

Online learning

CPSC 532S: Modern Statistical Learning Theory

4 April 2022

cs.ubc.ca/~dsuth/532S/22/

Admin

- **Project presentations** in class on Wednesday; see Piazza
 - Grade component will be based mainly on clarity
 - Zoom setup probably jankier than I was hoping; come in person if you can

Admin

- **Project presentations** in class on Wednesday; see Piazza
 - Grade component will be based mainly on clarity
 - Zoom setup probably jankier than I was hoping; come in person if you can
- **Project reports** due Friday, with usual late policy (-5 Saturday, -10 Sunday)
 - Grade a combination of presentation + “technical content”
 - Not going to be harsh, just looking to see you did *something*

Admin

- **Project presentations** in class on Wednesday; see Piazza
 - Grade component will be based mainly on clarity
 - Zoom setup probably jankier than I was hoping; come in person if you can
- **Project reports** due Friday, with usual late policy (-5 Saturday, -10 Sunday)
 - Grade a combination of presentation + “technical content”
 - Not going to be harsh, just looking to see you did *something*
- **A4** also due Friday, but **extended late policy**: only -1 point per day late
 - Hard deadline of next Friday the 15th

Admin

- **Project presentations** in class on Wednesday; see Piazza
 - Grade component will be based mainly on clarity
 - Zoom setup probably jankier than I was hoping; come in person if you can
- **Project reports** due Friday, with usual late policy (-5 Saturday, -10 Sunday)
 - Grade a combination of presentation + “technical content”
 - Not going to be harsh, just looking to see you did *something*
- **A4** also due Friday, but **extended late policy**: only -1 point per day late
 - Hard deadline of next Friday the 15th
- **Final**: take-home, available over most of the finals period
 - Available to “check out” when it’s done (probably ~April 12/13)
 - Will target a couple hours of work (\ll an assignment)
 - Will have at least 12 hrs to do, maybe 1-2 days (TBD); must hand in by April 25

Admin

- **Project presentations** in class on Wednesday; see Piazza
 - Grade component will be based mainly on clarity
 - Zoom setup probably jankier than I was hoping; come in person if you can
- **Project reports** due Friday, with usual late policy (-5 Saturday, -10 Sunday)
 - Grade a combination of presentation + “technical content”
 - Not going to be harsh, just looking to see you did *something*
- **A4** also due Friday, but **extended late policy**: only -1 point per day late
 - Hard deadline of next Friday the 15th
- **Final**: take-home, available over most of the finals period
 - Available to “check out” when it’s done (probably ~April 12/13)
 - Will target a couple hours of work (\ll an assignment)
 - Will have at least 12 hrs to do, maybe 1-2 days (TBD); must hand in by April 25
- Optional **bonus assignment**: when it’s done, probably ~April 12/13, due April 25

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat
 - We learn how to predict as we go

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat
 - We learn how to predict as we go
 - Focusing on binary classification to start

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat
 - We learn how to predict as we go
 - Focusing on binary classification to start
 - Usual analysis does *not* assume a fixed distribution \mathcal{D}

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat
 - We learn how to predict as we go
 - Focusing on binary classification to start
 - Usual analysis does *not* assume a fixed distribution \mathcal{D}
 - Labels can even be chosen **adversarially**

Online learning

- Class so far has been in the **(passive) batch setting**:
 - Observe training set $S \sim \mathcal{D}^n$, pick h , test on new examples from \mathcal{D}
- Today: the **online** setting
 - See an x_t , make a prediction \hat{y}_t , see true label y_t , repeat
 - We learn how to predict as we go
 - Focusing on binary classification to start
 - Usual analysis does *not* assume a fixed distribution \mathcal{D}
 - Labels can even be chosen **adversarially**

Hello Danica,

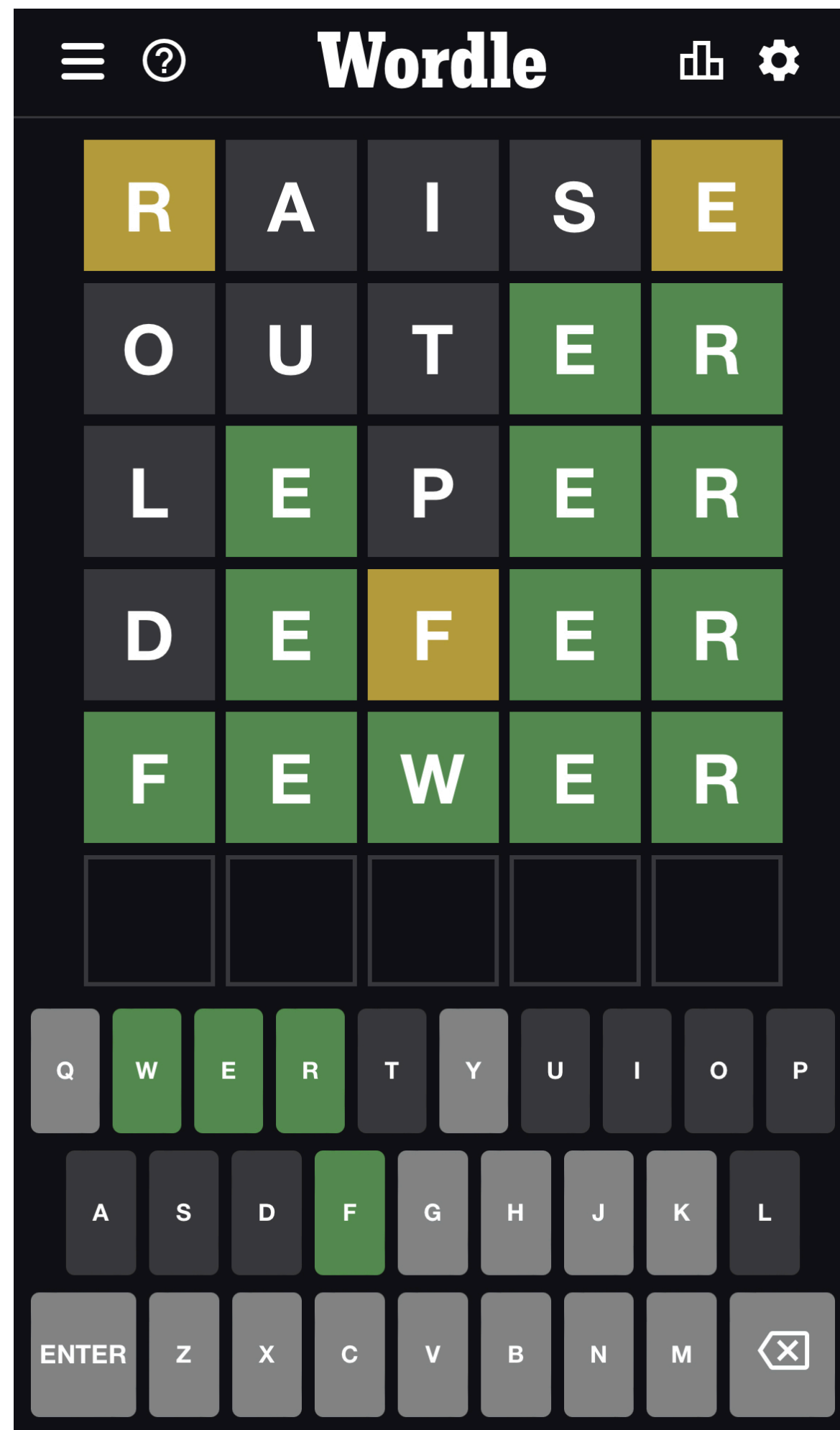
I am incredibly sorry about this! It looks like the earlier CMT emails went to my spam folder. I can do this review within the next 12 hours (i.e. by midnight

Realizable online setting

- **Realizable** setting: labels y_t have to be consistent with some $h^* \in \mathcal{H}$

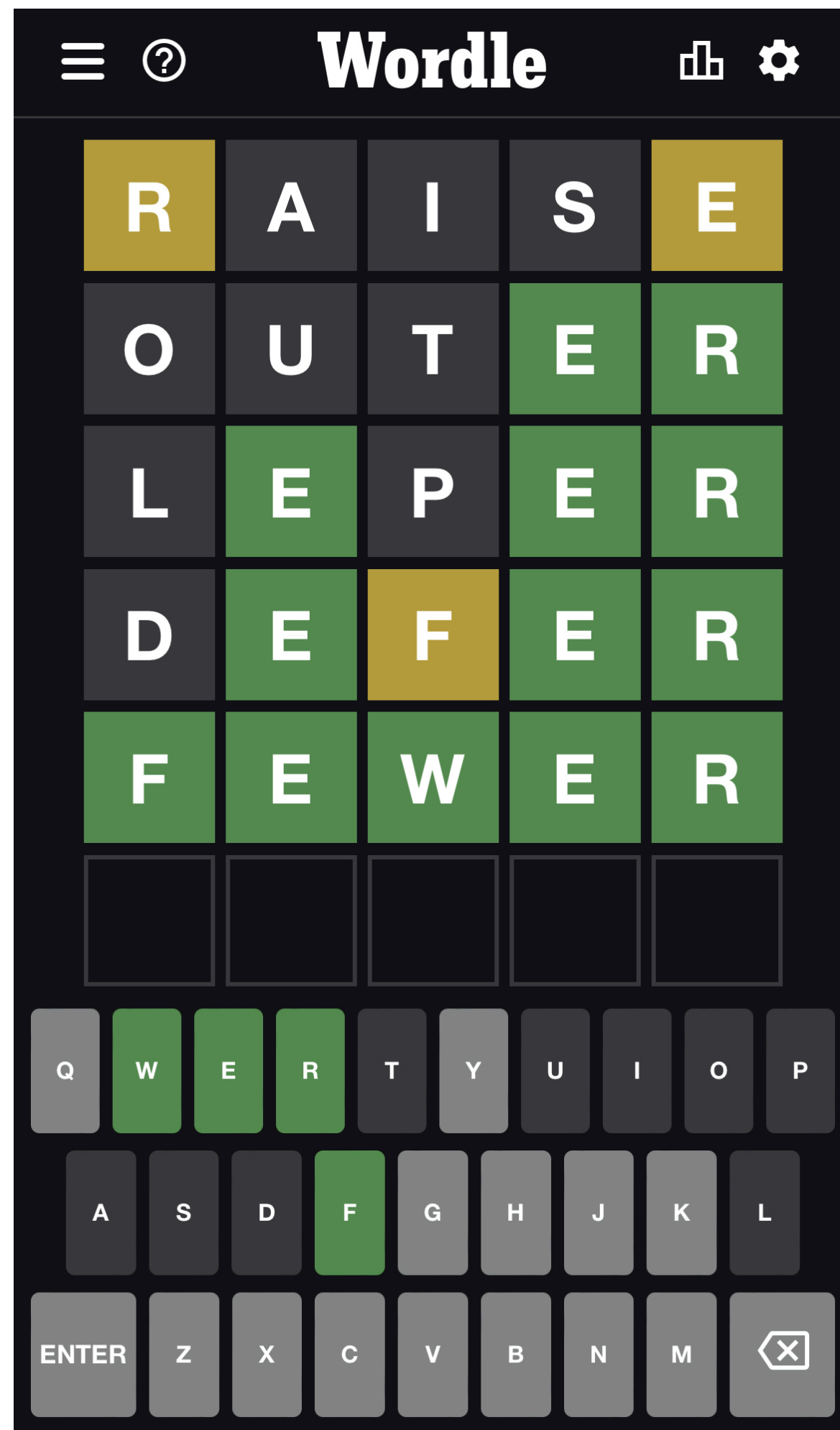
Realizable online setting

- **Realizable** setting: labels y_t have to be consistent with some $h^* \in \mathcal{H}$



Realizable online setting

- **Realizable** setting: labels y_t have to be consistent with some $h^* \in \mathcal{H}$

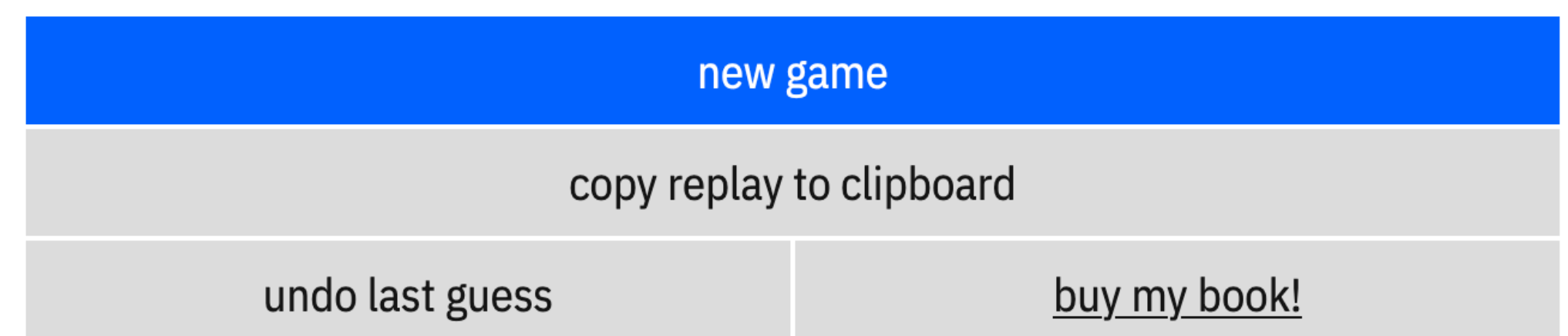


ABSURDLE by [gntm](#)



R	A	I	S	E
P	O	U	T	Y
W	O	O	L	Y
F	O	L	L	Y
J	O	L	L	Y
H	O	L	L	Y
D	O	L	L	Y
G	O	L	L	Y

You guessed successfully in 8 guesses!



Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$
- If \mathcal{H} is finite, consider the algorithm Consistent (basically ERM):

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$
- If \mathcal{H} is finite, consider the algorithm Consistent (basically ERM):
 - Start with the **version space** $V_1 = \mathcal{H}$

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$
- If \mathcal{H} is finite, consider the algorithm Consistent (basically ERM):
 - Start with the **version space** $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t = h(x_t)$ for any arbitrary $h \in V_t$

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$
- If \mathcal{H} is finite, consider the algorithm Consistent (basically ERM):
 - Start with the **version space** $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t = h(x_t)$ for any arbitrary $h \in V_t$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

Mistake bounds

- Take a sequence $S = ((x_1, h^*(x_1)), \dots, (x_T, h^*(x_T)))$
- $M_A(S)$ is the number of **mistakes** the algorithm A makes on S
- $M_A(\mathcal{H})$ is the **worst-case** number of mistakes for **any** S with labels in \mathcal{H}
- \mathcal{H} is **online learnable** if there's an A with $M_A(\mathcal{H}) < \infty$
- If \mathcal{H} is finite, consider the algorithm Consistent (basically ERM):
 - Start with the **version space** $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t = h(x_t)$ for any arbitrary $h \in V_t$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Have mistake bound $M_{\text{Consistent}}(\mathcal{H}) \leq |\mathcal{H}| - 1$

A smarter algorithm for finite, realizable \mathcal{H}

- If Consistent made a mistake, we might only remove **one** h from V_t
- Better algorithm can always either (a) be right or (b) make lots of progress

A smarter algorithm for finite, realizable \mathcal{H}

- If Consistent made a mistake, we might only remove **one** h from V_t
- Better algorithm can always either (a) be right or (b) make lots of progress
- Halving:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \left| \{h \in V_t : h(x_t) = r\} \right|$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

A smarter algorithm for finite, realizable \mathcal{H}

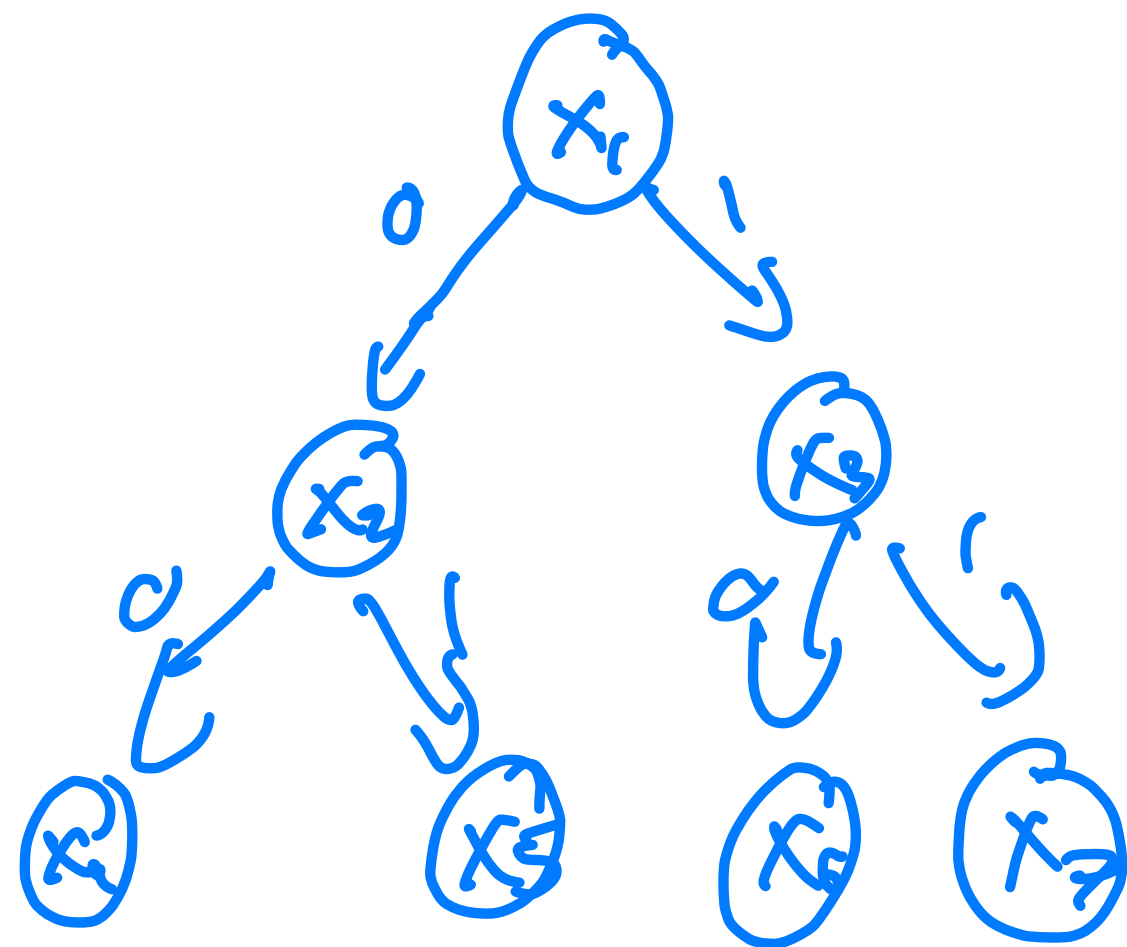
- If Consistent made a mistake, we might only remove **one** h from V_t
- Better algorithm can always either (a) be right or (b) make lots of progress
- Halving:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \left| \{h \in V_t : h(x_t) = r\} \right|$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
 - If we were wrong, we removed **at least half** of V_t

A smarter algorithm for finite, realizable \mathcal{H}

- If Consistent made a mistake, we might only remove **one** h from V_t
- Better algorithm can always either (a) be right or (b) make lots of progress
- Halving:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \left| \{h \in V_t : h(x_t) = r\} \right|$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- If we were wrong, we removed **at least half** of V_t
- $M_{\text{Halving}}(\mathcal{H}) \leq \log_2 |\mathcal{H}|$ – way better bound

Online learnability

- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1



Online learnability

- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1
- \mathcal{H} **shatters a tree** if everywhere in the tree is reached by some $h \in \mathcal{H}$

Online learnability

- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1
- \mathcal{H} **shatters a tree** if everywhere in the tree is reached by some $h \in \mathcal{H}$
- The **Littlestone dimension** $\text{Ldim}(\mathcal{H})$ is the max depth of any tree \mathcal{H} shatters

Online learnability

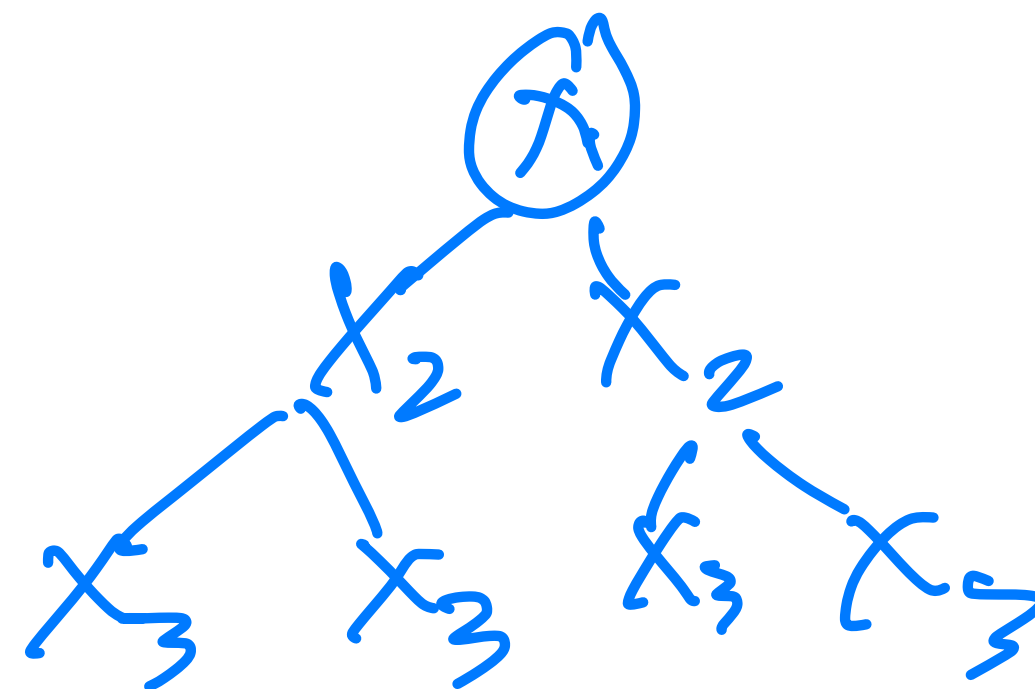
- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1
- \mathcal{H} **shatters a tree** if everywhere in the tree is reached by some $h \in \mathcal{H}$
- The **Littlestone dimension** $\text{Ldim}(\mathcal{H})$ is the max depth of any tree \mathcal{H} shatters
- Any algorithm A must have $M_A(\mathcal{H}) \geq \text{Ldim}(\mathcal{H})$

Online learnability

- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1
- \mathcal{H} **shatters a tree** if everywhere in the tree is reached by some $h \in \mathcal{H}$
- The **Littlestone dimension** $\text{Ldim}(\mathcal{H})$ is the max depth of any tree \mathcal{H} shatters
- Any algorithm A must have $M_A(\mathcal{H}) \geq \text{Ldim}(\mathcal{H})$
- If \mathcal{H} can shatter a set, it can shatter any tree from that set

Online learnability

- Think about the **game tree** for the learner and the **adversary**
 - Put points $x_t \in \mathcal{X}$ into a full binary tree
 - Start at the root, move left if learner predicts 0, right if it predicts 1
- \mathcal{H} **shatters a tree** if everywhere in the tree is reached by some $h \in \mathcal{H}$
- The **Littlestone dimension** $\text{Ldim}(\mathcal{H})$ is the max depth of any tree \mathcal{H} shatters
- Any algorithm A must have $M_A(\mathcal{H}) \geq \text{Ldim}(\mathcal{H})$
- If \mathcal{H} can shatter a set, it can shatter any tree from that set
 - $\text{VCdim}(\mathcal{H}) \leq \text{Ldim}(\mathcal{H})$



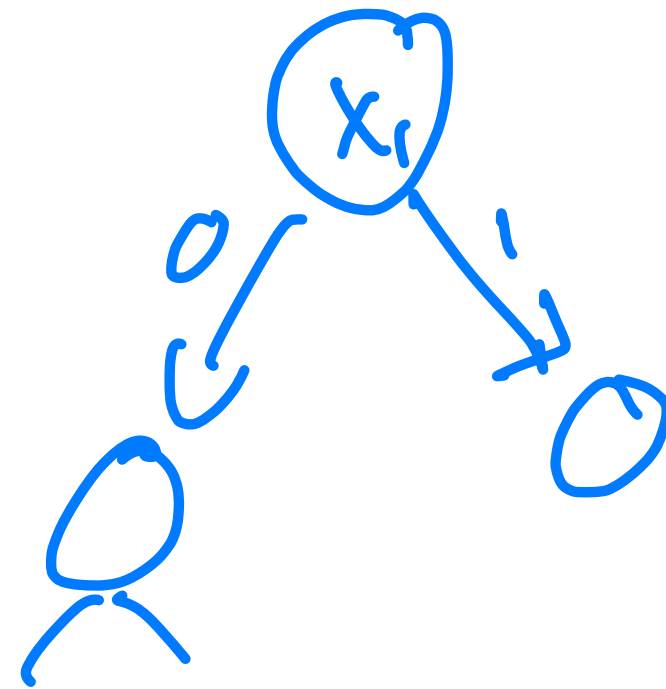
Littlestone dimension examples

- If \mathcal{H} is finite, can't shatter a full tree deeper than $\log_2 |\mathcal{H}|$

Littlestone dimension examples

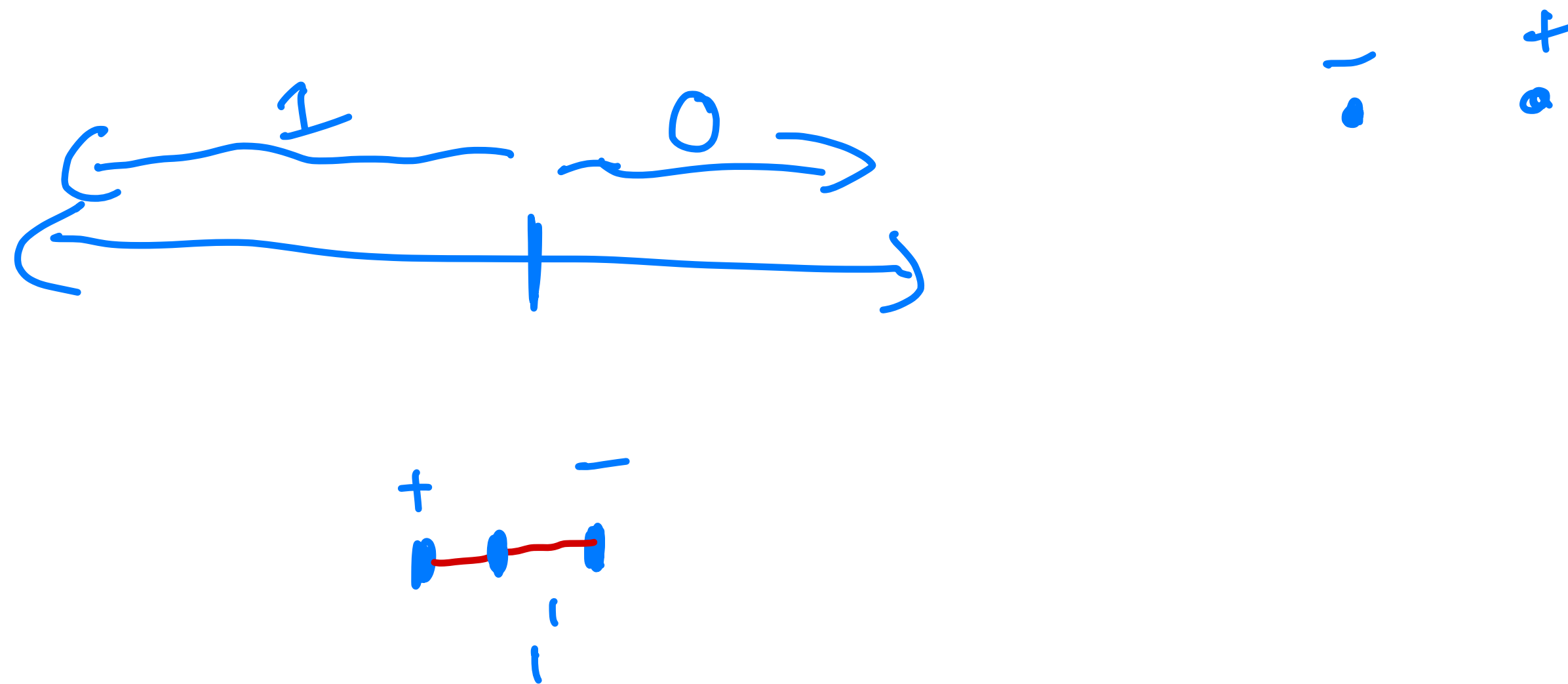
- If \mathcal{H} is finite, can't shatter a full tree deeper than $\log_2 |\mathcal{H}|$
- If $\mathcal{X} = [d]$, $\mathcal{H} = \{x \mapsto \mathbb{I}(x = i) : i \in [d]\}$, have $\text{Ldim}(\mathcal{H}) = 1$

$\{1, 2, 3\}$



Littlestone dimension examples

- If \mathcal{H} is finite, can't shatter a full tree deeper than $\log_2 |\mathcal{H}|$
- If $\mathcal{X} = [d]$, $\mathcal{H} = \{x \mapsto \mathbb{I}(x = i) : i \in [d]\}$, have $\text{Ldim}(\mathcal{H}) = 1$
- If $\mathcal{X} = [0,1]$ and $\mathcal{H} = \{x \mapsto \mathbb{I}(x \leq a) : a \in [0,1]\}$, have $\text{Ldim}(\mathcal{H}) = \infty$ (!)



Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$

Standard Optimal Algorithm

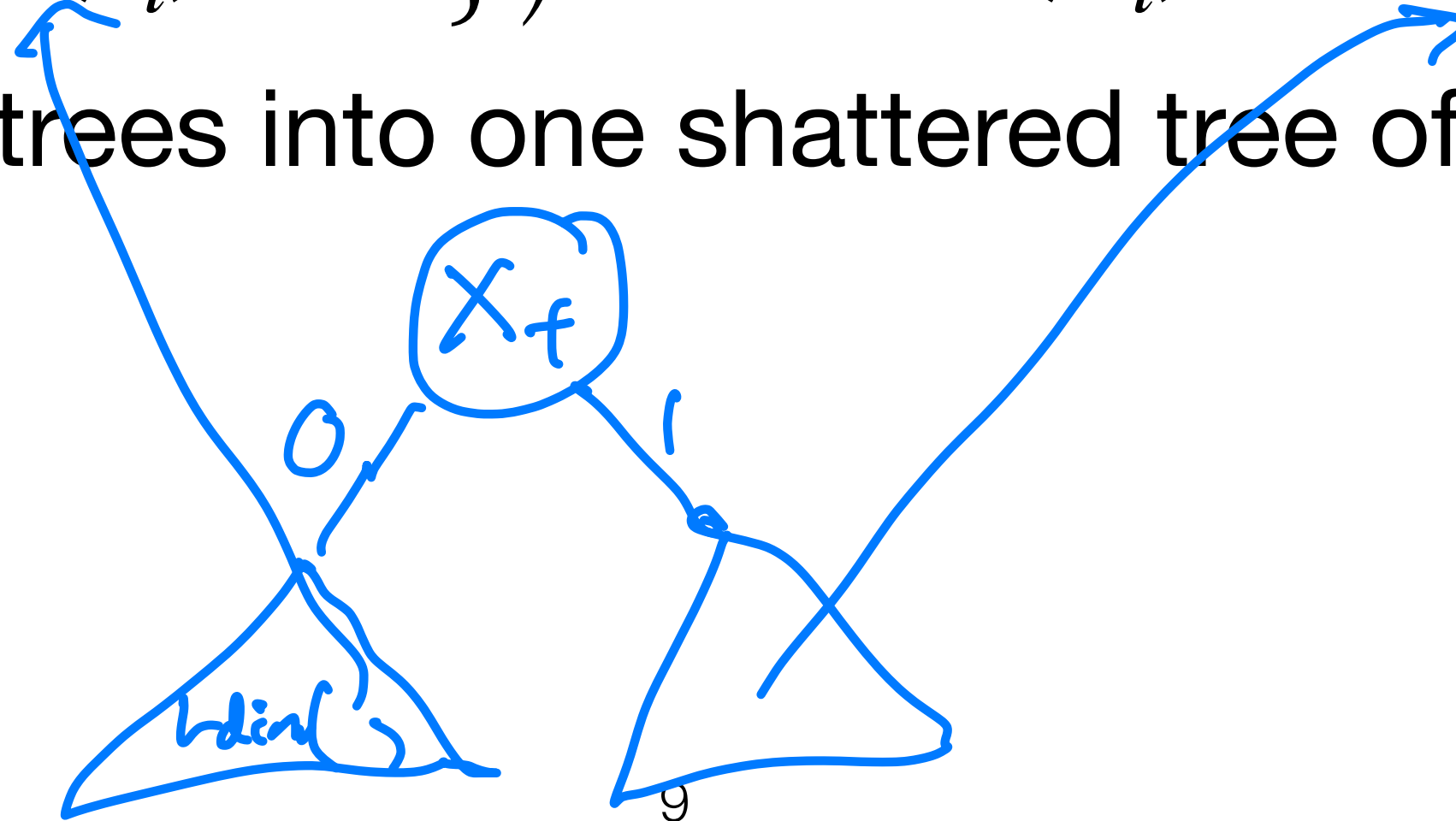
- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:

Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:
 - If not, $\operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 0\} \right) = \operatorname{Ldim}(V_t) = \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 1\} \right)$

Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:
 - If not, $\operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 0\} \right) = \operatorname{Ldim}(V_t) = \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 1\} \right)$
 - Then combine shattered trees into one shattered tree of depth $\operatorname{Ldim}(V_t) + 1$



Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:
 - If not, $\operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 0\} \right) = \operatorname{Ldim}(V_t) = \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 1\} \right)$
 - Then combine shattered trees into one shattered tree of depth $\operatorname{Ldim}(V_t) + 1$
 - But then $\operatorname{Ldim}(V_t) = \operatorname{Ldim}(V_t) + 1 \dots$ contradiction

Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:
 - If not, $\operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 0\} \right) = \operatorname{Ldim}(V_t) = \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 1\} \right)$
 - Then combine shattered trees into one shattered tree of depth $\operatorname{Ldim}(V_t) + 1$
 - But then $\operatorname{Ldim}(V_t) = \operatorname{Ldim}(V_t) + 1 \dots$ contradiction
- Thus $M_{\text{SOA}}(\mathcal{H}) = \operatorname{Ldim}(\mathcal{H})$, the best possible mistake bound

Standard Optimal Algorithm

- Like Halving, but tries to reduce Littlestone dimension instead of cardinality:
 - Start with the version space $V_1 = \mathcal{H}$
 - Given x_t , predict $\hat{y}_t \in \operatorname{argmax}_{r \in \{0,1\}} \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = r\} \right)$
 - Seeing y_t , update $V_{t+1} = \{h \in V_t : h(x_t) = y_t\}$
- Whenever we make a mistake, $\operatorname{Ldim}(V_{t+1}) \leq \operatorname{Ldim}(V_t) - 1$:
 - If not, $\operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 0\} \right) = \operatorname{Ldim}(V_t) = \operatorname{Ldim} \left(\{h \in V_t : h(x_t) = 1\} \right)$
 - Then combine shattered trees into one shattered tree of depth $\operatorname{Ldim}(V_t) + 1$
 - But then $\operatorname{Ldim}(V_t) = \operatorname{Ldim}(V_t) + 1 \dots$ contradiction
- Thus $M_{\text{SOA}}(\mathcal{H}) = \operatorname{Ldim}(\mathcal{H})$, the best possible mistake bound
- But SOA is not necessarily easy to compute!

(pause)

You could use this time to do
<https://seoi.ubc.ca/surveys>
(if you want)

Unrealizable online learning

- In the batch setting:
 - Realizable PAC assumes labels come from $h^* \in \mathcal{H}$
 - Agnostic PAC just has us **compete with** best $h^* \in \mathcal{H}$
- In the online setting:
 - Realizable assumes labels come from $h^* \in \mathcal{H}$

Unrealizable online learning

- In the batch setting:
 - Realizable PAC assumes labels come from $h^* \in \mathcal{H}$
 - Agnostic PAC just has us **compete with** best $h^* \in \mathcal{H}$
- In the online setting:
 - Realizable assumes labels come from $h^* \in \mathcal{H}$
 - Unrealizable has us compete with best $h^* \in \mathcal{H}$

Unrealizable online learning

- In the batch setting:
 - Realizable PAC assumes labels come from $h^* \in \mathcal{H}$
 - Agnostic PAC just has us **compete with** best $h^* \in \mathcal{H}$
- In the online setting:
 - Realizable assumes labels come from $h^* \in \mathcal{H}$
 - Unrealizable has us compete with best $h^* \in \mathcal{H}$

$$\text{Regret}_A(h, T) = \sup_{(x_1, y_1), \dots, (x_T, y_T)} \left[\sum_{t=1}^T |\hat{y}_t - y_t| - \sum_{t=1}^T |h(x_t) - y_t| \right]$$

$$\hat{y}_t = A((x_1, y_1), \dots, (x_{t-1}, y_{t-1}), x_t)$$

Unrealizable online learning

- In the batch setting:
 - Realizable PAC assumes labels come from $h^* \in \mathcal{H}$
 - Agnostic PAC just has us **compete with** best $h^* \in \mathcal{H}$
- In the online setting:
 - Realizable assumes labels come from $h^* \in \mathcal{H}$
 - Unrealizable has us compete with best $h^* \in \mathcal{H}$

$$\text{Regret}_A(h, T) = \sup_{(x_1, y_1), \dots, (x_T, y_T)} \left[\sum_{t=1}^T |\hat{y}_t - y_t| - \sum_{t=1}^T |h(x_t) - y_t| \right]$$

$$\text{Regret}_A(\mathcal{H}, T) = \sup_{h \in \mathcal{H}} \text{Regret}_A(h, T)$$

Unrealizable online learning

- In the batch setting:
 - Realizable PAC assumes labels come from $h^* \in \mathcal{H}$
 - Agnostic PAC just has us **compete with** best $h^* \in \mathcal{H}$
- In the online setting:
 - Realizable assumes labels come from $h^* \in \mathcal{H}$
 - Unrealizable has us compete with best $h^* \in \mathcal{H}$

$$\text{Regret}_A(h, T) = \sup_{(x_1, y_1), \dots, (x_T, y_T)} \left[\sum_{t=1}^T |\hat{y}_t - y_t| - \sum_{t=1}^T |h(x_t) - y_t| \right]$$

$$\text{Regret}_A(\mathcal{H}, T) = \sup_{h \in \mathcal{H}} \text{Regret}_A(h, T)$$

- Ideally, we want **sublinear regret**: $\frac{1}{T} \text{Regret}_A(\mathcal{H}, T) \xrightarrow{T \rightarrow \infty} 0$

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes



Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes



Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$



Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$
- To avoid this:

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$
- To avoid this:
 - Learner has **random** prediction, $\Pr(\hat{y}_t = 1) = p_t$

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$
- To avoid this:
 - Learner has **random** prediction, $\Pr(\hat{y}_t = 1) = p_t$
 - Adversary commits to y_t without knowing the roll

Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$
- To avoid this:
 - Learner has **random** prediction, $\Pr(\hat{y}_t = 1) = p_t$
 - Adversary commits to y_t without knowing the roll



Regret: impossible to avoid

- Regret: “how much better it would have been to just play $h(x_t)$ every time”
- Consider $\mathcal{H} = \{x \mapsto 0, x \mapsto 1\}$
 - Adversary could always just say “no, you’re wrong” and get T mistakes
 - For any sequence of true y_t , either $x \mapsto 0$ or $x \mapsto 1$ has $\leq \frac{T}{2}$ mistakes
 - So regret would be at least $T - \frac{T}{2} = \frac{T}{2}$
- To avoid this:
 - Learner has **random** prediction, $\Pr(\hat{y}_t = 1) = p_t$
 - Adversary commits to y_t without knowing the roll
 - Measure **expected** loss $\Pr(\hat{y}_t \neq y_t) = |p_t - y_t|$



Low regret for online classification

- For every \mathcal{H} , there's an algorithm with

$$\text{Regret}_A(\mathcal{H}, T) \leq \sqrt{2 \min \left(\log |\mathcal{H}|, (1 + \log T) \text{Ldim}(\mathcal{H}) \right) T}$$

- Also a lower bound of $\Omega \left(\sqrt{\text{Ldim}(\mathcal{H}) T} \right)$
- Based on Weighted-Majority algorithm for **learning with expert advice**

Learning from expert advice

- There are d available experts who make predictions



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots$



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots, T$
 - Follow with probabilities $w_t = \tilde{w}_t / \|\tilde{w}_t\|_1$



Learning from expert advice

- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots$
 - Follow with probabilities $w_t = \tilde{w}_t / \|\tilde{w}_t\|_1$
 - Update based on costs v_t as $\tilde{w}_{t+1} = \tilde{w}_t \exp(-\eta v_t)$ (exp is elementwise)



Learning from expert advice



- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots$
 - Follow with probabilities $w_t = \tilde{w}_t / \|\tilde{w}_t\|_1$
 - Update based on costs v_t as $\tilde{w}_{t+1} = \tilde{w}_t \exp(-\eta v_t)$ (exp is elementwise)
- **Theorem** (SSBD 21.11): $\sum_{t=1}^T \langle w_t, v_t \rangle - \min_{i \in [d]} \sum_{t=1}^T (v_t)_i \leq \sqrt{2 \log(d) T}$ if $T > 2 \log d$

Learning from expert advice



- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots$
 - Follow with probabilities $w_t = \tilde{w}_t / \|\tilde{w}_t\|_1$
 - Update based on costs v_t as $\tilde{w}_{t+1} = \tilde{w}_t \exp(-\eta v_t)$ (exp is elementwise)
- **Theorem** (SSBD 21.11): $\sum_{t=1}^T \langle w_t, v_t \rangle - \min_{i \in [d]} \sum_{t=1}^T (v_t)_i \leq \sqrt{2 \log(d) T}$ if $T > 2 \log d$
- Can avoid knowing T by *doubling trick*: run for $T = 1, T = 2, T = 4, \dots$ sequentially

Learning from expert advice



- There are d available experts who make predictions
- At time t , learner chooses to follow expert i with probability $(w_t)_i$
- Sees potential costs $v_t \in \mathbb{R}^d$; pays expectation $\langle w_t, v_t \rangle$
- Weighted-Majority algorithm:
 - Start with $\tilde{w}_1 = (1, \dots, 1)$; $\eta = \sqrt{2 \log(d) / T}$
 - For $t = 1, 2, \dots$
 - Follow with probabilities $w_t = \tilde{w}_t / \|\tilde{w}_t\|_1$
 - Update based on costs v_t as $\tilde{w}_{t+1} = \tilde{w}_t \exp(-\eta v_t)$ (exp is elementwise)
- **Theorem** (SSBD 21.11): $\sum_{t=1}^T \langle w_t, v_t \rangle - \min_{i \in [d]} \sum_{t=1}^T (v_t)_i \leq \sqrt{2 \log(d) T}$ if $T > 2 \log d$
- Can avoid knowing T by *doubling trick*: run for $T = 1, T = 2, T = 4, \dots$ sequentially
 - Only blows up regret by $< 3.5x$ (SSBD exercise 21.4)

Low regret for online classification

- For finite \mathcal{H} , we can just run Weighted-Majority with each $h \in \mathcal{H}$

Low regret for online classification

- For finite \mathcal{H} , we can just run Weighted-Majority with each $h \in \mathcal{H}$
 - Plugging in previous theorem, $\text{Regret}_{\text{WM}}(\mathcal{H}, T) \leq \sqrt{2 \log |\mathcal{H}| T}$

Low regret for online classification

- For finite \mathcal{H} , we can just run Weighted-Majority with each $h \in \mathcal{H}$
 - Plugging in previous theorem, $\text{Regret}_{\text{WM}}(\mathcal{H}, T) \leq \sqrt{2 \log |\mathcal{H}| T}$
- For infinite \mathcal{H} , we need a not-too-big set of experts where one is still good

Low regret for online classification

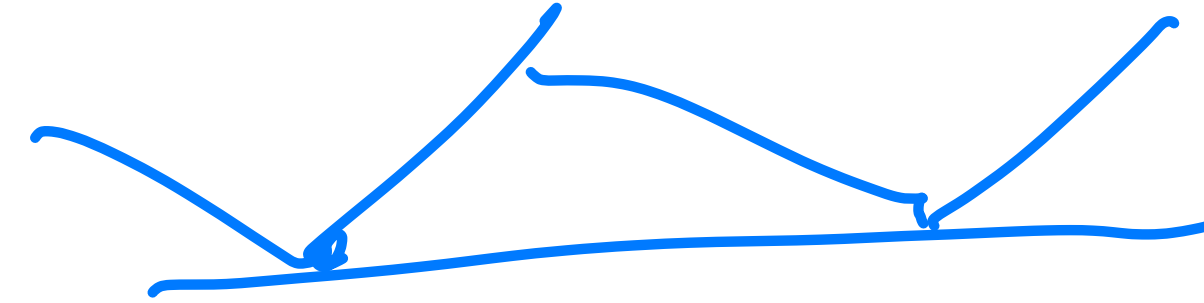
- For finite \mathcal{H} , we can just run Weighted-Majority with each $h \in \mathcal{H}$
 - Plugging in previous theorem, $\text{Regret}_{\text{WM}}(\mathcal{H}, T) \leq \sqrt{2 \log |\mathcal{H}| T}$
- For infinite \mathcal{H} , we need a not-too-big set of experts where one is still good
 - Expert(i_1, i_2, \dots, i_L) runs SOA on x_1, \dots, x_T ,
but takes choice with smaller Ldim on indices i_1, i_2, \dots, i_L

Low regret for online classification

- For finite \mathcal{H} , we can just run Weighted-Majority with each $h \in \mathcal{H}$
 - Plugging in previous theorem, $\text{Regret}_{\text{WM}}(\mathcal{H}, T) \leq \sqrt{2 \log |\mathcal{H}| T}$
- For infinite \mathcal{H} , we need a not-too-big set of experts where one is still good
 - Expert(i_1, i_2, \dots, i_L) runs SOA on x_1, \dots, x_T ,
but takes choice with smaller Ldim on indices i_1, i_2, \dots, i_L
 - Can show (21.13-14) that one expert is as good as the best $h \in \mathcal{H}$,
and there aren't too many of them,
giving $\text{Regret}_A(\mathcal{H}, T) \leq \sqrt{2(1 + \log T) \text{Ldim}(\mathcal{H}) T}$

Online convex optimization

- **Online convex optimization** is
 - Convex hypothesis class \mathcal{H}
 - At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$



Online convex optimization

- **Online convex optimization** is

- Convex hypothesis class \mathcal{H}

- At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$

- $$\text{Regret}(w^*, T) = \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*), \quad \text{Regret}(\mathcal{H}, T) = \sup_{w^* \in \mathcal{H}} \text{Regret}(w^*, T)$$

Online convex optimization

- **Online convex optimization** is
 - Convex hypothesis class \mathcal{H}
 - At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$
- $\text{Regret}(w^*, T) = \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*), \quad \text{Regret}(\mathcal{H}, T) = \sup_{w^* \in \mathcal{H}} \text{Regret}(w^*, T)$
- **Online gradient descent** (exactly like SGD) has:

Online convex optimization

- **Online convex optimization** is

- Convex hypothesis class \mathcal{H}

- At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$

- $$\text{Regret}(w^*, T) = \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*), \quad \text{Regret}(\mathcal{H}, T) = \sup_{w^* \in \mathcal{H}} \text{Regret}(w^*, T)$$

- **Online gradient descent** (exactly like SGD) has:

- $$\text{Regret}(w^*, T) \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \quad \text{where } v_t \in \partial \ell_t(w_t) \text{ are step directions}$$

Online convex optimization

- **Online convex optimization** is

- Convex hypothesis class \mathcal{H}

- At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$

- $$\text{Regret}(w^*, T) = \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*), \quad \text{Regret}(\mathcal{H}, T) = \sup_{w^* \in \mathcal{H}} \text{Regret}(w^*, T)$$

- **Online gradient descent** (exactly like SGD) has:

- $$\text{Regret}(w^*, T) \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \quad \text{where } v_t \in \partial \ell_t(w_t) \text{ are step directions}$$

- $$\text{Regret}(w^*, T) \leq \frac{1}{2} (\|w^*\|^2 + \rho^2) \sqrt{T} \quad \text{if } \ell_t \text{ are } \rho\text{-Lipschitz, } \eta = 1/\sqrt{T}$$

Online convex optimization

- **Online convex optimization** is

- Convex hypothesis class \mathcal{H}

- At each step: learner picks $w_t \in \mathcal{H}$, environment picks convex loss $\ell_t(w_t)$

- $$\text{Regret}(w^*, T) = \sum_{t=1}^T \ell_t(w_t) - \sum_{t=1}^T \ell_t(w^*), \quad \text{Regret}(\mathcal{H}, T) = \sup_{w^* \in \mathcal{H}} \text{Regret}(w^*, T)$$

- **Online gradient descent** (exactly like SGD) has:

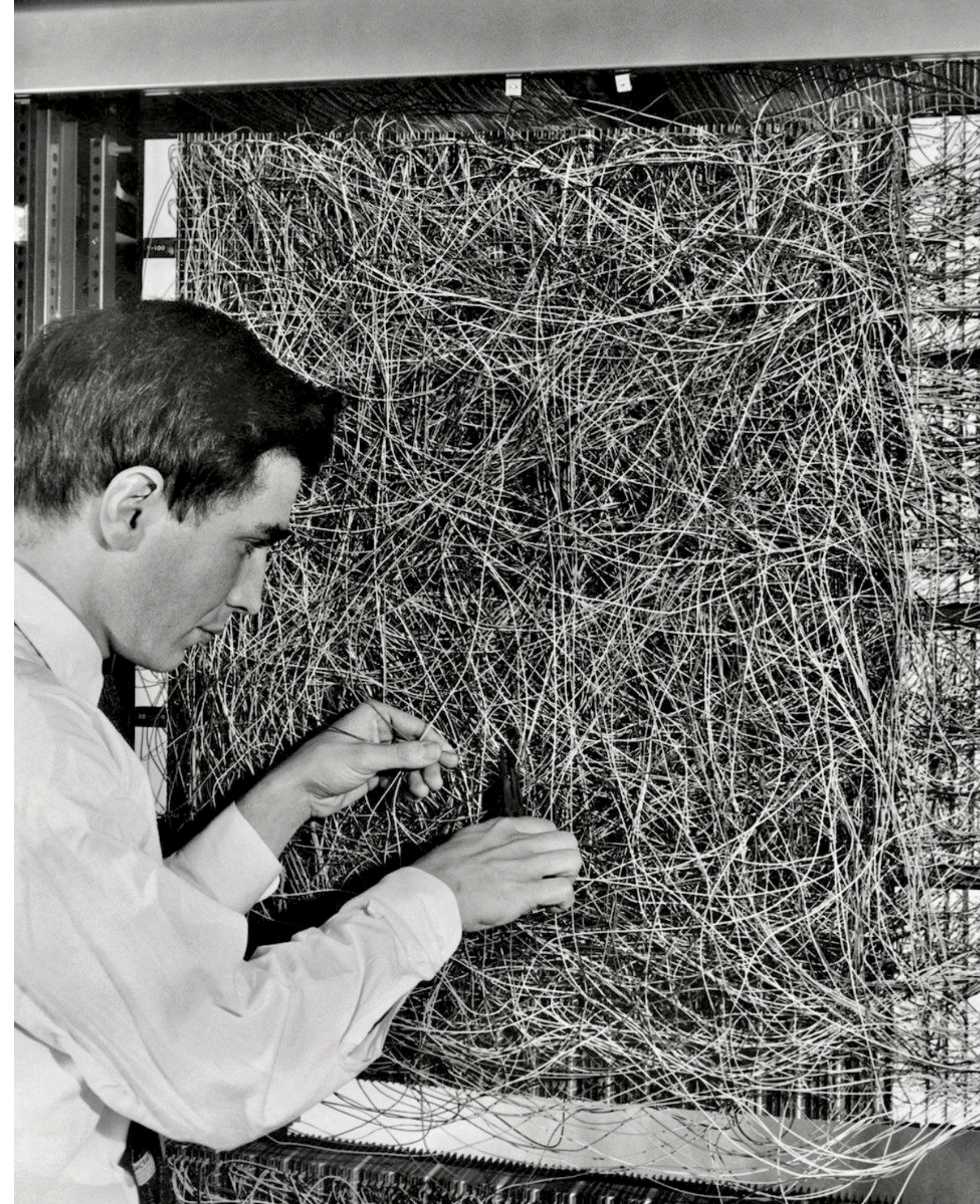
- $$\text{Regret}(w^*, T) \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \quad \text{where } v_t \in \partial \ell_t(w_t) \text{ are step directions}$$

- $$\text{Regret}(w^*, T) \leq \frac{1}{2} (\|w^*\|^2 + \rho^2) \sqrt{T} \quad \text{if } \ell_t \text{ are } \rho\text{-Lipschitz, } \eta = 1/\sqrt{T}$$

- $$\text{Regret}(w^*, T) \leq B\rho\sqrt{T} \quad \text{if } \ell_t \text{ are } \rho\text{-Lipschitz, } \mathcal{H} \text{ is } B\text{-bounded, } \eta = B/(\rho\sqrt{T})$$

Online Perceptron

- You learned about Batch Perceptron in HW3
- Original algorithm is online
- Essentially identical, just only update on mistake
- Corresponds to online gradient descent on hinge loss
- Get same $(R/\gamma)^2$ margin-based mistake bound
 - $L_{\text{dim}} = \infty$ without the margin condition



Online-to-batch conversion

- If we have a good online algorithm, we have a good batch algorithm:
just run it on the batch

Online-to-batch conversion

- If we have a good online algorithm, we have a good batch algorithm: just run it on the batch
- MRT **Lemma** 8.14: If $S \sim \mathcal{D}^T$ gives h_1, \dots, h_T for $0 \leq \ell(h, (x, y)) \leq M$,

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(h_t) \leq \frac{1}{T} \sum_{t=1}^T \ell(h_t(x_t), y_t) + M \sqrt{\frac{2}{T} \log \frac{1}{\delta}}$$

Online-to-batch conversion

- If we have a good online algorithm, we have a good batch algorithm: just run it on the batch
- MRT **Lemma** 8.14: If $S \sim \mathcal{D}^T$ gives h_1, \dots, h_T for $0 \leq \ell(h, (x, y)) \leq M$,

$$\frac{1}{T} \sum_{t=1}^T L_{\mathcal{D}}(h_t) \leq \frac{1}{T} \sum_{t=1}^T \ell(h_t(x_t), y_t) + M \sqrt{\frac{2}{T} \log \frac{1}{\delta}}$$

- MRT **Theorem** 8.15: if $\ell(\cdot, z)$ is also convex,

$$L_{\mathcal{D}} \left(\frac{1}{T} \sum_{t=1}^T h_t \right) \leq \inf_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \frac{1}{T} \text{Regret}_A(\mathcal{H}, T) + 2M \sqrt{\frac{2}{T} \log \frac{2}{\delta}}$$

(pause)

You could use this time to do
<https://seoi.ubc.ca/surveys>
(if you want)

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),
 - for all subsets $H \subseteq \mathcal{H}$, $\Pr(A(S_1) \in H) \leq \exp(\epsilon) \Pr(A(S_2) \in H) + \delta$

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),
 - for all subsets $H \subseteq \mathcal{H}$, $\Pr(A(S_1) \in H) \leq \exp(\epsilon) \Pr(A(S_2) \in H) + \delta$
- Called **pure** DP if $\delta = 0$

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),
 - for all subsets $H \subseteq \mathcal{H}$, $\Pr(A(S_1) \in H) \leq \exp(\epsilon) \Pr(A(S_2) \in H) + \delta$
- Called **pure** DP if $\delta = 0$
- Used in practice (US Census, Apple, ...), tons of work on algorithms

Differential privacy

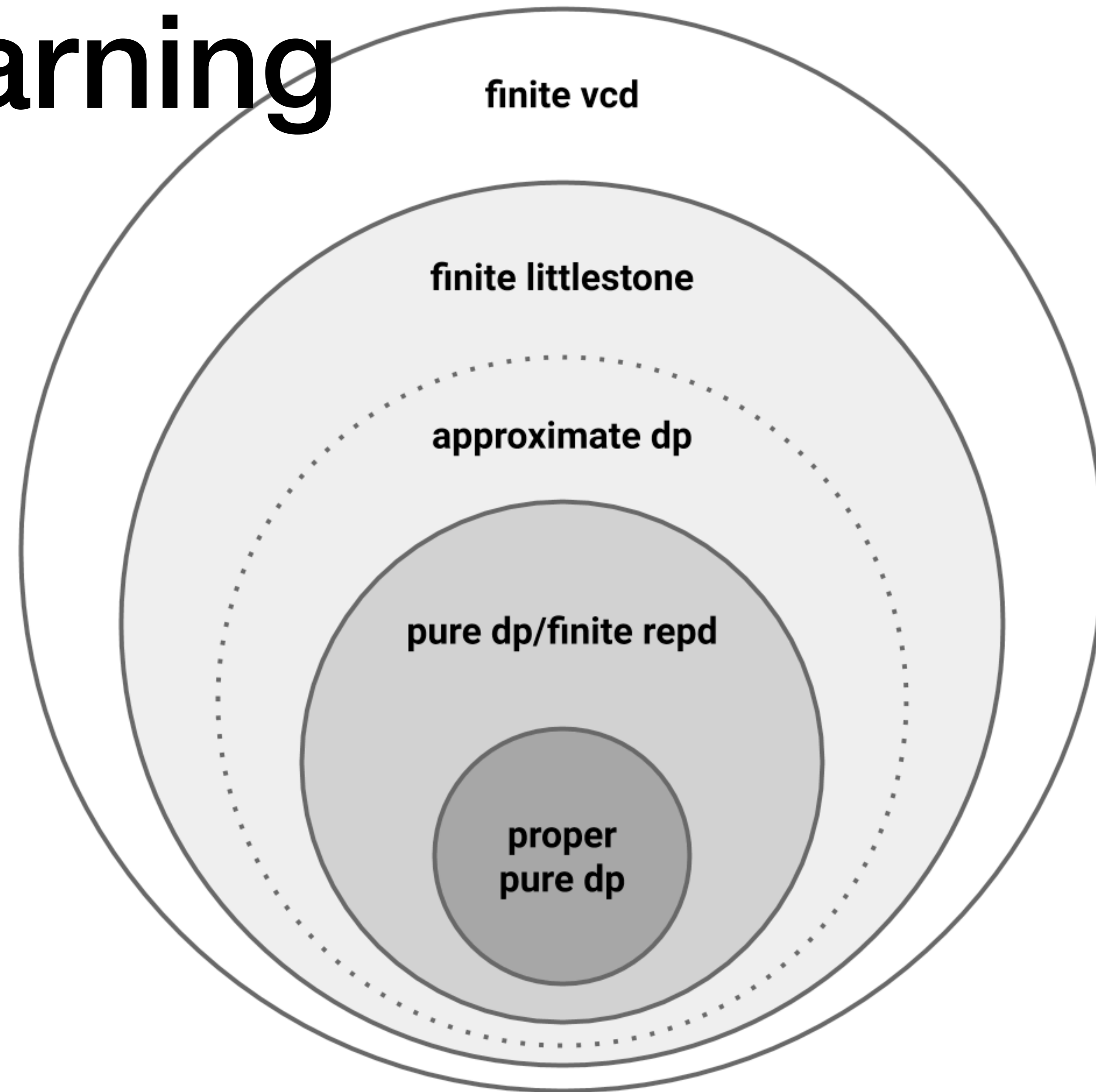
- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),
 - for all subsets $H \subseteq \mathcal{H}$, $\Pr(A(S_1) \in H) \leq \exp(\epsilon) \Pr(A(S_2) \in H) + \delta$
- Called **pure** DP if $\delta = 0$
- Used in practice (US Census, Apple, ...), tons of work on algorithms
 - Mijung Park and Mathias Lecuyer both work on this, teach courses (532P next fall, 538L now [but not next year])

Differential privacy

- Randomized learning algorithm $A(S)$ is called **(ϵ, δ) differentially private** if
 - for all S_1, S_2 that differ on a single element (i.e. one person's data),
 - for all subsets $H \subseteq \mathcal{H}$, $\Pr(A(S_1) \in H) \leq \exp(\epsilon) \Pr(A(S_2) \in H) + \delta$
- Called **pure** DP if $\delta = 0$
- Used in practice (US Census, Apple, ...), tons of work on algorithms
 - Mijung Park and Mathias Lecuyer both work on this, teach courses (532P next fall, 538L now [but not next year])
- Can be thought of as a particular **form of stability**

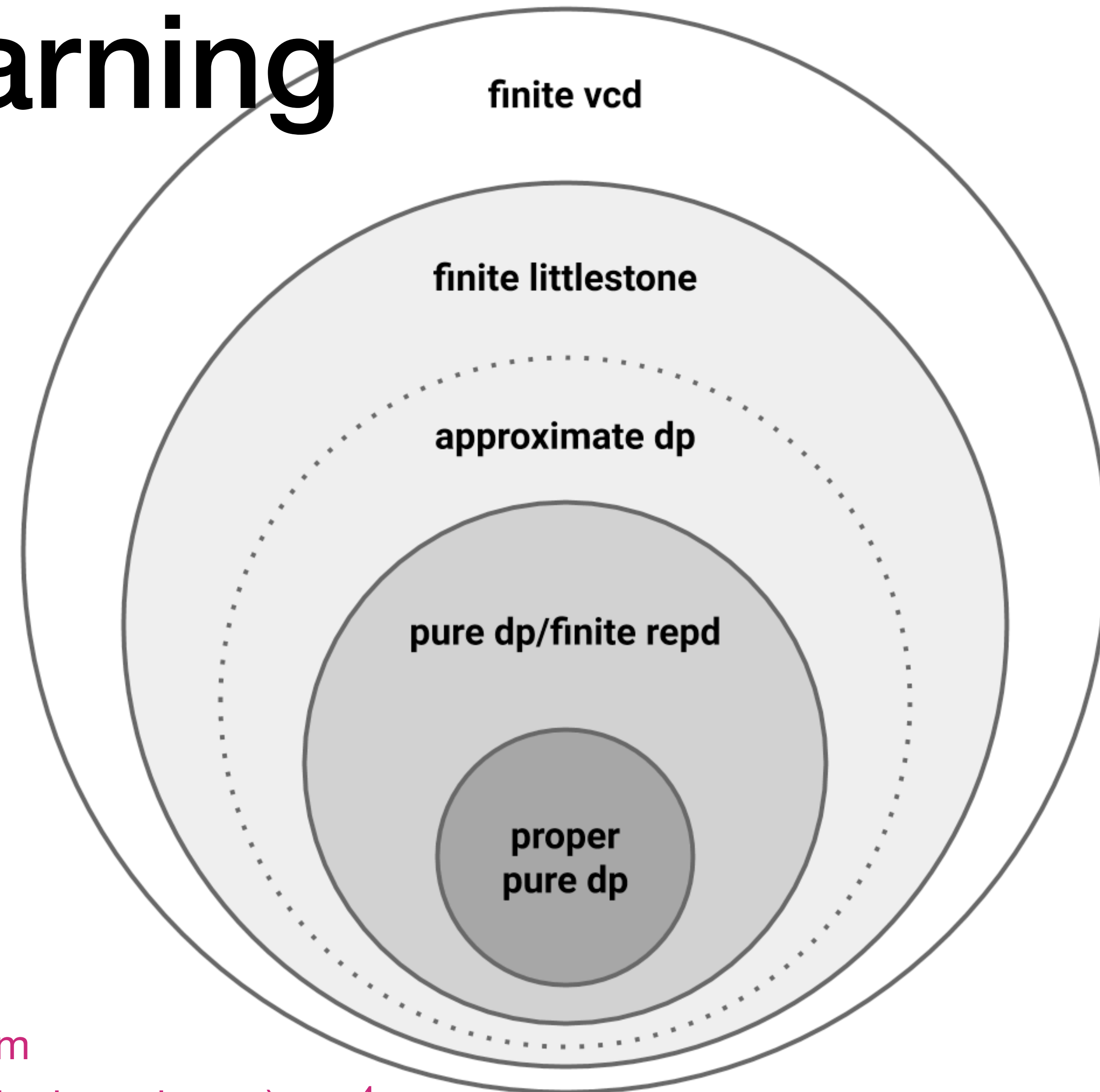
DP and online learning

- Feldman and Xiao 2014:
Pure private PAC learning takes $\Omega(\text{Ldim}(\mathcal{H}))$ samples
 - Related to communication complexity



DP and online learning

- Feldman and Xiao 2014:
Pure private PAC learning takes $\Omega(\text{Ldim}(\mathcal{H}))$ samples
 - Related to communication complexity
- Alon, Livni, Malliaris, Moran 2019:
Approximate private PAC learning takes $\Omega(\log^*(\text{Ldim}(\mathcal{H})))$ samples

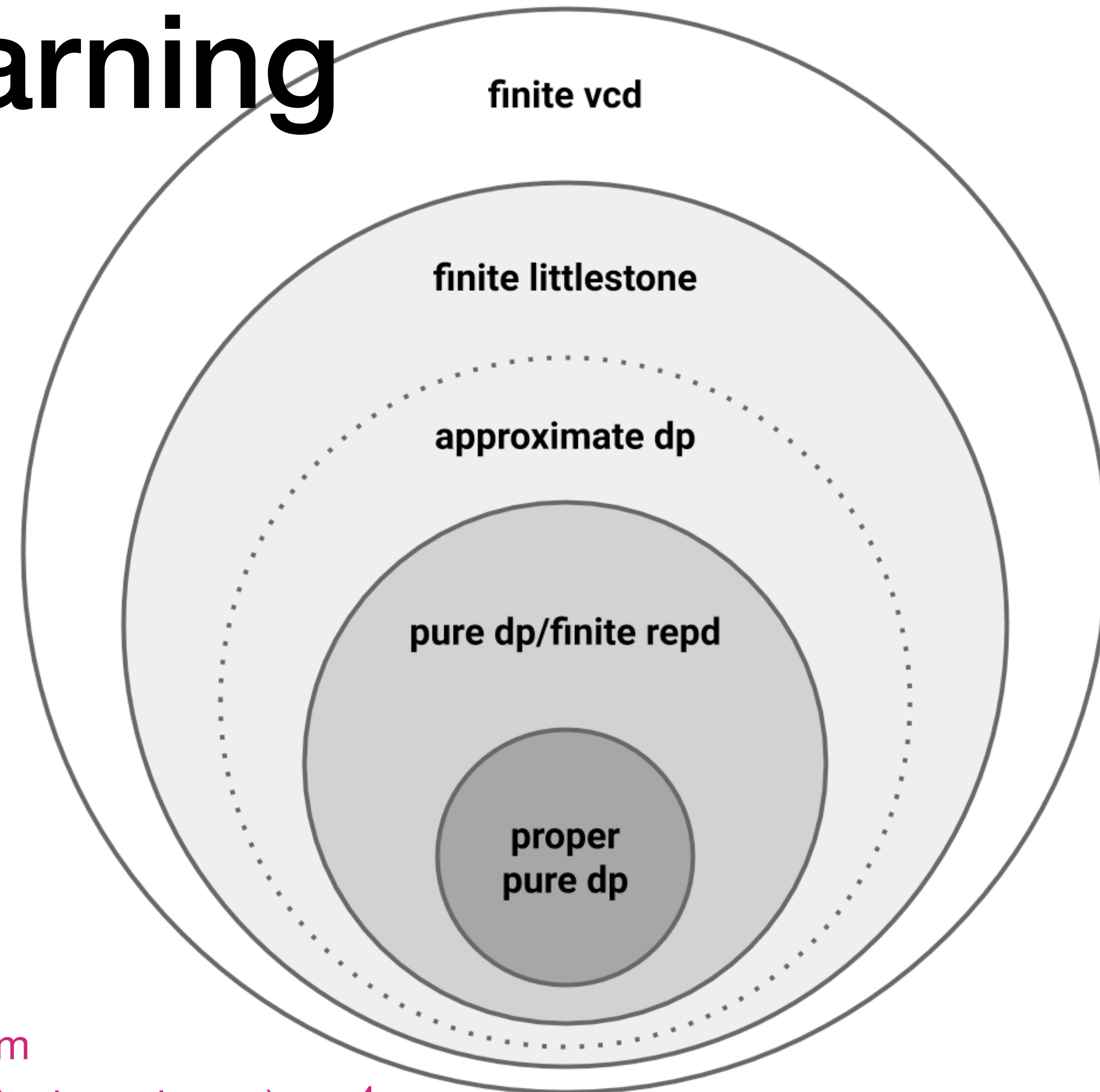


\log^* = iterated logarithm

$\log^*(\text{number of atoms in the universe}) \approx 4$

DP and online learning

- Feldman and Xiao 2014:
Pure private PAC learning takes $\Omega(\text{Ldim}(\mathcal{H}))$ samples
 - Related to communication complexity
- Alon, Livni, Malliaris, Moran 2019:
Approximate private PAC learning takes $\Omega(\log^*(\text{Ldim}(\mathcal{H})))$ samples
- Bun, Livni, Moran 2020:
Approximate private PAC learning in $2^{\mathcal{O}(\text{Ldim}(\mathcal{H}))}$ samples
 - analysis via “global stability”



\log^* = iterated logarithm

$\log^*(\text{number of atoms in the universe}) \approx 4$

DP and online learning

- Can learn differentially privately iff can learn online
 - Close connections via stability
 - But huge gap in sample and time complexity
 - Indications (Bun 2020) that converting one to the other isn't possible with polynomial time + sample complexity
 - Still a lot to understand here

Some of the stuff we didn't cover

- **Multiclass learning**: can use same techniques, need right loss
- **Ranking**: which search results are most relevant?
- **Boosting**: combine “weak learners” to a strong one (kind of like A3 Q3 b)
- **Transfer learning** / **out-of-domain generalization** / ...: train on \mathcal{D} , test on \mathcal{D}'
- Do ImageNet Classifiers Generalize to ImageNet? / The Ladder mechanism
- **Robustness**: what if we have some adversarially-corrupted training data?
- **Unsupervised learning** (just the PCA question on A1)
“How well can we ‘understand’ a data distribution?”
- **Semi-supervised learning** (just the algorithm from A4)
- **Active learning**: if x s are available but labeling them is expensive,
can we choose which to label?
- **Multi-armed bandits**: which action should I take?
- **Reinforcement learning**: interacting with an environment with hidden state
- ...