

Lecture 6: Regularization

Recall from previous lecture:

D is unknown distribution over \mathbb{R}^d , $f: \mathbb{R}^d \rightarrow \{0,1\}$ is ground truth, $f \in \mathcal{A}$ for some known concept class \mathcal{A} .

Given: labeled training set $(x_1, y_1), \dots, (x_n, y_n)$
where $x_i \sim D$ independent,
 $y_i = f(x_i)$.

Goal: Output $g \in \mathcal{A}$ with good performance:

perfect recovery: $g = f$

mostly focus on this \rightarrow PAC learning: $\Pr_{x \sim D} [g(x) \neq f(x)] \leq \epsilon$ w.p. $1-\delta$.

Thm: If $\mathcal{A} = \{f_1, \dots, f_m\}$, then $n \geq \frac{1}{\epsilon} (\log m + \log \frac{1}{\delta})$ suffices.

(aside: $\log = \ln$).

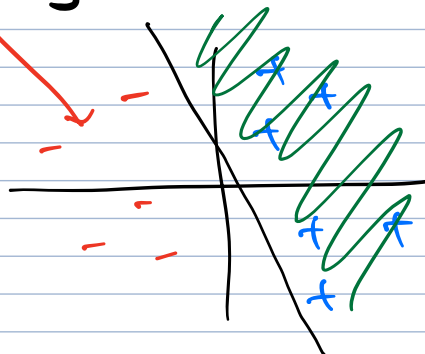
Thm: If $\mathcal{A} = \{\text{linear classifiers}\}$, then $n \geq \Omega(\frac{1}{\epsilon} (d + \log \frac{1}{\delta}))$ suffices.

Same rates in agnostic setting ($f \notin \mathcal{A}$) but w/ $\frac{1}{\epsilon^2}$, using ERM

Recall: linear classifiers are parametrized by $\theta \in \mathbb{R}^d$ unit vector,

$$f_\theta(x) = \text{sign}(\langle \theta, x \rangle)$$

$$= \begin{cases} 1 & \text{if } \langle \theta, x \rangle \geq 0 \\ 0 & \text{o.w.} \end{cases}$$



empirical risk minimizer

Q: What can you do if $d \gg n$?

A: **Nothing!** (in the worst case)

A: Regularization (for many settings)

(not classification!)

In this class, we will explore this through the lens of linear regression

Linear regression: $x_1, \dots, x_n \sim D$ in \mathbb{R}^d

$y_1, \dots, y_n \in \mathbb{R}$ \leftarrow not $\{0,1\}$.

ground truth $\theta^* \in \mathbb{R}^d$ \leftarrow not unit

Promise: $y_i \approx \langle \theta^*, x_i \rangle$, $\forall i=1, \dots, n$ \leftarrow some noise

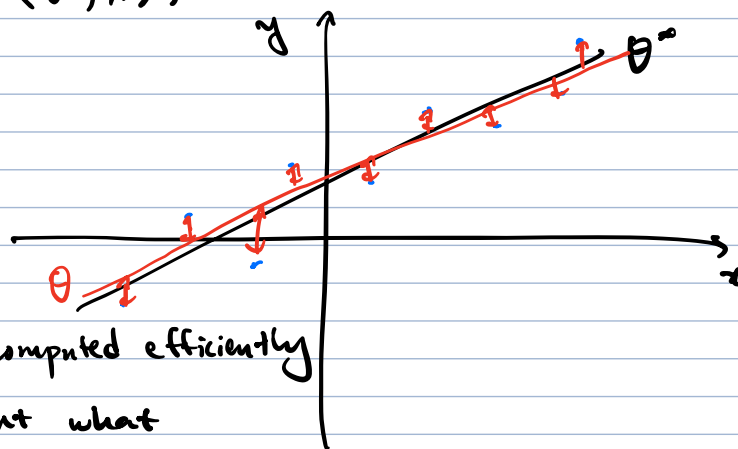
Goal: output θ that:
1) is close to θ^* , or (often assumed to be Gaussian)
2) acts like θ^* for typical data

$$\text{GenError: } \mathbb{E}_{x \sim D} (\langle \theta, x \rangle - \langle \theta^*, x \rangle)^2$$

ERM for linear regression:

$$L_{\text{train}}(\theta) := \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2$$

Output: $\arg\min_{\theta \in \mathbb{R}^d} L_{\text{train}}(\theta)$
 \uparrow can be computed efficiently



This works well when $n \gg d$, but what happens if $n \ll d$?

In this case, there are many θ with the same training loss.

Why? In this case, $\text{span}(\{x_1, \dots, x_n\})$ has $\dim \leq n \ll d$.

for any θ , and any $v \in (\text{span}\{x_1, \dots, x_n\})^\perp$, $\leftarrow \dim = d - n$.

$$\langle \theta, x_i \rangle = \langle \theta + v, x_i \rangle \quad \forall i=1, \dots, n \quad \langle v, x_i \rangle = 0 \quad \forall i$$

$$\text{so } L_{\text{train}}(\theta) = L_{\text{train}}(\theta + v)$$

How to choose good θ ?

Intuition: we should favor simple solutions

"simple" depends on the problem though

Common, useful notions of simple:

- small (low norm) solutions
- sparse solutions

Regularization is a method that lets you favor such "simple" solutions.

In this class: ℓ_2 and ℓ_1 regularization

\uparrow
small solutions
are simple

\uparrow
sparse solutions
are simple.

ℓ_2 -regularization aka "ridge regression", "Tikhonov regularization"
 [Hoerl, Kennard '70]

$$L_{\text{ridge}}(\theta) = \underbrace{\sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2}_{\text{error}} + \underbrace{\lambda \|\theta\|_2^2}_{\text{simplicity}}, \quad \lambda > 0 \quad \text{some parameter}$$

Output $\arg\min_{\theta} L_{\text{ridge}}(\theta)$ \leftarrow can be computed efficiently

example: $z_1, \dots, z_n \in \mathbb{R}$,

$$x_1 = (z_1, z_1, 0) \quad y_1 = 10z_1$$

$$x_2 = (z_2, z_2, 0) \quad y_2 = 10 \cdot z_2$$

$$\vdots$$
$$x_n = (z_n, z_n, 0) \quad y_n = 10z_n$$

Which $\theta \in \mathbb{R}^3$ is a solution w/ 0 error?

$$\langle \theta, x_i \rangle = \theta_1 z_i + \theta_2 z_i + \theta_3 \cdot 0$$

so θ is a perfect fit as long as $\theta_1 + \theta_2 = 10$,

θ_3 can be arbitrary

Ridge regression will choose θ w/ minimal l_2 -norm

$$\theta_3 = 0 \quad \leftarrow \text{removes spurious feature!}$$

$$\min \theta_1^2 + \theta_2^2 \quad \text{s.t.} \quad \theta_1 + \theta_2 = 10.$$

$$\theta_1 = \theta_2 = 5$$

ridge solution (for some suitable λ)

$$\theta \approx (5, 5, 0).$$

l_1 (and l_0) regularization

What if we want a sparse sol'n? eg. $(10, 0, 0)$

$$L_{l_0}(\theta) = \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2 + \lambda \|\theta\|_0$$

$$\|\theta\|_0 = |\{i : \theta_i \neq 0\}|$$

But this cannot be
computed efficiently

\leftarrow is limit
(in a certain sense)
of $\|\theta\|_p$ as $p \rightarrow 0$.

least absolute shrinkage & selection operator
 \downarrow

Instead: take the "convex proxy": l_1 (aka LASSO)

$$L_{\text{lasso}}(\theta) = \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2 + \lambda \|\theta\|_1$$

[Santosa, Symes '86]
[Tibshirani, '96]

$$\text{eg } x_1 = (z_1, z_1, 0) \quad y_1 = 10z_1$$

$$x_2 = (z_2, z_2, 0) \quad \vdots$$

$$\vdots$$
$$x_n = (z_n, z_n, 0) \quad y_n = 10z_n$$

What is the l_1 -minimizing solution?

$$\theta_1 z_i + \theta_2 z_i = 10 z_i \rightarrow z\theta_1 + \theta_2 = 10$$

Again $\theta_3 = 0$. Now $l_2: \min \theta_1 + \theta_2$ s.t.

$$2\theta_1 + \theta_2 = 10$$

What about l_1 ?

$$\min |\theta_1| + |\theta_2|$$

$$\text{s.t. } 2\theta_1 + \theta_2 = 10.$$

$$\theta_1 = 5, \theta_2 = 0. \rightarrow \text{only 1 nonzero!}$$

$$\theta_1 = 4, \theta_2 = 2$$

has nonzero on both.

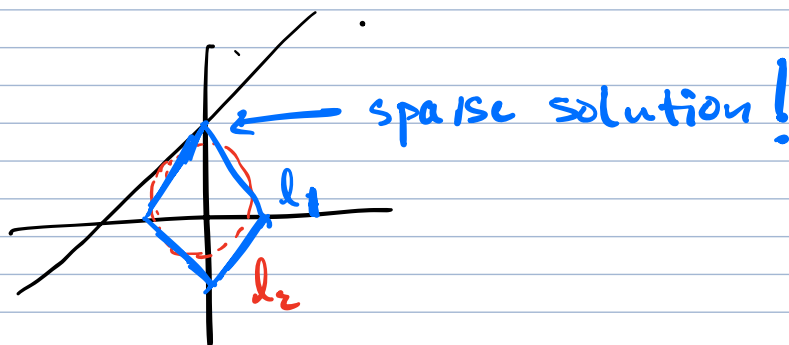
Why does l_1 enforce sparsity?

l_1 - regularization is the "soft version" of the following "hard" constraint

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$\min \|\theta\|_1$ such that

$\xrightarrow{\text{affine constraint}} X\theta = y \leftarrow \text{set of } \theta \text{ w/ 0 train error}$



Why sparsity?

- In practice, there are often many spurious features! Useful to prune them
- sparsity offers statistical advantages!

Suppose ground truth is sparse, i.e.

$$\mathcal{A} = \{ \langle \theta, x \rangle : \|\theta\|_0 \leq k \}.$$

Recall: Intuitively, generalization $\approx \log |\# \text{distinct } \theta \text{ in } \mathcal{A}|$.

How many free parameters for \mathcal{A} ?

A vector in A can be specified by:

1. Choose its nonzero coordinates S , $|S| \leq k$.

2. Choose a k -dim vector on this support

$$\binom{n}{k} \cdot C^k \approx n^k \cdot C^k$$

so generalization $\approx \log(n C^k)$

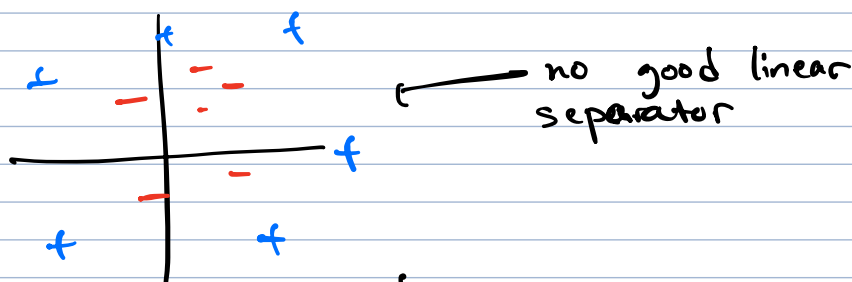
$$\approx k \log n \ll d \text{ if } k \text{ small}$$

rate for general vectors.

Kernelization

What if linear is insufficiently expressive?

recall from last lecture:



We can use a kernel $k: \mathbb{R}^d \rightarrow \mathbb{R}^m$, $m \gg d$.

hope there is a classifier in this higher space.

Polynomial kernel:

$$k(x_1, \dots, x_d) = (1, x_1, \dots, x_d, x_1^2, x_1 x_2, \dots, x_d^2)$$

(degree 2 polynomial kernel).

$$k: \mathbb{R}^d \rightarrow \mathbb{R}^{(d+1)^2}$$

For $\theta \in \mathbb{R}^{(d+1)^2}$

$$\langle \theta, k(x) \rangle = \sum_{i=0}^d \sum_{j=0}^d \theta_{ij} x_i x_j \quad (\text{set } x_0 = 1)$$

degree r kernel

$$k: \mathbb{R}^d \rightarrow \mathbb{R}^{(d+1)^r} \leftarrow \text{huge!}$$

Recipe: kernelize + regularize

Rule of thumb: want to be in regime where everything just barely works!