

$$\textcircled{1} f: \{0,1\}^d \rightarrow \{0,1\}$$

$$\text{AND}_S(x) = \bigwedge_{i \in S} x_i \quad S \subseteq [d]$$

$$(x^1, y_1), \dots, (x^m, y_m) \quad \text{where } x^l \in \{0,1\}^d \text{ \& } y_l = \text{AND}_{S^*}(x^l) \text{ for non-noise } S^*$$

Initialize set of features  $S$  to be  $[d]$ , full set of features

For each example  $(x^l, y_l)$ :

$$\textcircled{1} \hat{y}_l = \text{AND}_S(x^l) \quad \text{where } S \text{ is current set of features}$$

$\textcircled{2}$  IF  $\hat{y}_l \neq y_l$  UPDATE set of features. Remove all features  $x_i$  such that  $x_i = 0$  in  $x^l$ . Continue to next example

At end, output set of remaining features,  $\hat{S}$ .

Predict  $\hat{y}(x) = \text{AND}_{\hat{S}}(x)$  for any new input,  $x$

- Since we removed only features that we know cannot be in  $S^*$ , the above algorithm shows that the true set of features  $S^*$  will be a subset of remaining features at end.
- Since the algorithm removes at least one feature per mistake, it will make at most  $d$  mistakes in entire run-time, which is a polynomial bound
- Run-time is also polynomial in  $d$  since each example can be processed in  $O(d)$ . There are at most  $d$  examples seen before true set of features is identified. So, the overall run-time is  $O(d^2)$



② Consider the following scenario with  $Z$  experts  $A$  &  $B$ .

- $A$  predicts 1 for each day,  $t$
- $B$  predicts 0 for each day,  $t$
- Since the algorithm is deterministic, an adversary can fix all outcomes such that the predictions are always wrong.

Then, at least 1 of  $A$  &  $B$  will have an error rate of  $\leq 0.5$ , & the algorithm error rate is 1.

- Thus, showing that no deterministic algorithm can do better than a factor of  $Z$  compared to the best expert



③ There's some  $W^*$  & we see  $(x^1, y^1), \dots, (x^T, y^T)$ :

$$y^T = \text{sign}(\langle W^*, x^T \rangle), \quad W_i^* \geq 0, \quad \sum_{i=1}^d W_i^* = 1$$

MWM:

$$① \quad W(0, i) = 1 \quad \forall i \in [d], \quad u_1^0 = \frac{W(0, i)}{\sum_{j=1}^d W(0, j)}$$

② For each step  $t=1, \dots, T$

$$\text{Predict: } \text{sign}(\langle u^{t-1}, x^t \rangle)$$

$$\text{Update: } \begin{array}{l} \text{if correct} \quad W(t, i) = W(t-1, i) \quad \forall i \\ u^t = u^{t-1} \end{array}$$

$$\text{else} \quad L(t, i) \in [0, 1]$$

$$W(t, i) = (1 - \epsilon \cdot L(t, i)) W(t-1, i)$$

$$u^t = \text{normalized vector of } W(t)$$

from  
office  
hour?

$$L(t, i) = -y^t x_i^t$$

$$\text{Regret Bound: } \mathbb{E}[L(T)] \leq L^*(T) + O(\sqrt{T \ln d})$$



$$\sum_{t=1}^T \left[ \sum_{i=1}^d \Pr[\text{pick expert } i \text{ at step } t] \cdot L(t, i) \right] = \sum_{t=1}^T \sum_{i=1}^d -y^t u_i^t x_i^t = \sum_{t=1}^T -y^t \langle u^t, x^t \rangle$$

Loss of Best Expert: let  $L^t$  be vector of losses on step  $t$ ,  $L^t \in [-1, 1]^d$

$$\sum_{t=1}^T L^t \quad \text{vector of losses for each expert}$$

$$L^*(T) = \min_i \sum_{t=1}^T L^t \leq \langle W^*, \sum_{t=1}^T L^t \rangle$$

$$\Rightarrow \mathbb{E}[L(T)] \leq \sum_{t=1}^T \langle W^*, L^t \rangle + O(\sqrt{T \ln d}) \quad \begin{array}{l} \text{// sum term} = 0 \text{ if correct} \\ L_i^t = -y^t x_i^t \text{ if not} \end{array}$$

$$\mathbb{E}[L(T)] \leq (\# \text{ mistakes}) \cdot (-\gamma) + O(\sqrt{T \ln d})$$

$$\Rightarrow (\# \text{ mistakes}) \cdot \gamma + \mathbb{E}[L(T)] \leq O(\sqrt{T \ln d})$$



④ a)  $\{(0,1), (1,0)\}$

To determine the line of best fit through the origin we know that we want to maximize the following expression:

$$\max_V \sum_{x \in X} \|\text{Proj}_L(x)\|_2^2 \equiv \max_{V: \|V\|_2=1} \sum_{x \in X} \langle x, V \rangle^2$$

→ We know:  $\|\text{Proj}_V(x)\|_2 = \frac{\langle x, V \rangle}{\|V\|_2}$

So, the problem boils down to:

$$\max_{V: \|V\|_2=1} \langle (x_1, y_1), \vec{V} \rangle^2 + \langle (x_2, y_2), \vec{V} \rangle^2$$

For the points  $(0,1)$  &  $(1,0)$  we get

$$\max_{V: \|V\|_2=1} \langle (1,0), \vec{V} \rangle^2 + \langle (0,1), \vec{V} \rangle^2 \Rightarrow V_1^2 + V_2^2 = 1$$

- So, the best fit line is not unique for this problem.

Any line that passes through the origin will be the line of best fit.



$$b). \{ (0,1), (2,0) \}$$

- Following the same analysis as before, but with the new points we get:

$$\max_{\vec{v}: \|\vec{v}\|_2=1} \langle (0,1), \vec{v} \rangle^2 + \langle (2,0), \vec{v} \rangle^2 \Rightarrow v_2^2 + 4v_1^2$$

$$\Rightarrow \vec{v} = \langle 1, 0 \rangle$$

Thus the line of best fit is the X-axis.

In this case the best-fit line is unique. If we

think about it geometrically, this makes sense as moving away from the X-axis will have adverse effects because the  $(2,0)$  point has more "power" than the  $(0,1)$  point.



$$c) \cdot \{(0, -1), (2, 0)\}$$

- Following the same analysis as before, but with the new points we get:

$$\max_{V: \|V\|_2=1} \langle (0, -1), \vec{V} \rangle^2 + \langle (2, 0), \vec{V} \rangle^2 \Rightarrow V_2^2 + 4V_2^2$$

$$\Rightarrow \vec{V} = \langle 1, 0 \rangle$$

Again, the line of best-fit is the x-axis. The line of best-fit is unique for the same reasons as before



⑤ If we were given a dataset & asked to find the best-fit line that doesn't necessarily need to go through the origin, I would do the following:

- "Center" the data by subtracting the mean of each dimension from each data point, which basically shifts the origin to the center of the data
- Use the same approach as discussed in class to find the line of best-fit
- Re-shift the line to the correct position by using the means again