

# Computational Learning Theory (COLT)

**Goals:** Theoretical characterization of

1. Difficulty of machine learning **problems**
  - Under what conditions is learning possible and impossible?
2. Capabilities of machine learning **algorithms**
  - Under what conditions is a particular learning algorithm assured of learning successfully?

# Some easy to compute functions are not learnable

- Cryptography
  - $E_k$  where  $k$  specifies the key
  - Even if values of  $E_k$  are known for polynomially many dynamically chosen inputs
    - Computationally infeasible to deduce an algorithm for  $E_k$  or an approximation to it

---

## What General Laws govern machine and non-machine Learners?

1. Identify classes of learning problems
  - as difficult or easy, independent of learner
2. No. of training examples
  - necessary or sufficient to learn
  - How is it affected if learner can pose queries?
3. Characterize no. of mistakes learner will make before learning
4. Characterize inherent computational complexity of classes of learning problems

## Goal of COLT

- Inductively Learn a Target Function Given:
  - Only training examples of target function
  - Space of candidate hypotheses
- *Sample Complexity*: how many training examples are needed to converge with high probability to a successful hypothesis?
- *Computational Complexity*: how much computational effort is needed for learner to converge with high probability to a successful hypothesis?
- *Mistake Bound*: how many training examples will the learner misclassify before converging to a successful hypothesis?

# Two frameworks for analyzing learning algorithms

## 1. *Probably Approximately Correct (PAC)* framework

- Identify classes of hypotheses that can/cannot be learned from a polynomial number of training samples
  - Finite hypothesis space
  - Infinite hypotheses (*VC dimension*)
- Define natural measure of complexity for hypothesis spaces (*VC dimension*) that allows bounding the number of training examples required for inductive learning

## 2. *Mistake bound* framework

- Number of training errors made by a learner before it determines correct hypothesis

---

# Probably Learning an Approximately Correct Hypothesis

- Problem Setting of *Probably Approximately Correct (PAC)* learning
- Sample complexity of learning boolean valued concepts from noise-free training data

## Problem setting of PAC learning model

- $X$  = set of all possible instances over which target function may be defined  
= *set of all people*
- Each person described by a set of attributes  
*age ( young, old )*  
*height ( short, tall )*
- Target concept  $c$  in  $C$   
corresponds to some subset of  $X$   
 $c: X \rightarrow \{0, 1\}$  “people who are skiers”  
 $c(X) = 1$  if  $X$  is a positive example  
 $c(X) = 0$  if  $X$  is a negative example

## Distribution in PAC learning model

- Instances are generated at random from  $X$  according to some probability distribution  $D$
- Learner  $L$  considers some set of possible hypotheses when attempting to learn target concept
- After observing a sequence of training examples of target concept  $c$ ,  $L$  must output some hypothesis  $h$  from  $H$  which is an estimate of  $c$



## Error of a Hypothesis

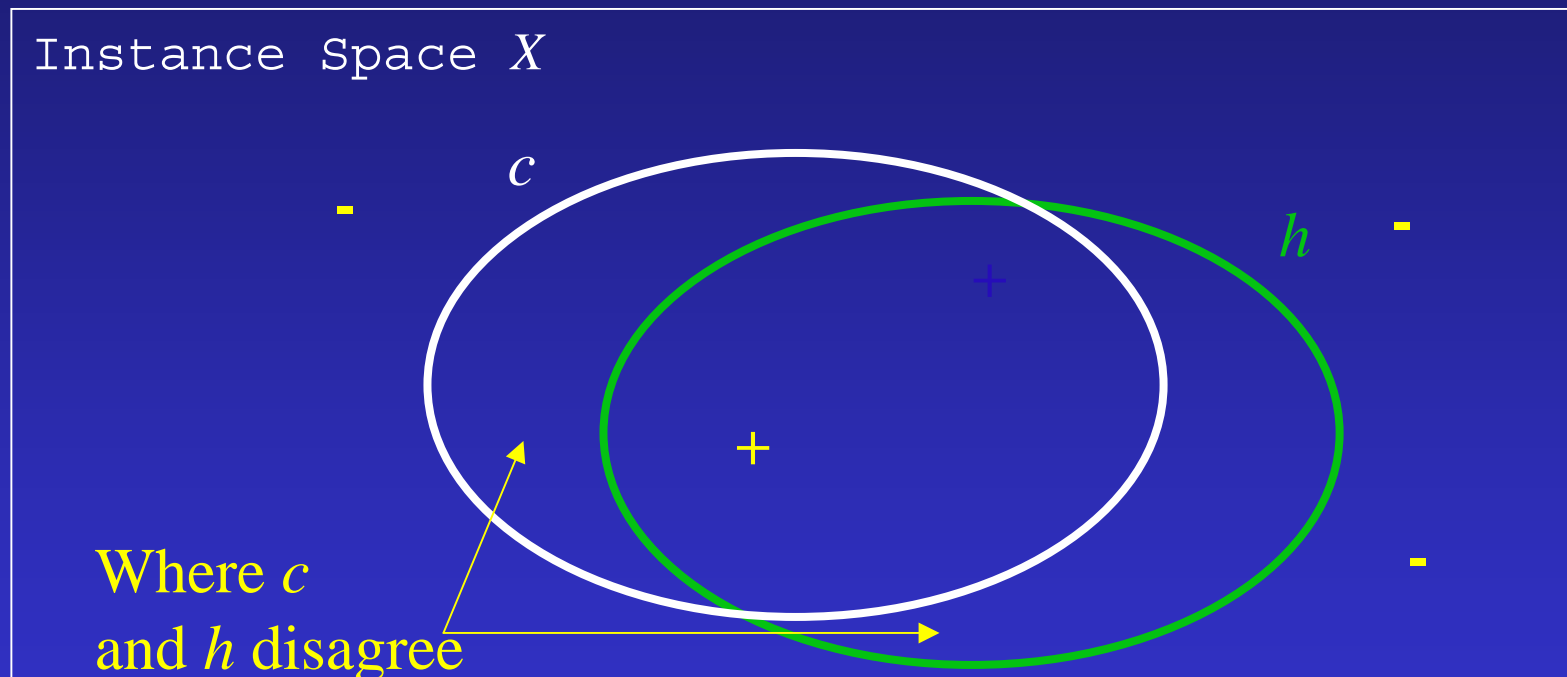
- The true error of a hypothesis  $h$

with respect to target concept  $c$  and distribution  $D$

is the probability that  $h$  will misclassify an instance drawn at random according to  $D$

$$\mathbf{error}_D(\mathbf{h}) = \mathbf{Pr}_{\mathbf{x} \in D}[\mathbf{c}(\mathbf{x}) \neq \mathbf{h}(\mathbf{x})]$$

## Error of hypothesis $h$



$h$  has a nonzero error with respect to  $c$  although they agree on all training samples

# Probably Approximately Correct (PAC) Learnability

- Characterize concepts learnable from
  - a reasonable number of randomly drawn training examples
  - a reasonable amount of computation
- Strong characterization is futile:
  - No of training examples needed to learn hypothesis  $h$  for which  $error_D(h) = 0$  cannot be determined because
    - *Unless there are training examples for every instance in  $X$  there may be multiple hypotheses consistent with training examples*
    - *Training examples are picked at random and therefore may be misleading*

## PAC Learnable

- Consider a concept class  $C$  defined over
  - a set of instances  $X$  of length  $n$  and
  - a learner  $L$  using hypothesis space  $H$
- $C$  is PAC-learnable by  $L$  using  $H$  if
  - for all  $c \in C$ ,
  - distributions  $D$  over  $X$ ,
  - such that  $0 < \epsilon < 1/2$
  - learner  $L$  will with *probability at least*  $(1-\delta)$
  - output a hypothesis  $h \in H$
  - such that  $\text{error}_D(h) \leq \epsilon$
  - in time that is polynomial in  $1/\epsilon$ ,  $1/\delta$ ,  $n$  and  $\text{size}(c)$

## PAC learnability requirements of Learner $L$

- With arbitrarily high probability  $(1-\delta)$ 
  - output a hypothesis having arbitrarily low error  $(\epsilon)$
- Do so efficiently
  - In time that grows at most polynomially with  $1/\epsilon$  and  $1/\delta$ , which define the strength of our demands on the output hypothesis
  - With  $n$  and  $size(c)$  that define the inherent complexity of the underlying instance space  $X$  and concept class  $C$
  - $n$  is the size of instances in  $X$ 
    - If instances are conjunctions of  $k$  Boolean variables,  $n=k$
  - $size(c)$  is the encoding length of  $c$  in  $C$ , eg, no of Boolean features actually used to describe  $c$

---

## Computational Resources *vs* No of Training Samples Required

- Two are closely related
- If  $L$  requires some minimum processing time per example, then for  $C$  to be PAC-learnable,  $L$  must learn from a polynomial number of training examples
  - To show some class  $C$  of target concepts is PAC-learnable is to first show that each target concept in  $C$  can be learned from a polynomial number of training examples

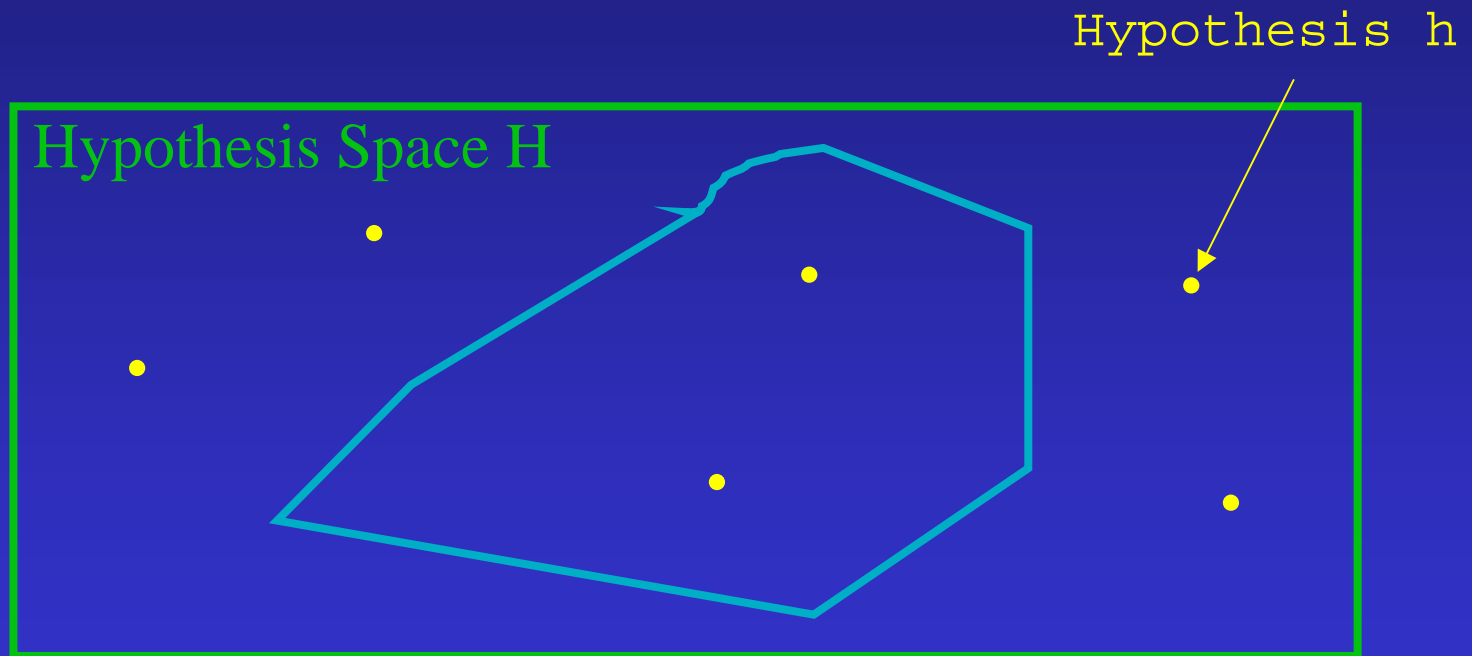
---

# Sample Complexity for Finite Hypothesis Spaces

- Sample Complexity
  - No. of training examples required
  - Growth with problem size
- Bound on no. of training samples needed for **consistent learners**— learners that perfectly fit the training data

# Version Space

- Contains all plausible versions of the target concept

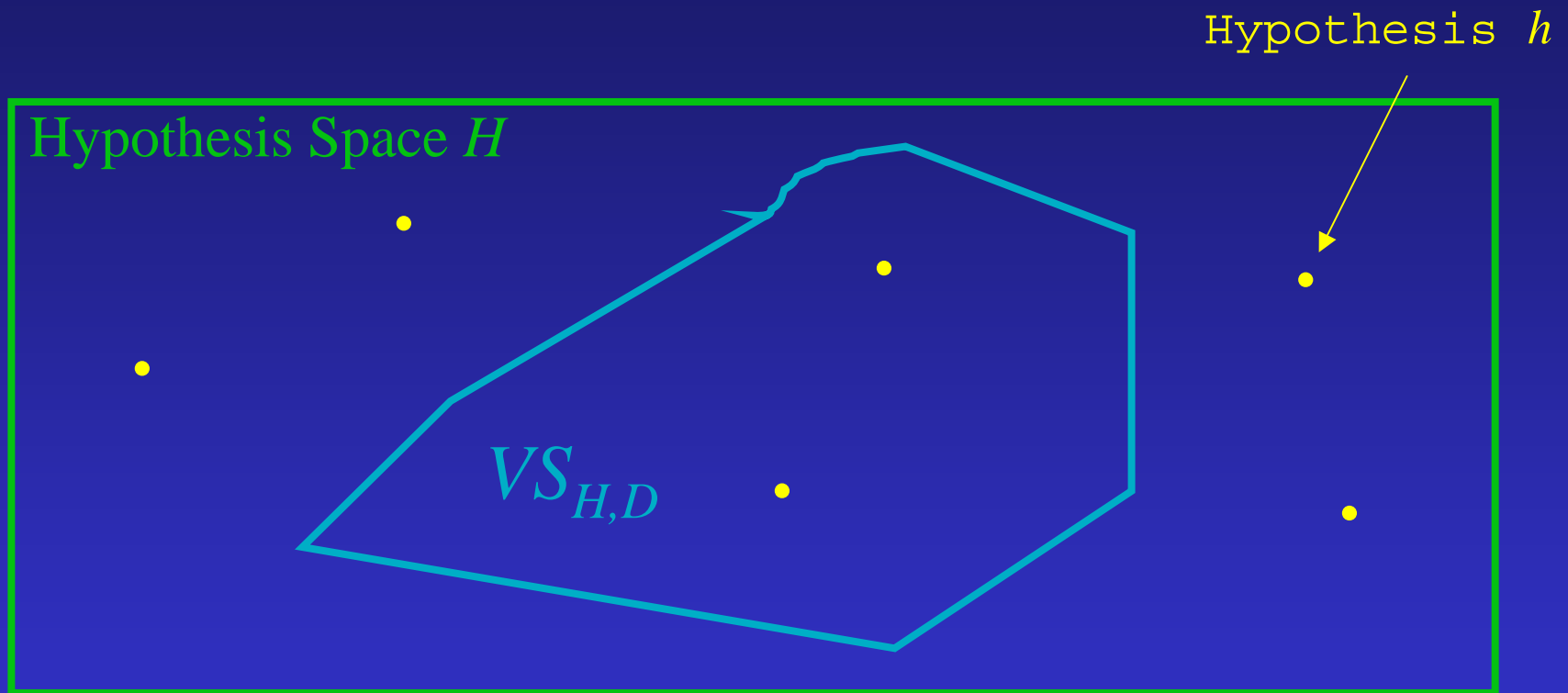




## Version Space

- A hypothesis  $h$  is consistent with training examples  $D$ 
  - iff  $h(x)=c(x)$  for each example  $\langle x, c(x) \rangle$  in  $D$
- **Version space** with respect to
  - hypothesis  $H$  and
  - training examples  $D$ ,
  - is a subset of hypotheses from  $H$  consistent with the training examples in  $D$

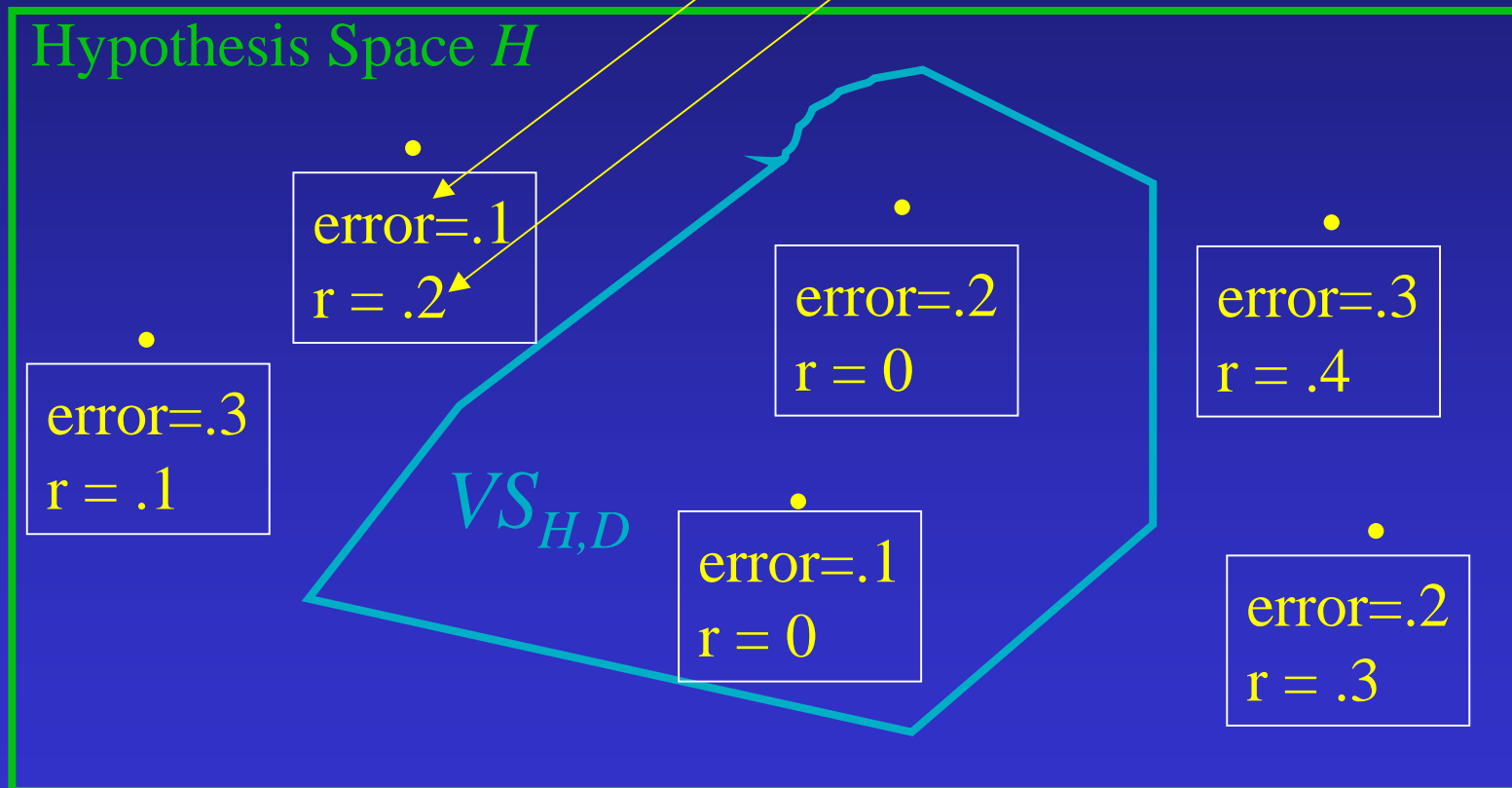
# Version Space



$$VS_{H,D} = \{h \in H \mid (\forall (\mathbf{x}, \mathbf{c}(\mathbf{x})) \in D)(h(\mathbf{x}) = \mathbf{c}(\mathbf{x}))\}$$

# Version Space with associated errors

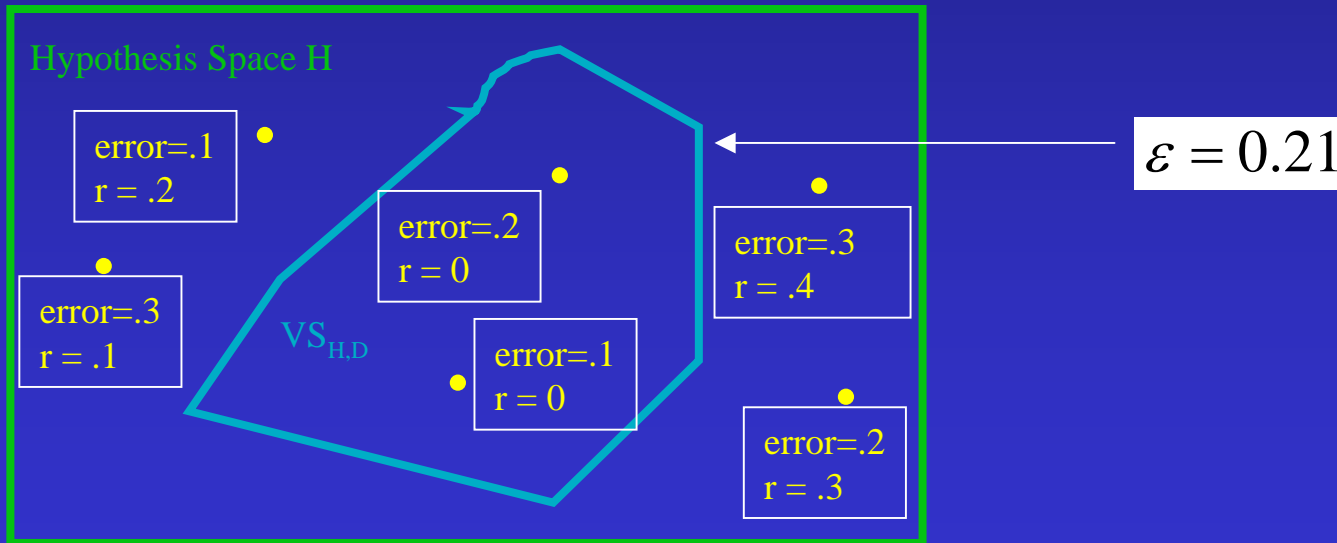
error is the true error,  
r is the training error



# Exhausting the Version Space: true error is less than $\varepsilon$

- The version space is  $\varepsilon$ -exhausted with respect to  $c$  and  $D$  if

$$(\forall h \in VS_{H,D}) \dots \text{error}_D(h) < \varepsilon$$



## Upper bound on probability of not $\varepsilon$ -exhausted

- **Theorem:**

- If the hypothesis space  $H$  is finite
- $D$  is a sequence of  $m \geq 1$  independent random samples of concept  $c$
- Then for any  $0 < \varepsilon < 1$
- Probability of not  $\varepsilon$ -exhausted (with respect to  $c$ ) is less than or equal to

$$|H| e^{-\varepsilon m}$$

- Bounds the probability that  $m$  training samples will fail to eliminate all “bad” hypotheses

## Number of training samples required

$$|\mathbf{H}| e^{-\varepsilon m} \leq \delta$$

Rearranging

$$m \geq \frac{1}{\varepsilon} (\ln |\mathbf{H}| + \ln(1/\delta))$$

Probability of failure  
is below some desired level

- Provides general bound on the no. of training samples
  - sufficient for any consistent learner to learn any target concept in  $H$  for
  - any desired values of  $\delta$  and  $\varepsilon$

## Generalization to non-zero training error

$$\mathbf{m} \geq \frac{1}{2\epsilon^2} (\ln |\mathbf{H}| + \ln(1/\delta))$$

- $m$  grows as the square of  $1/\epsilon$  rather than linearly
- Called agnostic learning

# Conjunctions of Boolean literals are PAC-learnable

- Sample complexity

$$\mathbf{m} \geq \frac{1}{\varepsilon} (\mathbf{n} \ln 3 + \ln(1/\delta))$$



## K-term DNF and CNF concepts

$$\mathbf{m} \geq \frac{1}{\varepsilon} (\mathbf{nk} \ln 3 + \ln(1/\delta))$$