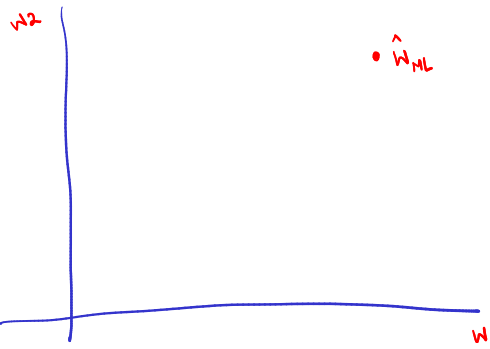


Linear regression / Ridge regression

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2 + \underbrace{\lambda \|w\|^2}_{\text{REGULARIZATION}} = \sum_{i=1}^d w_i^2 = \|w\|^2$$

PARAMETER SPACE

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$



Where is  $\hat{w}_R \rightarrow$  solution of the ridge regression problem?

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2 + \lambda \|w\|^2$$

↑  
Ⓐ

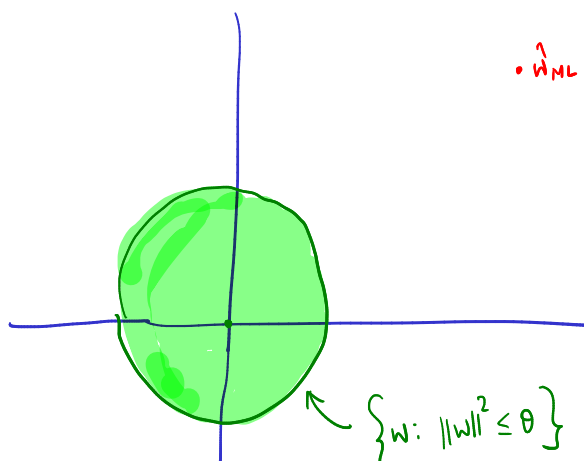
≡ is equivalent to

$$\begin{array}{l} \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2 \\ \text{s.t. } \|w\|^2 \leq \theta \end{array}$$

← Ⓑ

depends on  $\lambda$

For every choice of  $\lambda > 0$ ,  $\exists \theta$  s.t the optimal solutions of problems Ⓐ and Ⓑ coincide.



$$\|w\|^2 \leq \theta$$

$$\Rightarrow w_1^2 + w_2^2 \leq \theta$$

$$\{w: \|w\|^2 \leq \theta\}$$

- What is the loss / error / objective function value of linear regression of  $\hat{w}_{ML}$

$$\sum_{i=1}^n (\hat{w}_{ML}^T x_i - y_i)^2 = f(\hat{w}_{ML})$$

Consider the set of all  $w$  s.t

$$f(w) = f(\hat{w}_{ML}) + c$$

$c \geq 0$

$$S_c = \left\{ w : f(w) = f(\hat{w}_{ML}) + c \right\}$$

i.e., every  $w \in S_c$  satisfies

$$\underbrace{\|X^T w - y\|^2}_{f(w)} = \underbrace{\|X^T \hat{w}_{ML} - y\|^2}_{f(\hat{w}_{ML})} + c$$

on simplification [Please do this],

one gets

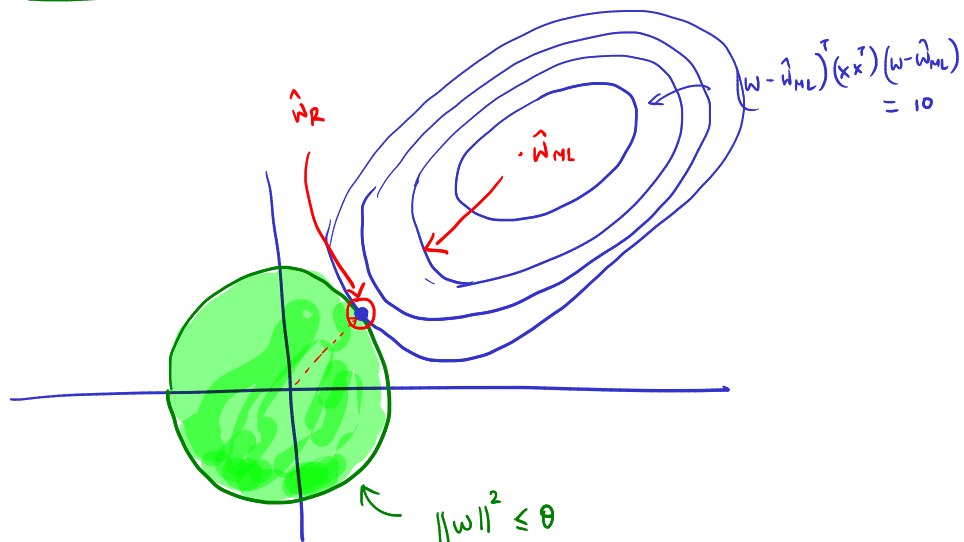
$$(w - \hat{w}_{ML})^T (X^T X) (w - \hat{w}_{ML}) = c'$$

Some constant that depends on  $c, (X^T X), \hat{w}_{ML}$  and not on  $w$

If  $X^T X = I$ ,

$$(w - \hat{w}_{ML})^T I (w - \hat{w}_{ML}) = c'$$

$$\|w - \hat{w}_{ML}\|^2 = c'$$



Conclusions: • Ridge regression pushes feature values towards 0. But does not necessarily make it 0.

• An alternate way to regularize would then be using  $\|\cdot\|_1$  norm instead of  $\|\cdot\|_2^2$  norm

$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

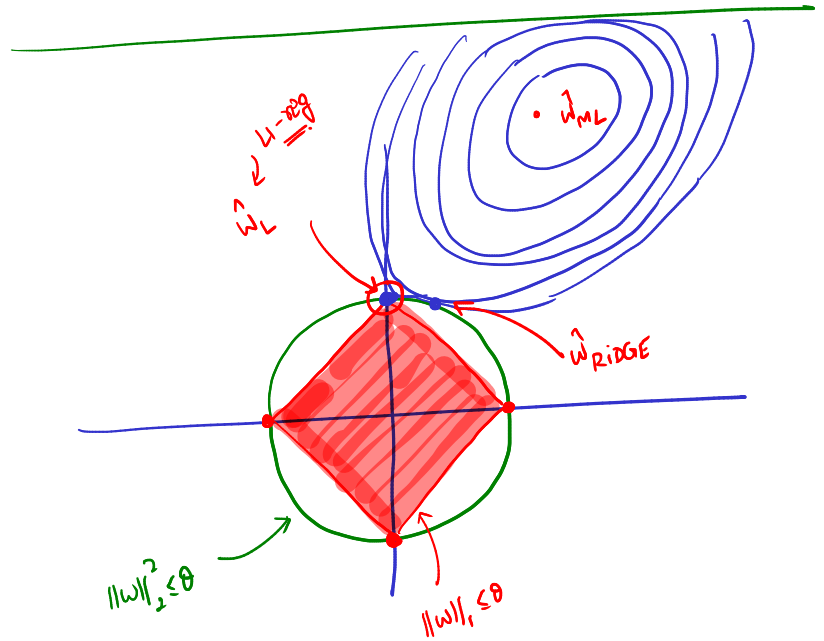
$L_1$  Regularization

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2 + \lambda \|w\|_1$$

$$=$$

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (\tilde{w}^T x_i - y_i)^2$$

s.t.  $\|w\|_1 \leq \theta$

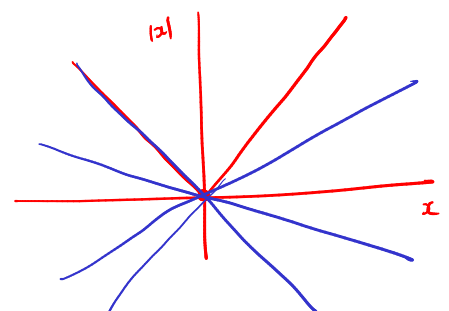
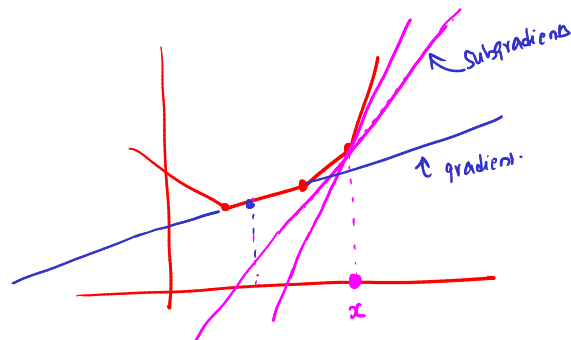


$L_1$  Reg: LASSO - LEAST ABSOLUTE SHRINKAGE and SELECTION OPERATOR.

Points

- LASSO does not have a closed form solution
- Sub-gradient methods are usually used to solve LASSO.

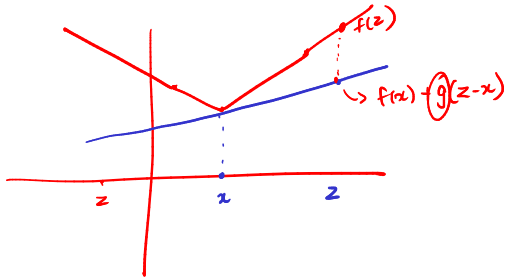
Sub-grad at 0  
=  $[-1, 1]$



## Subgradient

A vector  $g \in \mathbb{R}^d$  is a sub-gradient of  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  at a point  $x \in \mathbb{R}^d$  if

$$\forall z \quad f(z) \geq f(x) + g^T(z-x)$$



## Why Subgradients?

- If function  $f$  to minimize is a Convex function, then sub-gradient descent Converges!
  - There are other special purpose methods for LASSO  $\rightarrow$  (S1) IRLS [Iterative Reweighted Least Squares].
-