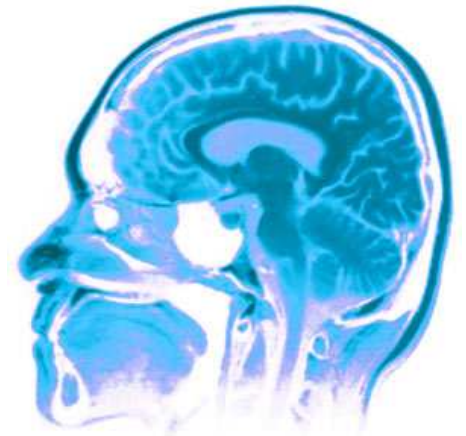




CPS C340



Sparse regularization and feature selection



Nando de Freitas

October, 2012

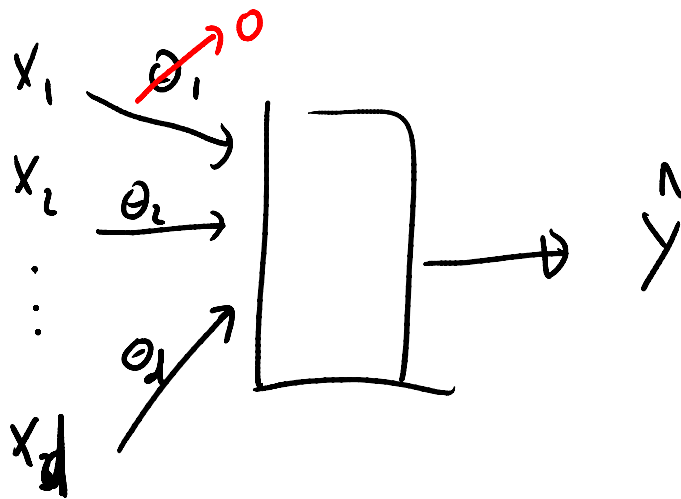
University of British Columbia

Outline of the lecture

This lecture introduces one of the most popular modern techniques for regression and variable selection: Sparse regularization. The goal is for you to:

- ❑ Learn regularization with the L1 norm.
- ❑ Understand how regularizers can be used to automatically select input features.
- ❑ Understand the concept of sub-gradients as part of the derivation of an algorithm.
- ❑ Understand the pseudo-code of a coordinate descent algorithm to estimate the parameters of linear sparse models.

Selecting features for prediction



$$\hat{y} = x_1 \theta_1 + x_2 \theta_2 + \dots + x_d \theta_d$$

A red arrow points from the θ_1 term in the equation to a red circle containing the number 0.

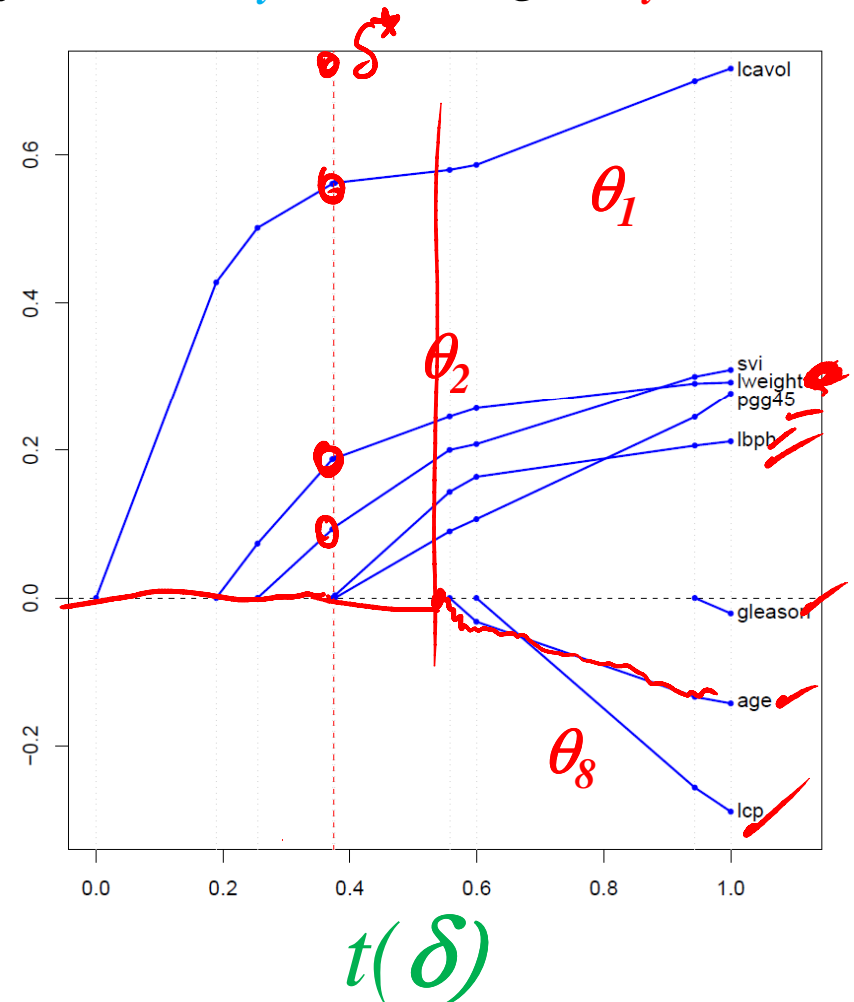
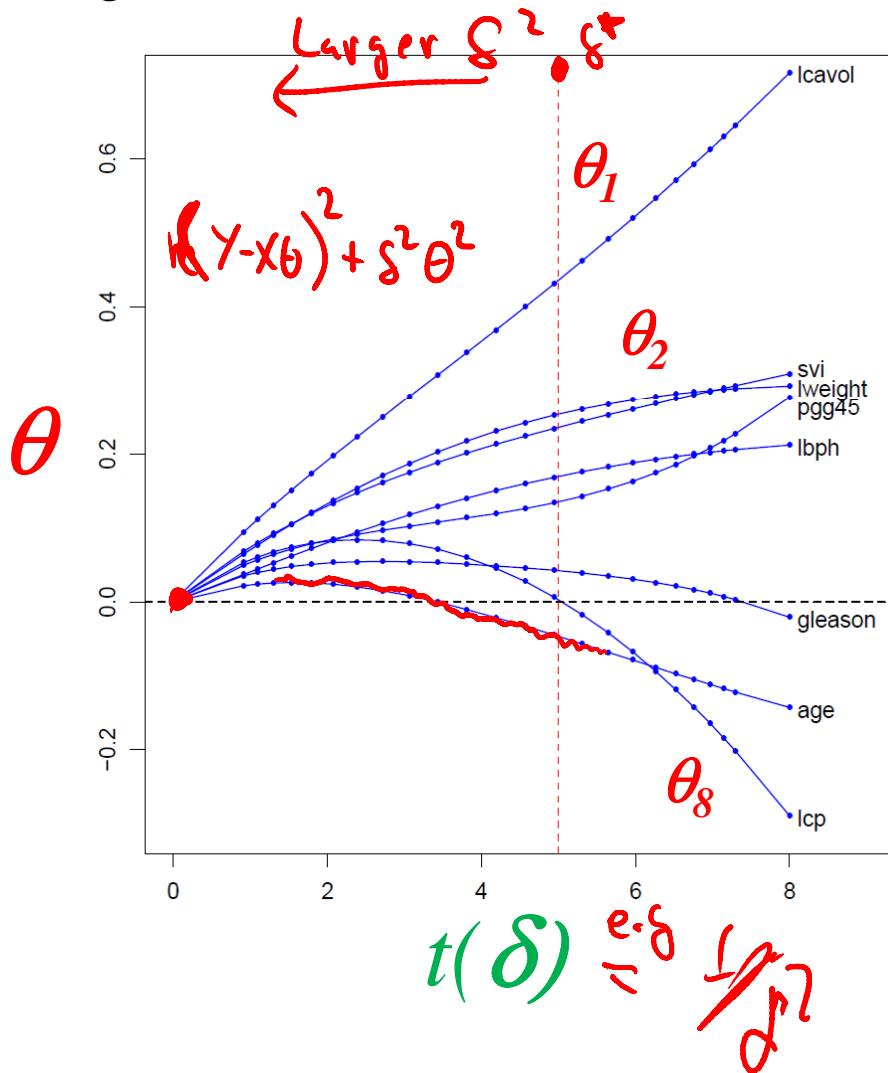
x_1 is expensive

x_1 does not contribute to good predictions \hat{y}

Then we want $\theta_1 \rightarrow 0$

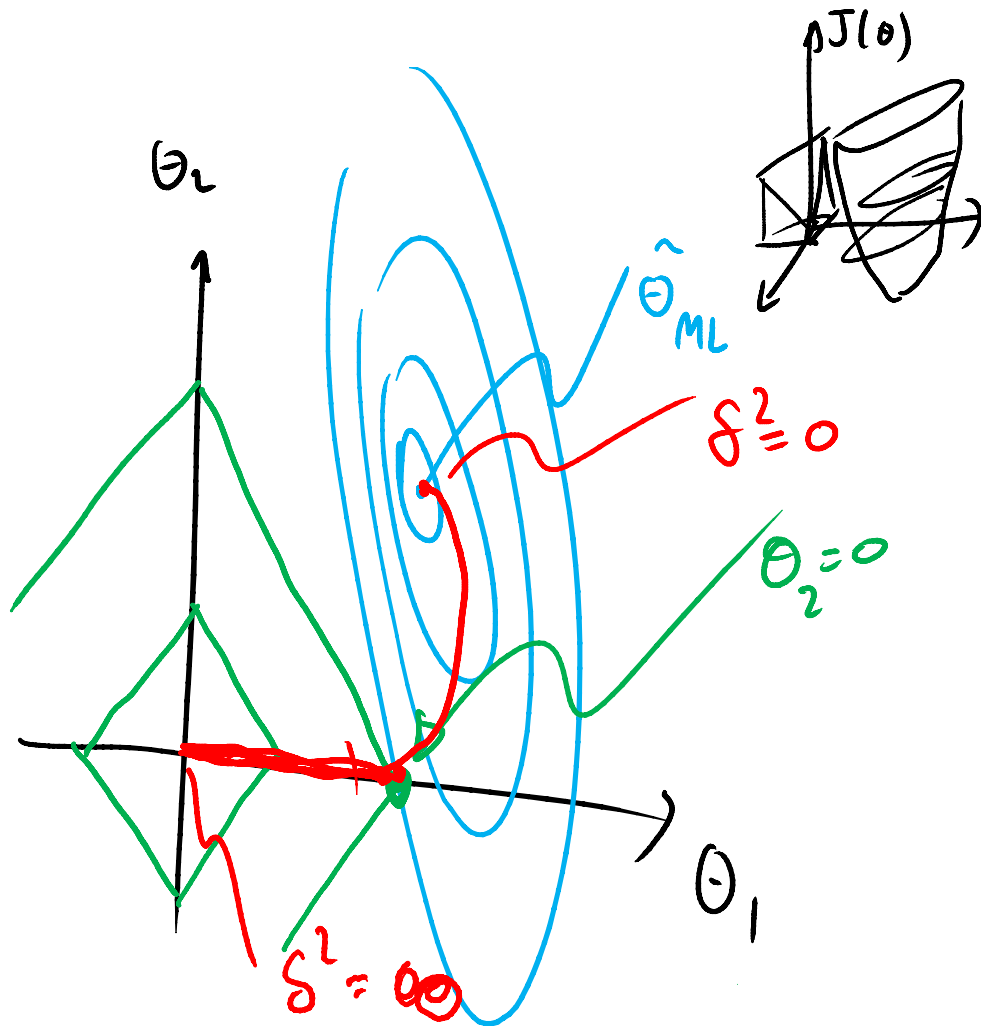
Selecting features for prediction

As δ increases, $t(\delta)$ decreases and each θ_i goes to zero, but too slowly for ridge. Lasso will ensure that irrelevant features x_i have weight $\theta_i = 0$.



The Lasso: least absolute selection and shrinkage operator

$$J(\theta) = \underbrace{(Y - X\theta)^T (Y - X\theta)}_{\text{Least Squares}} + \delta^2 \underbrace{\sum_{j=1}^d |\theta_j|}_{L_1 \text{ Norm}}$$

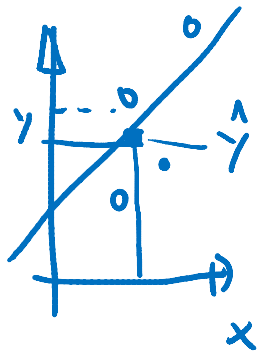


in 2D $\theta = (\theta_1, \theta_2)$

$$|\theta_1| + |\theta_2| = \text{const}$$

$$\begin{aligned} \theta_1 + \theta_2 &= \text{const} \\ \theta_1 - \theta_2 &= \text{const} \\ -\theta_1 - \theta_2 &= \text{const} \\ -\theta_1 + \theta_2 &= \text{const} \end{aligned}$$

Differentiating the objective function



$$J(\theta) = (y - \overset{n \times d}{\mathbf{X}}\theta)^T (y - \overset{d \times 1}{\mathbf{X}}\theta) + \delta^2 \sum_{j'=1}^d |\theta_{j'}| \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$= \sum_{i=1}^n (y_i - \underbrace{x_i^T \theta}_{\hat{y}_i})^2 + \delta^2 \sum_{j'=1}^d |\theta_{j'}|$$

$$= \left[\sum_{i=1}^n (y_i - x_{ij} \theta_j - x_{i-j}^T \theta_{-j})^2 \right] + \delta^2 \sum_{j'=1}^d |\theta_{j'}|$$

$$\underline{\theta} = (\theta_1, \theta_2, \theta_3)$$

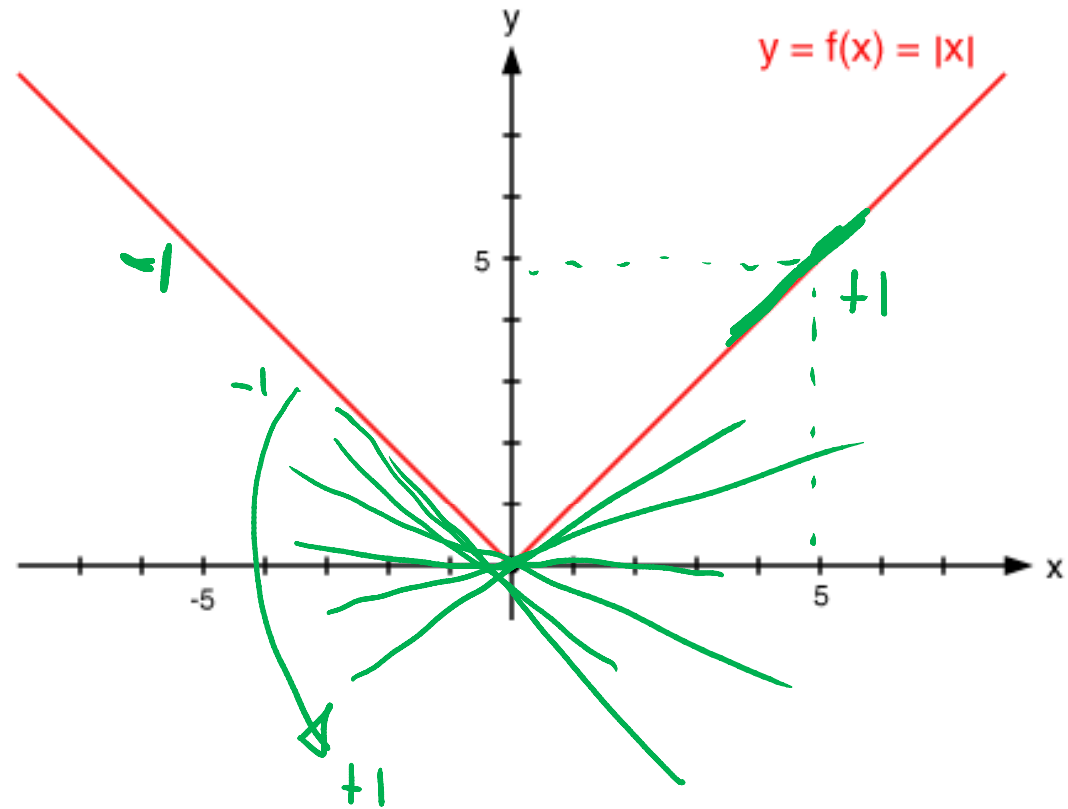
$$\underline{\theta}_2 = (\theta_1, \theta_3)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^n 2 (y_i - \overline{x_{ij} \theta_j} - \overline{x_{i-j}^T \theta_{-j}}) (\overline{-x_{ij}}) + \delta^2 \frac{\partial}{\partial \theta_j} (|\theta_1| + \dots + |\theta_d|)$$

$$= 2 \sum_{i=1}^n x_{ij}^2 \theta_j - 2 \sum_{i=1}^n (y_i - x_{i-j}^T \theta_{-j}) x_{ij} + \delta^2 \frac{\partial}{\partial \theta_j} |\theta_j|$$

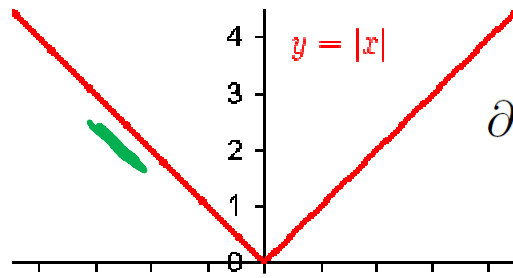
Differentiating the objective function

Subdifferentials



$$f(x) = |x|$$

$$\partial f(x) = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ \{+1\} & \text{if } x > 0 \end{cases}$$



$$\partial f(x) = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ \{+1\} & \text{if } x > 0 \end{cases}$$

$$\partial_{\theta_j} J(\boldsymbol{\theta}) = a_j \theta_j - c_j + \delta^2 \partial_{\theta_j} |\theta_j|$$

$$= \begin{cases} \{a_j \theta_j - c_j - \delta^2\} & \text{if } \theta_j < 0 \leftarrow \frac{\partial |\theta_j|}{\partial \theta_j} = -1 \\ [-c_j - \delta^2, -c_j + \delta^2] & \text{if } \theta_j = 0 \\ \{a_j \theta_j - c_j + \delta^2\} & \text{if } \theta_j > 0 \leftarrow \frac{\partial |\theta_j|}{\partial \theta_j} = +1 \end{cases}$$

Hence, the estimate of the j -th parameter, **given the other parameters**, is

$$\hat{\theta}_j = \begin{cases} (c_j + \delta^2)/a_j & \text{if } c_j < -\delta^2 \text{ when } \theta_j < 0 \\ 0 & \text{if } c_j \in [-\delta^2, \delta^2] \\ (c_j - \delta^2)/a_j & \text{if } c_j > \delta^2 \end{cases}$$

$$a_j \theta_j - c_j - \delta^2 = 0$$

$$a_j \theta_j = c_j + \delta^2$$

$$\hat{\theta}_j = \frac{c_j + \delta^2}{a_j}$$

Coordinate descent algorithm for sparse prediction

1. Initialize $\underline{\Theta}$, e.g. $\underline{\Theta} = (\underline{X}^T \underline{X} + \delta^2 \mathbf{I})^{-1} \underline{X}^T \underline{y}$ (ridge)

2. REPEAT UNTIL CONVERGED

% \underline{X} is n by d

% \underline{y} is n by 1

3. FOR $j=1, 2, \dots, d$ DO

4. $a_j = 2 \sum_{i=1}^n x_{ij}^2$ ✓

5. $c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \underline{x}_i^T \underline{\Theta} + x_{ij} \Theta_j)$ ✓

6. IF $c_j < -\delta^2$

7. $\Theta_j = (c_j + \delta^2) / a_j$

8. ELSEIF $c_j > \delta^2$

9. $\Theta_j = (c_j - \delta^2) / a_j$

10. ELSE

11. $\Theta_j = 0$

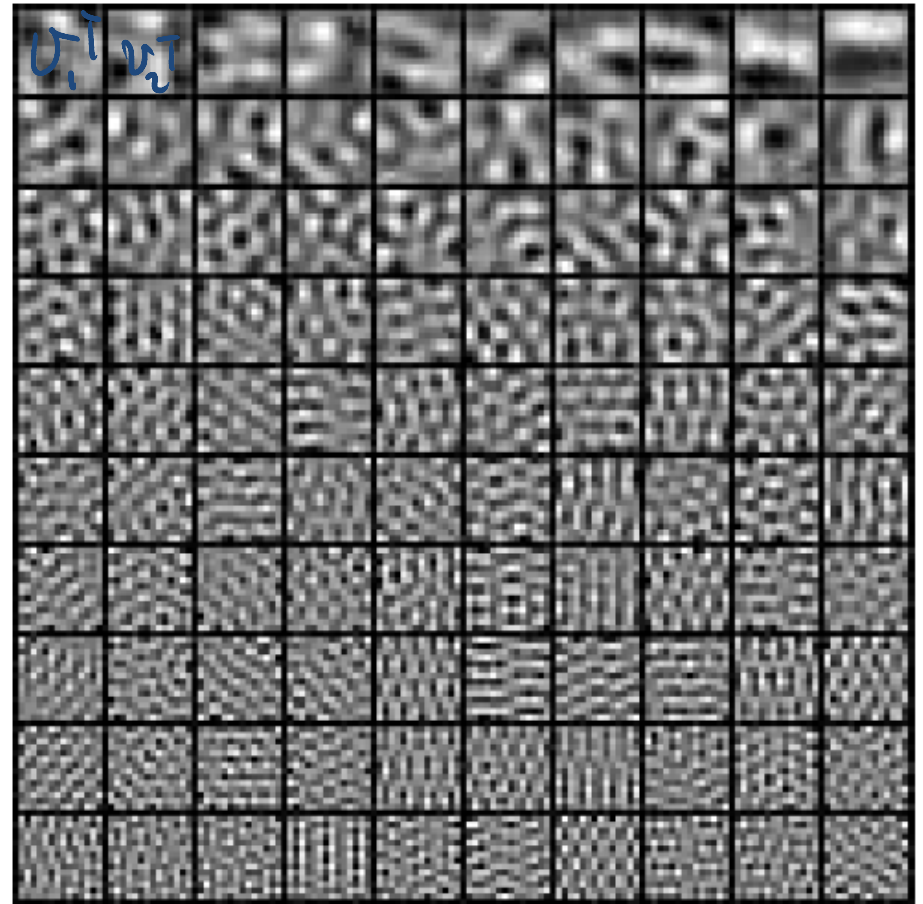
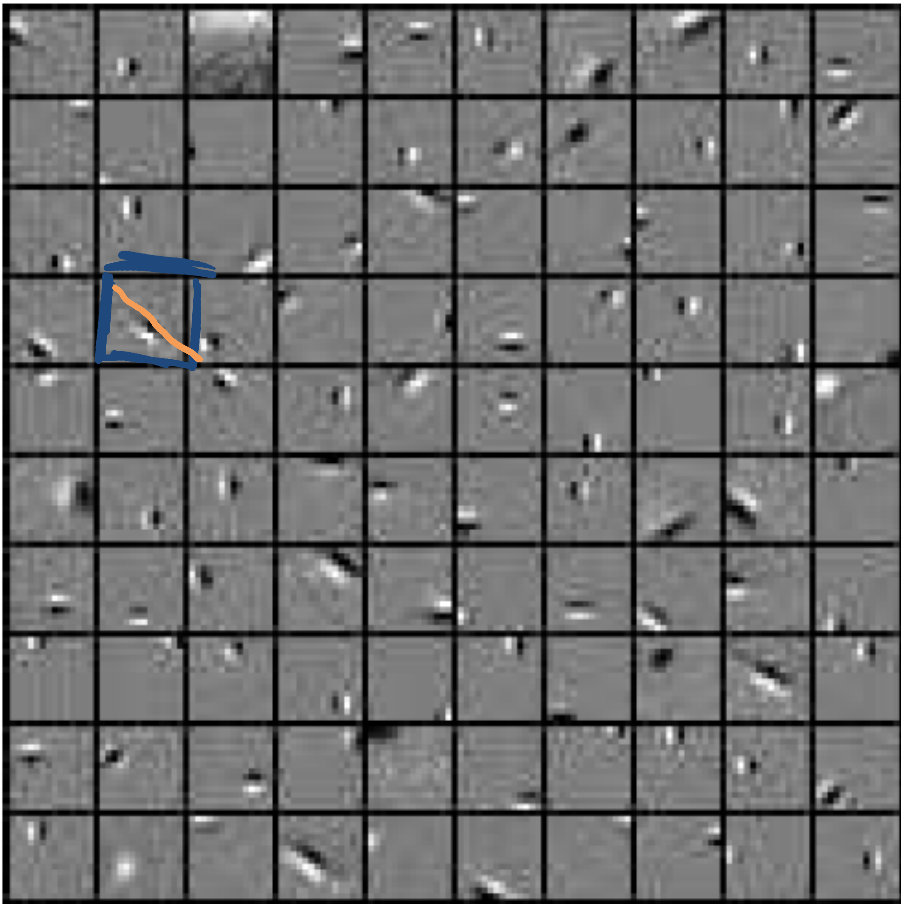
$$\underline{X}^* \hat{\underline{\Theta}} = \hat{\underline{y}}$$

The effect of L1 regularization on PCA

$$\mathbf{B}^*, \mathbf{C}^* = \underset{\mathbf{B}, \mathbf{C}}{\operatorname{arg\,min}} \|\mathbf{X} - \mathbf{BC}\|_2^2 + \lambda \|\mathbf{C}\|_1$$

s.t. $\|\mathbf{b}_j\|_2^2 = 1, \forall j.$

$\mathbf{B} = \mathbf{U}\Sigma$
 $\mathbf{C} = \mathbf{V}^T$



Next lecture

In the next lecture, we go back to probability so as to get enough background to understand classification.