

# Lasso Regression

---

Nipun Batra

February 5, 2020

IIT Gandhinagar

- LASSO  $\rightarrow$  Least absolute shrinkage and selection operator

# Lasso Regression

- LASSO  $\rightarrow$  Least absolute shrinkage and selection operator
- Popular as it leads to a sparse solution.

# Constructing the Objective Function

- Find a  $\theta_{opt}$  such that

$$\theta_{opt} = \arg \min_{\theta} (Y - X\theta)^T (Y - X\theta) : \|\theta\|_1 < s \quad (1)$$

# Constructing the Objective Function

- Find a  $\theta_{opt}$  such that

$$\theta_{opt} = \arg \min_{\theta} (Y - X\theta)^T (Y - X\theta) : \|\theta\|_1 < s \quad (1)$$

- Using KKT conditions

$$\theta_{opt} = \arg \min_{\theta} \underbrace{(Y - X\theta)^T (Y - X\theta) + \delta^2 \|\theta\|_1}_{\text{convex function}} \quad (2)$$

- Since  $|\theta|$  is not differentiable, we cannot solve,

$$\frac{\partial(Y - X\theta)^T(Y - X\theta) + \delta^2 \|\theta\|_1}{\partial\theta} = 0 \quad (3)$$

## Solving the Objective

- Since  $|\theta|$  is not differentiable, we cannot solve,

$$\frac{\partial(Y - X\theta)^T(Y - X\theta) + \delta^2 \|\theta\|_1}{\partial\theta} = 0 \quad (3)$$

- How to Solve? Use Coordinate descent!

# Sample Dataset

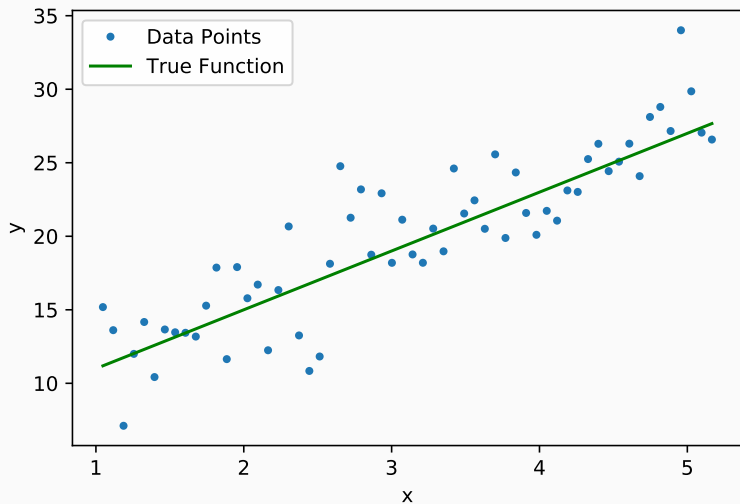


Figure 1:  $y = 4x + 7$



# Geometric Interpretation

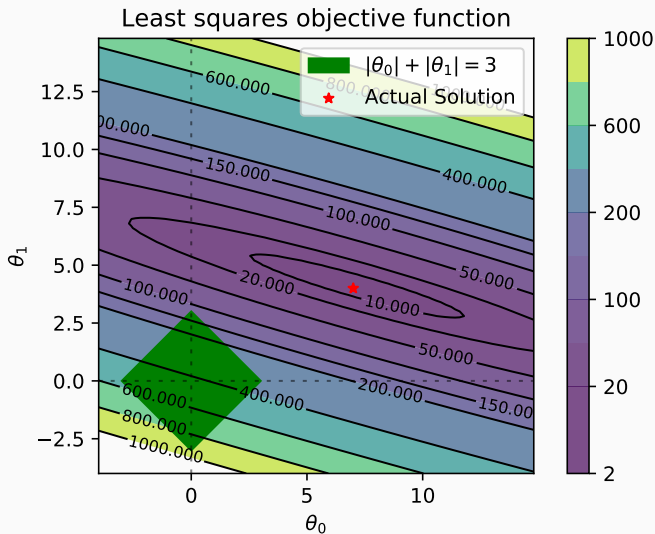
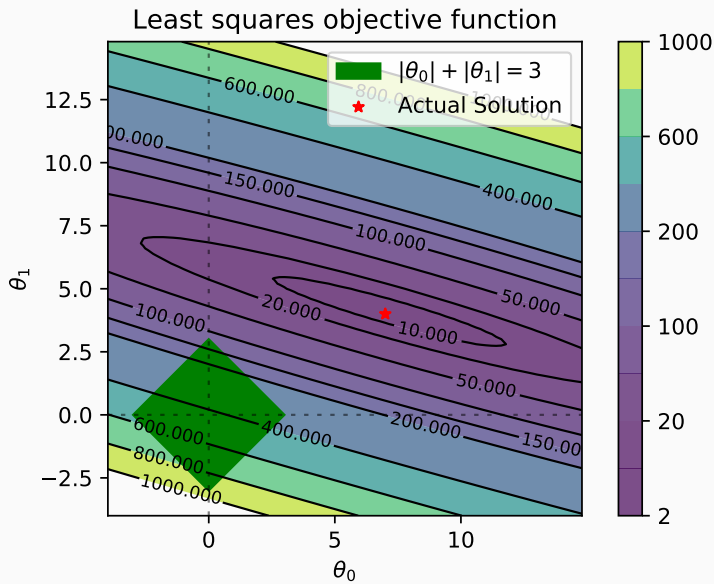


Figure 2: Lasso regression

# Effect of $\mu$ - Regularization of Parameters



# Effect of $\mu$ - Regularization of Parameters

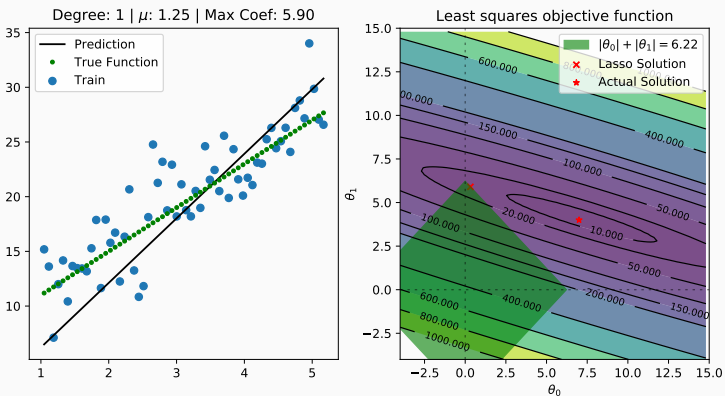


Figure 4:  $\mu = 1.25$   
(on the *Sample Dataset*)

# Effect of $\mu$ - Regularization of Parameters

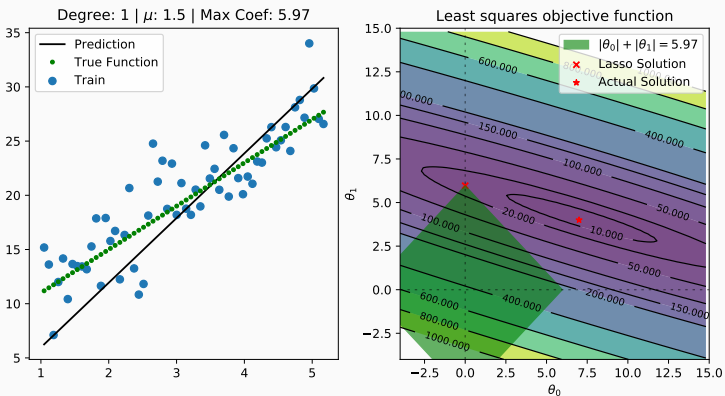


Figure 5:  $\mu = 1.5$   
(on the *Sample Dataset*)

# Effect of $\mu$ - Regularization of Parameters

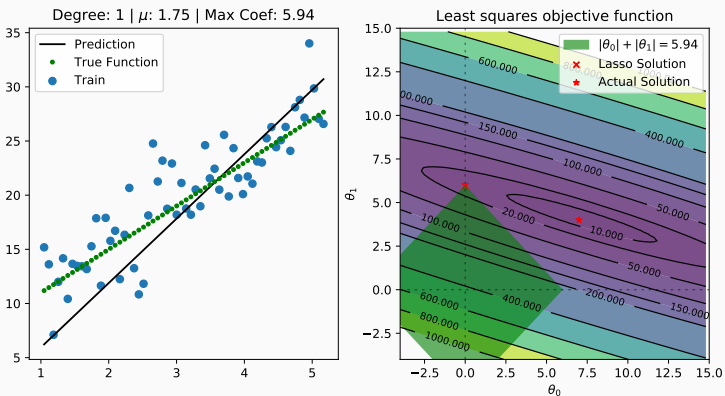


Figure 6:  $\mu = 1.75$   
(on the *Sample Dataset*)

# Effect of $\mu$ - Regularization of Parameters

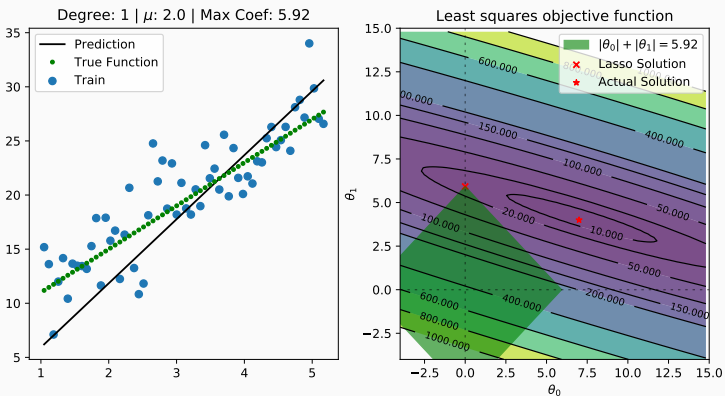
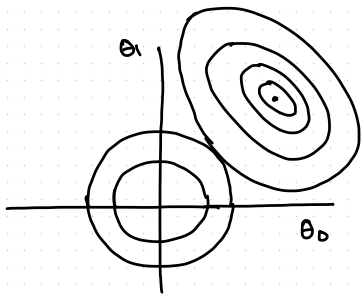


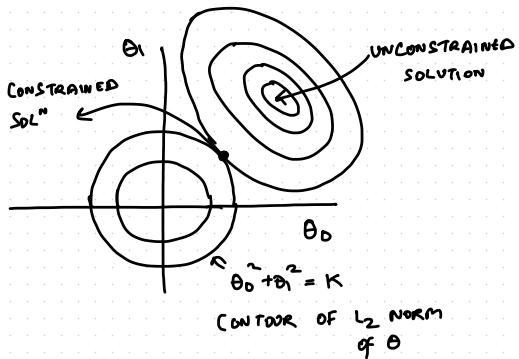
Figure 7:  $\mu = 2.0$   
(on the Sample Dataset)

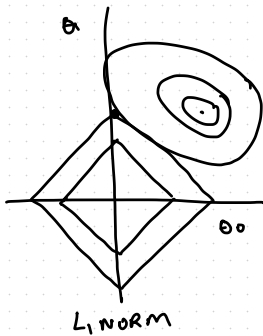
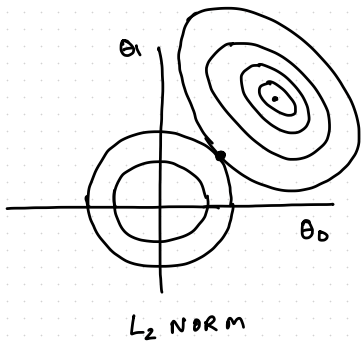
# WHY LASSO GIVES SPARSITY

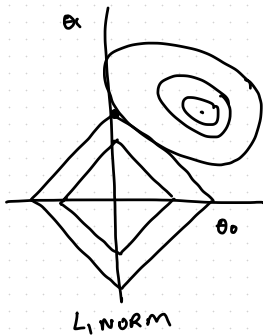
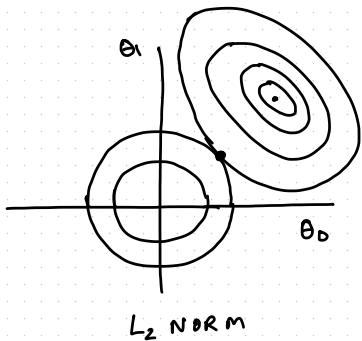
- ① GEOMETRIC INTERPRETATION
- ② G.D. BASED INTERPRETATION



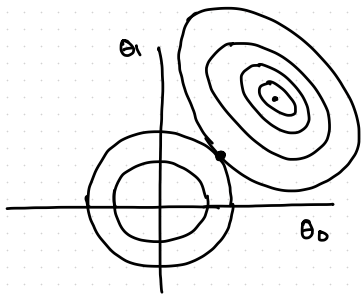




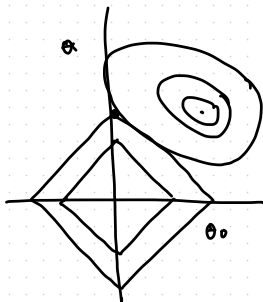




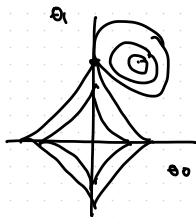
$L_p$  NORM  
( $0 < p < 1$ )



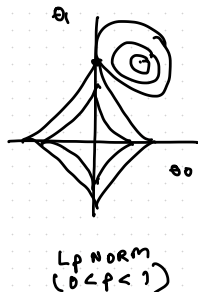
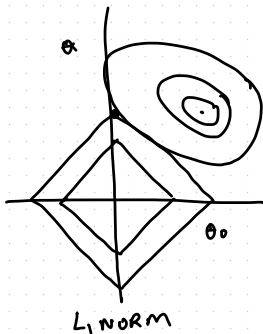
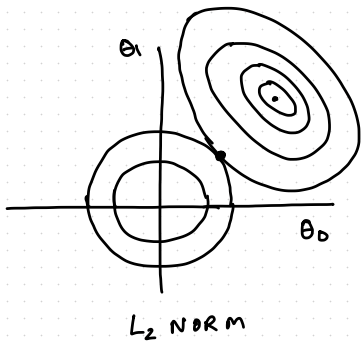
$L_2$  NORM



$L_1$  NORM



$L_p$  NORM  
( $0 < p < 1$ )

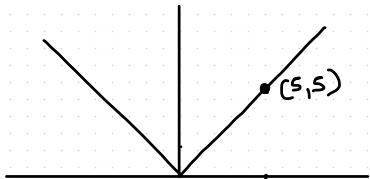


SPARSITY  $\longrightarrow$   
 PROB. OF INTERSECTING AXIS  $\longrightarrow$   
 DIFFICULTY OF SOLVING  $\longrightarrow$

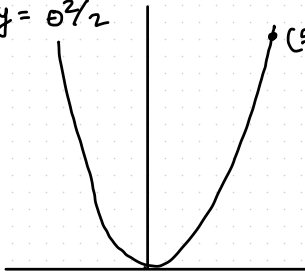
$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

$$y = \theta^2/2$$

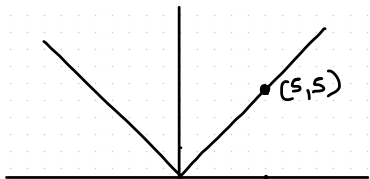
$$y = |x| \quad (\text{FOR NOW ASSUME } x > 0)$$



$$y = x^2/2$$

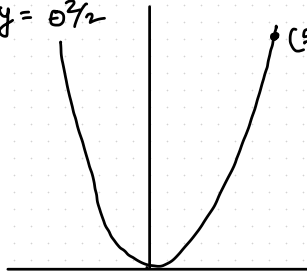


$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$



$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{ASSUME } \theta > 0)$$

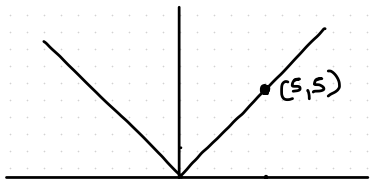
$$y = \theta^2/2$$



$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$



$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

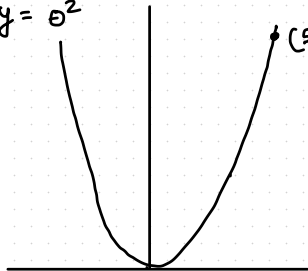


$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

$$\theta_0^1 = \theta_0^0 - 0.5 * 1 = 4.5$$

$\text{LET } \alpha = 0.5$

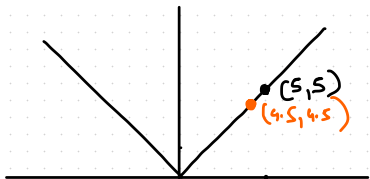
$$y = \theta^2$$



$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\theta_0^1 = \theta_0^0 - 0.5 * 5 = 2.5$$

$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

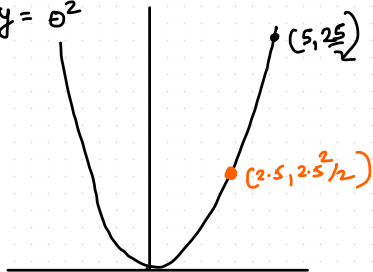


$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

$$\theta'_0 = \theta_0^0 - 0.5 * 1 = 4.5$$

LET  $\alpha = 0.5$

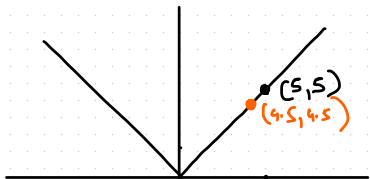
$$y = \theta^2$$



$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\theta'_0 = \theta_0^0 - 0.5 * 5 = 2.5$$

$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

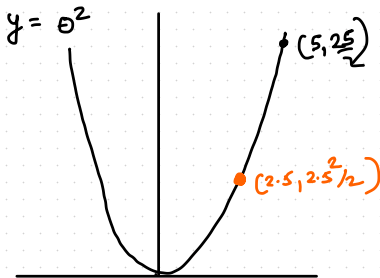


$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

$$\boxed{\text{LET } \alpha = 0.5}$$

$$\theta_0^1 = \theta_0^0 - 0.5 \times 1 = 4.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 \times 1 = 4.0$$

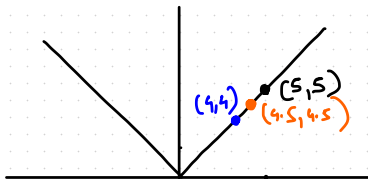


$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\theta_0^1 = \theta_0^0 - 0.5 \times 5 = 2.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 \times 2.5 = 1.25$$

$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

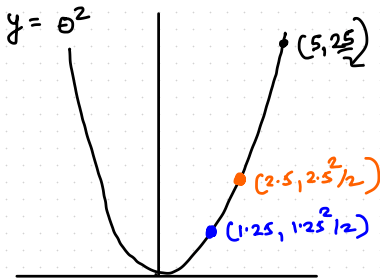


$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

LET  $\alpha = 0.5$

$$\theta_0^1 = \theta_0^0 - 0.5 * 1 = 4.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 * 1 = 4.0$$

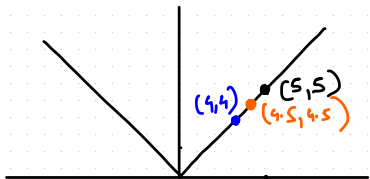


$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\theta_0^1 = \theta_0^0 - 0.5 * 5 = 2.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 * 2.5 = 1.25$$

$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$



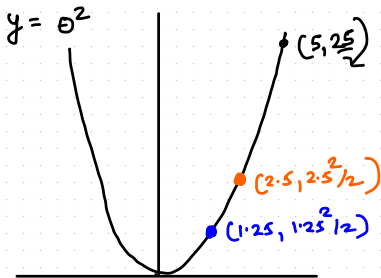
$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

LET  $\alpha = 0.5$

$$\theta_0^1 = \theta_0^0 - 0.5 \times 1 = 4.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 \times 1 = 4.0$$

$\theta_0^t = \theta_0^{t-1} - 0.5$



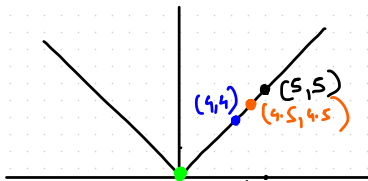
$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\theta_0^1 = \theta_0^0 - 0.5 \times 5 = 2.5$$

$$\theta_0^2 = \theta_0^1 - 0.5 \times 2.5 = 1.25$$

$\theta_0^t = \theta_0^{t-1} - 0.5 \theta_0^{t-1} = 0.5 \theta_0^{t-1}$

$$y = |\theta| \quad (\text{FOR NOW ASSUME } \theta > 0)$$

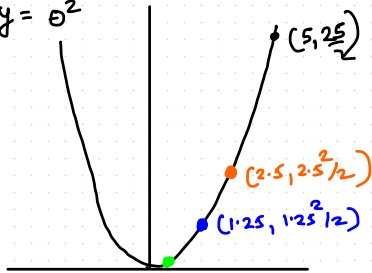


$$\frac{\partial y}{\partial \theta} = 1 \quad (\text{Assume } \theta > 0)$$

$$\theta_0^{10} = 0$$

$\text{LET } \alpha = 0.5$

$$y = \theta^2$$



$$\frac{\partial y}{\partial \theta} = \frac{2\theta}{2} = \theta$$

$$\begin{aligned} \theta_0^{10} &= 5 * (0.5)^{10} \\ &= 0.0048 \end{aligned}$$

(Approaching 0  
but not exactly  
zero)

# Regularization path of lasso regression

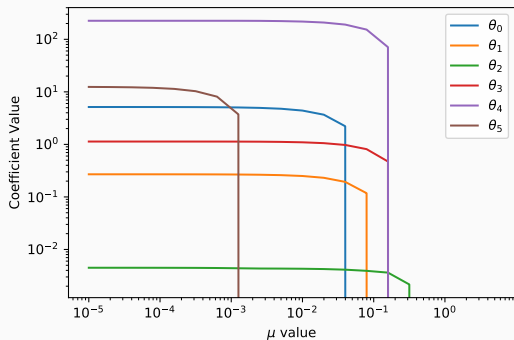


Figure 8: Regularization path of  $\theta_i$

- LASSO inherently does feature selection!



## LASSO and feature selection

- LASSO inherently does feature selection!
- Sets coefficients of “less important” features to zero.

# LASSO and feature selection

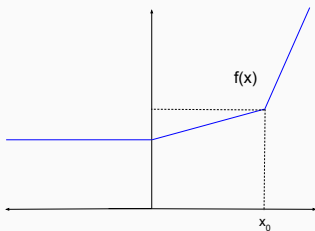
- LASSO inherently does feature selection!
- Sets coefficients of “less important” features to zero.
- Sparse and memory efficient and often more interpretable models.

# Subgradient

- Generalizes gradient to convex but non-differentiable problems
- Examples:
  - $f(x) = |x|$

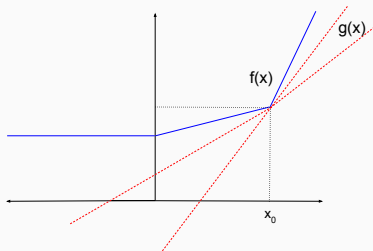
## Task at hand

- TASK: find derivative of  $f(x)$  at  $x = x_0$



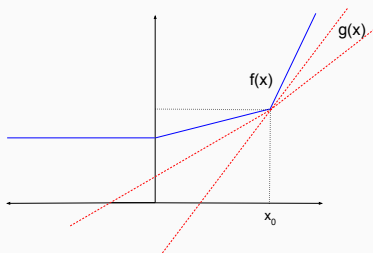
# Solution

- Construct a differentiable  $g(x)$ 
  - Intersecting  $f(x)$  at  $x = x_0$
  - Below or on  $f(x)$  for all  $x$



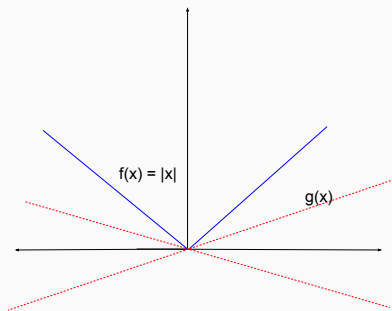
# Solution

- Compute slope of  $g(x)$  at  $x = x_0$



## Another Example: $f(x) = |x|$

- Subgradient of  $f(x)$  belongs to  $[-1, 1]$



- Another optimisation method (akin to gradient descent)



# Coordinate Descent

- Another optimisation method (akin to gradient descent)
- Objective:  $\text{Min}_{\theta} f(\theta)$

# Coordinate Descent

- Another optimisation method (akin to gradient descent)
- Objective:  $\text{Min}_{\theta} f(\theta)$
- Key idea: Sometimes difficult to find minimum for all coordinates

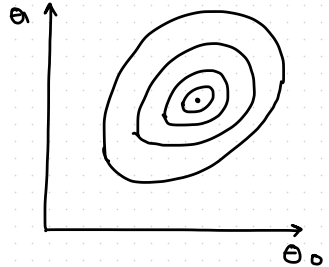
# Coordinate Descent

- Another optimisation method (akin to gradient descent)
- Objective:  $\text{Min}_{\theta} f(\theta)$
- Key idea: Sometimes difficult to find minimum for all coordinates
- ..., but, easy for each coordinate

# Coordinate Descent

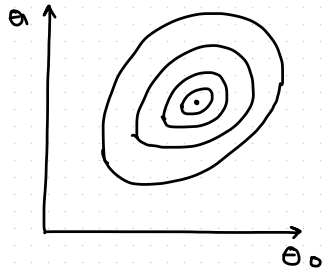
- Another optimisation method (akin to gradient descent)
- Objective:  $\text{Min}_{\theta} f(\theta)$
- Key idea: Sometimes difficult to find minimum for all coordinates
- ..., but, easy for each coordinate
- turns into a  $1D$  optimisation problem

# COORDINATE DESCENT ALGORITHM



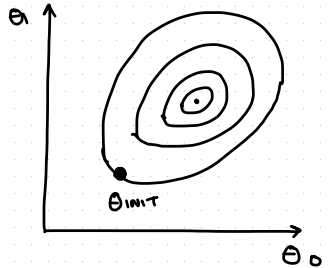
# COORDINATE DESCENT ALGORITHM

GOAL:  $\min_{\theta} f(\theta)$



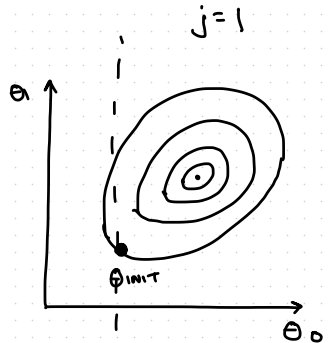
# COORDINATE DESCENT ALGORITHM

1) INIT  $\theta$



# COORDINATE DESCENT ALGORITHM

- 1) INIT  $\theta$
- 2) WHILE NOT CONVERGED
  - 2.1) PICK COORDINATE 'j'





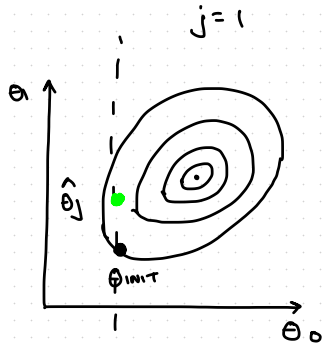
# COORDINATE DESCENT ALGORITHM

1) INIT  $\theta$

2) WHILE NOT CONVERGED

2.1) PICK COORDINATE 'j'

2.2)  $\hat{\theta}_j = \min_{\phi} f(\theta_0, \phi)$



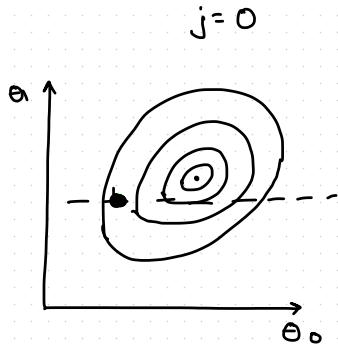
# COORDINATE DESCENT ALGORITHM

1) INIT  $\theta$

2) WHILE NOT CONVERGED

✓ 2.1) PICK COORDINATE 'j'

$$2.2) \hat{\theta}_j = \min_{\phi} f(\phi, \theta_1)$$



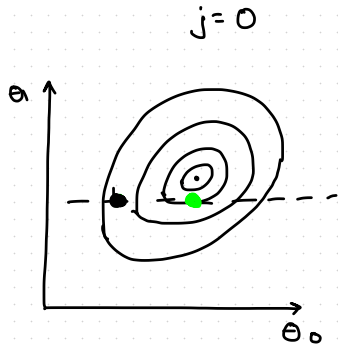
# COORDINATE DESCENT ALGORITHM

1) INIT  $\theta$

2) WHILE NOT CONVERGED

2.1) PICK COORDINATE 'j'

✓ 2.2)  $\hat{\theta}_j = \min_{\phi} f(\theta_0, \phi)$



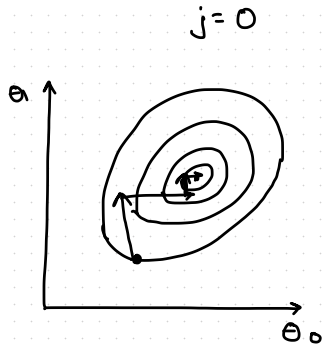
# COORDINATE DESCENT ALGORITHM

1) INIT  $\theta$

2) WHILE NOT CONVERGED

2.1) PICK COORDINATE 'j'

2.2)  $\hat{\theta}_j = \min_{\phi} f(\theta_0, \phi)$



- Picking next coordinate:

- Picking next coordinate:

# Coordinate Descent

- Picking next coordinate: random, round-robin
- No step-size to choose!

# Coordinate Descent

- Picking next coordinate: random, round-robin
- No step-size to choose!
- Converges for Lasso objective



## Coordinate Descent : Example

Learn  $y = \theta_0 + \theta_1 x$  on following dataset, using coordinate descent where initially  $(\theta_0, \theta_1) = (2, 3)$  for 2 iterations.

x	y
1	1
2	2
3	3

## Coordinate Descent : Example

Our predictor,  $\hat{y} = \theta_0 + \theta_1 x$

Error for  $i^{\text{th}}$  datapoint,  $\epsilon_i = y_i - \hat{y}_i$

$$\epsilon_1 = 1 - \theta_0 - \theta_1$$

$$\epsilon_2 = 2 - \theta_0 - 2\theta_1$$

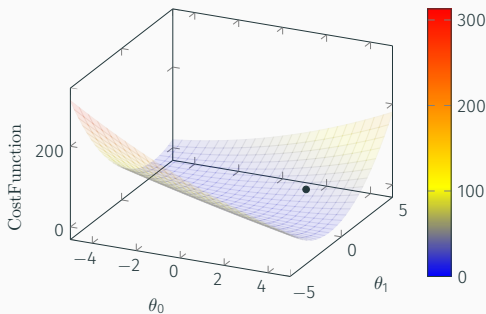
$$\epsilon_3 = 3 - \theta_0 - 3\theta_1$$

$$\text{MSE} = \frac{\epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2}{3} = \frac{14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1}{3}$$

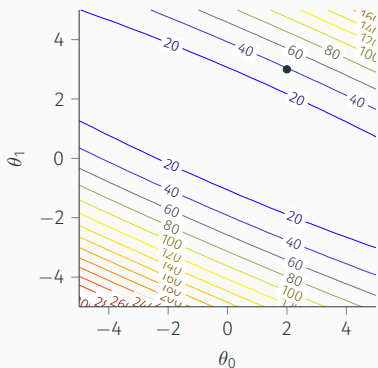
# Iteration 0

$$\text{MSE} = \frac{1}{3}(14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1)$$

Surface Plot



Contour plot, view from top



# Coordinate Descent : Example

## Iteration 1

INIT:  $\theta_0 = 2$  and  $\theta_1 = 3$

$\theta_1 = 3$  optimize for  $\theta_0$

# Coordinate Descent : Example

## Iteration 1

INIT:  $\theta_0 = 2$  and  $\theta_1 = 3$

$\theta_1 = 3$  optimize for  $\theta_0$

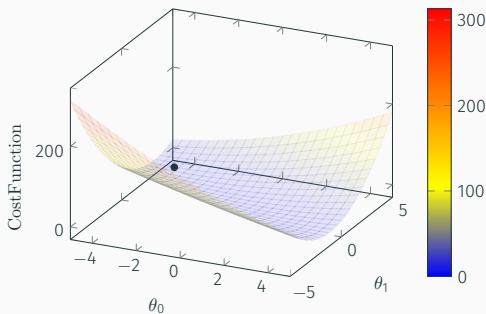
$$\frac{\partial \text{MSE}}{\partial \theta_0} = 6\theta_0 + 24 = 0$$

$$\theta_0 = -4$$

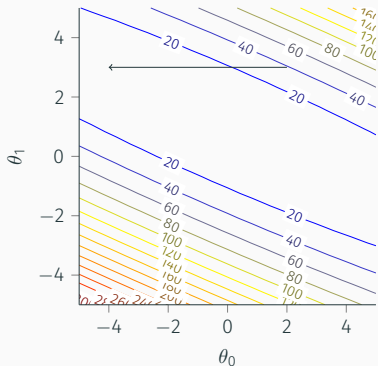
# Iteration 1

$$\text{MSE} = \frac{1}{3}(14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1)$$

Surface Plot



Contour plot, view from top



# Coordinate Descent : Example

## Iteration 2

INIT:  $\theta_0 = -4$  and  $\theta_1 = 3$

$\theta_0 = -4$  optimize for  $\theta_1$

# Coordinate Descent : Example

## Iteration 2

INIT:  $\theta_0 = -4$  and  $\theta_1 = 3$

$\theta_0 = -4$  optimize for  $\theta_1$

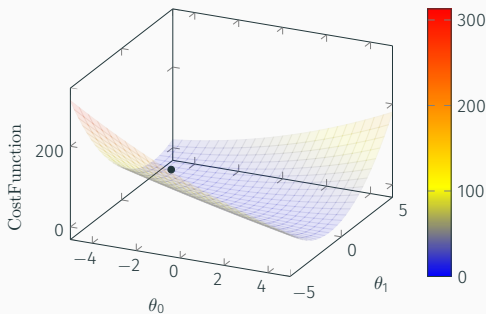
$\theta_1 = 2.7$



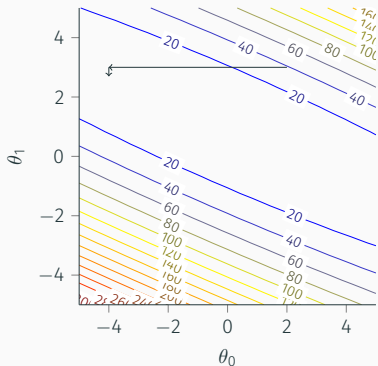
# Iteration 2

$$\text{MSE} = \frac{1}{3}(14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1)$$

Surface Plot



Contour plot, view from top



## Iteration 3

INIT:  $\theta_0 = -4$  and  $\theta_1 = 2.7$

$\theta_1 = 2.7$  optimize for  $\theta_0$

# Coordinate Descent : Example

## Iteration 3

INIT:  $\theta_0 = -4$  and  $\theta_1 = 2.7$

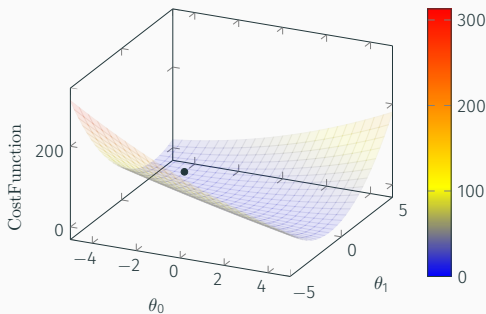
$\theta_1 = 2.7$  optimize for  $\theta_0$

$\theta_0 = -3.4$

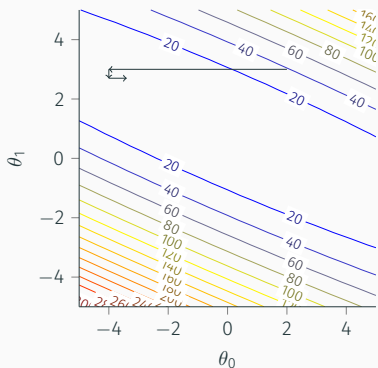
# Iteration 3

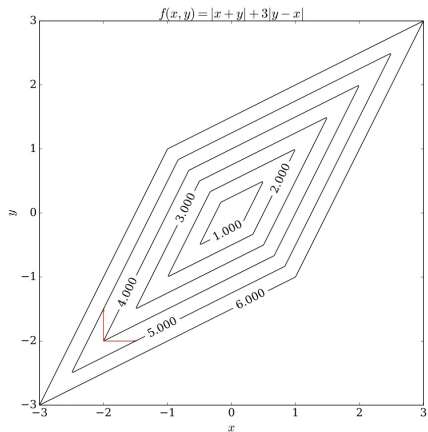
$$\text{MSE} = \frac{1}{3}(14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1)$$

Surface Plot

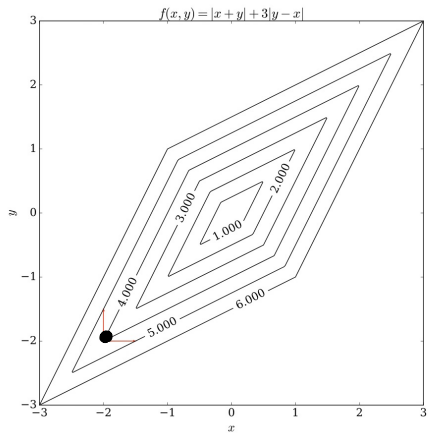


Contour plot, view from top



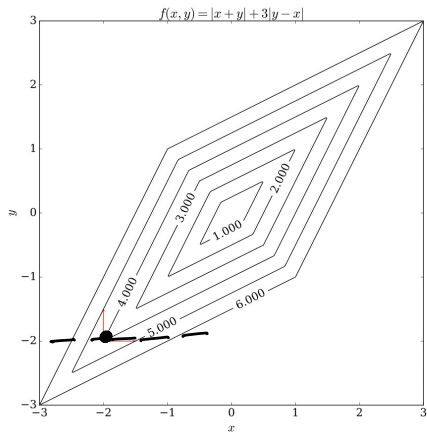


FAILURE OF COORDINATE  
DESCENT



FAILURE OF COORDINATE  
DESCENT

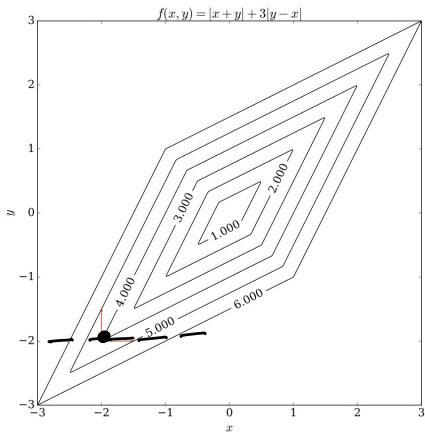
START WITH  $(x, y) = (-2, -2)$



FAILURE OF COORDINATE  
DESCENT

START WITH  $(x, y) = (-2, -2)$

FIX  $y = -2$ , OPTIMIZE  
ABOUT  $x$ .



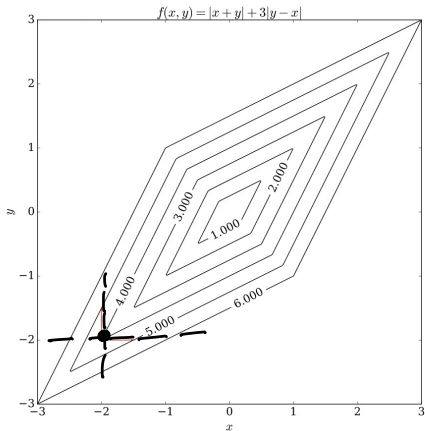
FAILURE OF COORDINATE  
DESCENT

START WITH  $(x, y) = (-2, -2)$

FIX  $y = -2$ , OPTIMIZE  
ABOUT  $x$ .

OBJECTIVE INCREASES  
IN BOTH DIRECTIONS





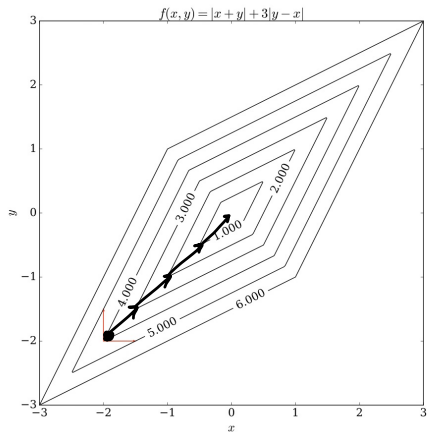
FAILURE OF COORDINATE  
DESCENT

START WITH  $(x, y) = (-2, -2)$

FIX  $y = -2$ , OPTIMIZE  
ABOUT  $x$ .

OBJECTIVE INCREASES  
IN BOTH DIRECTIONS

SIMILAR IF WE FIX  
 $x$  and OPTIMIZE ABOUT  $y$ .



GRADIENT DESCENT  
WILL WORK!

- NEED SIMULTANEOUS  
UPDATE IN BOTH  
COORDINATES

# Coordinate Descent for Unregularized Regression

- Express error as a difference of  $y_i$  and  $\hat{y}_i$

$$\hat{y}_i = \sum_{j=0}^d \theta_j x_i^j = \theta_0 x_i^0 + \theta_1 x_i^1 + \theta_2 x_i^2 + \dots + \theta_d x_i^d \quad (4)$$

$$\epsilon_i = y_i - \hat{y}_i \quad (5)$$

$$= y_i - \theta_0 x_i^0 + \theta_1 x_i^1 + \dots + \theta_d x_i^d \quad (6)$$

$$= y_i - \sum_{j=0}^d \theta_j x_i^j \quad (7)$$

## Coordinate Descent for Unregularized regression

$$\sum_{i=1}^N \epsilon^2 = \text{RSS} = \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \theta_d x_i^d \right) \right)^2$$

## Coordinate Descent for Unregularized regression

$$\sum_{i=1}^N \epsilon^2 = \text{RSS} = \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \theta_d x_i^d \right) \right)^2$$
$$\frac{\partial \text{RSS}(\theta_j)}{\partial \theta_j} = 2 \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \dots \right) \right) \left( -x_i^j \right)$$

## Coordinate Descent for Unregularized regression

$$\sum_{i=1}^N \epsilon^2 = \text{RSS} = \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \theta_d x_i^d \right) \right)^2$$

$$\begin{aligned} \frac{\partial \text{RSS}(\theta_j)}{\partial \theta_j} &= 2 \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \dots \right) \right) (-x_i^j) \\ &= 2 \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_d x_i^d \right) \right) (-x_i^j) + 2 \sum_{i=1}^N \theta_j (x_i^j)^2 \end{aligned}$$

## Coordinate Descent for Unregularized regression

$$\begin{aligned}\sum_{i=1}^N \epsilon^2 = \text{RSS} &= \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \theta_d x_i^d \right) \right)^2 \\ \frac{\partial \text{RSS}(\theta_j)}{\partial \theta_j} &= 2 \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_j x_i^j + \dots \right) \right) (-x_i^j) \\ &= 2 \sum_{i=1}^N \left( y_i - \left( \theta_0 x_i^0 + \dots + \theta_d x_i^d \right) \right) (-x_i^j) + 2 \sum_{i=1}^N \theta_j (x_i^j)^2\end{aligned}$$

where:

$$\hat{y}_i^{(-j)} = \theta_0 x_i^0 + \dots + \theta_d x_i^d$$

is  $\hat{y}_i$  without  $\theta_j$

## Coordinate Descent for Unregularized regression

$$\text{Set } \frac{\partial \text{RSS}(\theta_j)}{\partial \theta_j} = 0$$

$$\theta_j = \sum_{i=1}^N \frac{(y_i - (\theta_0 x_i^0 + \dots + \dots + \theta_d x_i^d)) (x_i^j)}{(x_i^j)^2} = \frac{\rho_j}{z_j}$$

$$\rho_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(-j)})$$

$$z_j = \sum_{i=1}^N (x_i^j)^2$$

$z_j$  is the squared of  $\ell_2$  norm of the  $j^{\text{th}}$  feature



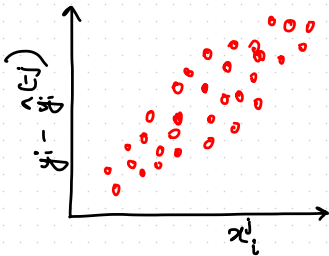
UNDERSTANDING  $\beta_j$  IN COORDINATE DESCENT

$$\beta_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(-j)})$$

# UNDERSTANDING $\beta_j$ IN COORDINATE DESCENT

$$\beta_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(j)})$$

CASE 1

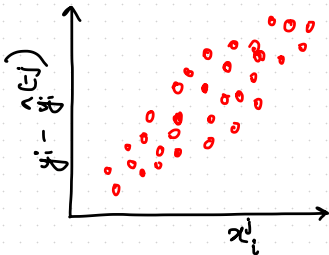


$x_i^j$  STRONG +VE CORR.  
WITH  $y_i - \hat{y}_i^{(j)}$

# UNDERSTANDING $\beta_j$ IN COORDINATE DESCENT

$$\beta_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(j)})$$

CASE 1



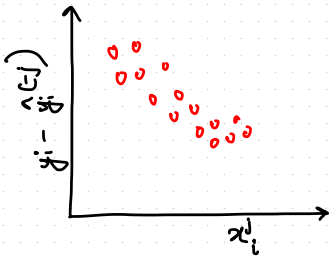
$x_i^j$  STRONG +VE CORR.  
WITH  $y_i - \hat{y}_i^{(j)}$

⇓  
 $j^{\text{th}}$  FEATURE IS IMP.T.  
AND ITS COEFFICIENT  
+VE

# UNDERSTANDING $\beta_j$ IN COORDINATE DESCENT

$$\beta_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(j)})$$

## CASE II



$x_i^j$  STRONG -VE CORR.  
WITH  $y_i - \hat{y}_i^{(j)}$

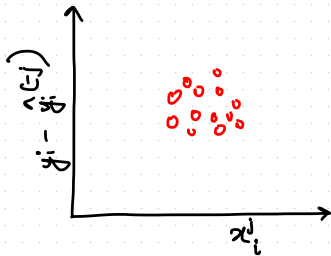


$j^{\text{th}}$  FEATURE IS IMP.  
AND ITS COEF. -VE

# UNDERSTANDING $\beta_j$ IN COORDINATE DESCENT

$$\beta_j = \sum_{i=1}^N x_i^j (y_i - \hat{y}_i^{(j)})$$

## CASE II



$x_i^j$  WEAK WITH  $y_i - \hat{y}_i^{(j)}$  CORR.  
↓  
 $j^{\text{th}}$  FEATURE IS **NOT** IMP.  
AND ITS COEF.  $\rightarrow 0$

# Coordinate Descent for Lasso Regression

$$\text{Minimize } \underbrace{\sum_{i=1}^N \epsilon^2 + \delta^2 \{|\theta_0| + |\theta_1| + \dots + |\theta_j| + \dots + |\theta_d|\}}_{\text{LASSO OBJECTIVE}}$$

$$\frac{\partial}{\partial \theta_j} (\text{LASSO OBJECTIVE}) = -2\rho_j + 2\theta_j z_j + \delta^2 \frac{\partial}{\partial \theta_j} |\theta_j|$$

$$\frac{\partial}{\partial \theta_j} |\theta_j| = \begin{cases} 1 & \theta_j > 0 \\ [-1, 1] & \theta_j = 0 \\ -1 & \theta_j < 0 \end{cases}$$

## Coordinate Descent for Lasso Regression

- Case 1:  $\theta_j > 0$

$$2\rho_j + 2\theta_j z_j + \delta^2 = 0$$

$$\theta_j = \frac{\rho_j - \frac{\delta^2}{2}}{z_j}$$

$$\rho_j > \frac{\delta^2}{2} \Rightarrow \theta_j = \frac{\rho_j - \frac{\delta^2}{2}}{z_j}$$

# Coordinate Descent for Lasso Regression

- Case 1:  $\theta_j > 0$

$$2\rho_j + 2\theta_j z_j + \delta^2 = 0$$

$$\theta_j = \frac{\rho_j - \frac{\delta^2}{2}}{z_j}$$

$$\rho_j > \frac{\delta^2}{2} \Rightarrow \theta_j = \frac{\rho_j - \frac{\delta^2}{2}}{z_j}$$

- Case 2:  $\theta_j < 0$

$$\rho_j < \frac{\delta^2}{2} \Rightarrow \theta_j = \frac{\rho_j + \delta^2/2}{z_j} \tag{8}$$



# Coordinate Descent for Lasso Regression

- Case 3:  $\theta_j = 0$

$$\frac{\partial}{\partial \theta_j}(\text{LASSO OBJECTIVE}) = -2\rho_j + 2\theta_j z_j + \underbrace{\delta^2 \frac{\partial}{\partial \theta_j} |\theta_j|}_{[-1,1]}$$
$$\in \underbrace{[-2\rho_j - \delta^2, -2\rho_j + \delta^2]}_{\{0\} \text{ lies in this range}}$$

$$-2\rho_j - \delta^2 \leq 0 \text{ and } -2\rho_j + \delta^2 \leq 0$$

$$-\frac{\delta^2}{2} \leq \rho_j \leq \frac{\delta^2}{2} \Rightarrow \theta_j = 0$$

## Summary of Lasso Regression

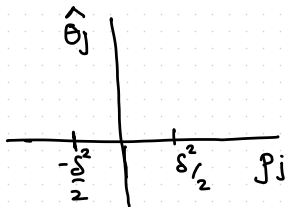
$$\theta_j = \left[ \begin{array}{ll} \frac{\rho_j + \frac{\delta^2}{2}}{z_j} & \text{if } \rho_j < -\frac{\delta^2}{2} \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \rho_j \leq \frac{\delta^2}{2} \\ \frac{\rho_j - \frac{\delta^2}{2}}{z_j} & \text{if } \rho_j > \frac{\delta^2}{2} \end{array} \right] \quad (9)$$

# LASSO (SOFT) THRESHOLDING

$$\theta_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \beta_j \leq \frac{\delta^2}{2} \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$

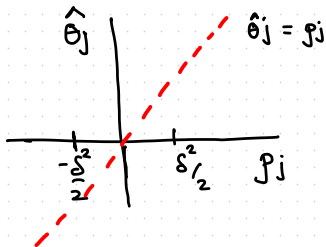
# LASSO (SOFT) THRESHOLDING

$$\hat{\theta}_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \beta_j \leq \frac{\delta^2}{2} \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$



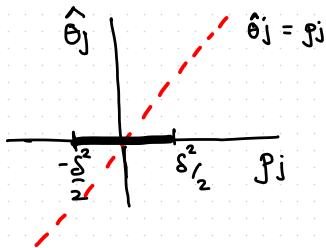
# LASSO (SOFT) THRESHOLDING

$$\hat{\theta}_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \beta_j \leq \frac{\delta^2}{2} \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$



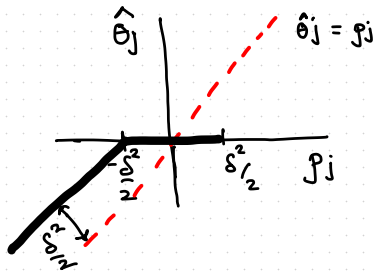
# LASSO (SOFT) THRESHOLDING

$$\hat{\theta}_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \beta_j \leq \frac{\delta^2}{2} \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$



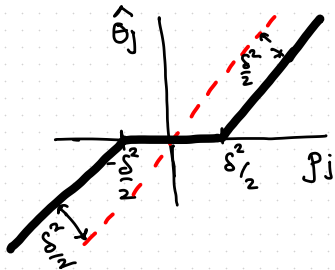
# LASSO (SOFT) THRESHOLDING

$$\hat{\theta}_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\delta^2/2 \leq \beta_j \leq \delta^2/2 \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$



# LASSO (SOFT) THRESHOLDING

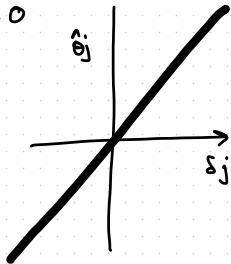
$$\hat{\theta}_j = \begin{cases} \frac{\beta_j + \delta^2/2}{z_j} & \text{if } \beta_j < -\delta^2/2 \\ 0 & \text{if } -\frac{\delta^2}{2} \leq \beta_j \leq \frac{\delta^2}{2} \\ \frac{\beta_j - \delta^2/2}{z_j} & \text{if } \beta_j > \delta^2/2 \end{cases}$$





# LASSO (SOFT) THRESHOLDING

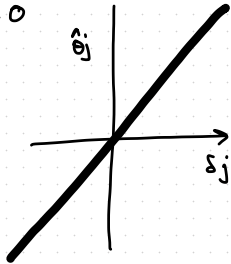
$$s^2 = 0$$



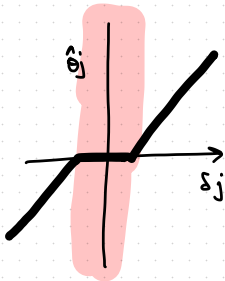
NO REGULARISATION

# LASSO (SOFT) THRESHOLDING

$$s^2 = 0$$



NO REGULARISATION



REGULARISATION  
↓  
SPARSITY