

19. Linear Predictors (part 2: regression)

Monday, February 19, 2024 5:01 PM

§9.2 Linear Regression

$\mathcal{H} = \mathcal{L}_d$, $X \in \mathbb{R}^d$, $Y = \mathbb{R}$. Usually use squared loss $\ell(h, (x, y)) := (h(x) - y)^2$
using ℓ' loss $|h(x) - y|$ is a LP

so $\hat{L}_S(h) = \frac{1}{n} \sum_{i=1}^n (h(x_i) - y_i)^2$ is mean-squared-error, MSE. Very traditional.

ERM is least-squares problem

argmin $\frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2$ let $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$ $\begin{matrix} \leftarrow d \rightarrow \\ \uparrow n \downarrow \end{matrix}$
 $w \in \mathbb{R}^d$
(Assume $n \geq d$ for now)
 $= \frac{1}{n} \|X \cdot w - y\|_2^2$
solve in 1000's of ways.
prefer matrix/vector format numerically

Some major classes of sol'n method:

① dense methods.

Solution satisfies normal equations,

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

$$X^T (Xw - y) = 0, \text{ i.e. } w = (X^T X)^{-1} X^T y$$

Do not call $\text{inv}(X^T X)$! Don't do cholesky on $X^T X$ either! (slightly better) } these are inaccurate though fast

Instead: QR

$$n \times \begin{bmatrix} d \\ x \end{bmatrix} = \begin{bmatrix} Q \\ R \end{bmatrix} \begin{bmatrix} d \\ R \end{bmatrix}$$

$Q^T Q = I$, R is upper triangular.
eg. Gram-Schmidt, $O(nd^2)$

*even better: QR \rightarrow pivoting

$$w = (X^T X)^{-1} X^T y = (R^T Q^T Q R)^{-1} R^T Q^T y = (R^T R)^{-1} R^T Q^T y$$

cancel by hand

$$= R^{-1} R^T R^{-1} Q^T y$$

cancel by hand

and solve $w = R^{-1} Q^T y$ as $Rw = Q^T y$

using backward substitution since R is triangular
 $O(d^2)$ time instead of $O(d^3)$

(Matlab's backslash $w = X \backslash y$)
(will do a good job)

② If large, use Krylov subspace method like conjugate gradient (CG)

or better yet a variant optimized for least squares (LSQR...)

Can be combined w/ preconditioners

Always beats gradient descent

③ Stochastic Gradient Descent (SGD) if huge

④ New state-of-the-art methods (randomized, hybrid, approximate, ...)

if extremely large and/or very ill-conditioned.

19a. Linear Predictors (part 2: regression)

Monday, February 19, 2024

5:16 PM

What if I want to do **quadratic regression** or other **polynomial regression**?

① Sure... but it can be recast as linear regression by adding features!

eg. $X = \mathbb{R}$, $p(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3$

embed $x \mapsto \phi(x) = (1, x, x^2, x^3) \in \mathbb{R}^4$ and do

(homogeneous) linear classification here.

eg. $X = \mathbb{R}^2$, $p(x) = a_0 + a_1 x_1 + a_2 x_2 + a_3 x_1 x_2 + a_4 x_1^2 + a_5 x_2^2$
 $\phi(x) = (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2)$

② In high d, think twice... it's a lot of parameters

Sometimes we tame this by not considering cross-terms.

Learning theory for linear regression (ch. 11 in Mohri)

via Rademacher complexity

generic result: Thm 11.3 If the range of the loss l is bounded in $[0, M]$ and

$\hat{y} \mapsto l(\hat{y}, y)$ is μ -Lipschitz continuous, then

$$L_D(h) \leq \hat{L}_S(h) + \begin{cases} 2\mu R_m(\mathcal{H}) + M \cdot \sqrt{\log(1/\delta)/2m} \\ 2\mu \hat{R}_S(\mathcal{H}) + 3M \sqrt{\log(2/\delta)/2m} \end{cases}$$

The problem?

If $l(\hat{y}, y) = (\hat{y} - y)^2$, this is neither **bounded** nor **Lipschitz**

unless y and $\hat{y} (= \langle w, x \rangle)$ are bounded, which is an unusual

(though not unheard of) assumption. eg., bound $\|w\|$ (as we saw

in HW 3) and assume $\|x\|$ is bounded.

Alternatives: pseudo-dimension (§11.2.3 Mohri)

"Shattering" (and hence VC dim.) made sense when Y was finite,

but if $Y = \mathbb{R}$, what to do?

19b. Linear Predictors (part 2: regression)

Wednesday, February 21, 2024

2:38 PM

a new notion of "shattering"

Def 11.4 Shattering of real-valued function families

Let \mathcal{F} be a family of functions $f: Z \rightarrow \mathbb{R}$ (eg. $Z = X \times Y$,
 $f = l \circ h$, $h \in \mathcal{H}$)

A set $C = (z_1, \dots, z_m) \subseteq Z$ is

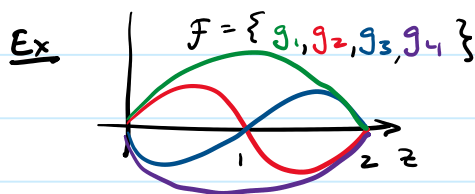
shattered by \mathcal{F} if $\exists (t_1, \dots, t_m) \in \mathbb{R}^m$

s.t.

$$\left| \left\{ \begin{bmatrix} \text{sign}(g(z_1) - t_1) \\ \vdots \\ \text{sign}(g(z_m) - t_m) \end{bmatrix} : g \in \mathcal{F} \right\} \right| = 2^m$$

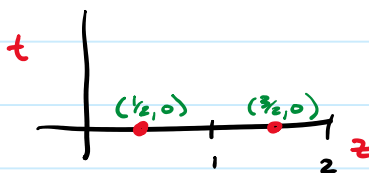
← "witnesses" to the shattering
← cardinality, not abs. value

You pick z_1, \dots, z_m and t_1, \dots, t_m once, then
look at all $g \in \mathcal{F}$



$$z_1 = 1/2, t_1 = 0$$

then $z_2 = 3/2, t_2 = 0$ are shattered



Def 11.5 The pseudo-dimension of \mathcal{F} is the size of the largest set that
can be shattered by \mathcal{F} . aka "pdim"

Facts (see Mohri)

- $\text{pdim}(\mathcal{L}_d) = d+1$ ($= \text{Vcdim}(\text{sign} \circ \mathcal{L}_d)$) via reducing to Vcdim w/ thresholds
- There are generalization error bounds involving pdim
(... but we won't cover. Not as nice/natural as Vcdim)