

Math / Comp 562 Lecture 1

2022/01/10



Lecture 1.

Formal Model.

Batch Supervised Machine Learning

X Domain $x \in X$
 Y Labels $y \in Y$ $Y = \{0, 1\}$ or $Y = \{-1, 1\}$ Class
 S^n training data $Y = \mathbb{R}$ regression
 $S = S^n = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$f: X \rightarrow Y$ hypothesis / predictor
prediction rule (classifier)

P data generation model
statistical data distribution
(unknown)
arbitrary

Noisy labels

$$p = p(x, y)$$

so on x , $P_Y(y) = P(y|x)$

Det. labels

$$y = f^*(x) \text{ and } p = p_x(x).$$

Defn Losses $\ell(y, z) : Y \times Y \rightarrow \mathbb{R}^+$

ℓ is a loss if $\ell(y, z) \geq 0$
with $\ell(y, z) = 0$ iff $y = z$

Ex Class:

$$\text{Given } y, y' \in Y$$
$$\ell_{0-1}(y, y') = \begin{cases} 0 & \text{if } y = y' \\ 1 & \text{otherwise} \end{cases}$$

Ex Regression

$$\ell_2(y, z) = (y - z)^2$$

Def Case: Global / generalization loss/error (risk)
of function $f : X \rightarrow Y$

$$R(f) = R_{\ell, \rho}(f) = \mathbb{E}[\ell(f(x), f^*(x))]$$

Ex $R(f) = \mathbb{E} \ell_{0-1}(f(x), f^*(x)) = \mathbb{P}(f(x) \neq f^*(x))$

Defn 2.1 Expected Risk / loss. noisy label case

Given $\ell: Y \times Y \rightarrow \mathbb{R}^+$ loss fn
 p prob dist on $X \times Y$
 $f: X \rightarrow Y$

the expected risk/loss of f is

$$\begin{aligned} \mathcal{R}(f) &= \mathcal{R}_{p,\ell}(f) = \mathbb{E}[\ell(y, f(x))] \\ &= \int_{X \times Y} \ell(y, f(x)) dp(x, y) \end{aligned}$$

Defn 2.2 Empirical Risk

Given $\ell: Y \times Y \rightarrow \mathbb{R}^+$ loss fn
data $S = S^n = \{(x_1, y_1), \dots, (x_n, y_n)\}$
 $f: X \rightarrow Y$

$$\hat{\mathcal{R}}(f) = \hat{\mathcal{R}}_{S^n, \ell}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Random ML

when f is nonrandom

$\mathcal{R}(f)$ is number

when $S^n =$ random iid
then

$f = A(S^n)$ is a
random
funct.

$\Rightarrow \mathcal{R}(f)$ random
variable

and

$\hat{\mathcal{R}}(f)$ r.v.

Dealing with noisy labels:

Given x' $y \sim p(z|x')$ $p_{x,y}(x,y)$

define conditional risk for any $z \in \mathcal{Y}$

$$r(z|x') = \mathbb{E}[\ell(y,z) | x=x'] = \int \ell(y,z) dp(y|x')$$

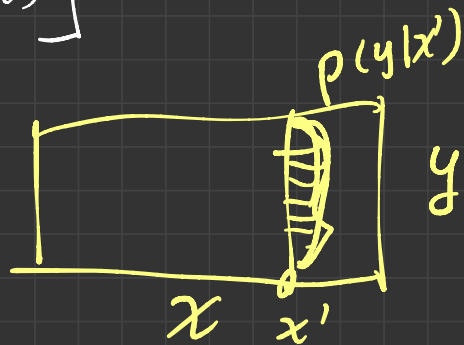
EX $\begin{cases} y=1 & p=0.8 \\ y=-1 & p=0.2 \end{cases}$
 x'

$$r(1|x') = 0.2$$

$$r(-1|x') = 0.8$$

$$\ell = \ell_{0,1}$$

$$\mathcal{R}(f) = \mathbb{E}_{x,y \sim p_{x,y}}[\ell(y, f(x))] = \mathbb{E}_{x' \sim p_x}[r(f(x')|x')]$$



Defn $f^*: X \rightarrow Y$ is a Bayes predictor (fn

$$\text{if } f^*(x') \in \underset{z \in Y}{\operatorname{argmin}} \mathbb{E}[l(y, z) | x=x'] \\ = \underset{z \in Y}{\operatorname{argmin}} r(z|x')$$

prop 2.1

All Bayes predictors have risk equal to

$$R^* = \mathbb{E}_{x' \sim p_x} \left[\inf_{z \in Y} r(z|x') \right]$$

Note: can have more than one
 $f^*(x') = +1$ or -1

$+1$	$p=0.5$
-1	$p=0.5$
x'	

Defn Excess Risk

The excess risk of f is $R(f) - R^*$
(always ≥ 0).

Cases

Binary Class: $f^*(x) = \arg \max_{z \in \{0,1\}} \underbrace{P(y=z | x=x')}_{\eta(x)}$

$$\mathcal{R}^* = \mathbb{E}[\min(\eta(x), 1 - \eta(x))]$$

Regression: $y \in \mathbb{R}$ $\ell(y, z) = (y - z)^2$

$$f^*(x) = \mathbb{E}[y | x=x']$$

2.3.2

Empirical Risk Minimization (ERM)

choose f to minimize $\hat{R}(f)$ empirical Risk / Loss

even though really care about

$\hat{R}(f)$ general Loss / Exp-Risk

Defn Hyp Class

$$\mathcal{H} = \{ f_{\theta}: \mathcal{X} \rightarrow \mathcal{Y} \mid \theta \in \Theta \}$$

EX Feature $\varphi: \mathcal{X} \rightarrow \mathbb{R}^d$ $\varphi(x) = (\varphi_1(x), \dots, \varphi_d(x))^T$

$$f_{\theta}(x) = \theta^T \varphi(x), \quad \theta \in \mathbb{R}^d = \Theta$$

$$\text{Do: } \min_{f \in \mathcal{H}} \hat{R}(f) \iff \min_{\theta \in \Theta} \hat{R}(f_{\theta})$$

functional opt parametric opt

PROS

- can optimise, using GD / SGD
- can be applied in high dim

CONS

- harder opt in non-convex case (NN)
- ⇒ need a good feature vector
 - by experts
 - learned by nn (no theory)
- need to control overfitting
- in classification case, need to deal with optimization on real values
 - round to class. (solved using score based losses)

2.3.1 Compare to "Local Averaging" methods

instead of learning f_{θ} , interpolate (average) directly from data.

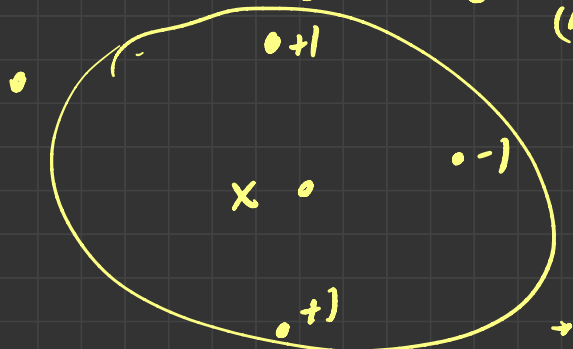
EX KNN, decision trees,

K-nearest neighbors:

Given $S^n = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$X = \mathbb{R}^d$ $Y = \{0, 1\}$

Define: $f(x)$ by (i) find k nearest nbrs of x
(ii) $f(x) =$ majority vote of labels of nbrs.


$$f(x) = \text{maj} \{+1, +1, -1\} \\ = +1$$

pros :

- no opt. / train
- easy to implement
- works well (often) in low dim

Cons :

- slow at query time: need to search
- bad in high dimensions
because of curse of dim.
- the distance fn is crucial
- need to tune k hyperparam-
can underfit / overfit with wrong k

Risk decomposition

Defn Given $\cdot \mathcal{L}$

$$\cdot \mathcal{H} = \{f_\theta \mid \theta \in \Theta\}$$

$$\cdot \triangleright S^n = \{(x_i, y_i)\}_{i=1}^n$$

Way to get rid of randomness coming from f_θ^1 . Still random.

Defn Approximation error: $E_A(\mathcal{H}) = \inf_{\theta \in \Theta} \{ \mathcal{R}(f_\theta) - \mathcal{R}^* \}$

↙ deterministic

how well can approx f^* using \mathcal{H} .

Estimation Error $E_E(f_\theta^1, \mathcal{H}) = \mathcal{R}(f_\theta^1) - \inf_{\theta \in \Theta} \mathcal{R}(f_\theta)$

↙ random

Then

$$\text{Excess Risk} = E\mathcal{R}(f_\theta^1) - \mathcal{R}^*$$

$$= E_E + E_A$$

Uniform Deviation $\leq U(\mathcal{H}) = \sup_{\theta \in \Theta} |\mathcal{R}(f_\theta) - \hat{\mathcal{R}}(f_\theta)|$ random by S^n

by defn: $\mathcal{R}(f_{\hat{\theta}}) - \hat{\mathcal{R}}(f_{\hat{\theta}}) \leq U(\mathcal{H})$

$$E_{\mathcal{E}}(f_{\hat{\theta}}, \mathcal{H}) = \mathcal{R}(f_{\hat{\theta}}) - \inf_{\theta \in \Theta} \mathcal{R}(f_\theta)$$

$$\begin{aligned} &= \underbrace{-\{\mathcal{R}(f_{\hat{\theta}}) - \hat{\mathcal{R}}(f_{\hat{\theta}})\}}_{\leq 2U(\mathcal{H})} + \underbrace{\{\hat{\mathcal{R}}(f_{\hat{\theta}}) - \hat{\mathcal{R}}(f_{\theta^*})\}}_{E_{\text{opt}}(f_{\hat{\theta}}, \mathcal{H})} + \underbrace{\{\hat{\mathcal{R}}(f_{\theta^*}) - \mathcal{R}(f_{\theta^*})\}}_{\leq 2U(\mathcal{H})} \\ &\leq 2U(\mathcal{H}) + E_{\text{opt}} \end{aligned}$$

Thus

$$\text{Excess Risk} = E_A + E_{\mathcal{E}}$$

$$\leq \underbrace{E_A}_{\uparrow \text{fit}} + \underbrace{E_{\text{opt}}}_{\uparrow \text{opt}} + \underbrace{2U(\mathcal{H})}_{\uparrow \text{capacity to overfit}}$$

Diff NN/ML:

NN: control fit & opt, not cap. ML: cap \checkmark opt not E_A .