

Statistical learning (for regression/classification)

Data (x, y) drawn iid from unknown
features \uparrow labels \downarrow distribution D

Loss function $\ell(\hat{y}, y) \geq 0$ e.g. $(\hat{y} - y)^2$

Predictor $h: X \rightarrow Y$

Statistical risk $L_D(h) = \mathbb{E}[\ell(h(X), Y)]$

Training set S : pairs (X, Y) drawn iid from D

Learning algorithm A maps S to h

H_A : ~~set~~ set of all predictors output by A
(e.g., linear predictors)

In stat learning we care about variance error

$$L_D(A(S)) - \inf_{h \in H_A} L_D(h)$$

\uparrow
random var.

Online learning asks: can we study machine learning without using statistical assumptions?

Online learning

- Data points arrive one by one in a stream
- No stat assumptions on the stream

Motivations:

- Sensor data, user interaction data, financial data
- Large datasets can only be accessed sequentially

Online learning is incremental

Start from initial predictor $h_1 \in \mathcal{H}$

For $t = 1, 2, \dots$ (stream)

- observe next data point (x_t, y_t) in stream
- test current predictor h_t and measure loss $\ell(h(x_t), y_t)$
- update $h_t \rightarrow h_{t+1} \in \mathcal{H}$

Note: update is local as it depends on $h_t, (x_t, y_t)$

Notation: $\ell(h_t(x_t), y_t) \rightarrow \ell_t(h_t) \in [0, 1]$

h_1, h_2, \dots stream of predictors generated

Sequential risk of h_1, h_2, \dots $\sum_{t=1}^T l_t(h_t)$

Regret (cfr. variance error):

$$R_T = \sum_{t=1}^T l_t(h_t) - \inf_{h \in \mathcal{H}} \sum_{t=1}^T l_t(h)$$

Since $l_t \in [0, 1]$ $R_T = O(T)$

Learning takes place when $R_T = o(T)$

Learning w/ expert advice: a notable special case

K actions (e.g. stocks), $l_t(i) \in [0, 1]$ is loss of action i at time t

We consider randomized predictors

Model at time t is distribution p_t over $\{1, \dots, K\}$ actions

Regret $R_T = \mathbb{E} \left[\sum_{t=1}^T l_t(I_t) \right] - \min_{i=1, \dots, K} \sum_{t=1}^T l_t(i)$

For $t = 1, 2, \dots$

- Draw $I_t \sim p_t$ and pay $l_t(I_t)$

- Observe l_t and update $p_t \rightarrow p_{t+1}$

$$L_{t-1}(i) = \sum_{s=1}^{t-1} \ell_s(i) \quad \text{total loss of } i \text{ up to time } t-1$$

$$\text{Hedge algorithm: } p_t(i) \propto \exp(-\eta L_{t-1}(i))$$

$\eta > 0$

Regret analysis

$$w_t(i) = e^{-\eta L_{t-1}(i)}, \quad w_1(i) = 1 \text{ as } L_0(i) = 0$$

$$W_t = \sum_j w_t(j) \quad p_t(i) = w_t(i) / W_t$$

Note: $w_{t+1}(i) = w_t(i) e^{-\eta \ell_t(i)}$ update rule

$$\ln \frac{W_{T+1}}{W_1} = \sum_{t=1}^T \ln \frac{W_{t+1}}{W_t}$$

$$\ln \frac{W_{t+1}}{W_t} = \ln \sum_i \frac{w_{t+1}(i)}{W_t} = \ln \sum_i \boxed{\frac{w_t(i)}{W_t}} e^{-\eta \ell_t(i)} \quad \begin{matrix} \nwarrow p_t(i) \end{matrix}$$

$$\leq \ln \sum_i p_t(i) \left(1 - \eta \ell_t(i) + \frac{\eta^2}{2} \ell_t(i)^2 \right) \quad e^{-x} \leq 1 - x + \frac{x^2}{2} \quad (x > 0)$$

$$= \ln \left(1 - \eta \sum_i p_t(i) \ell_t(i) + \frac{\eta^2}{2} \sum_i p_t(i) \ell_t(i)^2 \right)$$

$$\leq -\eta \underbrace{\sum_i p_t(i) \ell_t(i)}_{\mathbb{E}[\ell_t(I_t)]} + \frac{\eta^2}{2} \sum_i p_t(i) \ell_t(i)^2 \quad \ln(1-x) \leq -x$$

$$\mathbb{E}[\ell_t(I_t)]$$

Summing over $t = 1, \dots, T$:

$$\textcircled{1} \ln \frac{\bar{W}_{T+1}}{W_1} \leq -\eta \mathbb{E} \left[\sum_t \ell_t(I_t) \right] + \frac{\eta^2}{2} \sum_t \sum_i \ell_t(i)^2 p_t(i)$$

$$K_T^* = \underset{i=1, \dots, K}{\operatorname{argmin}} \sum_{t=1}^T \ell_t(i) \quad \text{best action in first } T \text{ steps}$$

$$\begin{aligned} \textcircled{2} \ln \frac{\bar{W}_{T+1}}{W_1} &= \ln \sum_i e^{-\eta L_T(i)} - \ln K \geq \ln e^{-\eta L_T(K_T^*)} - \ln K \\ &= -\eta L_T(K_T^*) - \ln K \end{aligned}$$

Divide $\textcircled{1}$ and $\textcircled{2}$ by $\eta > 0$ and rearrange:

$$R_T \leq \frac{\ln K}{\eta} + \frac{\eta}{2} \sum_t \sum_i \ell_t(i)^2 p_t(i)$$

$\hookrightarrow \leq 1 \quad \hookrightarrow \sum_i p_t(i) = 1$

$$\text{Hence } R_T \leq \frac{\ln K}{\eta} + \frac{\eta}{2} T$$

$$\text{choosing } \eta \approx \sqrt{\frac{\ln K}{T}} \text{ gives } R_T \leq \sqrt{T \ln K}$$

⊛ Not improvable! (except for constants)

Dynamic tuning $\eta_t \approx \sqrt{\frac{\ln K}{t}}$ gives

$$R_T \leq \sqrt{T \ln K} \quad \text{simultaneously for all } T$$