# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...
- Augment  $p$  with 200 fake (noise) covariates
  - 21,321 total parameters



# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...
- Augment  $p$  with 200 fake (noise) covariates
  - 21,321 total parameters
- Key: (# original effects): (# of fake effects)
  - **No order** to blue: **lower** red is better

# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...
- Augment  $p$  with 200 fake (noise) covariates
  - 21,321 total parameters
- Key: (# original effects): (# of fake effects)
  - **No order** to blue: **lower** red is better

METHOD	#MAIN	#PAIR
PLASSO	4 : 0	2 : 78

# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...
- Augment  $p$  with 200 fake (noise) covariates
  - 21,321 total parameters
- Key: (# original effects): (# of fake effects)
  - **No order** to blue: **lower** red is better

METHOD	#MAIN	#PAIR
PLASSO	4 : 0	2 : 78
HLASSO	6 : 46	4 : 38

# Experiments: Real data

- Covariates and response: Auto MPG
- $N = 398$ ,  $p = 6$  (real-valued), but...
- Augment  $p$  with 200 fake (noise) covariates
  - 21,321 total parameters
- Key: (# original effects): (# of fake effects)
  - **No order** to blue: **lower** red is better

METHOD	#MAIN	#PAIR
Our method	3 : 0	1 : 0
PLASSO	4 : 0	2 : 78
HLASSO	6 : 46	4 : 38

# Conclusions

**We provide:** fast, accurate detection of pairwise interactions

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML* 2019. ArXiv:1905.06501**

# Conclusions

**We provide:** fast, accurate detection of pairwise interactions

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:**

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

- Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*
- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)



# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption; (3) Improve scaling in  $N$

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

• In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption; (3) Improve scaling in  $N$

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

• In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

JH Huggins, T Campbell, M Kasprzak, and T Broderick. Scalable Gaussian process inference with finite-data mean and variance guarantees. *AISTATS 2019*.

R Agrawal, T Campbell, JH Huggins, and T Broderick. Data-dependent compression of random features for large-scale kernel approximation. *AISTATS 2019*.

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption; (3) Improve scaling in  $N$ ; (4) Applications!

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

• In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

JH Huggins, T Campbell, M Kasprzak, and T Broderick. Scalable Gaussian process inference with finite-data mean and variance guarantees. *AISTATS 2019*.

R Agrawal, T Campbell, JH Huggins, and T Broderick. Data-dependent compression of random features for large-scale kernel approximation. *AISTATS 2019*.

# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption; (3) Improve scaling in  $N$ ; (4) Applications!

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

• In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

JH Huggins, T Campbell, M Kasprzak, and T Broderick. Scalable Gaussian process inference with finite-data mean and variance guarantees. *AISTATS 2019*.

R Agrawal, T Campbell, JH Huggins, and T Broderick. Data-dependent compression of random features for large-scale kernel approximation. *AISTATS 2019*.



# Conclusions

**We provide:** fast, accurate detection of pairwise (and higher-order) interactions

**Up next:** (1) Nonlinearity; (2) Remove strong hierarchy assumption; (3) Improve scaling in  $N$ ; (4) Applications!

**R Agrawal, BL Trippe, JH Huggins, and T Broderick. The Kernel interaction trick: Fast Bayesian discovery of pairwise interactions in high dimensions. *ICML 2019*. ArXiv:1905.06501**

*Thanks to Pyro contributors! Martin Jankowiak, Du Phan, Neeraj Pradhan*

- In Pyro: [https://pyro.ai/examples/sparse\\_regression.html](https://pyro.ai/examples/sparse_regression.html)

**R Agrawal, T Broderick. The SKIM-FA Kernel: High-Dimensional Variable Selection and Nonlinear Interaction Discovery in Linear Time. <https://arxiv.org/abs/2106.12408>**

JH Huggins, T Campbell, M Kasprzak, and T Broderick. Scalable Gaussian process inference with finite-data mean and variance guarantees. *AISTATS 2019*.

R Agrawal, T Campbell, JH Huggins, and T Broderick. Data-dependent compression of random features for large-scale kernel approximation. *AISTATS 2019*.

WT Stephenson, S Ghosh, TD Nguyen, M Yurochkin, SK Deshpande, and T Broderick. Measuring the sensitivity of Gaussian processes to kernel

15 choice. *AISTATS 2022*, to appear.