# Advanced Introduction to Machine Learning, CMU-10715

## Vapnik–Chervonenkis Theory

Barnabás Póczos

**ML** MACHINE LEARNING DEPARTMENT

**Carnegie Mellon.**
School of Computer Science

# Learning Theory

We have explored many ways of learning from data
But…

- How good is our classifier, really?

- How much data do we need to make it "good enough"?

# Review of what we have learned so far

# Notation

$$R(f) = \Pr[Y \neq f(X)] \qquad R^* = R(f^*) = \inf_{f:\mathcal{X}\to\mathbb{R}} R(f) \qquad f^* = \arg\inf_{f:\mathcal{X}\to\mathbb{R}} R(f)$$

$$R^*_{\mathcal{F}} = R(f^*_{\mathcal{F}}) = \inf_{f\in\mathcal{F}} R(f) \qquad f^*_{\mathcal{F}} = \arg\inf_{f\in\mathcal{F}} R(f)$$

$$\widehat{R}_n(f) = \frac{1}{n}\sum_{i=1}^{n} 1_{\{Y_i \neq f(X_i)\}} \qquad \widehat{R}^*_{n,\mathcal{F}} = \inf_{f\in\mathcal{F}} \widehat{R}_n(f) \qquad f^*_{n,\mathcal{F}} = \arg\min_{f\in\mathcal{F}} \widehat{R}_n(f)$$

This is what the learning algorithm produces

We will need these definitions, please copy it!

$R(f) = $ Risk $\qquad\qquad R^* = $ Bayes risk

$\widehat{R}_n(f) = $ Empricial risk $\qquad f^* = $ Bayes classifier

$f^*_{n,\mathcal{F}} = $ the classifier that the learning algorithm produces

# Big Picture

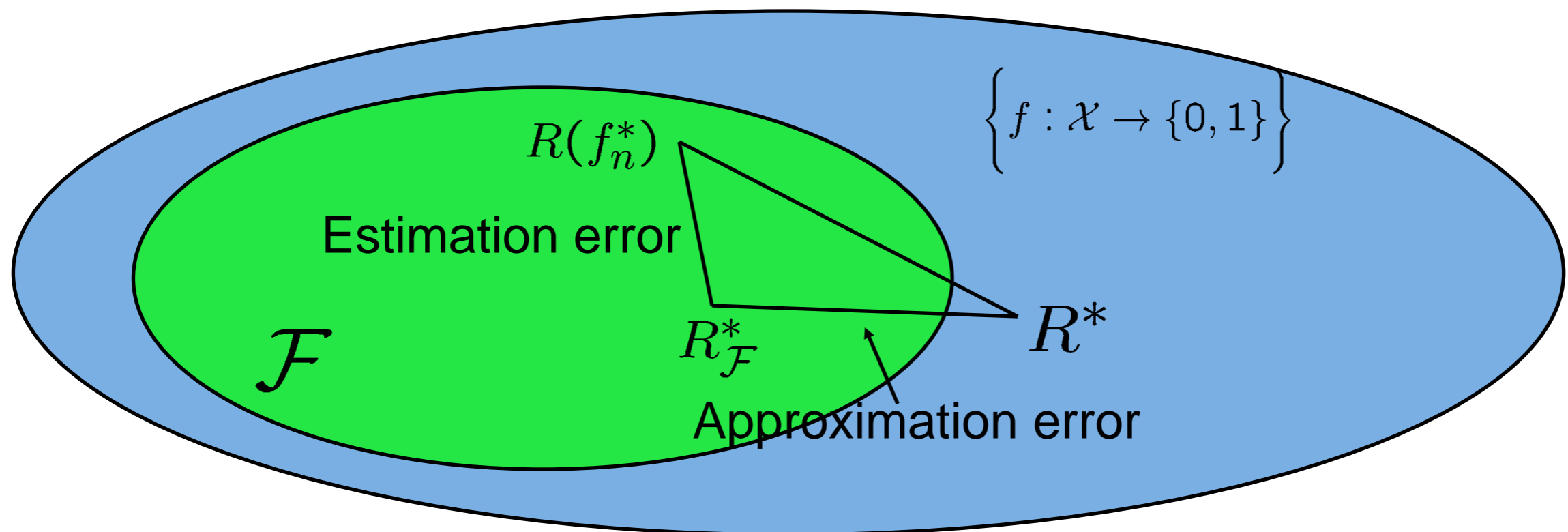**Ultimate goal:** $R(f_n^*) - R^* = 0$

ERM: $f_n^* = f_{n,\mathcal{F}}^* = \arg\min_{f \in \mathcal{F}} \widehat{R}_n(f) = \arg\min_{f \in \mathcal{F}} \frac{1}{n}\sum_{i=1}^n L(Y_i, f(X_i))$

Risk of the classifier $f_n^*$    Estimation error    Approximation error

$$R(f_n^*) - R^* = \overbrace{R(f_n^*) - R_{\mathcal{F}}^*} + \overbrace{R_{\mathcal{F}}^* - R^*}$$

Bayes risk      Bayes risk

$$R_{\mathcal{F}}^* = \inf_{g \in \mathcal{F}} R(g)$$    Best classifier in $\mathcal{F}$



$\left\{ f : \mathcal{X} \to \{0,1\} \right\}$

$R(f_n^*)$

Estimation error

$R_{\mathcal{F}}^*$

$R^*$

$\mathcal{F}$

Approximation error

5

# Big Picture: Illustration of Risks

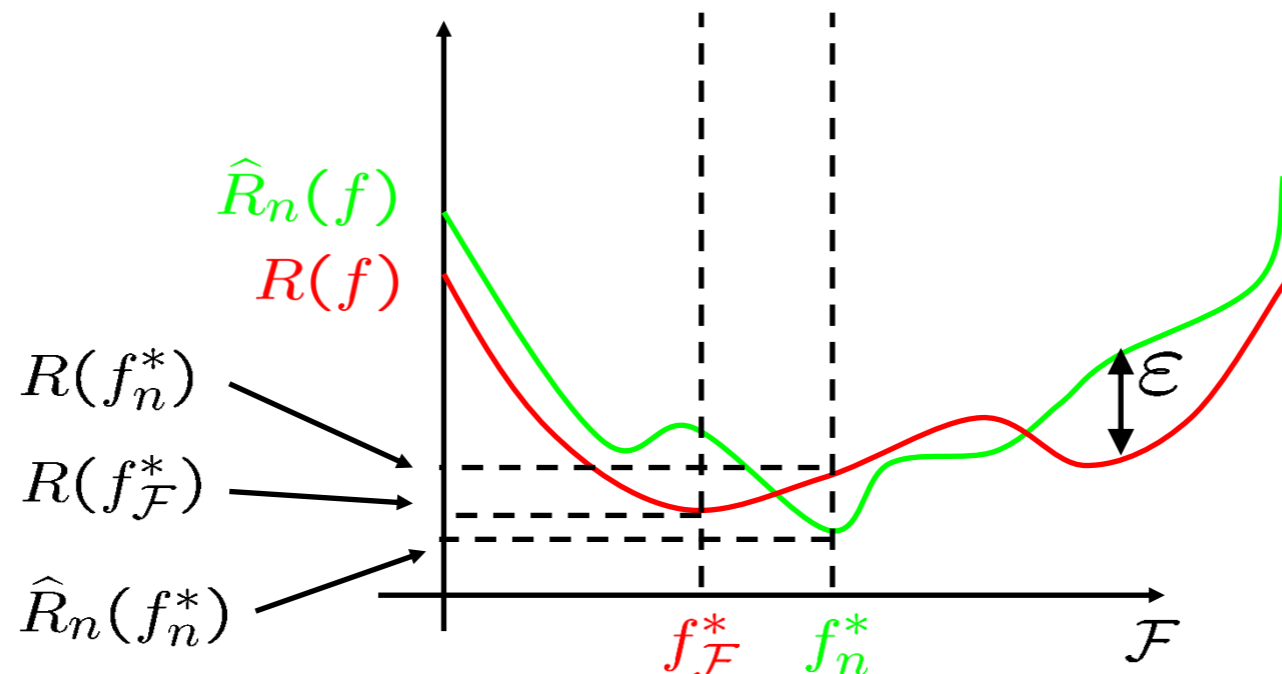$$|\widehat{R}_n(f_n^*) - R(f_n^*)| \leq \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| = \varepsilon$$

$$|R(f_n^*) - R(f_\mathcal{F}^*)| \leq 2\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| = 2\varepsilon$$

$$|\widehat{R}_n(f_n^*) - R(f_\mathcal{F}^*)| \leq 3\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| = 3\varepsilon$$

**Upper bound**
$$\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$$

**Goal of Learning:**
For a fixed $\mathcal{F}$, make the $|R(f_n^*) - R(f_\mathcal{F}^*)|$ estimation error small

# Learning Theory

# Outline

From Hoeffding's inequality, we have seen that

Theorem: Let $\mathcal{F} = \{f : \mathcal{X} \to \{0, 1\}\}$, and $|\mathcal{F}| \leq N$

$$\begin{cases} \Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > \varepsilon\right) \leq 2N \exp\left(-2n\varepsilon^2\right) \\[2em] \Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| < \sqrt{\dfrac{\log(N) + \log(2/\delta)}{2n}}\right) \geq 1 - \delta \end{cases}$$

These results are useless if N is big, or infinite. (e.g. all possible hyper-planes)

Today we will see how to fix this with the Shattering coefficient and VC dimension

# Outline

From Hoeffding's inequality, we have seen that

**Theorem:** Let $\mathcal{F} = \{f : \mathcal{X} \to \{0, 1\}\}$, and $|\mathcal{F}| \leq N$

$$\Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > \varepsilon\right) \leq 2N \exp\left(-2n\varepsilon^2\right)$$

$$\Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| < \sqrt{\frac{\log(N) + \log(2/\delta)}{2n}}\right) \geq 1 - \delta$$

After this fix, we can say something meaningful about this too:

$$|R(f_n^*) - R(f_{\mathcal{F}}^*)| \leq 2 \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| = 2\varepsilon$$

Best true risk in $\mathcal{F}$

This is what the learning algorithm produces and its true risk

# Hoeffding inequality

**Theorem:** Let $\mathcal{F} = \{f : \mathcal{X} \to \{0, 1\}\}$, and $|\mathcal{F}| \leq N$

$$\Rightarrow \Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > \varepsilon\right) \leq 2N \exp\left(-2n\varepsilon^2\right)$$

$$\Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| < \sqrt{\frac{\log(N) + \log(2/\delta)}{2n}}\right) \geq 1 - \delta$$

**Definition:** $\displaystyle \widehat{R}_n(f) = \frac{1}{n} \sum_{i=1}^{n} 1_{\{Y_i \neq f(X_i)\}}$

**Observation!**

It does not matter how many elements $\mathcal{F}$ has. All that matters in the union bound is how many elements

$$\{[f(X_1), \ldots, f(X_n)] \; f \in \mathcal{F}\}$$

has. (The effective size of $\mathcal{F}$). It can't even be more than $2^n$.

# McDiarmid's Bounded Difference Inequality

Suppose $X_1, X_2, \ldots, X_n$ are independent and assume that

$$\sup_{x_1, x_2, \ldots, x_n, \widehat{x}_i} |f(x_1, x_2, \ldots, x_n) - f(x_1, x_2, \ldots, x_{i-1}, \widehat{x}_i, x_{i+1}, \ldots, x_n)| \leq c_i$$
$$\text{for } 1 \leq i \leq n$$

(**Bounded Difference Assumption**: replacing the $i$-th coordinate $x_i$ changes the value of $f$ by at most $c_i$.)

## It follows that

$$\Pr\{f(X_1, X_2, \ldots, X_n) - E[f(X_1, X_2, \ldots, X_n)] \geq \varepsilon\} \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}\right)$$

$$\Pr\{E[f(X_1, X_2, \ldots, X_n)] - f(X_1, X_2, \ldots, X_n) \geq \varepsilon\} \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}\right)$$

$$\Pr\{|E[f(X_1, X_2, \ldots, X_n)] - f(X_1, X_2, \ldots, X_n)| \geq \varepsilon\} \leq 2\exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n c_i^2}\right).$$

# Bounded Difference Condition

Our main goal is to bound $\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

**Lemma:**

The **"bounded difference condition"** is satisfied for $\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

**Proof:**

$$\widehat{R}_n(f) = \frac{1}{n} \sum_{i=1}^{n} 1_{\{f(X_i) \neq Y_i\}}$$

Let $g$ denote the following function:

$$g(Z_1, \ldots, Z_n) = g((X_1, Y_1), \ldots, (X_n, Y_n)) = \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$$

**Observation:**

If we change $Z_i = (X_i, Y_i)$, then $g$ can change $c_i = 1/n$ at most.

(Look at how much $\widehat{R}_n(f)$ can change if we change either $X_i$ or $Y_i$!)

$\Rightarrow$ McDiarmid can be applied for *g!*

# Bounded Difference Condition

The **"bounded difference condition"** is satisfied for $\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

**Corollary:**

for $g = \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

$$\Pr\{g - \mathbb{E}[g] \geq \varepsilon\} \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^{n} c_i^2}\right) \qquad c_i = 1/n$$

$$\Pr\left\{\left|\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| - \mathbb{E}[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|]\right| \geq \varepsilon\right\} \leq 2\exp\left(-2\varepsilon^2 n\right)$$

$\Rightarrow \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$ is concentrated around its mean!

Therefore, it is enough to study how $\mathbb{E}[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|]$ behaves.

The Vapnik-Chervonenkis inequality does that with the ***shatter coefficient*** (and ***VC dimension)!***

# Concentration and Expected Value

$$Z_n = \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$$

# Vapnik-Chervonenkis inequality

**Our main goal is to bound** $\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

We already know:

$$\Pr\left\{ \left| \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| - \mathbb{E}[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|] \right| \geq \varepsilon \right\} \leq 2\exp\left(-2\varepsilon^2 n\right)$$

**Vapnik-Chervonenkis inequality:**

$$\mathbb{E}\left[ \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| \right] \leq 2\sqrt{\frac{\log(2S_{\mathcal{F}}(n))}{n}}$$

**Corollary: Vapnik-Chervonenkis theorem:**

$$\Pr\left( \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > t \right) \leq 4S_{\mathcal{F}}^2(n)\exp(-nt^2/8)$$

We will define $S_{\mathcal{F}}(n)$ later.

# Shattering

# How many points can a linear boundary classify exactly in 1D?

2 pts

3 pts

There exists placement s.t. all labelings can be classified
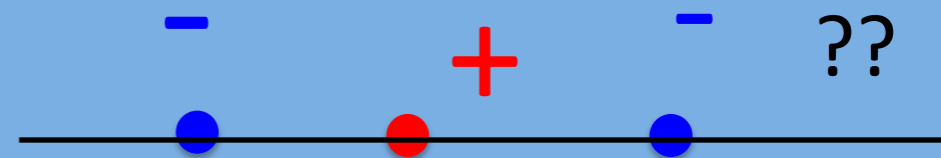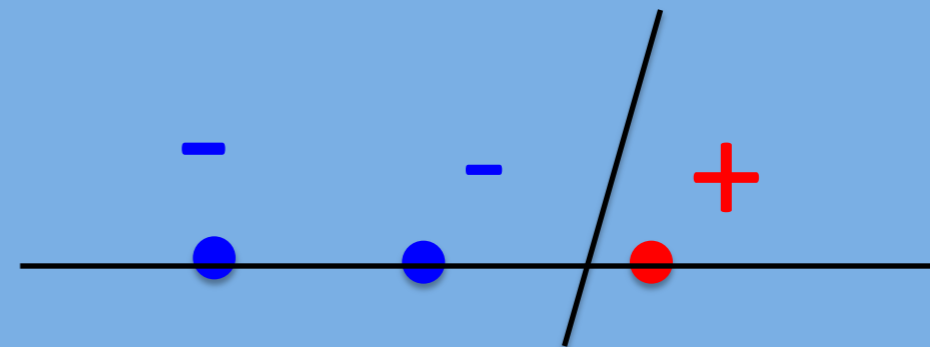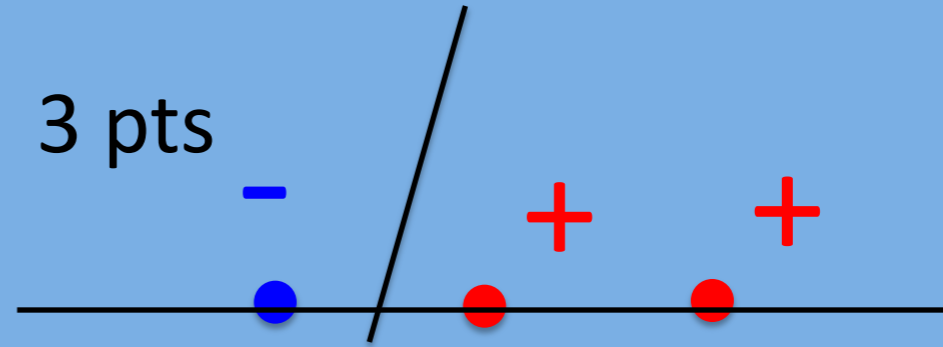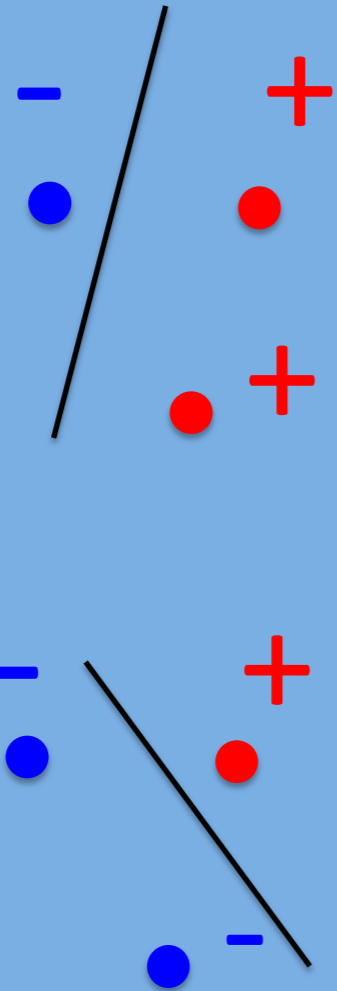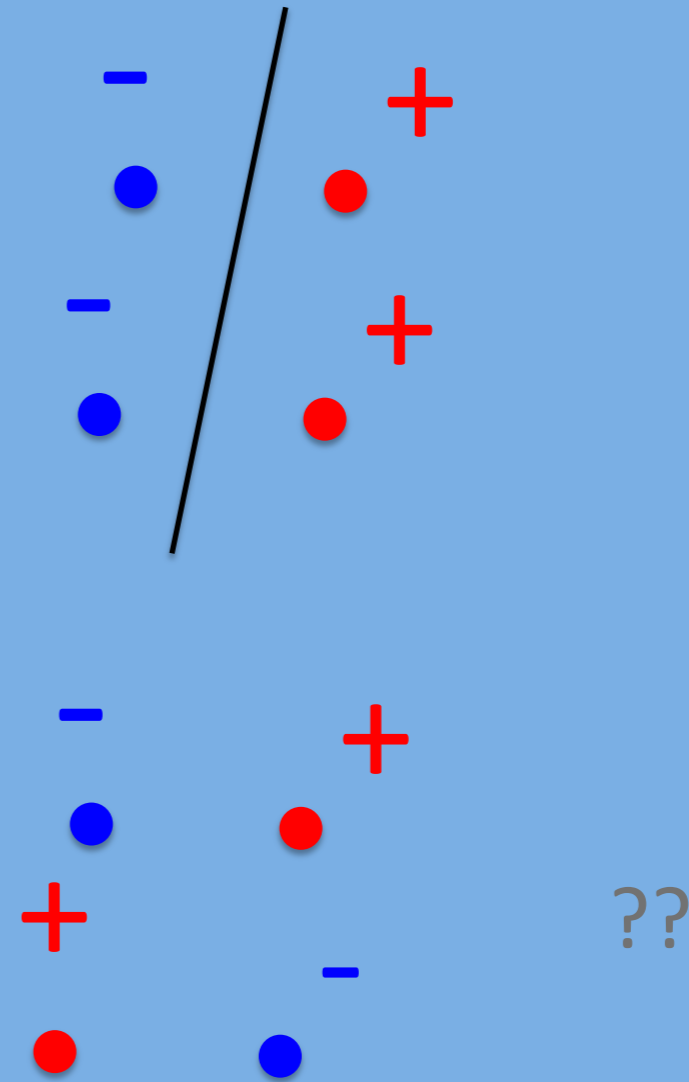
??

The answer is 2

# How many points can a linear boundary classify exactly in 2D?

3 pts

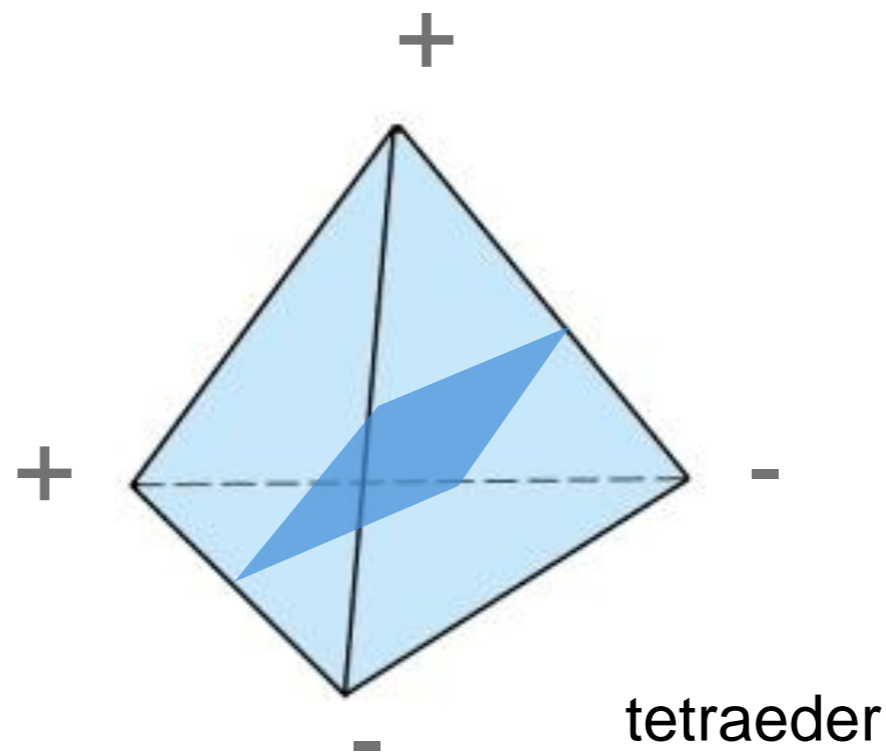There exists placement s.t. all labelings can be classified

4 pts

??

The answer is 3

# How many points can a linear boundary classify exactly in 3D?

The answer is 4



tetraeder

# How many points can a linear boundary classify exactly in d-dim?

The answer is d+1

# Growth function, Shatter coefficient

Let $\mathcal{F} = \mathcal{X} \rightarrow \{0, 1\}$

How many different behaviour can we get with $[f(x_1), \ldots, f(x_n)]$, $f \in \mathcal{F}$?

**Definition**

$$S_{\mathcal{F}}(x_1, \ldots, x_n) = |\{f(x_1), \ldots, f(x_n)\}; f \in \mathcal{F}|$$

(=5 in this example)

**Growth function, Shatter coefficient**

$$S_{\mathcal{F}}(n) = \max_{x_1, \ldots, x_n} |\{f(x_1), \ldots, f(x_n)\}; f \in \mathcal{F}|$$

maximum number of behaviors on *n* points

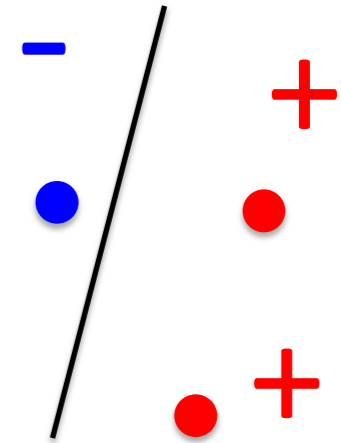| $|\mathcal{F}| = 7$ | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $f_1$ | 0 | 0 | 0 |
| $f_2$ | 0 | 1 | 0 |
| $f_3$ | 1 | 1 | 1 |
| $f_4$ | 1 | 0 | 0 |
| $f_5$ | 0 | 1 | 1 |
| $f_6$ | 0 | 1 | 0 |
| $f_7$ | 1 | 1 | 1 |

# Growth function, Shatter coefficient

**Definition**

$$S_{\mathcal{F}}(x_1, \ldots, x_n) = |\{f(x_1), \ldots, f(x_n)\}; f \in \mathcal{F}|$$

**Growth function, Shatter coefficient**

$$S_{\mathcal{F}}(n) = \max_{x_1, \ldots, x_n} |\{f(x_1), \ldots, f(x_n)\}; f \in \mathcal{F}|$$

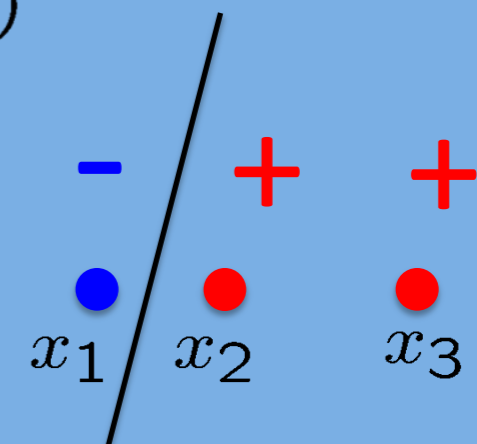maximum number of behaviors on *n* points

**Example:** Half spaces in 2D $\Rightarrow S_{\mathcal{F}}(3) = 2^3 = 8$

(Although $\exists x_1, x_2, x_3$ such that $S_{\mathcal{F}}(x_1, x_2, x_3) = 6 < 8$)

$\{\emptyset\}, \{x_1\}, \{x_3\}, \{x_1, x_2\}, \{x_2, x_3\}, \{x_1, x_2, x_3\}$
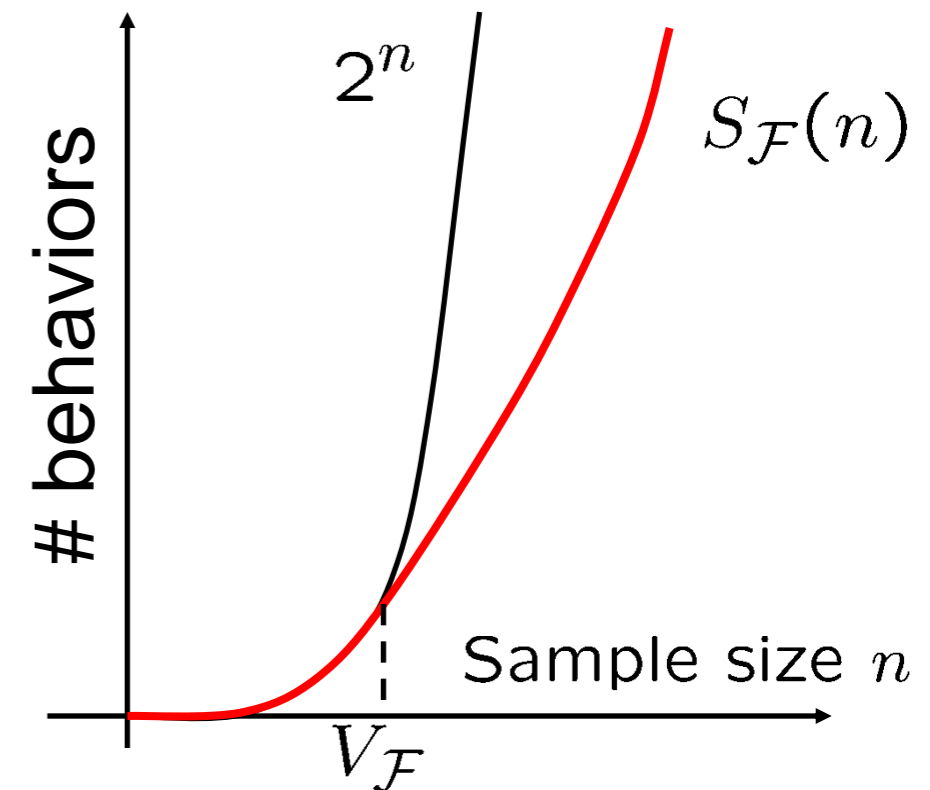We can't get $\{x_2\}$ and $\{x_1, x_3\}$

# VC-dimension

**Definition**

$$S_{\mathcal{F}}(x_1, \dots, x_n) = |\{f(x_1), \dots, f(x_n)\}; f \in \mathcal{F}|$$

**Growth function, Shatter coefficient**

$$S_{\mathcal{F}}(n) = \max_{x_1, \dots, x_n} |\{f(x_1), \dots, f(x_n)\}; f \in \mathcal{F}|$$

maximum number of behaviors on *n* points



**Definition: VC-dimension**

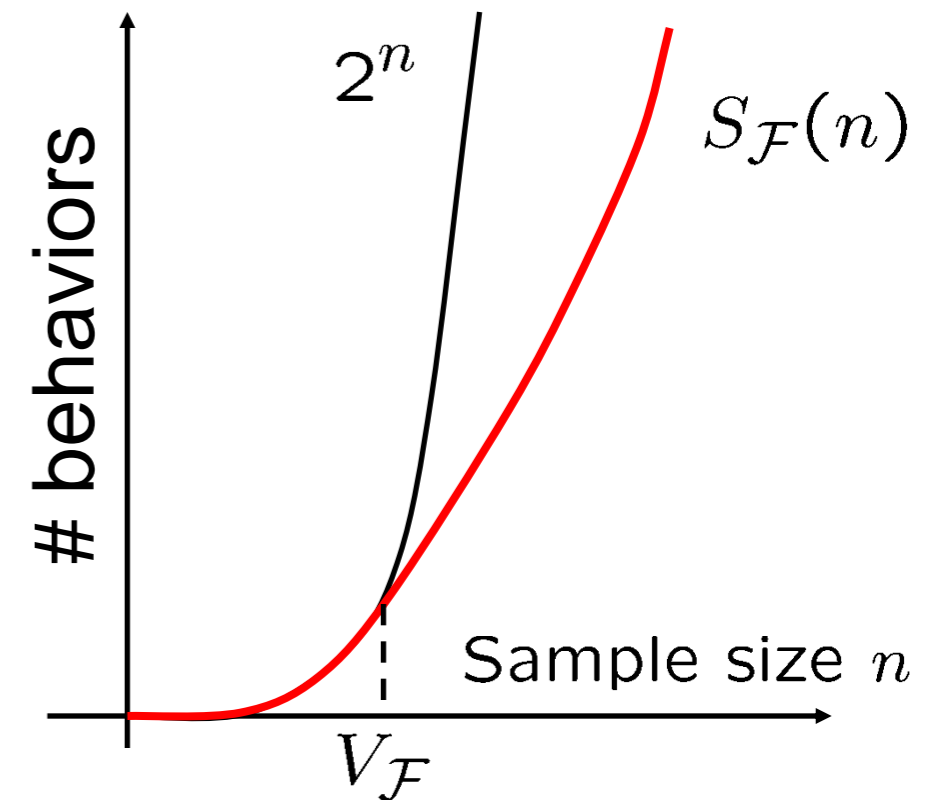$$V_{\mathcal{F}} = \max\{n : S_{\mathcal{F}}(n) = 2^n\}$$

**Definition: Shattering**

$\mathcal{F}$ shatters the sample $x_1, \dots, x_n$ iff $\mathcal{F}$ has all the $2^n$ behaviors on the sample.

**Note:** $V_{\mathcal{F}}$ is the size of largest shattered sample

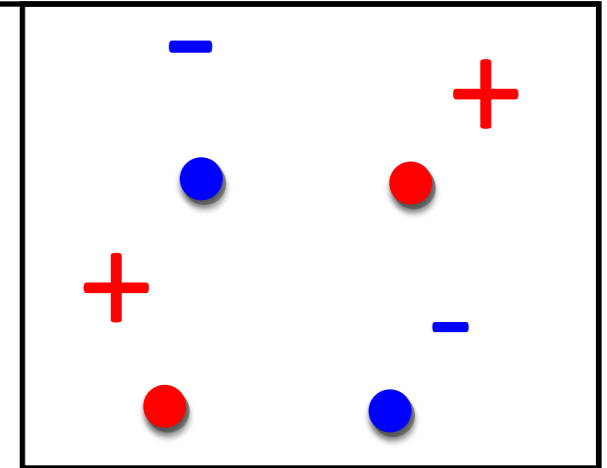# VC-dimension

**Definition** $V_{\mathcal{F}} = \max\{n : S_{\mathcal{F}}(n) = 2^n\}$



- If the VC dimension is $n$, then we can find $n$ points that can be shattered, i.e. show $2^n$ behaviours.

- $n+1$ points never show $2^{n+1}$ behaviours.

# VC-dimension

- You pick set of points $x_1, \ldots, x_n$

- Adversary assigns labels $y_1, \ldots, y_n$

- If $VC_{\mathcal{F}} \geq n$, then you find a hypothesis $f$ in $\mathcal{F}$ consistent with the labels, i.e. $f(x_i) = y_i \quad (1 \leq i \leq n)$

- If $VC_{\mathcal{F}} = n$, then for any $n+1$ points, there exists a labeling that cannot be shattered (can't find a hypothesis $f$ in $\mathcal{F}$ consistent with it)

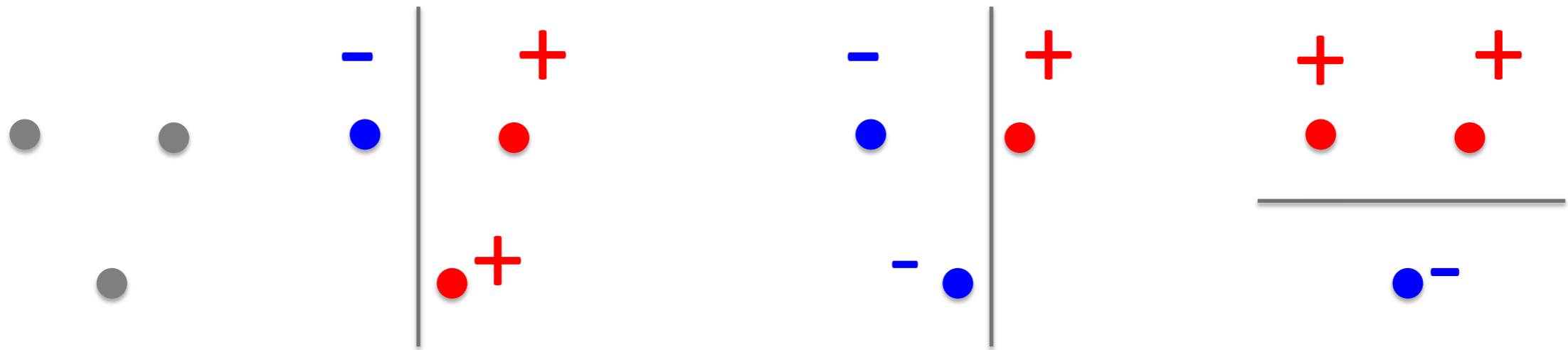The VC dimension measures how rich $\mathcal{F}$ is.

If the VC dimension is high, e.g. $\infty$, then it is easy to overfit!

# Examples

What's the VC dim. of decision stumps in 2d?



There is a placement of 3 pts that can be shattered $\Rightarrow$ VC dim $\geq 3$

## What's the VC dim. of decision stumps in 2d?

If VC dim = 3, then for all placements of 4 pts, there exists a labeling that can't be shattered

3 collinear

1 in convex hull of other 3

quadrilateral

# VC dim. of axis parallel rectangles in 2d

What's the VC dim. of axis parallel rectangles in 2d?

$$f(x) = \text{sign}(1 - 2 \cdot 1_{\{x \in \text{ rectangle}\}})$$



There is a placement of 3 pts that can be shattered $\Rightarrow$ VC dim $\geq$ 3

There is a placement of 4 pts that can be shattered $\Rightarrow$ VC dim $\geq 4$
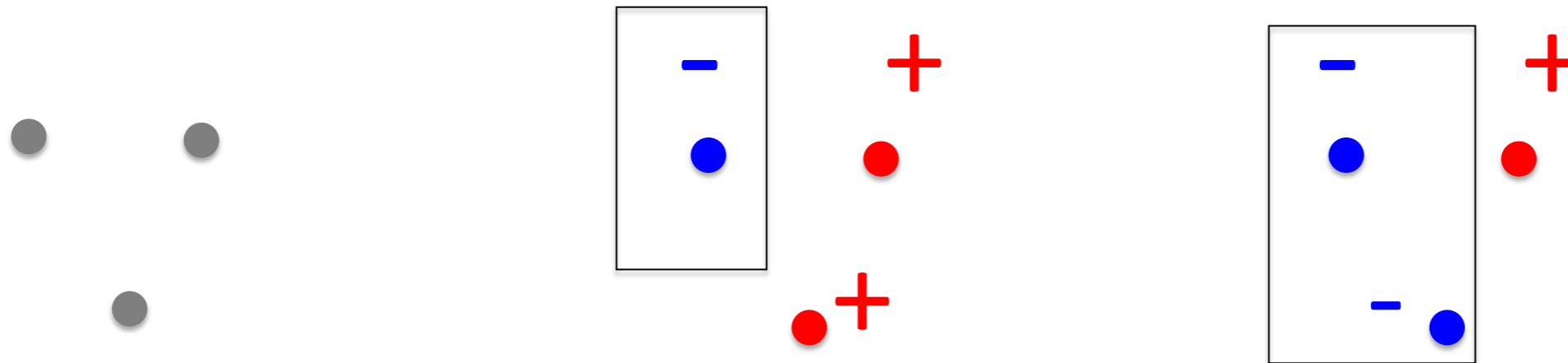
# VC dim. of axis parallel rectangles in 2d

What's the VC dim. of axis parallel rectangles in 2d?

$$f(x) = \text{sign}(1 - 2 \cdot 1_{\{x \in \text{ rectangle}\}})$$

If VC dim = 4, then for all placements of 5 pts, there exists a labeling that can't be shattered



4 collinear

2 in convex hull

1 in convex hull

pentagon

# Sauer's Lemma

We already know that $S_{\mathcal{F}}(n) \le 2^n$   [Exponential in *n*]

> Sauer's lemma:
> $$S_{\mathcal{F}}(n) \le \sum_{k=0}^{VC_{\mathcal{F}}} \binom{n}{k}$$

The VC dimension can be used to upper bound the shattering coefficient.

Corollary:   $S_{\mathcal{F}(n)} \le (n+1)^{VC_{\mathcal{F}}}$   [Polynomial in *n*]

$$S_{\mathcal{F}}(n) \le \left(\frac{ne}{VC_{\mathcal{F}}}\right)^{VC_{\mathcal{F}}}$$

# Proof of Sauer's Lemma

Write all different behaviors on a sample $(x_1, x_2, \ldots x_n)$ in a matrix:

|  | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $f_1$ | 0 | 0 | 0 |
| $f_2$ | 0 | 1 | 0 |
| $f_3$ | 1 | 1 | 1 |
| $f_4$ | 1 | 0 | 0 |
| $f_5$ | 0 | 1 | 0 |
| $f_6$ | 1 | 1 | 1 |
| $f_7$ | 0 | 1 | 1 |

$|\mathcal{F}| = 7$

|  | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $f_1$ | 0 | 0 | 0 |
| $f_2$ | 0 | 1 | 0 |
| $f_3$ | 1 | 1 | 1 |
| $f_4$ | 1 | 0 | 0 |
| $f_7$ | 0 | 1 | 1 |

$|\mathcal{F}| = 7$

# Proof of Sauer's Lemma

$|\mathcal{F}| = 7$

|       | $x_1$ | $x_2$ | $x_3$ |
|-------|-------|-------|-------|
| $f_1$ | 0     | 0     | 0     |
| $f_2$ | 0     | 1     | 0     |
| $f_3$ | 1     | 1     | 1     |
| $f_4$ | 1     | 0     | 0     |
| $f_7$ | 0     | 1     | 1     |

$= A$

Shattered subsets of columns:

$$\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_3\}, \{x_1, x_2\}, \{x_1, x_3\}$$

We will prove that

$$S_\mathcal{F}(x_1, \ldots, x_n) = \# \text{ rows}(A) \leq \# \text{ shattered subsets of columns of } A \leq \sum_{k=0}^{VC_\mathcal{F}} \binom{n}{k}$$

Therefore,

$$S_{\mathcal{F}(n)} = \max_{x_1, \ldots, x_n} S_\mathcal{F}(x_1, \ldots, x_n) \leq \sum_{k=0}^{VC_\mathcal{F}} \binom{n}{k}$$

# Proof of Sauer's Lemma

$|\mathcal{F}| = 7$

|       | $x_1$ | $x_2$ | $x_3$ |
|-------|-------|-------|-------|
| $f_1$ | 0     | 0     | 0     |
| $f_2$ | 0     | 1     | 0     |
| $f_3$ | 1     | 1     | 1     |
| $f_4$ | 1     | 0     | 0     |
| $f_7$ | 0     | 1     | 1     |

$= A$

Shattered subsets of columns:

$$\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_3\}, \{x_1, x_2\}, \{x_1, x_3\}$$

Lemma 1    # shattered subsets of columns of $A \leq \sum\limits_{k=0}^{VC_{\mathcal{F}}} \binom{n}{k}$

In this example: $6 \leq 1+3+3=7$

Lemma 2    # rows$(A) \leq$ # shattered subsets of columns of $A$

for any binary matrix with no repeated rows.
In this example: $5 \leq 6$

# **Proof of Lemma 1**

$|\mathcal{F}| = 7$

| | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $f_1$ | 0 | 0 | 0 |
| $f_2$ | 0 | 1 | 0 |
| $f_3$ | 1 | 1 | 1 |
| $f_4$ | 1 | 0 | 0 |
| $f_7$ | 0 | 1 | 1 |

$= A$

Shattered subsets of columns:

$$\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_3\}, \{x_1, x_2\}, \{x_1, x_3\}$$

In this example: $6 \leq 1+3+3=7$

Lemma 1   # shattered subsets of columns of $A \leq \sum_{k=0}^{VC_\mathcal{F}} \binom{n}{k}$

Proof

$VC_\mathcal{F}$ is the size of largest imaginable shattered sample. $VC_\mathcal{F} = \max\{n : S_\mathcal{F}(n) = 2^n\}$

If a shattered subsets of columns has $d$ elements, then $VC_\mathcal{F} \geq d$

For example if $\{x_1, x_3\}$ are shattered in $A$, then $VC_\mathcal{F} \geq 2$.

Lemma 2 $\quad\#\ \mathrm{rows}(A) \leq \#\ \text{shattered subsets of columns of }A$

for any binary matrix with no repeated rows.

Proof Induction on the number of columns

**Base case:** A has one column. There are three cases:

$A = (0) \quad \Rightarrow \quad 1 \leq 1 \qquad$ shattered subsets of columns: $\{\emptyset\}$

$A = (1) \quad \Rightarrow \quad 1 \leq 1 \qquad$ shattered subsets of columns: $\{\emptyset\}$

$A = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad \Rightarrow \quad 2 \leq 2 \qquad$ shattered subsets of columns: $\{\emptyset\}, \{x_1\}$

**Inductive case:** A has at least two columns. $x_m$

Let $A'$ be $A$ minus its last column $x_m$ removed

In $A'$ each row can occure once or twice.

If "twice" $\Rightarrow$ move one of them to $B$ the other to $C$

If "once" $\Rightarrow$ move them to $C$

$$\begin{array}{|c|c|} \hline C & \\ \hline A' & \\ \hline B & \\ \hline \end{array} = A$$

**We have,**

$$\# \text{ rows}(A) = \# \text{ rows}(B) + \# \text{ rows}(C)$$

$$\leq \# \text{ shattered subsets of columns of } (B)$$
$$+ \# \text{ shattered subsets of columns of } (C)$$

**By induction (less columns)**

| 0 | 0 | 0 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 1 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 1 |

$$\{\emptyset\} \qquad\qquad \{\emptyset\}, \{x_1\}, \{x_2\}\{x_1, x_2\}$$

# shattered subsets of columns of $(B)$ + # shattered subsets of columns of $(C)$

$\leq$ # shattered subsets of columns of $(A)$

$$\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_3\}, \{x_1, x_2\}, \{x_1, x_3\}$$

because

$x_m$
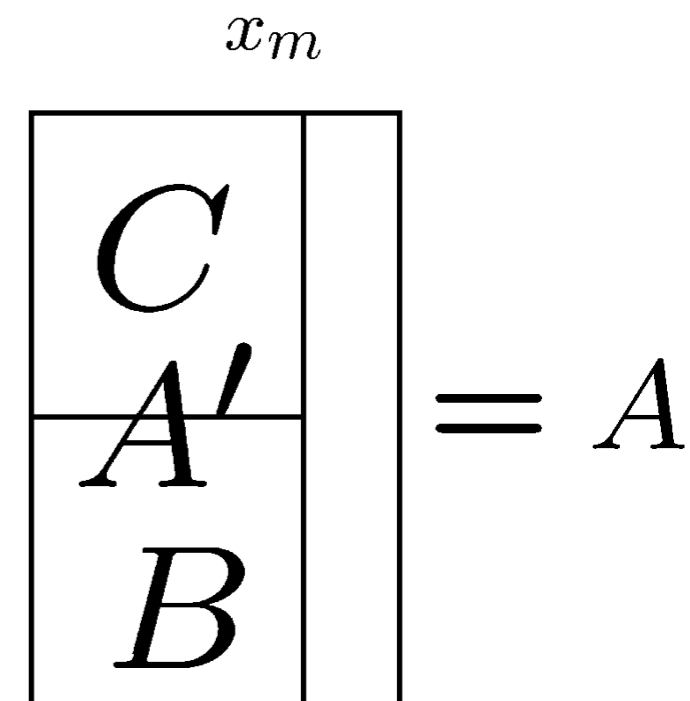
"once" $\Rightarrow$ move them to $C$

Therefore, if $C$ shatters $S$ e.g. $\{x_1, x_2\}$, then $A$ shatters $S$.

"twice" $\Rightarrow$ move one of them to $B$ the other to $C$

Therefore, if $B$ shatters $S$, then $A$ shatters $S \cup x_m$.

$$\begin{array}{c} C \\ A' \\ B \end{array} = A$$

| | | |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 1 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 1 |

# Vapnik-Chervonenkis inequality

When $|\mathcal{F}| = N < \infty$, we already know $\mathbb{E}\left[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|\right] \leq \sqrt{\frac{\log(2N)}{2n}}$

**Vapnik-Chervonenkis inequality:** [We don't prove this]

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|\right] \leq 2\sqrt{\frac{\log(2S_{\mathcal{F}}(n))}{n}}$$

## From Sauer's lemma:

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|\right] \leq 2\sqrt{\frac{\log(2S_{\mathcal{F}}(n))}{n}} \leq 2\sqrt{\frac{VC_{\mathcal{F}}\log(n+1) + \log 2}{n}}$$

Since $\quad |R(f_n^*) - R(f_{\mathcal{F}}^*)| \leq 2\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

Therefore, $\mathbb{E}[|R(f_n^*) - R(f_{\mathcal{F}}^*)|] \leq 4\sqrt{\frac{VC_{\mathcal{F}}\log(n+1) + \log 2}{n}}$   Yayyy! ☺

Estimation error

# Linear (hyperplane) classifiers

We already know that

$$\mathbb{E}[|R(f_n^*) - R(f_{\mathcal{F}}^*)|] \leq 4\sqrt{\frac{VC_{\mathcal{F}}\log(n+1) + \log 2}{n}}$$

Estimation error

For linear classifiers in dimension when $\mathcal{X} = \mathbb{R}^d$: $VC_{\mathcal{F}} = d + 1$.

$$\Rightarrow \mathbb{E}[|R(f_n^*) - R(f_{\mathcal{F}}^*)|] \leq 4\sqrt{\frac{(d+1)\log(n+1) + \log 2}{n}}$$

Estimation error

If we do feature map first, $x = \phi(x) \in \mathbb{R}^{d'}$, then linear separation (SVM) $\Rightarrow VC_{\mathcal{F}} = d' + 1$.

Estimation error
$$\Rightarrow \mathbb{E}[|R(f_n^*) - R(f_{\mathcal{F}}^*)|] \leq 4\sqrt{\frac{(d'+1)\log(n+1) + \log 2}{n}}$$

# Vapnik-Chervonenkis Theorem

We already know from McDiarmid:

$$\Pr\left\{ \left| \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| - \mathbb{E}[\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|] \right| \geq \varepsilon \right\} \leq 2\exp\left(-2\varepsilon^2 n\right)$$

Vapnik-Chervonenkis inequality: $\mathbb{E}\left[ \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| \right] \leq 2\sqrt{\dfrac{\log(2S_{\mathcal{F}}(n))}{n}}$

Corollary: Vapnik-Chervonenkis theorem:  [We don't prove them]

$$\Pr\left( \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > t \right) \leq 4 S_{\mathcal{F}}(2n) \exp(-nt^2/8)$$

$$\Pr\left( \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > t \right) \leq 8 S_{\mathcal{F}}(n) \exp(-nt^2/32)$$

Hoeffding + Union bound for finite function class:

When $|\mathcal{F}| = N < \infty$, $\Rightarrow \Pr\left( \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > t \right) \leq 2N \exp\left(-2nt^2\right)$

# PAC Bound for the Estimation Error

**VC theorem:**

$$\Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| > t\right) \le 8 S_{\mathcal{F}}(n) \exp(-nt^2/32)$$

**Inversion:**

$$8 S_{\mathcal{F}}(n) \exp(-nt^2/32) \le \delta \qquad \Rightarrow t^2 \ge \frac{32}{n} \log\left(\frac{8 S_{\mathcal{F}}(n)}{\delta}\right)$$

$$\Rightarrow \Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| \le 8\sqrt{\frac{\log(S_{\mathcal{F}}(n)) + \log\left(\frac{8}{\delta}\right)}{2n}}\right) \ge 1 - \delta$$

$$S_{\mathcal{F}}(n) \le \left(\frac{ne}{VC_{\mathcal{F}}}\right)^{VC_{\mathcal{F}}} \Rightarrow \Pr\left(\sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)| \le 8\sqrt{\frac{VC_{\mathcal{F}} \log\left(\frac{ne}{VC_{\mathcal{F}}}\right) + \log\left(\frac{8}{\delta}\right)}{2n}}\right) \ge 1 - \delta$$

Don't forget that $|R(f_n^*) - R(f_{\mathcal{F}}^*)| \le 2 \sup_{f \in \mathcal{F}} |\widehat{R}_n(f) - R(f)|$

**Estimation error**

$$\Rightarrow \Pr\left(|R(f_n^*) - R(f_{\mathcal{F}}^*)| \le 16\sqrt{\frac{\log(VC_{\mathcal{F}} \log\left(\frac{ne}{VC_{\mathcal{F}}}\right) + \log\left(\frac{8}{\delta}\right)}{2n}}\right) \ge 1 - \delta$$
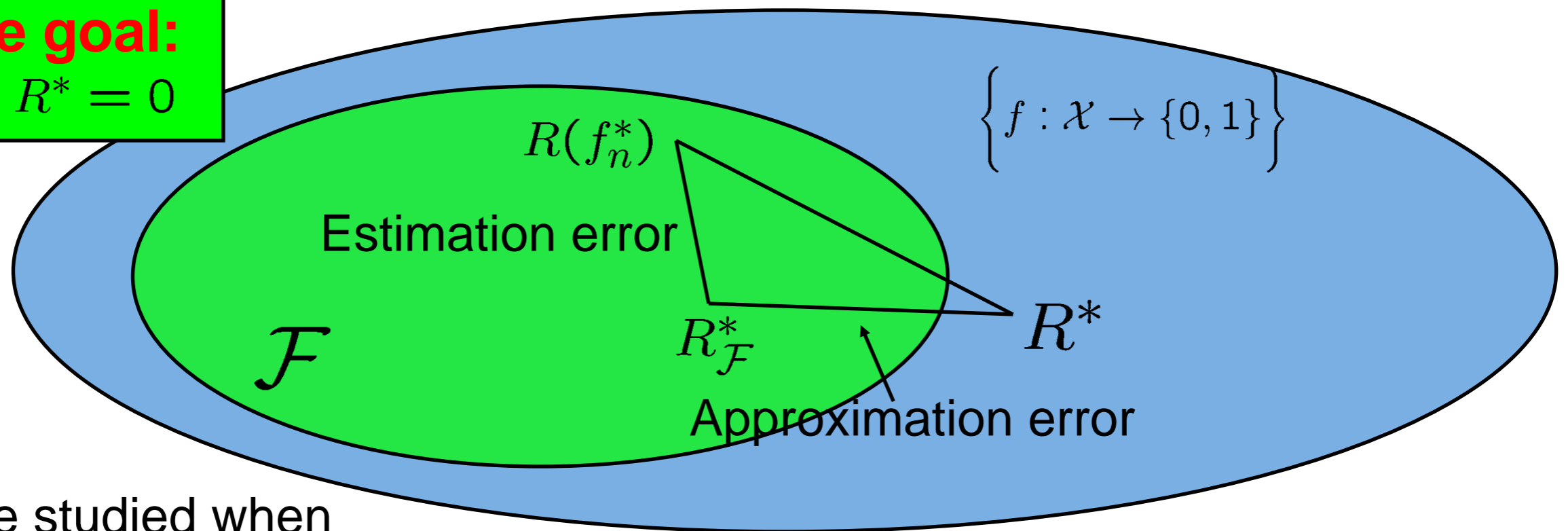
# Structoral Risk Minimization

Risk of the classifier $f_n^*$    Estimation error    Approximation error

$$R(f_n^*) - R^* = \overbrace{R(f_n^*) - R_{\mathcal{F}}^*} + \overbrace{R_{\mathcal{F}}^* - R^*}$$

Bayes risk

**Ultimate goal:**
$R(f_n^*) - R^* = 0$

$R(f_n^*)$

$\left\{ f : \mathcal{X} \to \{0,1\} \right\}$

Estimation error

$\mathcal{F}$

$R_{\mathcal{F}}^*$    $R^*$

Approximation error

So far we studied when
estimation error $\to 0$, but we also want approximation error $\to 0$

Let $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq ... \subseteq \mathcal{F}_n \subseteq ...$ such that $VC_{\mathcal{F}_1} \leq VC_{\mathcal{F}_2} \leq ... \leq VC_{\mathcal{F}_n} \leq ...$

Many different variants…
penalize too complex models to avoid overfitting

# What you need to know

Complexity of the classifier depends on number of points that can be classified exactly

Finite case – Number of hypothesis
Infinite case – Shattering coefficient, VC dimension

PAC bounds on true error in terms of empirical/training error and complexity of hypothesis space

Empirical and Structural Risk Minimization

Thanks for your attention ☺