# Lecture 6: Regularization

Recall from previous lecture:

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

<u>Given:</u> labeled training set $(X_1, y_1), \dots, (X_n, y_n)$
where $X_i \sim D$ independent,
$y_i = f(X_i)$.

<u>Goal:</u> Output $g \in \mathcal{A}$ with good performance:

perfect recovery: $g = f$

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

<u>Thm:</u> If $\mathcal{A} = \{f_1, \dots, f_m\}$, then $n \gg \frac{1}{\varepsilon}(\log m + \log 1/\delta)$ suffices.

(aside: $\log = \ln$).
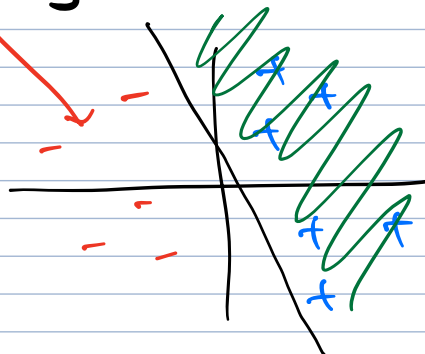
<u>Thm:</u> If $\mathcal{A} = \{$ linear classifiers $\}$, then $n \geq \Omega(\frac{1}{\varepsilon}(d + \log 1/\delta))$ suffices.
Same rates in <u>agnostic</u> setting $(f \notin \mathcal{A})$ but w/ $\frac{1}{\varepsilon^2}$, using **ERM**

empirical risk minimizer



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

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

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

<u>Q:</u> What can you do if $d \gg n$?

<u>A:</u> Nothing! (in the worst case)

<u>A:</u> 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

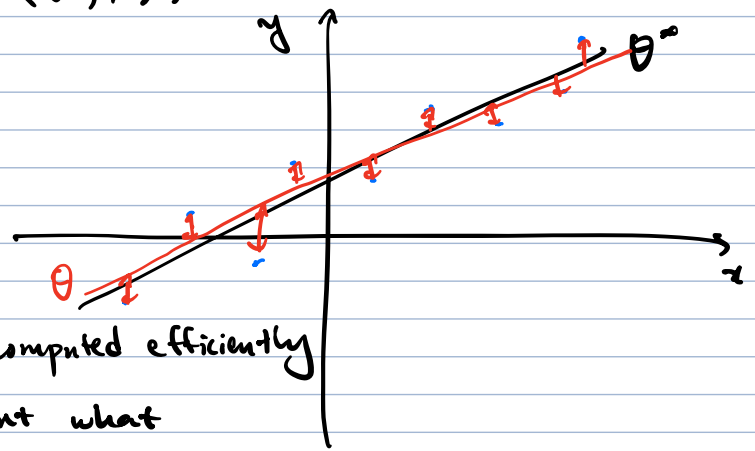<u>Goal:</u> output $\theta$ that: 1) is close to $\theta^*$, or (often assumed to be Gaussian)
2) acts like $\theta^*$ for typical data

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

ERM for linear regression:

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

Output: $\underset{\theta \in \mathbb{R}^d}{argmin} \ L_{train}(\theta)$

$\hookrightarrow$ can be computed efficiently

This works well when $n >> d$, but what happens if $n << d$?

In this case, there are many $\theta$ with the same training loss.

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

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

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

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

How to choose good $\theta$?

Intuition: we should favor <u>simple</u> 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: $\ell_2$ and $\ell_1$ regularization

small solutions are simple     sparse solutions are simple.

$\ell_2$ - regularization aka "ridge regression", "Tikhonov regularization"
[Hoerl, Kennard '70]

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

Output $argmin_\theta \ L_{ridge}(\theta)$ $\leftarrow$ can be computed efficiently

example: $z_1, \ldots, z_n \in \mathbb{R}, \eta_1, \ldots, \eta_n \in \mathbb{R}$

$$x_1 = (z_1, z_1, \eta_1) \qquad y_1 = 10 z_1$$
$$x_2 = (z_2, z_2, \eta_2) \qquad y_2 = 10 \cdot z_2$$
$$\vdots \qquad\qquad\qquad \vdots$$
$$x_n = (z_n, z_n, \eta_n) \qquad y_n = 10 z_n$$

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

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

so $\theta$ is a perfect fit as long as $\theta_1 + \theta_2 = 10$, $\theta_3$ can be arbitrary.

Ridge regression will choose $\theta$ w/ <u>minimal $\ell_2$ - norm</u>

$\theta_3 = 0$ $\leftarrow$ 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 $\lambda$)
$$\theta \approx (5, 5, 0).$$

---

## $\ell_1$ (and $\ell_0$) regularization

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

$$\rightarrow L_{\ell_0}(\theta) = \sum_{i=1}^n \left( \langle \theta, x_i \rangle - y_i \right)^2 + \lambda \|\theta\|_0$$

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

<span style="color:red">But this cannot be computed efficiently</span>

<span style="color:red">is limit.
(in a certain sense)
of $\|\theta\|_p^p$ as $p \to 0$.</span>

least absolute shrinkage & selection operator
$\downarrow$

Instead: take the "convex proxy": $\ell_1$ (aka LASSO)

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

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

eg
$$x_1 = (2z_1, z_1, \eta_1) \qquad y_1 = 10 z_1$$
$$x_2 = (2z_2, z_2, \eta_2) \qquad \vdots$$
$$\vdots$$
$$x_n = (2z_n, z_n, \eta_n) \qquad y_n = 10 z_n$$

What is the $\ell_1$ - minimizing solution?

$$\theta_1 2 z_i + \theta_2 z_i = 10 z_i \quad \rightarrow 2\theta_1 + \theta_2 = 10$$

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

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

$$\underline{\theta_1 = 4}, \quad \underline{\theta_2 = 2}$$

has nonzero on both.

What about $l_1$?

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

s.t. $2\theta_1 + \theta_2 = 10.$

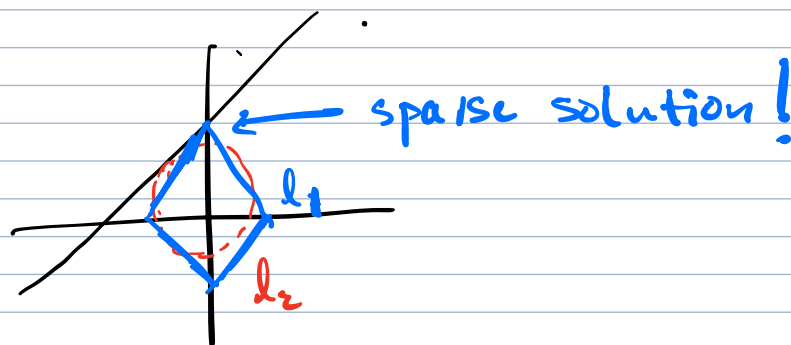$\theta_1 = 5, \boxed{\theta_2 = 0}. \longrightarrow$ only 1 nonzero!

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

affine constraint $\longrightarrow$ $X\theta = y \longleftarrow$ set of $\theta$ w/ 0 train error



sparse solution!

$l_1$

$l_2$

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.

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

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

How many free parameters for $A$?

A vector in $A$ can be specified by:

   1. Choose its nonzero coordinates $S$, <span style="color:red">$|S| \leq k$.</span>

   2. Choose a $k$-dim vector on this support

$$\binom{n}{k} \cdot C^k \approx n^k \cdot C^k$$
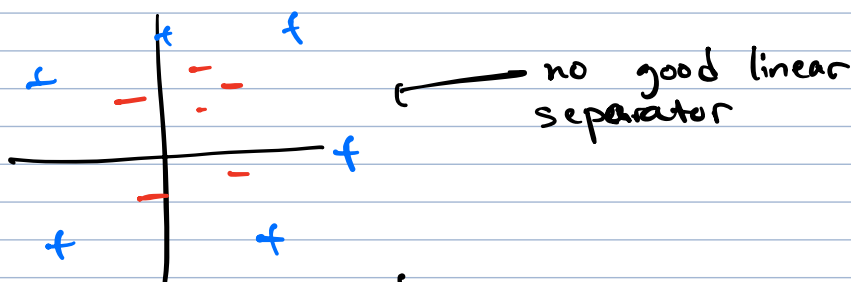
<span style="color:red">rate for general vectors.</span>

so    generalization $\approx \log\left((nC)^k\right)$

$$\approx k \log n \ll d \quad \text{.f}$$
$$k \text{ small}$$

---

## Kernelization

what if linear is insufficiently expressive?

recall from last lecture:



— no good linear separator

we can use a <u>kernel</u> $k: \mathbb{R}^d \to \mathbb{R}^m$, $m \gg d$.

hope there is a classifier in this higher space.

### Polynomial kernel:

$$k(x_1, \cdots, x_d) = \left(1, x_1, \cdots, x_d, x_1^2, x_1 x_2, \cdots, x_d^2\right)$$

(degree 2 polynomial kernel).

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

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

$$k: \mathbb{R}^d \to \mathbb{R}^{\boxed{(d+1)^r}} \quad \textcolor{red}{\Leftarrow \text{ huge!}}$$

Recipe: <u>kernelize</u> + <u>regularize</u>

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