

NIPS 2010

# **A Dirty Model for Multitask Learning**

Ali Jalali, Pradeep Ravikumar, Sujay Sanghavi and Chao Ruan  
UT Austin



# Motivation

---

- Modern Settings: **High-Dimensional** Problems
  - number of observations  $n \ll$  number of variables  $p$
  - Biology, Vision, Nanotechnology, Financial Analysis, ...
- **Low-Dimensional** Structure only hope for consistency?
  - Sparsity, Block Sparsity, Low-Rank, Graphical model Structure

## This talk:

- Superposition of structures: also low-dimensional but can model surprisingly “dirty” data.



# Multitask Learning

---

- Multiple tasks with some “shared” structure

## Problem:

- Learn tasks **jointly** (as opposed to **separately**)
  - e.g. Optical Character Recognition (OCR)

Writer 1

A	A	A	A	A	A
A	A	A	A	A	A

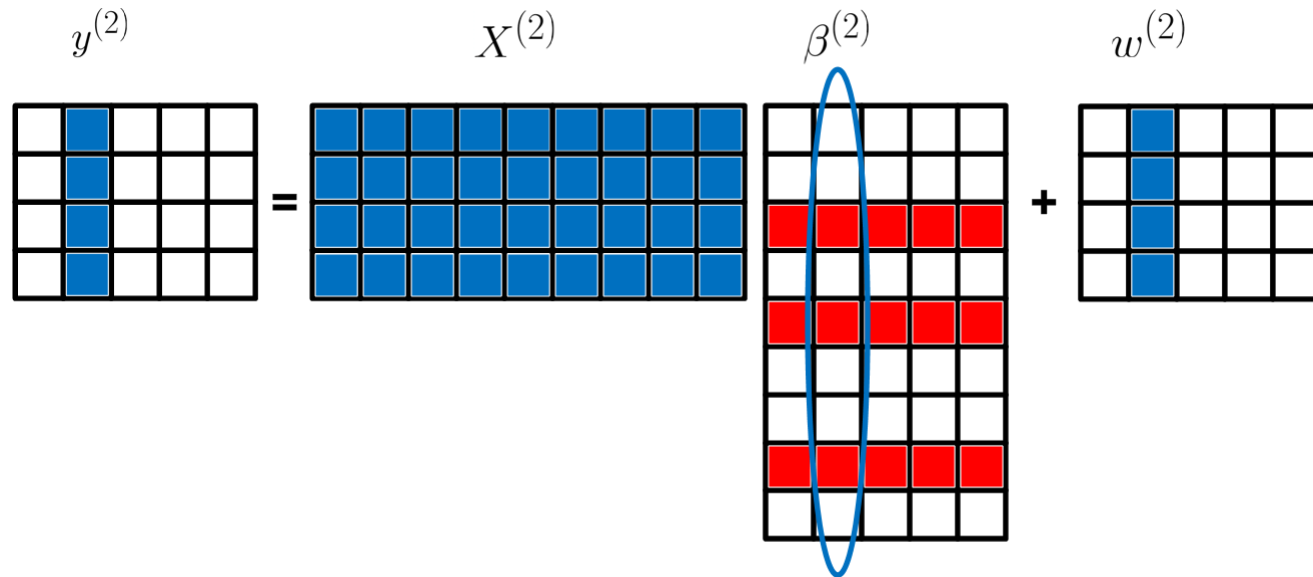
Writer 2

A	A	A	A	A	A
A	A	A	A	A	A



# Multiple Linear Regressions

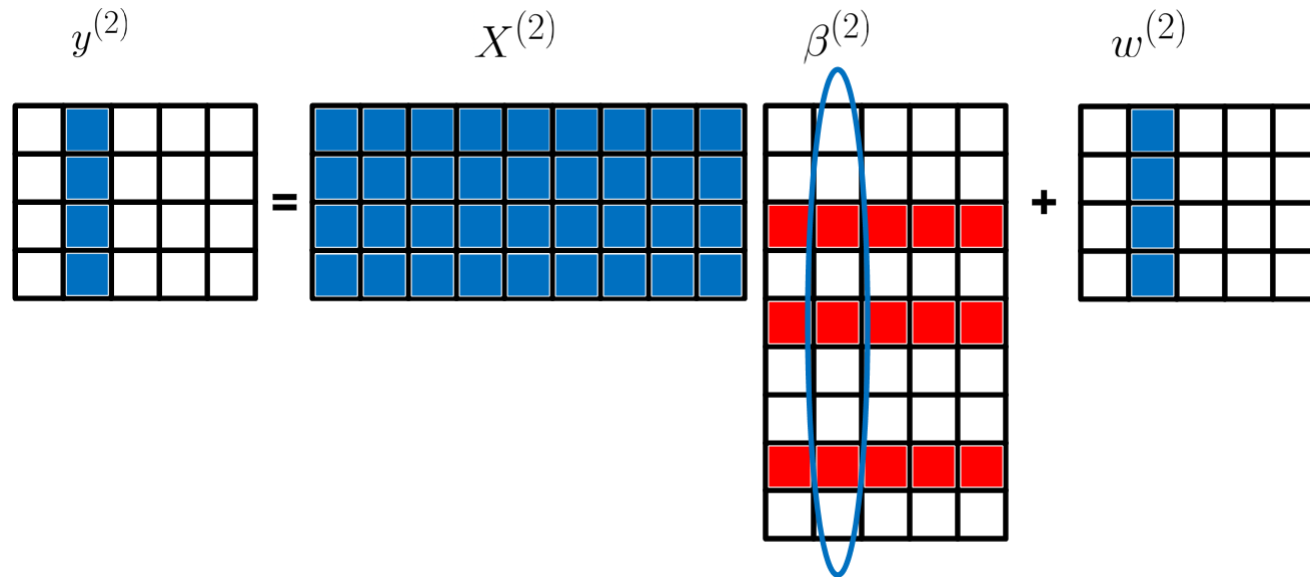
---



- Linear Model:  $y^{(k)} = X^{(k)}\beta^{(k)} + w^{(k)}$  for all tasks  $1 \leq k \leq r$
- **Problem:** Estimate  $\beta$  given  $n_k$  samples of  $X_i^{(k)}$ ,  $y_i^{(k)}$



# Leveraging Sparsity

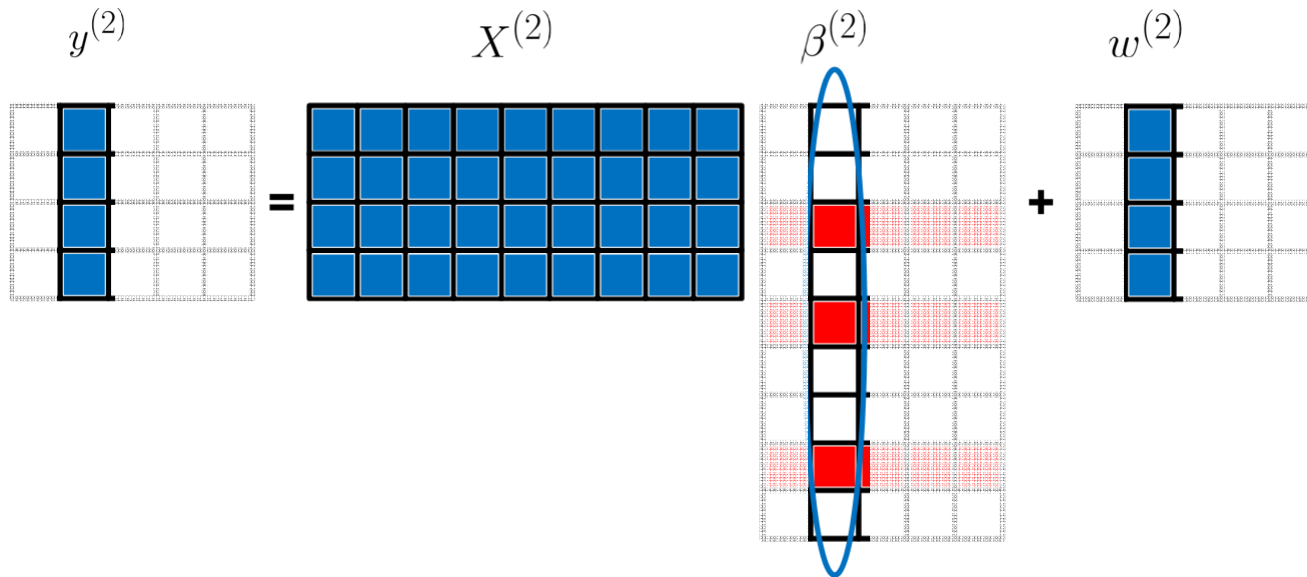


- Sparse structure: few features are relevant

- LASSO 
$$\min_{\beta^{(k)}} \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| y_i^{(k)} - X_i^{(k)} \beta^{(k)} \right\|_2^2 + \lambda_k \left\| \beta^{(k)} \right\|_1$$



# Leveraging Sparsity

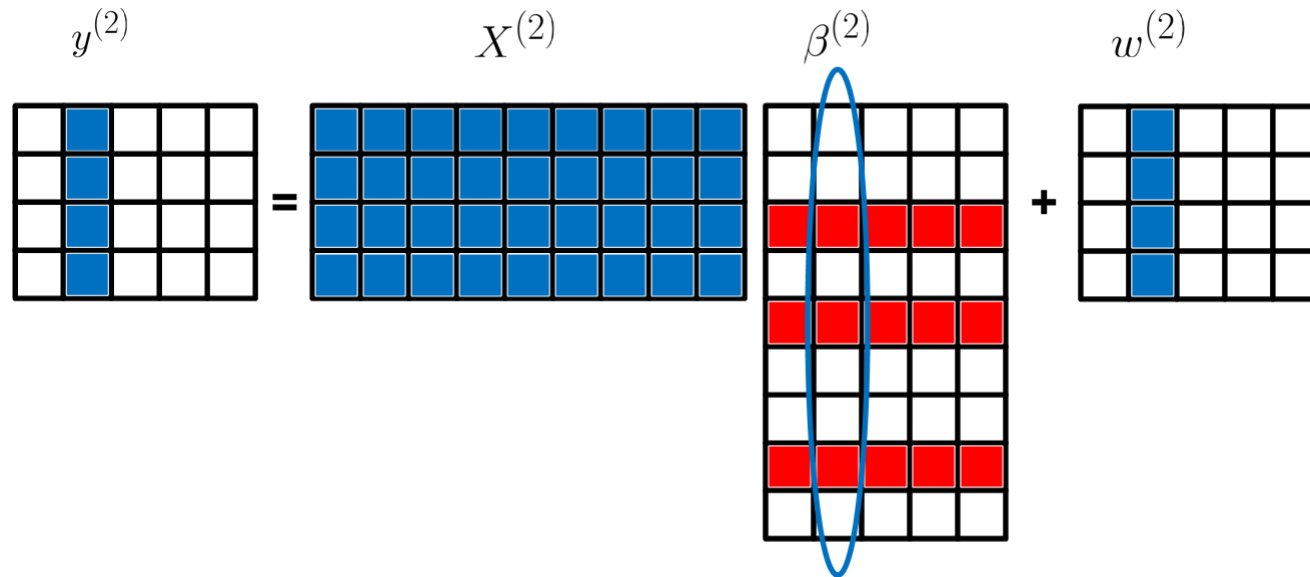


- Sparse structure: few features are relevant

- LASSO 
$$\min_{\beta^{(k)}} \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| y_i^{(k)} - X_i^{(k)} \beta^{(k)} \right\|_2^2 + \lambda_k \left\| \beta^{(k)} \right\|_1$$



# Leveraging block-sparsity



- Block-sparse structure: shared sparsity

- Group LASSO 
$$\min_{\beta} \sum_{k=1}^r \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| y_i^{(k)} - X_i^{(k)} \beta^{(k)} \right\|_2^2 + \lambda \|\beta\|_{1,\infty}$$

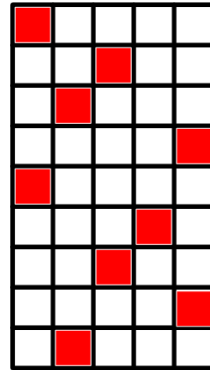
where  $\|\beta\|_{1,\infty} = \sum_j \max_k |\beta_j^{(k)}|$  (sum of maximum of rows)



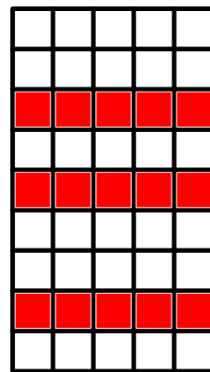
# Existing Methods Performance

---

- LASSO
  - Does not use shared sparsity



- Group LASSO
  - Does not model individual sparsity



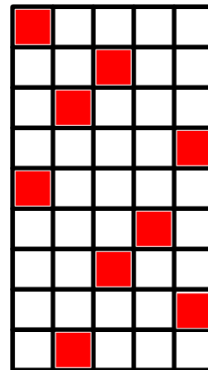




# Existing Methods Performance

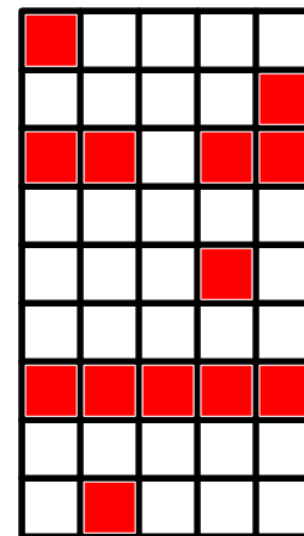
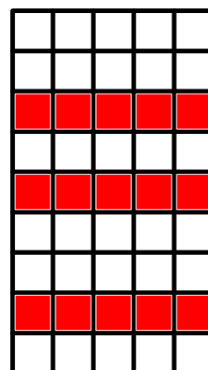
---

- LASSO
  - Does not use shared sparsity



How about more realistic cases?

- Group LASSO
  - Does not model individual sparsity

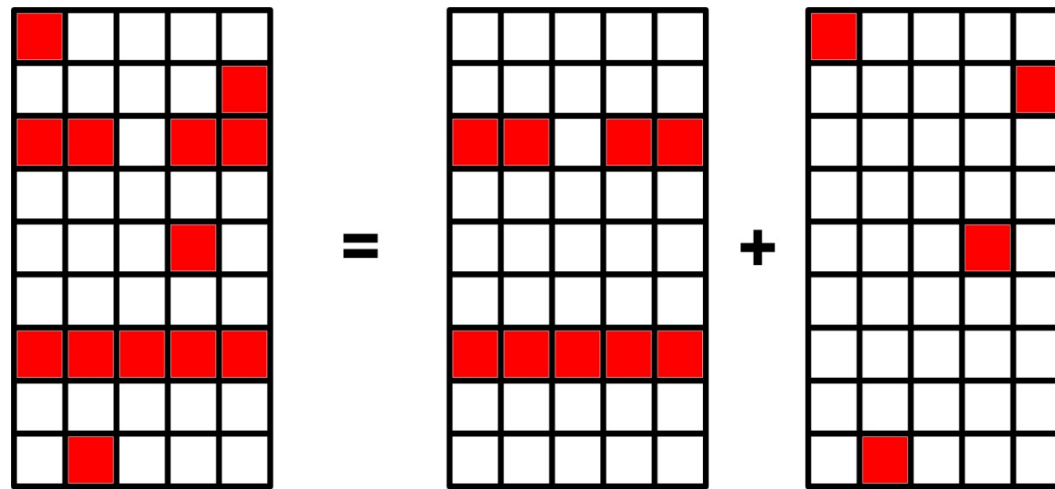


$\beta$



# Dirty Statistical Model

- Superposition of parameters with diff. structures



$$\beta = B + S$$

$$\|B\|_{1,\infty}$$

$$\|S\|_{1,1}$$

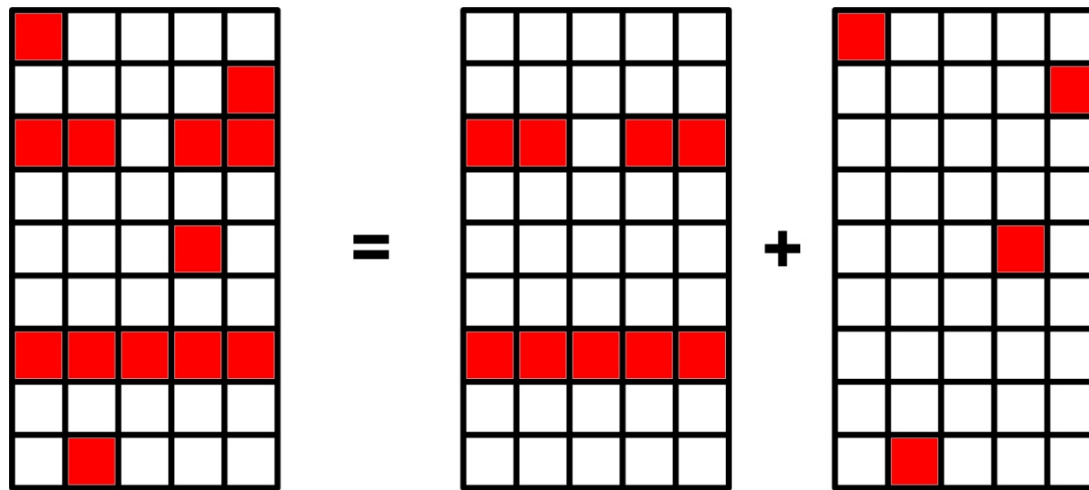
# shared features

# non-shared features



# Dirty Statistical Model

---



$$\beta = \mathbf{B} + \mathbf{S}$$

Algorithm:

$$\min_{B,S} \sum_{k=1}^r \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| y_i^{(k)} - X_i^{(k)} \left( B^{(k)} + S^{(k)} \right) \right\|_2^2 + \lambda_b \|B\|_{1,\infty} + \lambda_s \|S\|_{1,1}$$

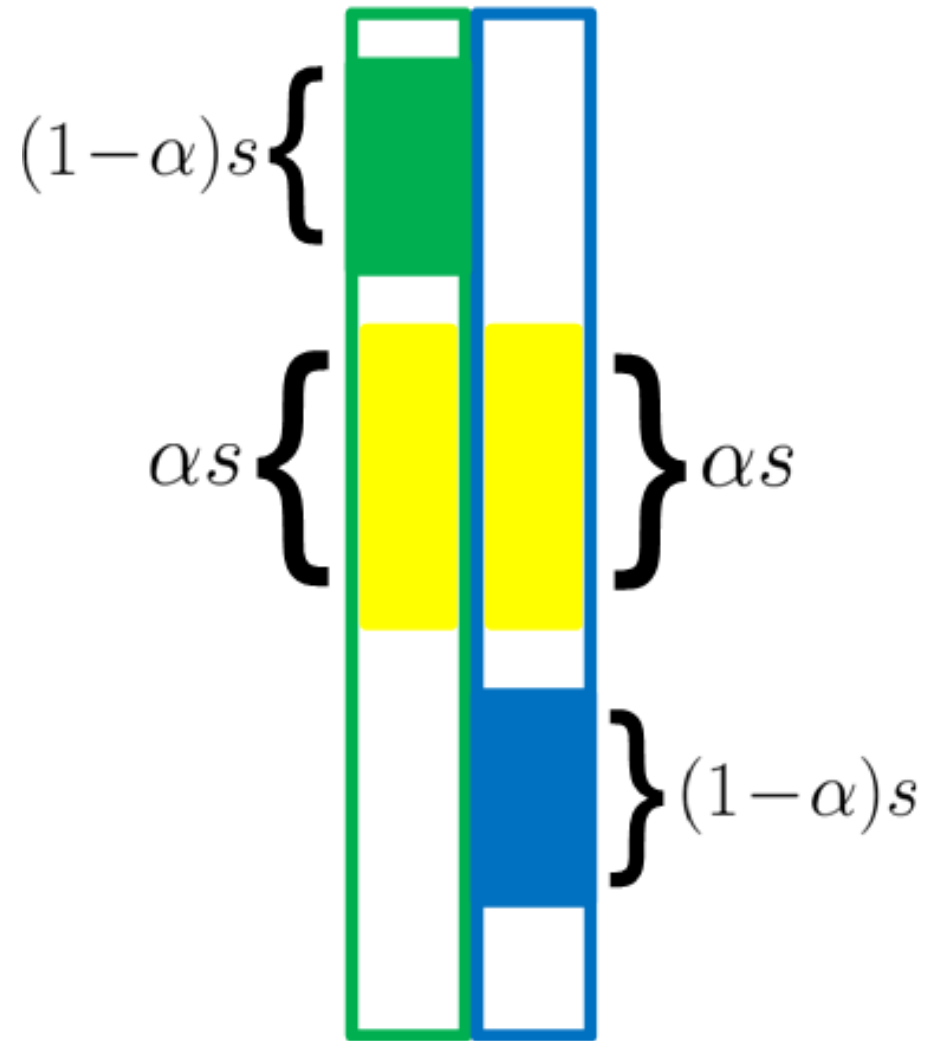
$$\text{output } \hat{\beta} = \hat{B} + \hat{S}$$



# Two Tasks Case

---

- Each task depends on “s” features
- $\alpha$ -portion of the features overlaps











# Phase Transition for Two Tasks

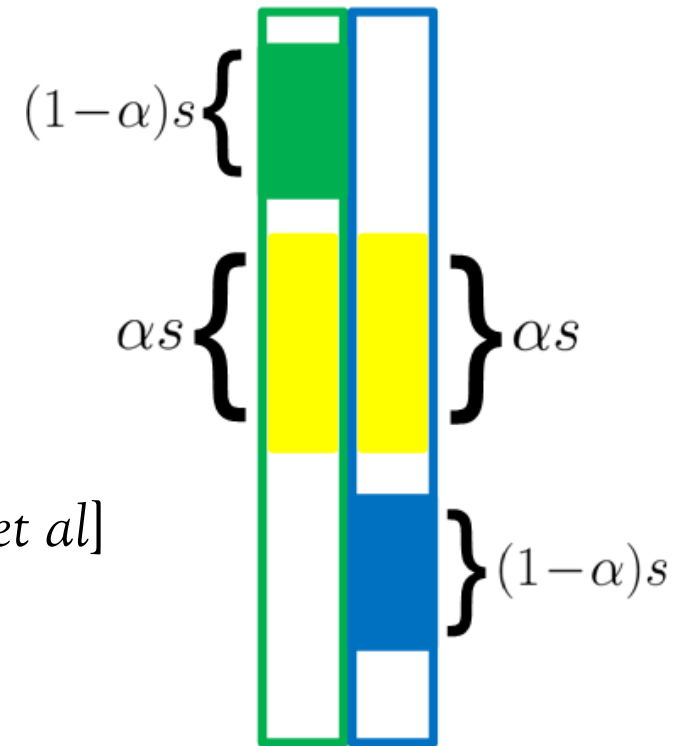
---

– LASSO:

$$\frac{n}{s \log(p)} \approx 2 \quad [\text{Wainwright}]$$

– Group LASSO:

$$\frac{n}{s \log(p)} \approx 4 - 3\alpha \quad [\text{Negahban et al}]$$







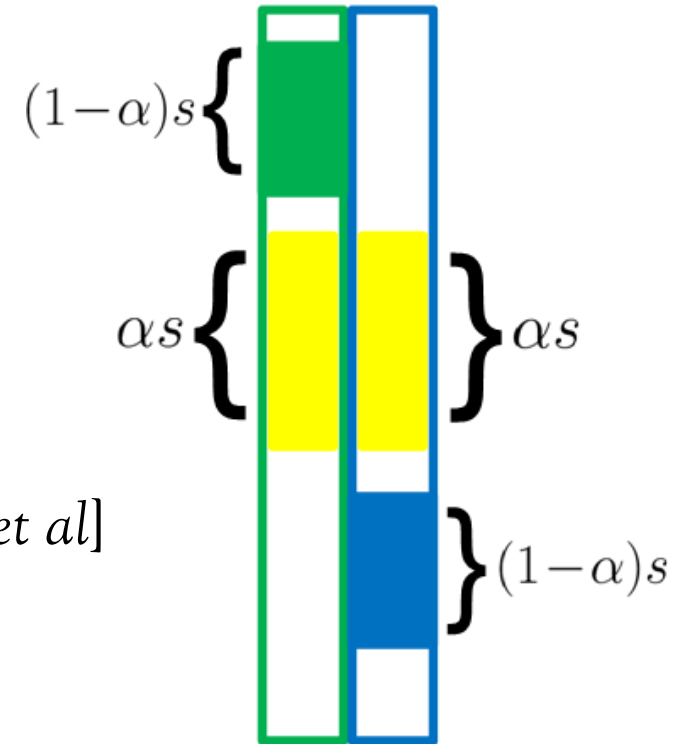
# Phase Transition for Two Tasks

– LASSO:

$$\frac{n}{s \log(p)} \approx 2 \quad [\text{Wainwright}]$$

– Group LASSO:

$$\frac{n}{s \log(p)} \approx 4 - 3\alpha \quad [\text{Negahban et al}]$$



## Consequences:

$\alpha < 2/3$  :: Lasso is better

$\alpha > 2/3$  :: Group-Lasso is better



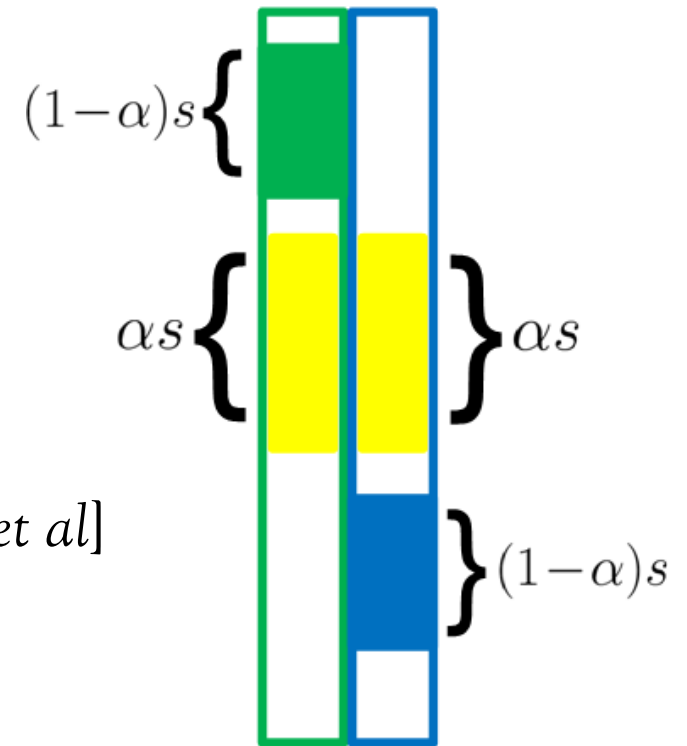
# Phase Transition for Two Tasks

– LASSO:

$$\frac{n}{s \log(p)} \approx 2 \quad [\text{Wainwright}]$$

– Group LASSO:

$$\frac{n}{s \log(p)} \approx 4 - 3\alpha \quad [\text{Negahban et al}]$$

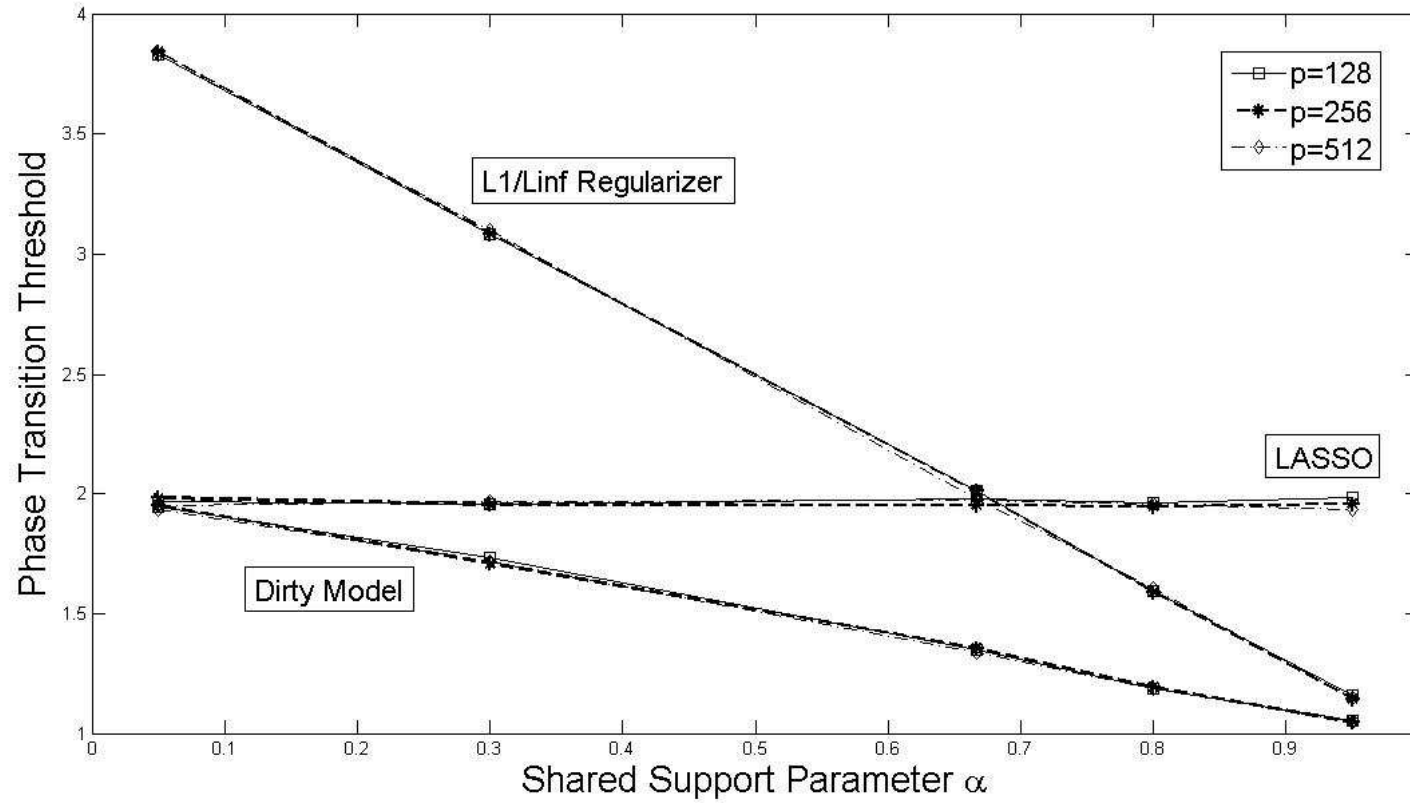


– Dirty Model:

$$\text{Theorem. } \frac{n}{s \log(p)} \approx 2 - \alpha$$



# Comparison



$$\text{Theorem. } \frac{n}{s \log(p)} \approx 2 - \alpha$$



# General $r$ -Task Case

---

- Model:  $y^{(k)} = X^{(k)}\beta^{(k)} + w^{(k)}$
- Algorithm:  $\min_{B,S} \sum_{k=1}^r \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| y_i^{(k)} - X_i^{(k)} \left( B^{(k)} + S^{(k)} \right) \right\|_2^2 + \lambda_b \|B\|_{1,\infty} + \lambda_s \|S\|_{1,1}$

**Theorem (Gaussian Design):** There exists a suitable choice of  $(\lambda_b, \lambda_s)$  for which our algorithm recovers the **exact sign support** of  $\beta$  with high probability provided that

$$n \geq K s \log(p)$$

where  $K$  only depends on  $X^{(k)}$  and  $r$ .



# General Dirty Models

---

$$\beta = \mathbf{B} + \mathbf{S}$$

- **Dirty Models: “Superposition of Simple Structures”**
  - Sparse + Low-Rank
    - Latent Variables, Graph Clustering, PCA with corruptions, etc
  - Block-Sparse + Low-Rank
    - Collaborative Filtering, PCA with Outliers, etc



# Summary

---

- Multi-task learning is challenging when there is partial overlap across tasks
  - relevant structure is neither sparsity nor block-sparsity
- A superposition of simple structures, i.e., dirty model surprisingly useful for modeling such “dirty” structure.
- For the multi-task learning problem, dirty model outperforms solo-structured lasso and group-lasso.