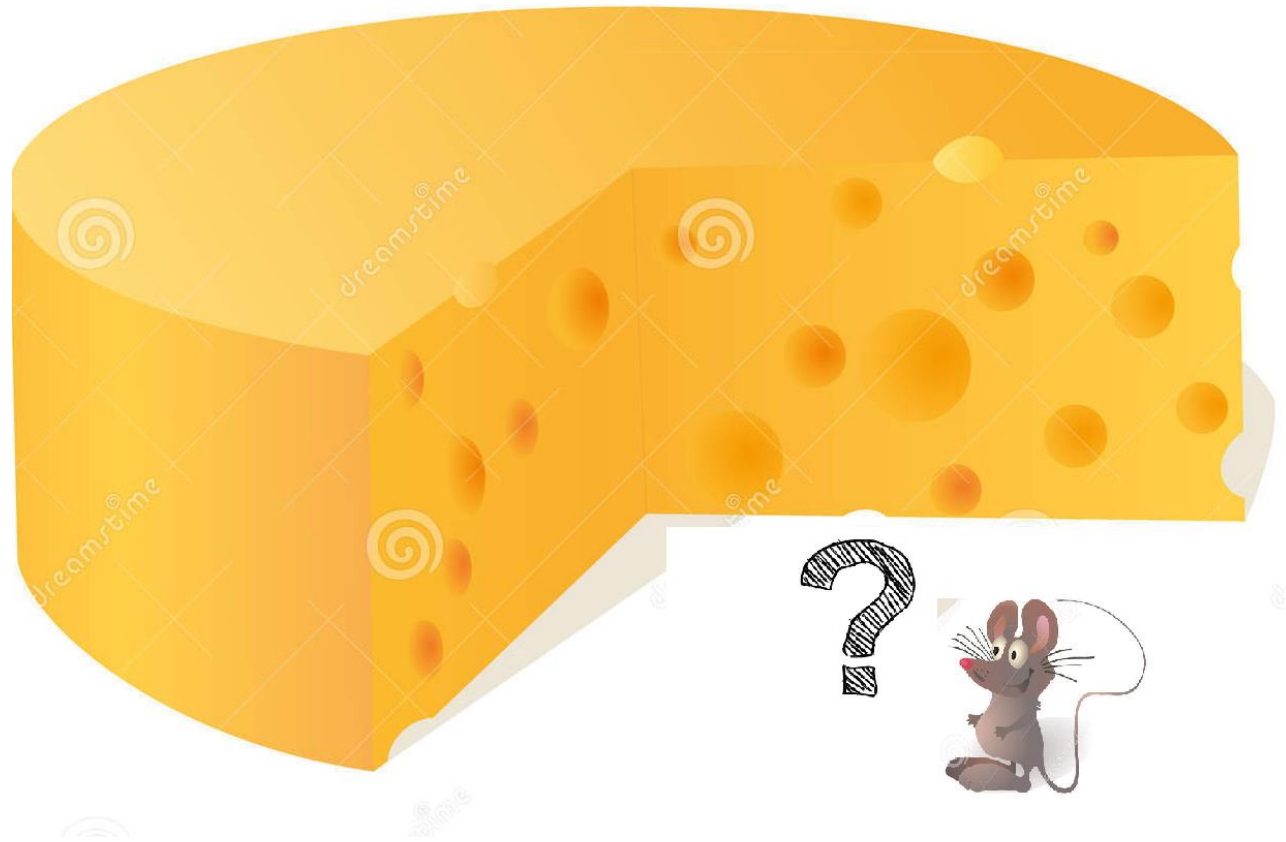


Big Data Class



LECTURER: DAN FELDMAN

TEACHING ASSISTANTS:

IBRAHIM JUBRAN

ALAA MAALOUF



Probability

Hoeffding's bound

(Simplified version – from Young'95 paper)

- For n independent random variables X_1, \dots, X_n where $X_i \in [0,1]$, with $E(X_i) \leq \mu_i$
- Let $X = \sum_{i=1}^n \frac{X_i}{n}$ and $\mu = \sum_{i=1}^n \mu_i$.
- Then: $\Pr[X \geq \mu + n\epsilon] < \frac{1}{e^{2n\epsilon^2}}$

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

- Let $\alpha = e^{4\epsilon} - 1$.

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

- Let $\alpha = e^{4\epsilon} - 1$.

$$\begin{aligned} & \Pr[\sum_{i=1}^n X_i \geq \mu + n\epsilon] \\ &= \Pr \left[\prod_{i=1}^n \frac{(1+\alpha)^{X_i}}{(1+\alpha)^{\mu_i + \epsilon}} \geq 1 \right] \end{aligned}$$

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

• Let $\alpha = e^{4\epsilon} - 1$.

Follows from:

1. For $0 \leq z \leq 1$, $(1 + \alpha)^z \leq 1 + \alpha z$
2. Markov's inequality : $\Pr[X \geq a] \leq \frac{E(X)}{a}$

$$\begin{aligned} & \Pr[\sum_{i=1}^n X_i \geq \mu + n\epsilon] \\ &= \Pr \left[\prod_{i=1}^n \frac{(1+\alpha)^{X_i}}{(1+\alpha)^{\mu_i + \epsilon}} \geq 1 \right] \\ &\leq \mathbb{E} \left[\prod_{i=1}^n \frac{1+\alpha X_i}{(1+\alpha)^{\mu_i + \epsilon}} \right] \end{aligned}$$

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

• Let $\alpha = e^{4\epsilon} - 1$.

Follows from:

1. For $0 \leq z \leq 1$, $(1 + \alpha)^z \leq 1 + \alpha z$
2. Markov's inequality : $\Pr[X \geq a] \leq \frac{E(X)}{a}$

$$\begin{aligned} & \Pr[\sum_{i=1}^n X_i \geq \mu + n\epsilon] \\ &= \Pr \left[\prod_{i=1}^n \frac{(1+\alpha)^{X_i}}{(1+\alpha)^{\mu_i + \epsilon}} \geq 1 \right] \\ &\leq \mathbb{E} \left[\prod_{i=1}^n \frac{1 + \alpha X_i}{(1+\alpha)^{\mu_i + \epsilon}} \right] \\ &= \prod_{i=1}^n \frac{1 + \alpha E(X_i)}{(1+\alpha)^{\mu_i + \epsilon}} \end{aligned}$$

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

• Let $\alpha = e^{4\epsilon} - 1$.

Follows from:

1. For $0 \leq z \leq 1$, $(1 + \alpha)^z \leq 1 + \alpha z$
2. Markov's inequality : $\Pr[X \geq a] \leq \frac{E(X)}{a}$

$$\begin{aligned} & \Pr[\sum_{i=1}^n X_i \geq \mu + n\epsilon] \\ &= \Pr \left[\prod_{i=1}^n \frac{(1+\alpha)^{X_i}}{(1+\alpha)^{\mu_i + \epsilon}} \geq 1 \right] \\ &\leq \mathbb{E} \left[\prod_{i=1}^n \frac{1 + \alpha X_i}{(1+\alpha)^{\mu_i + \epsilon}} \right] \\ &= \prod_{i=1}^n \frac{1 + \alpha E(X_i)}{(1+\alpha)^{\mu_i + \epsilon}} \\ &\leq \prod_{i=1}^n \frac{1 + \alpha \mu_i}{(1+\alpha)^{\mu_i + \epsilon}} \end{aligned}$$

Hoeffding's bound - Proof

(Simplified version – from Young'95 paper)

• Let $\alpha = e^{4\epsilon} - 1$.

Follows from:

1. For $0 \leq z \leq 1$, $(1 + \alpha)^z \leq 1 + \alpha z$
2. Markov's inequality : $\Pr[X \geq a] \leq \frac{E(X)}{a}$

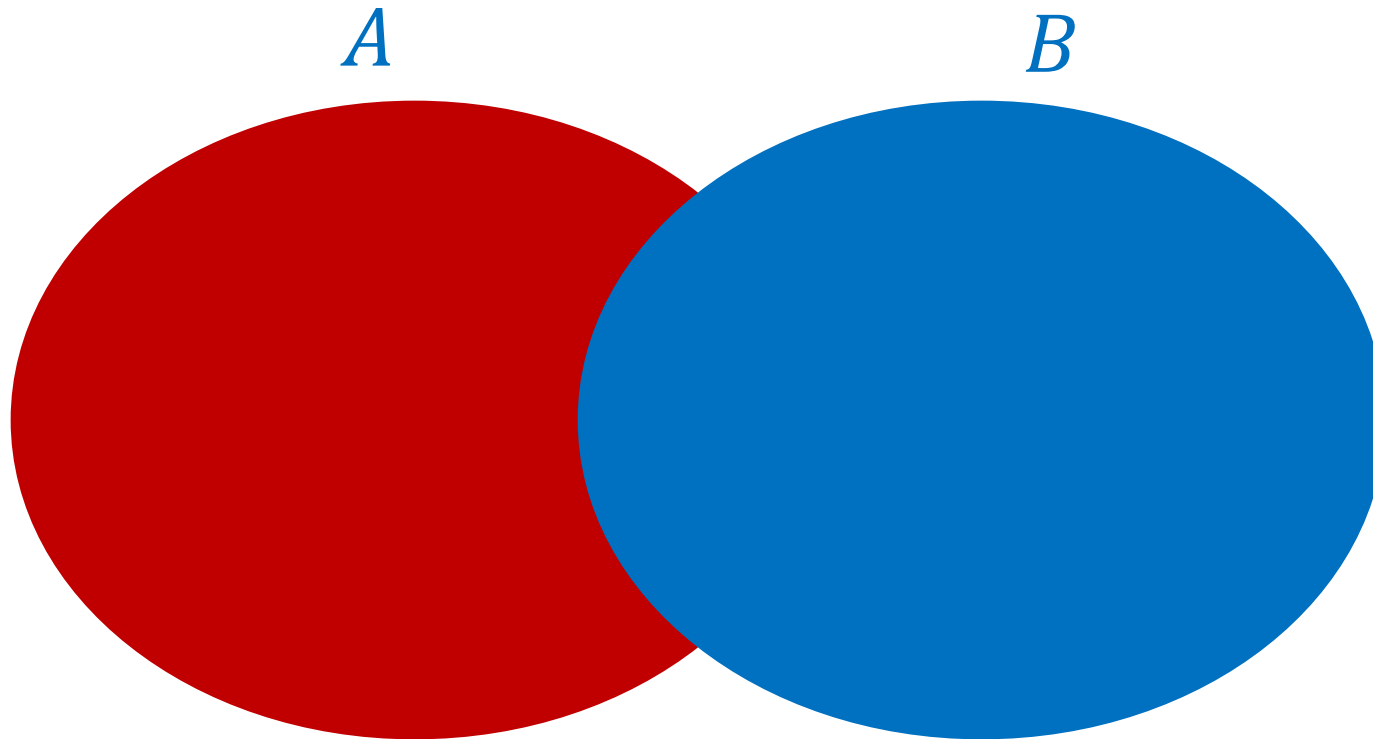
For $\epsilon > 0$, $\alpha = e^{4\epsilon} - 1$, $z \geq 0$:

$$1 + \alpha z < \frac{(1 + \alpha)^{z+\epsilon}}{e^{2\epsilon^2}}$$

$$\begin{aligned} & \Pr[\sum_{i=1}^n X_i \geq \mu + n\epsilon] \\ &= \Pr \left[\prod_{i=1}^n \frac{(1+\alpha)^{X_i}}{(1+\alpha)^{\mu_i+\epsilon}} \geq 1 \right] \\ &\leq \mathbb{E} \left[\prod_{i=1}^n \frac{1+\alpha X_i}{(1+\alpha)^{\mu_i+\epsilon}} \right] \\ &= \prod_{i=1}^n \frac{1+\alpha E(X_i)}{(1+\alpha)^{\mu_i+\epsilon}} \\ &\leq \prod_{i=1}^n \frac{1+\alpha \mu_i}{(1+\alpha)^{\mu_i+\epsilon}} \\ &< \prod_{i=1}^n \frac{1}{e^{2\epsilon^2}} \end{aligned}$$

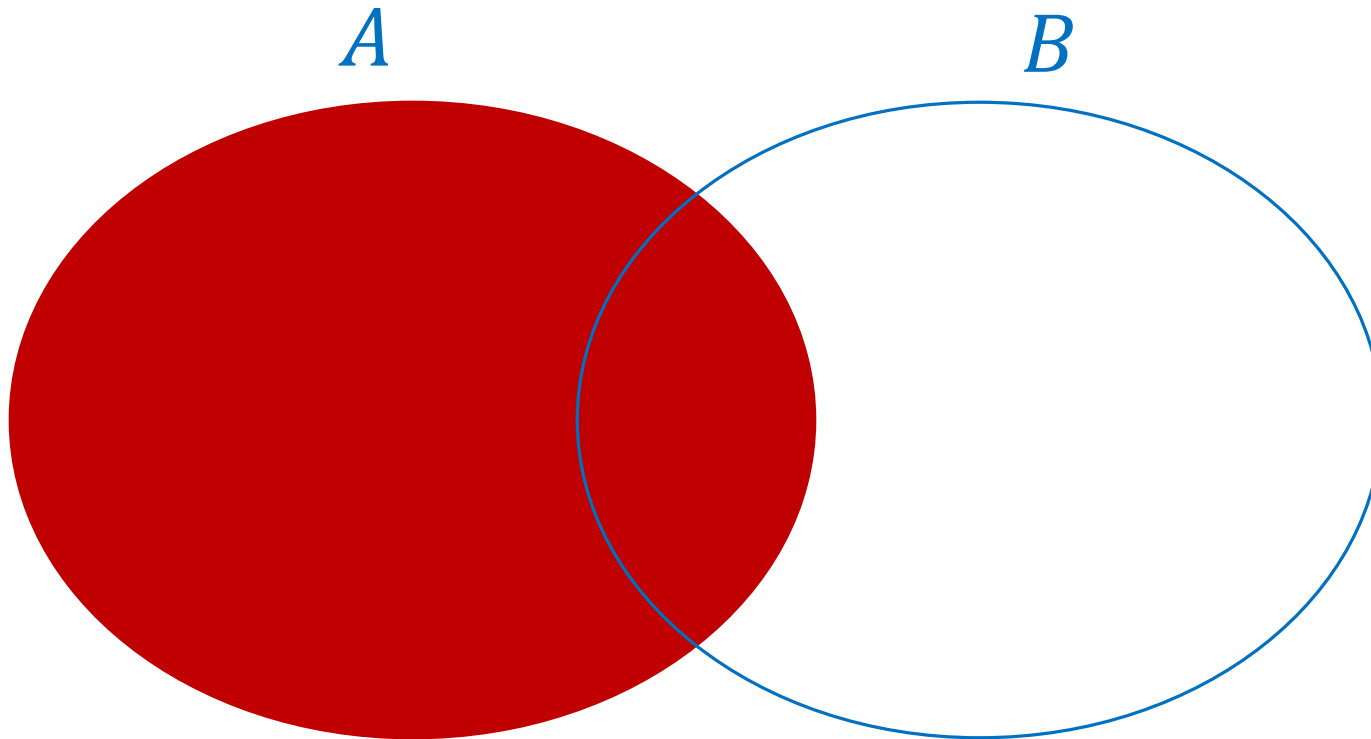
Union Bound (Boole's inequality)

$$A \cup B =$$



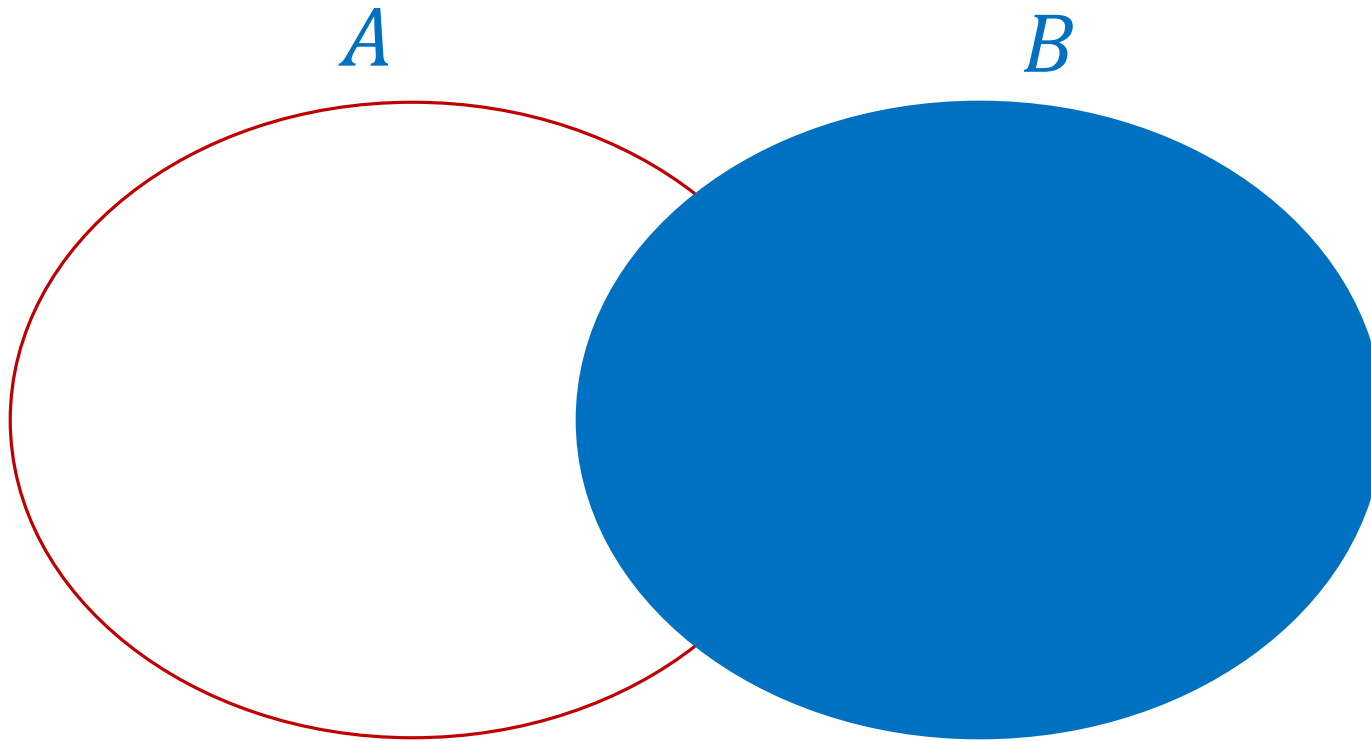
Union Bound (Boole's inequality)

$$A \cup B = A$$



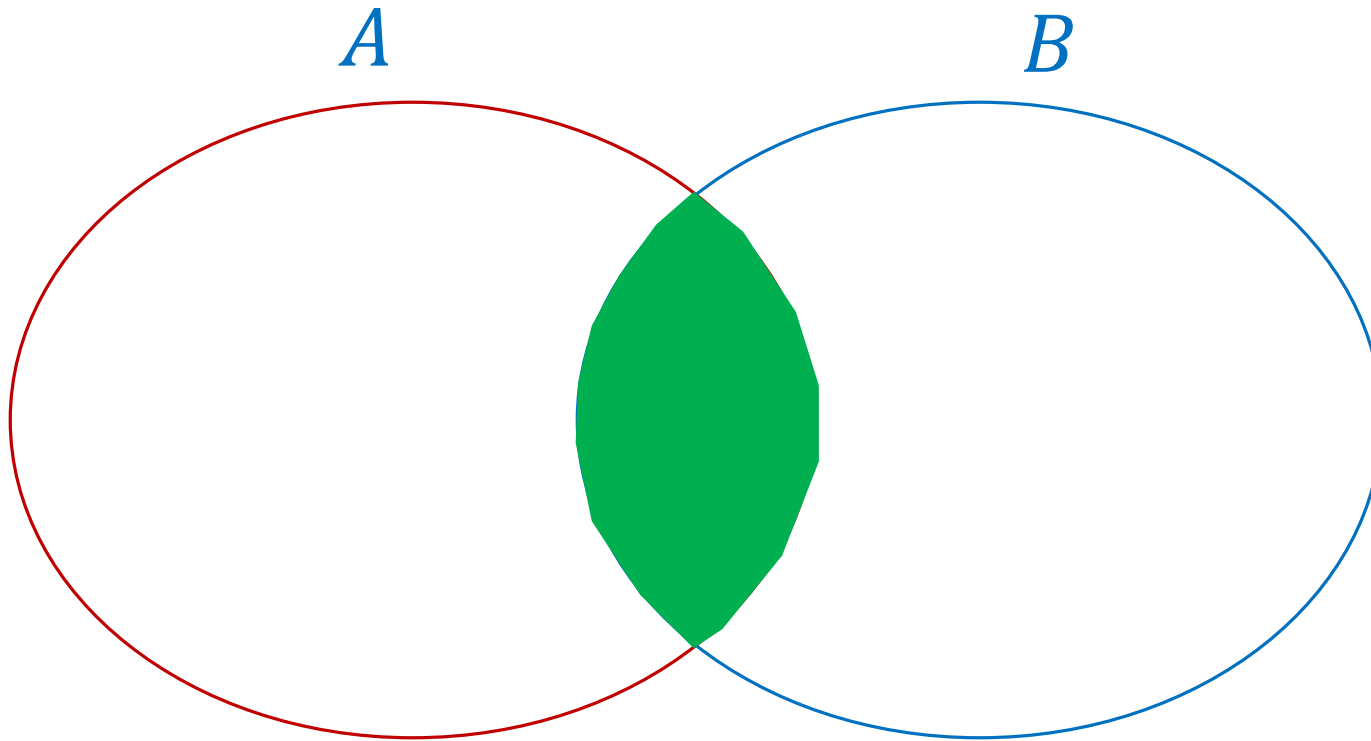
Union Bound (Boole's inequality)

$$A \cup B = A + B$$



Union Bound (Boole's inequality)

$$A \cup B = A + B - A \cap B$$



Union Bound (Boole's inequality)

$$\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Union Bound (Boole's inequality)

$$\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

→ For any events A_1, A_2, \dots, A_n :

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$

Union Bound (Boole's inequality)

$$\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

→ For any events A_1, A_2, \dots, A_n :

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$
$$P\left(\bigcup_{i=1}^n A_i\right) \geq \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j)$$

Union Bound (Boole's inequality)

$$\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

→ For any events A_1, A_2, \dots, A_n :

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$

$$P\left(\bigcup_{i=1}^n A_i\right) \geq \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j)$$

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) + \sum_{i < j < k} P(A_i \cap A_j \cap A_k)$$



Hoeffding's bound $\rightarrow \epsilon$ -sample

Pick a random sample S of F , it holds that:

$$Pr_1 = Pr \left(\left| \frac{|F \cap range_1|}{|F|} - \frac{|S \cap range_1|}{|S|} \right| > \epsilon \right)$$


Probability of failure for $range_1$

S is an ϵ -sample

Hoeffding's bound $\rightarrow \epsilon$ -sample

Pick a random sample S of F , it holds that:

$$Pr_1 = Pr \left(\left| \frac{|F \cap range_1|}{|F|} - \frac{|S \cap range_1|}{|S|} \right| > \epsilon \right) < \frac{1}{e^{2|S|\epsilon^2}}$$



Probability of failure for $range_1$ Hoeffding

Hoeffding's bound $\rightarrow \epsilon$ -sample

Pick a random sample S of F , it holds that:

$$Pr_1 = Pr \left(\left| \frac{|F \cap range_1|}{|F|} - \frac{|S \cap range_1|}{|S|} \right| > \epsilon \right) < \frac{1}{e^{2|S|\epsilon^2}}$$

Probability of
failure for $range_1$

Hoeffding

$$Pr_{\text{bad}} = Pr_1 \cup Pr_2 \cup \dots \cup Pr_m \leq m \cdot \frac{1}{e^{2|S|\epsilon^2}}$$

Union Bound

Hoeffding's bound $\rightarrow \epsilon$ -sample

The probability for failure should be small:

$$\Pr_{\text{bad}} \leq m \cdot \frac{1}{e^{2|S|\epsilon^2}} \leq \delta$$

$$\rightarrow \frac{m}{\delta} \leq e^{2|S|\epsilon^2}$$

$$\rightarrow |S| \geq \frac{1}{\epsilon^2} \left(\log m + \log \frac{1}{\delta} \right)$$

Hoeffding's bound $\rightarrow \epsilon$ -sample

The probability for failure should be small:

$$\Pr_{\text{bad}} \leq m \cdot \frac{1}{e^{2|S|\epsilon^2}} \leq \delta$$

$$\rightarrow \frac{m}{\delta} \leq e^{2|S|\epsilon^2}$$

$$\rightarrow |S| \geq \frac{1}{\epsilon^2} \left(\log m + \log \frac{1}{\delta} \right)$$

Might be infinite!

Handling $m \rightarrow \infty$

Example:

$Q = \text{circles in } R^d \rightarrow \text{range}(F, q, r) = \{f \in F \mid f \text{ inside the circle with center } q \text{ and radius } r\}$

\rightarrow number of different circles: $|Q| = m = \infty$.

However, the number of different equivalence classes is $\underline{n^{O(d)}}$ since: A sphere in R^d is determined by $d + 1$ points.

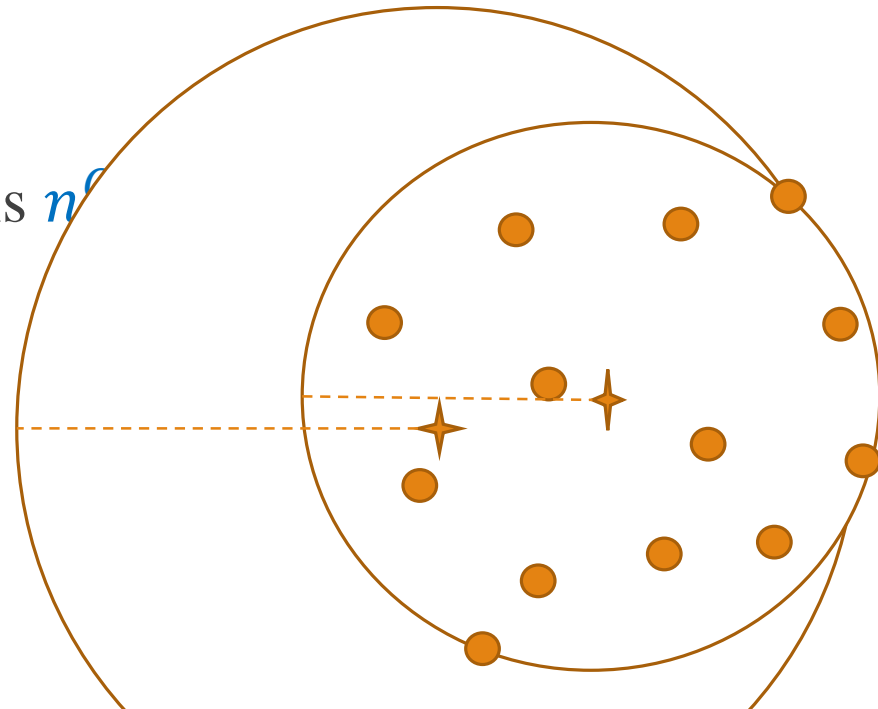
Handling $m \rightarrow \infty$

Example:

$Q = \text{circles in } R^d \rightarrow \text{range}(F, q, r) = \{f \in F \mid f \text{ inside the circle with center } q \text{ and radius } r\}$

\rightarrow number of different circles: $|Q| = m = \infty$.

However, the number of different equivalence classes is n^d .
 R^d is determined by $d + 1$ points.



Handling $m \rightarrow \infty$

Example:

$Q = \text{circles in } R^d \rightarrow \text{range}(F, q, r) = \{f \in F \mid f \text{ inside the circle with center } q \text{ and radius } r\}$

\rightarrow number of different circles: $|Q| = m = \infty$.

However, the number of different equivalence classes is $\underline{n^{O(d)}}$ since: A sphere in R^d is determined by $d + 1$ points.

$$\rightarrow |S| \geq \frac{1}{\epsilon^2} \left(\log n^d + \log \frac{1}{\delta} \right) = \frac{1}{\epsilon^2} \left(d \log n + \log \frac{1}{\delta} \right)$$

VC-dimension

Definition: (Range Space)

A **range space** is a pair (F, ranges) where F is a set, called ground set and ranges is a family (set) of subsets of F .

Definition: (VC-dimension)

The **VC-dimension** of a range space (F, ranges) is the size $|G|$ of the largest subset $G \subseteq F$ such that

$$|\{G \cap \text{range} \mid \text{range} \in \text{ranges}\}| \leq 2^{|G|}$$

VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{m}{d}\right)^d\right)$.

VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

Proof: for which values of m it holds that $\left(\frac{em}{d}\right)^d \leq 2^m$?

VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

Proof:

$$\left(\frac{em}{d}\right)^d \leq 2^m$$

VC-dimension

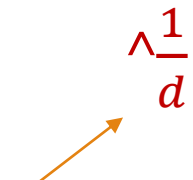
Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

Proof:

$$\left(\frac{em}{d}\right)^d \leq 2^m \rightarrow \frac{em}{d} \leq 2^{\frac{m}{d}}$$


VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

Proof:

$$ex \leq 2^x \text{ for } x \geq 4$$

$$\left(\frac{em}{d}\right)^d \leq 2^m \rightarrow \frac{em}{d} \leq 2^{\frac{m}{d}} \rightarrow \frac{m}{d} \geq 4$$

VC-dimension

Theorem:

Let f_1, \dots, f_m be real polynomials in $d \leq m$ variables, each of constant degree. Then the number of sign sequences $(\text{sign}(f_1(x)), \dots, \text{sign}(f_m(x)))$, $x \in R^d$, that consist of 1 and -1 is at most $O\left(\left(\frac{em}{d}\right)^d\right)$.

Corollary:

If $m \geq O(d)$, then the number of distinct sequences as in the above theorem is less than 2^m .

Proof:

$$\left(\frac{em}{d}\right)^d \leq 2^m \rightarrow \frac{em}{d} \leq 2^{\frac{m}{d}} \rightarrow \frac{m}{d} \geq 4 \rightarrow m = O(d) \quad \blacksquare$$

VC-dimension

Let $Q = \{q_1, \dots, q_k\} \subseteq R^d$ and $R = \{r_1, \dots, r_k\} \subseteq R$.

Then:

$$\text{range}(P, Q, R) = \{p \in P \mid \forall_i (\text{dist}^2(p, q_i) \leq r_i^2)\}$$

Consider the following polynomials:

$$\text{Poly}(P, Q, R) = \{ \|p_i - q_j\|^2 - r_j^2 \mid i \in [n], j \in [k] \}$$

→ $|\text{Poly}(P, Q, R)| = nk$ such polynomials in dk variables.

VC-dimension

Let $Q_1, Q_2 \subseteq R^d$ and $R_1, R_2 \subseteq R$.

Lemma:

If $Poly(P, Q_1, R_1)$ and $Poly(P, Q_2, R_2)$ have the **same sign sequence** for the nk polynomials, then

$$range(P, Q_1, R_1) = range(P, Q_2, R_2)$$

VC-dimension

Let $Q_1, Q_2 \subseteq R^d$ and $R_1, R_2 \subseteq R$.

Lemma:


If $Poly(P, Q_1, R_1)$ and $Poly(P, Q_2, R_2)$ have the **same sign sequence** for the nk polynomials, then

$$range(P, Q_1, R_1) = range(P, Q_2, R_2)$$

Proof:

Let $p \in range(P, Q_1, R_1)$

→ exists $j \in k$ such that $\|p - q_{j1}\|^2 - r_{j1}^2 \leq 0$



VC-dimension

Let $Q_1, Q_2 \subseteq R^d$ and $R_1, R_2 \subseteq R$.

Lemma:

If $Poly(P, Q_1, R_1)$ and $Poly(P, Q_2, R_2)$ have the **same sign sequence** for the nk polynomials, then

$$range(P, Q_1, R_1) = range(P, Q_2, R_2)$$

Proof:

Let $p \in range(P, Q_1, R_1)$

\rightarrow exists $j \in k$ such that $\|p - q_{j1}\|^2 - r_{j1}^2 \leq 0 \rightarrow \|p - q_{j2}\|^2 - r_{j2}^2 \leq 0$

VC-dimension

Let $Q_1, Q_2 \subseteq R^d$ and $R_1, R_2 \subseteq R$.

Lemma:

If $Poly(P, Q_1, R_1)$ and $Poly(P, Q_2, R_2)$ have the **same sign sequence** for the nk polynomials, then

$$range(P, Q_1, R_1) = range(P, Q_2, R_2)$$

Proof:

Let $p \in range(P, Q_1, R_1)$

→ exists $j \in k$ such that $\|p - q_{j1}\|^2 - r_{j1}^2 \leq 0 \rightarrow \|p - q_{j2}\|^2 - r_{j2}^2 \leq 0$

→ $p \in range(P, Q_2, R_2)$



VC-dimension

Conclusion:

#different ranges \leq #different sign sequences $\leq 2^m$



Last Lemma



If $m > O(d)$