# AGENDA

**01**  **Preliminaries**

Risk minimization, union bound, and Hoeffding inequality.

**02**  **The case of finite $\mathcal{H}$**

Uniform convergence, sample complexity, error bound, bias-variance tradeoff.

**03**  **The case of infinite $\mathcal{H}$**

VC dimension

# LEARNING THEORY
## PRELIMINARIES

We want to answer three main questions:

1. **Can we make formal the bias/variance tradeoff?**

2. **Why should doing well on the training set tell us anything about generalization error?**

3. **Are there conditions under which we can actually prove that learning algorithms will work well?**

We are going to **define** a **binary classifier**:

$$h_w(x) = g(w^T x)$$

$$g(z) = 1\{z \geq 0\}$$

We establish our **training set as**:

$$S = \left\{\left(x^{(i)}, y^{(i)}\right)\right\}_{i=1}^{m}, S \sim_{iid} D$$

We are going to define the **training error** $\widehat{\varepsilon}_S$ of a hypothesis $h_w$ in a **simple way:**

$$\widehat{\varepsilon}_S(h_w) = \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}\{h_w(x^{(i)}) \neq y^{(i)}\}$$

**Fraction of data points where the hypothesis is wrong.**

The **training error** is also called **RISK.**

As always, our **objective** consists in **minimizing** the **risk** (**training error**):

$$\widehat{w} = \underset{w}{\textbf{argmin}}\ \widehat{\mathcal{E}}_S(h_w)$$

**Minimizing** this **expression** deals with a **non-convex optimization problem**. **Logistic regression** and **SVM** are **convex approximations** to this **problem**.

We are going to **change** the **problem.** Now, the **objective** will reside in **choosing** the **hypothesis** function $h_w$ **instead** of the **parameters** $w$.

Thus, we define the **hypothesis class** $\mathcal{H}$ as the **class** of all **linear classifiers** that the **algorithm** is **choosing from**.

$$\mathcal{H} = \left\{ h_w : h_w(x) = 1\{w^T x \geq 0\}, w \in \mathbb{R}^{n+1} \right\}$$

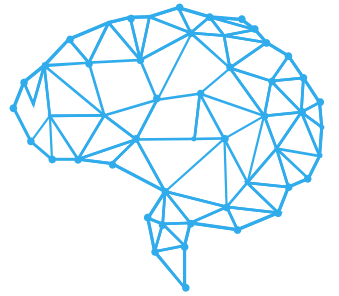Therefore, **empirical risk minimization** is **redefined** as:

$$\widehat{h_w} = \underset{h_w \in \mathcal{H}}{\textbf{argmin}} \, \widehat{\mathcal{E}}_S(h_w)$$

**NOTE**: the **hypothesis class** $\mathcal{H}$ can represent any set of functions.

Let us remember that the **main goal** resides in the **generalization error not** in the **training error**. The **generalization error** would be **defined as**:
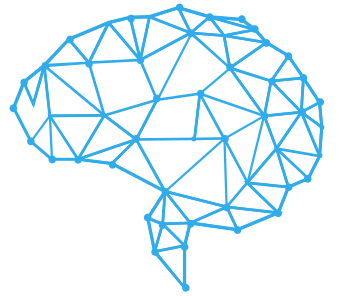
$$\varepsilon(h_w) = P_{(x,y) \sim D}(h_w(x) \neq y)$$

**Probability** that, if we now **draw** a **new example** $(x, y)$ **from** the **distribution D,** $h$ will **misclassify** it.

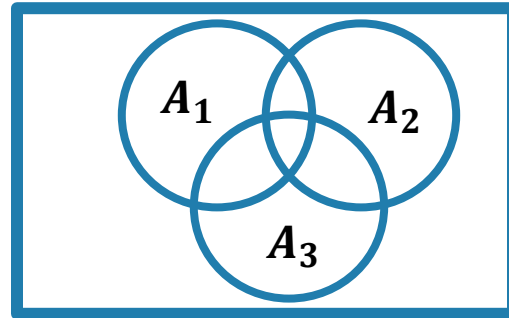We want to **make** an **estimation** $\widehat{\varepsilon}_S$ **(training error)** to get **close** to the **generalization error** $\varepsilon(h_w)$.

To **reduce** the **generalization error** we will **need two lemmas**:

1.  **Union bound:** Let $A_1, A_2, \ldots, A_k$ be $k$ **different events** (that may not be independent). Then $P(A_1 \cup \cdots \cup A_k) \leq P(A_1) + \cdots + P(A_k)$.



2.  **Hoeffding inequality:** let $Z1, \ldots, Zn$ be $n$ **independent** and **identically distributed** (**iid**) **random variables** drawn from a **Bernoulli**(φ) **distribution** with **mean** $\phi$. Therefore, $P(Z_i = 1) = \phi$ and $P(Z_i = 0) = 1 - \phi$. Let $\hat{\phi} = \frac{1}{m} \sum_{i=1}^{m} z_i$ and let any $\gamma > 0$ be fixed. Then
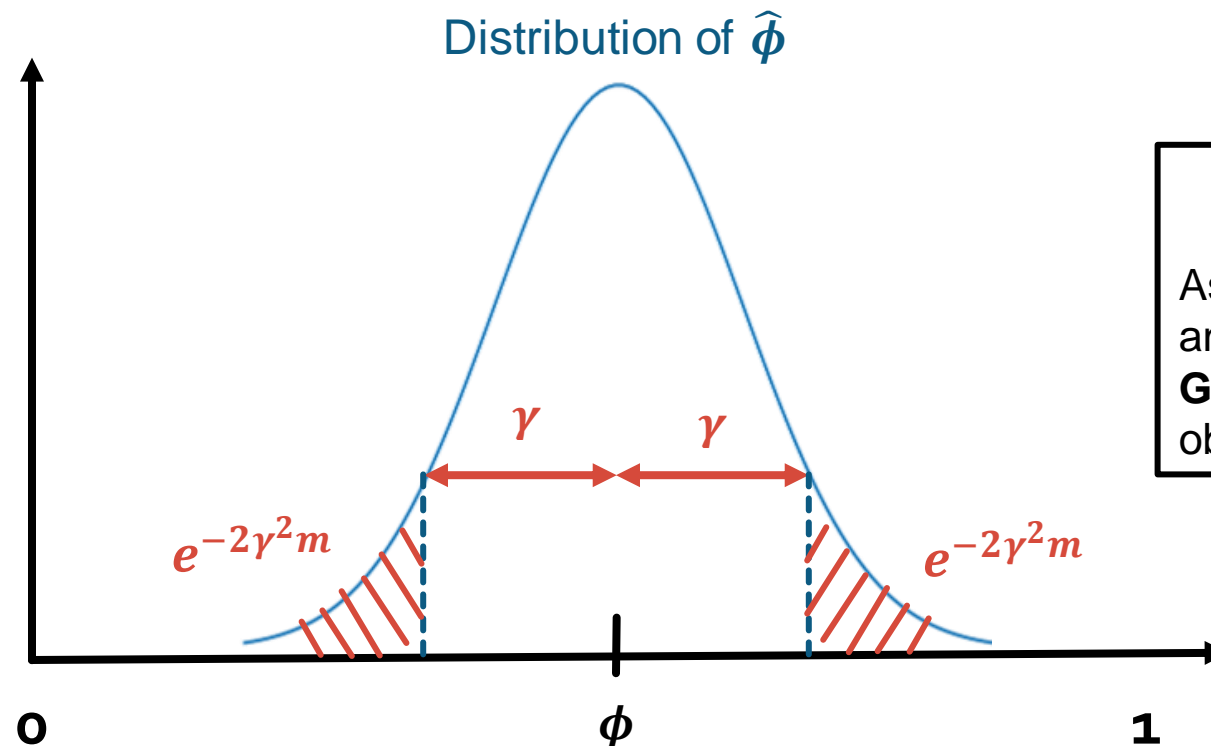
$$P\left(\left|\phi - \hat{\phi}\right| > \gamma\right) \leq 2e^{-2\gamma^2 m}$$

The **Hoeffding inequality** says that if we take $\widehat{\phi}$ (the **average** of $m$ **Bernoulli($\phi$) random variables**) to be our **estimate** of $\phi$, then the **probability** of **our being far** from the **true value** is **small,** so long as $m$ is **large.**



Distribution of $\widehat{\phi}$

$\gamma$ $\gamma$

$e^{-2\gamma^2 m}$ $e^{-2\gamma^2 m}$

0 $\phi$ 1

**Central Limit Theorem**

As you take more $m$ samples and **average them**, a **Gaussian distribution** is obtained.

THE CASE OF FINITE $\mathcal{H}$

Let us consider that we have a **finite hypothesis class** $\mathcal{H} = \{h_1, \dots, h_k\}$ consisting of $\boldsymbol{k}$ **hypotheses** or **functions** mapping from $\boldsymbol{\chi}$ to $\{0, 1\}$.

**Risk** minimization will **choose** the **hypothesis** with the **lowest training error**.

We are going to **prove** that:

1. $\widehat{\varepsilon}_S \approx \varepsilon$.

2. There is an upper-bound to $\widehat{\varepsilon}_S$.

Therefore, if we **minimize** the **training error,** the **generalization** error **will decrease** as well.

We are going to take a **fixed hypothesis** $h_j \in \mathcal{H}$ and will consider a **Bernoulli** random variable $Z_i \underset{iid}{\sim} D$ which **misclassifies** an **example** $Z_i = \mathbf{1}\{h_j(x^{(i)}) \neq y^{(i)}\} \in \{0, 1\}$.

The **probability** that, from a **fixed hypothesis** $h_j$, we **misclassify** an **example** is the **expected value (mean of the distribution)** or **generalization error**:

$$P(Z_i = 1) = \varepsilon(h_j)$$

On the other hand, we know that the **training error** is computed as the **fraction** of **misclassified examples (mean of the sample)**:

$$\widehat{\varepsilon}_S(h_j) = \frac{1}{m}\sum_{i=1}^{m} Z_i = \frac{1}{m}\sum_{i=1}^{m} \mathbf{1}\{h_j(x^{(i)}) \neq y^{(i)}\}$$

We will use the **Hoeffding inequality** to look at the **difference between** the **generalization** and **training errors**:

$$P\left(\left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) \leq 2e^{-2\gamma^2 m}$$

We have **proved** that for a **fixed hypothesis** $h_j$ and a **large training set**, the **training error** will **approximate** the **generalization error** with a **high probability**.

Now, let us **prove** this statement for **all** $h \in \mathcal{H}$.

Let us think of $A_j$ as an event $\left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma$, thus $P(A_j) \leq 2e^{-2\gamma^2 m}$. Using the union bound lemma we have:

$$P\left(\exists h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) = P\left(A_1 \bigcup \cdots \bigcup A_k\right)$$

$$P\left(\exists h \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) \leq \sum_{i=1}^{k} P(A_i)$$

$$P\left(\exists h \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) \leq \sum_{i=1}^{k} 2e^{-2\gamma^2 m}$$

$$P\left(\exists h \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) \leq 2ke^{-2\gamma^2 m}$$

The **probability** that such **hypothesis does not exist** is defined as:

$$P\left(\neg\exists h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) = 1 - P\left(\exists h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right)$$

$$P\left(\neg\exists h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) = P\left(\forall h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| \leq \gamma\right)$$

$$P\left(\neg\exists h_j \in \mathcal{H} / \left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| > \gamma\right) \geq 1 - 2ke^{-2\gamma^2 m}$$

We have proved that with **probability at least** $1 - 2ke^{-2\gamma^2 m}$, the **generalization error** $\mathcal{E}(h_j)$ will be **within** a **distance** $\gamma$ of the **training error** $\widehat{\mathcal{E}}(h_j)$ for all $h \in \mathcal{H}$.

**UNIFORM CONVERGENCE**
**(Holds for all hypotheses)**

There are **three quantities** of **interest**: $m, \gamma$, and the **probability** of **error** $\delta$. We can **bound one** in terms of the other two.

Given $\delta$ and $\gamma$, we can obtain the **size** of the **training set** $m$ for which the **training error** will be **within** $\gamma$ of the **generalization error** with at least probability $1 - \delta$.

$$\delta = 2ke^{-2\gamma^2 m}$$

$$m \geq -\frac{1}{2\gamma^2}\log\left(\frac{\delta}{2k}\right)$$

$$m \geq \frac{1}{2\gamma^2}\log\left(\left(\frac{\delta}{2k}\right)^{-1}\right)$$

$$m \geq \frac{1}{2\gamma^2}\log\left(\frac{2k}{\delta}\right)$$

Therefore, we need a **training set size** of $m \geq \frac{1}{2\gamma^2} log\left(\frac{2k}{\delta}\right)$ to guarantee that with **probability** of **at least** $1 - \delta$, we have that the **training error** $\widehat{\mathcal{E}}(h_j)$ is within $\gamma$ of the **generalization error** $\mathcal{E}(h_j)$. Formally

$$\left|\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)\right| \leq \gamma, \forall h \in \mathcal{H}$$

**SAMPLE COMPLEXITY**
**(Number of training examples needed to achieve a certain bounding error)**

**NOTE:** we can see that **even** if we **increase** the **number** of **hypotheses** $k$ in the class $\mathcal{H}$, the **number** of **training examples** $m$ needed will **remain small.**

If we hold $m$ and $\delta$ **fixed**, we can get the **following**:

$$\delta = 2ke^{-2\gamma^2 m}$$

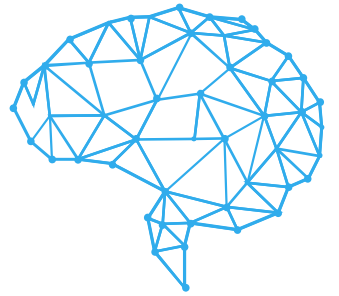$$\gamma = \sqrt{\frac{1}{2m}\log\left(\frac{2k}{\delta}\right)}$$

We want to make $\gamma$ the **upper bound error,** thus

$$\left|\varepsilon(h_j) - \widehat{\varepsilon}(h_j)\right| \leq \gamma$$

$$\left|\varepsilon(h_j) - \widehat{\varepsilon}(h_j)\right| \leq \sqrt{\frac{1}{2m}\log\left(\frac{2k}{\delta}\right)}$$

Let us assume that the **uniform convergence** $\forall h \in \mathcal{H}, |\mathcal{E}(h_j) - \widehat{\mathcal{E}}(h_j)| \leq \gamma$ **holds.**

Is there something that we can **prove** about the **generalization error** $\mathcal{E}$ using the **estimated hypothesis** $\widehat{h}$ with **empirical risk minimization**?  Remembering that:
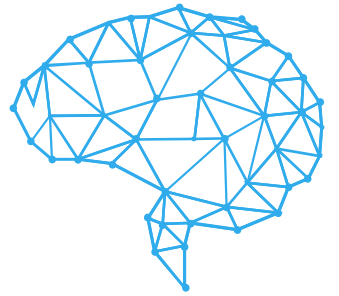
$$\widehat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \, \widehat{\mathcal{E}}(h)$$

Now, we are going to define the **best hypothesis** as the hypothesis that **minimizes** the **generalization error**:

$$h^* = \underset{h \in \mathcal{H}}{argmin} \, \mathcal{E}(h)$$

Starting with the **uniform convergence assumption** we have:

$$\left|\varepsilon(\boldsymbol{h_j}) - \widehat{\varepsilon}(\boldsymbol{h_j})\right| \leq \boldsymbol{\gamma}$$

$$\varepsilon(\widehat{\boldsymbol{h}}) - \widehat{\varepsilon}(\widehat{\boldsymbol{h}}) \leq \boldsymbol{\gamma}$$

$$\varepsilon(\widehat{\boldsymbol{h}}) \leq \widehat{\varepsilon}(\widehat{\boldsymbol{h}}) + \boldsymbol{\gamma}$$

Because we obtained $\widehat{\boldsymbol{h}}$ with **empirical risk minimization**, there is **no** other **hypothesis** with **less training error** than $\widehat{\boldsymbol{h}}$, thus $\widehat{\varepsilon}(\widehat{\boldsymbol{h}}) \leq \widehat{\varepsilon}(h^*)$ and the **inequality remains true:**

$$\varepsilon(\widehat{\boldsymbol{h}}) \leq \widehat{\varepsilon}(h^*) + \boldsymbol{\gamma}$$

$$\varepsilon(\widehat{\boldsymbol{h}}) \leq \widehat{\varepsilon}(h^*) + \boldsymbol{2\gamma}$$

**Theorem:**

Let $|H| = k$, and let any $n, \delta$ be fixed. Then with **probability at least $1 - \delta$**, we have that
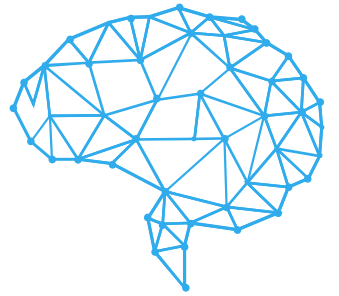
$$\varepsilon(\widehat{h}) \leq \widehat{\varepsilon}(h^*) + 2\gamma$$

$$\varepsilon(\widehat{h}) \leq \min_{h \in \mathcal{H}} \varepsilon(h) + 2\sqrt{\frac{1}{2m} log\left(\frac{2k}{\delta}\right)}$$

Thus, our **generalization error** of the **hypothesis** obtained with **ERM** $\varepsilon(\widehat{h})$ will be at most $2\gamma$ **higher** than the **training error** of the **best possible hypothesis** $\widehat{\varepsilon}(h^*)$.

By **analyzing** the **theorem** we can see the following:

$$\mathcal{E}(\widehat{h}) \leq \min_{h \in \mathcal{H}} \mathcal{E}(h) + 2\sqrt{\frac{1}{2m} log\left(\frac{2k}{\delta}\right)}$$
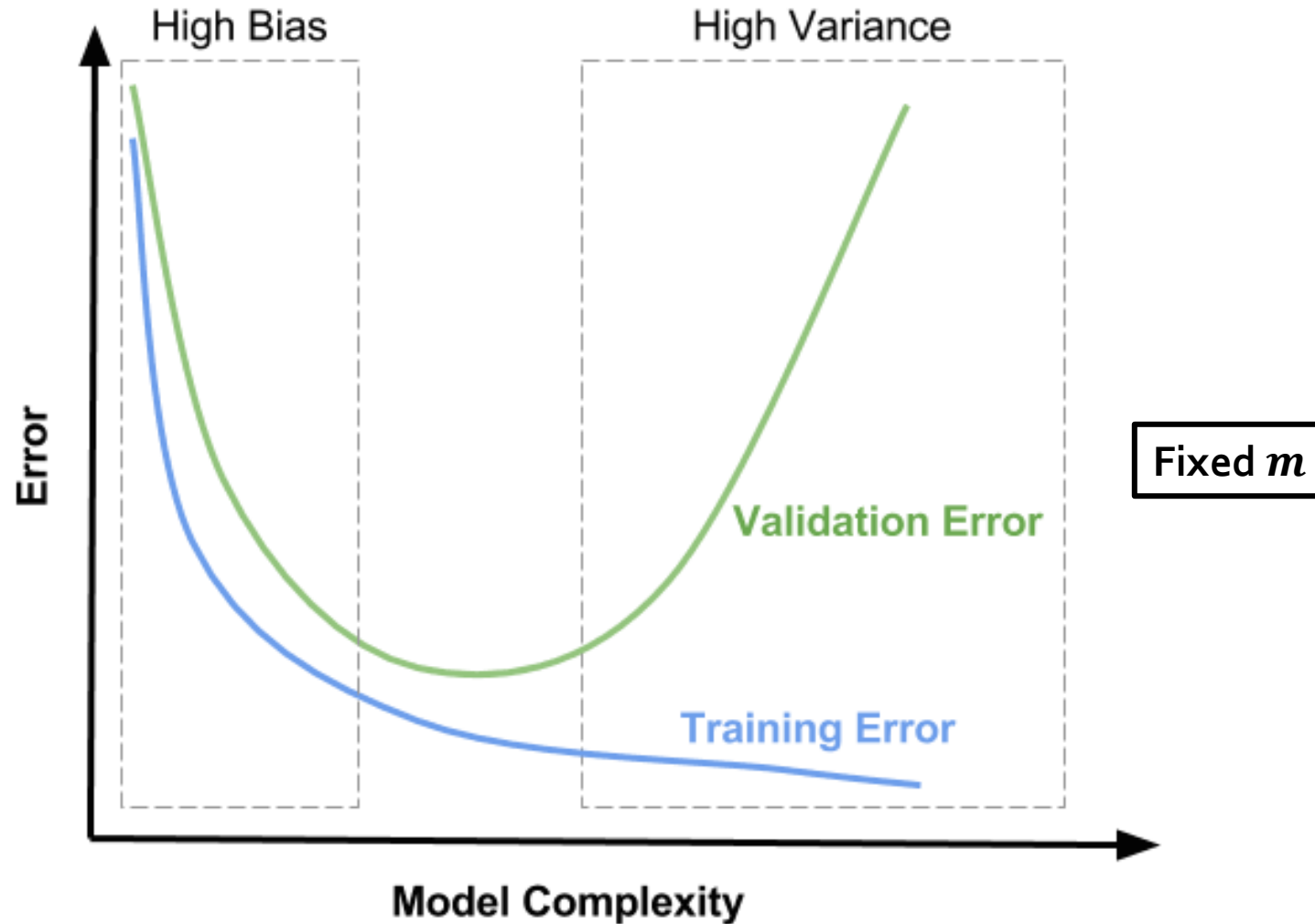
If we switch to a **larger hypothesis class** $\mathcal{H}' \supseteq \mathcal{H}$ (i.e. quadratic), then $\min_{h \in \mathcal{H}} \mathcal{E}(h)$ will **decrease** because we have a larger set of hypothesis for which we can obtain the minimum, thus we **reduce the bias**.

On the other hand, $k$ will **become larger** resulting in an **increase** of the **second term**, thus **increasing the variance**.
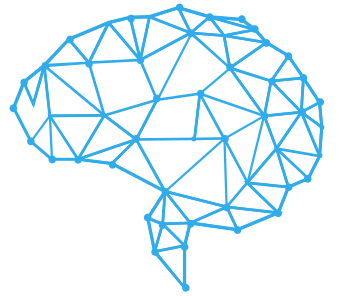
Fixed $m$

By looking at the **training** set **size $m$,** we can obtain the following **complexity bound**:

**Corollary:** Let $|\mathcal{H}| = k$, and let any $\delta, \gamma$ **be fixed**. Then for $\varepsilon(\widehat{h}) \leq \min\limits_{h \in H} \varepsilon(h) + 2\gamma$ to hold with **probability at least $1 - \delta$**, it suffices that:

$$m > \frac{1}{2\gamma^2} \log\left(\frac{2k}{\delta}\right)$$

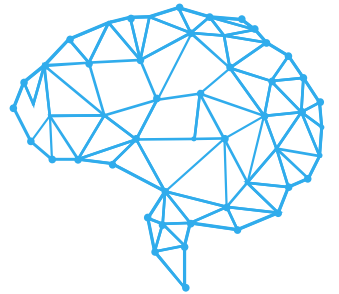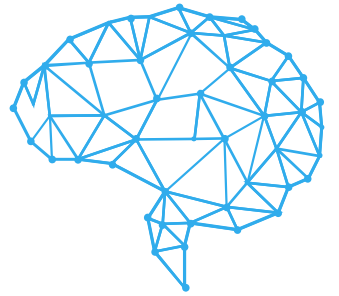$$m = O\left(\frac{1}{\gamma^2} \log\left(\frac{k}{\delta}\right)\right)$$

Many **hypothesis classes** contain an **infinite** number of **functions** (i.e. **any function parametrized** by **real** numbers). We want to **prove** the **previous results** for this **infinite space** of **functions**.

We are going to make some **statements** that are **not correct** at all but **will help** with the **understanding** of the **proof.**

Let us say that the class $\mathcal{H}$ is parametrized by $d$ **real** numbers. Because we are constrained by computers that **use 64 bits** to represent floating-point numbers, we have **at most** $k = 2^{64d}$ **different hypotheses**.

To **hold** the **theorem** $\mathcal{E}(\widehat{h}) \leq \widehat{\mathcal{E}}(h^*) + 2\gamma$ as **valid** with at least **probability** $1 - \delta$ , we need to suffice the corollary:
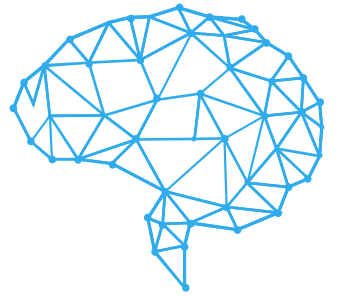
$$m = O\left(\frac{1}{\gamma^2}\log\left(\frac{k}{\delta}\right)\right)$$

$$m = O\left(\frac{1}{\gamma^2}\log\left(\frac{2^{64d}}{\delta}\right)\right) = O\left(\frac{d}{\gamma^2}\log\left(\frac{1}{\delta}\right)\right)$$

Therefore, the **number** of **training examples needed** is **at most linear** in the **parameters** of the **model** $d$.

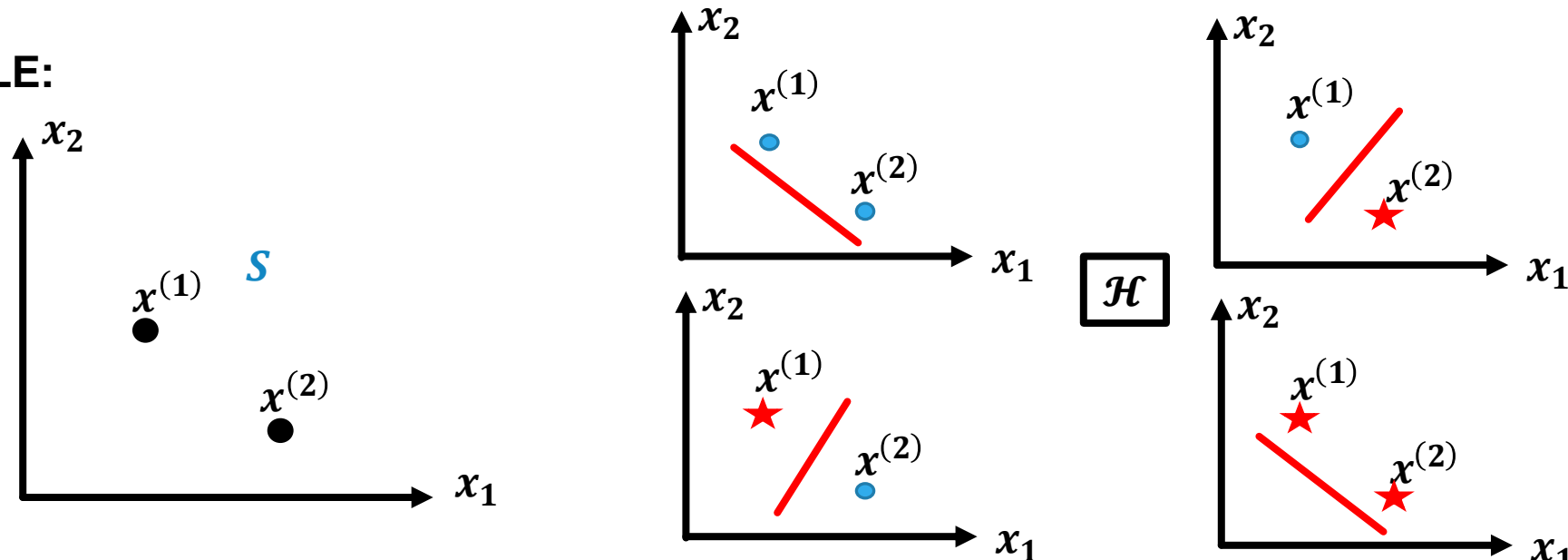We will **introduce** a **definition** to **make** the **proof** for **infinite classes** $\mathcal{H}$.

**DEFINITION:** given a set $S = \{x^{(1)}, ..., x^{(D)}\}$ of **points** $x^{(i)} \in \chi$, we say that $\mathcal{H}$ **SHATTERS** $S$ if $\mathcal{H}$ can realize any labeling on $S$. That is, if **for any set** of **labels** $\{y^{(1)}..., y^{(D)}\}$, there **exists some** $h \in \mathcal{H}$ so that $h(x^{(i)}) = y^{(i)}$ **for all** $i = 1,...D$.
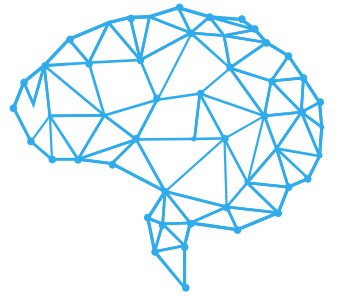
**EXAMPLE:**



SHATTER: for any possible labeling of these points, we can find a linear classifier that obtains "zero training error" on them.

**DEFINITION:**

The **Vapnik-Chervonenkis dimension** of $\mathcal{H}$, $\left(VC(\mathcal{H})\right)$ is the **size** of the **largest set** shattered by $\mathcal{H}$.

**EXAMPLE:**

If $\mathcal{H} = \{\textbf{linear classifiers in 2D}\}$, therefore $VC(\mathcal{H}) = 3.$ There is **no set of size 4** that it could **shatter**.

In a **general form** we have that if $\mathcal{H} = \{\textbf{linear classifiers in n Dimensions}\}$, therefore $VC(\mathcal{H}) = n + 1$.

**THEOREM:**

Let $\mathcal{H}$ be given and let $\boldsymbol{D = VC(\mathcal{H})}$. Then with **probability** at **least** $\boldsymbol{1 - \delta}$, we have that for all $\boldsymbol{h} \in \mathcal{H}$,

$$\left|\mathcal{E}(\boldsymbol{h}) - \widehat{\mathcal{E}}(\boldsymbol{h})\right| \leq \boldsymbol{O}\left(\sqrt{\frac{\boldsymbol{D}}{\boldsymbol{m}}\log\left(\frac{\boldsymbol{m}}{\boldsymbol{D}}\right) + \frac{\boldsymbol{1}}{\boldsymbol{m}}\log\left(\frac{1}{\delta}\right)}\right)$$

With probability at least $1 - \delta$, we also have that:

$$\widehat{\mathcal{E}}(\boldsymbol{h}) \leq \mathcal{E}(\boldsymbol{h}^*) + \boldsymbol{O}\left(\sqrt{\frac{\boldsymbol{D}}{\boldsymbol{m}}\log\left(\frac{\boldsymbol{m}}{\boldsymbol{D}}\right) + \frac{\boldsymbol{1}}{\boldsymbol{m}}\log\left(\frac{1}{\delta}\right)}\right)$$

**If a hypothesis class has finite VC dimension, then uniform convergence occurs as $m$ becomes large.**

**COROLLARY:**

For $|\mathcal{E}(h) - \widehat{\mathcal{E}}(h) \leq)| \leq \gamma$ to hold for all $h \in \mathcal{H}$ (and hence $\mathcal{E}(\widehat{h}) \leq \mathcal{E}(h^*) + 2\gamma$) with probability at least $1 - \delta$, it suffices that $n = O_{\gamma,\delta}(D)$.

Thus, **sample complexity** is **upper-bounded** by the **VC dimension**. Also, for "**most**" **hypothesis** classes, the **VC dimension** is also **roughly linear** in the **number** of **parameters**.

We **conclude** that for a **given hypothesis class** $\mathcal{H}$, the **number** of **training examples needed** to **achieve generalization error close** to that of the **optimal classifier** is usually **roughly linear** in the **number** of **parameters** of $\mathcal{H}$.