Sparse classifiers/models give insight

$(x_i, d_i)$, $i = 1 \cdots N$    $x_i^T w \approx d_i$    $Aw = [a_1 \; a_2 \cdots a_M] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} = \sum_{i=1}^{M} w_i a_i$

features, labels    $a_\ell$: $\ell^{th}$ feature component

Suppose $w_\ell \approx 0 \Rightarrow a_\ell$ is unimportant

If a small number of $w_\ell$ are nonzero, then only these few features matter! $w$ is sparse.
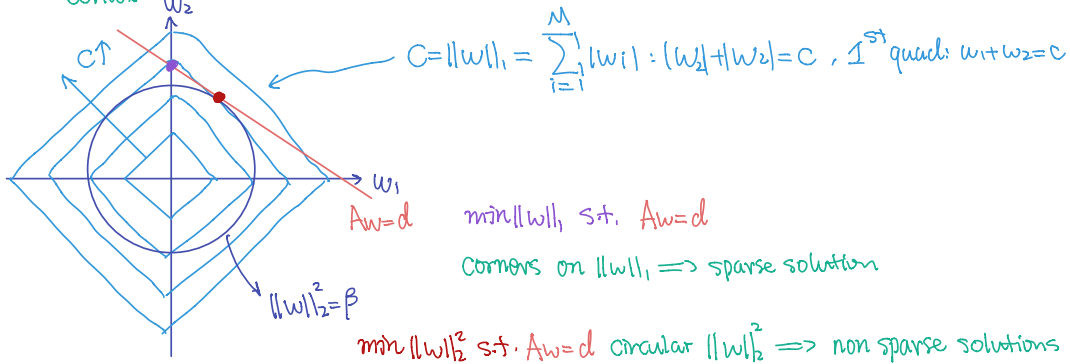
$\|w_0\| = \sum_{i=1}^{M} \mathbb{1}_{\{w_i \neq 0\}}$ (number of non-zero elements)

$\ell_0$-norm

$\|aw\|_0 \neq a\|w\|_0$    consider $\min_w \|w\|_0$ s.t. $\|Aw - d\|_2^2 < \varepsilon$    non-convex, intractable

Convex relaxation gives tractable problem

$\min_w \|w\|_1$ s.t. $\|Aw - d\|_2^2 < \varepsilon$ LASSO: Least Absolute Selection and Shrinkage Operator

convex



$C = \|w\|_1 = \sum_{i=1}^{M} |w_i| : (|w_1| + |w_2| = C$, $1^{st}$ quad: $w_1 + w_2 = C$

$Aw = d$    $\min \|w\|_1$ s.t. $Aw = d$

corners on $\|w\|_1 \Rightarrow$ sparse solution

$\|w\|_2^2 = \beta$

$\min \|w\|_2^2$ s.t. $Aw = d$ circular $\|w\|_2^2 \Rightarrow$ non sparse solutions

LASSO is a regularized least-squares problem

$\min_w \|w\|_1$ s.t. $\|Aw - d\|_2^2 < \varepsilon$ is equivalent to $\min_w \|Aw - d\|_2^2 + \lambda\|w\|_1$ for some $\lambda, \varepsilon$

Note: $\min_w \|w\|_1 + \frac{1}{\lambda}\|Aw - d\|_2^2$

LASSO

$w_L = \arg\min_w \|Aw - d\|_2^2 + \lambda\|w\|_1$

Sparse $w_L$

Can have small model error $w_{op} - w_L$

iterative solution

Ridge Regression

$w_R = \arg\min_w \|Aw - d\|_2^2 + \lambda\|w\|_2^2$

non-sparse $w_R$

great prediction error $\|Aw_{op} - Aw_R\|_2^2$

closed form solution

LASSO maybe used for model/feature selection

$w_L = \arg\min_w \|Aw - d\|_2^2 + \lambda\|w\|_1$    $S_L = \{i : [w_L]_i \neq 0\}$ selected features

$Aw_L = \sum_{i=1}^{M} a_i [w_L]_i = \sum_{i \in S_L} a_i [w_L]_i$

Debiasing    $A_L = \{a_i : i \in S_L\}$

$\hat{w}_L = \underset{w}{\arg\min} \|A_L w - d\|_2^2 = (A_L^T A_L)^{-1} A_L^T d$    avoids shrinkage due to $\|w\|_1$