

Lecture 12: Computational Learning Theory & Kernel Winter 2018

Kai-Wei Chang
CS @ UCLA

kw+cm146@kwchang.net

The instructor gratefully acknowledges Dan Roth, Vivek Srikumar, Sriram Sankararaman, Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Recap: Computational Learning Theory

- ❖ The Theory of Generalization

 - ❖ Using training instance to rule out incorrect hypotheses

- ❖ Probably Approximately Correct (PAC) learning

 - ❖ How many examples you need to see to obtain a learned function with error $\leq \epsilon$ with high probability

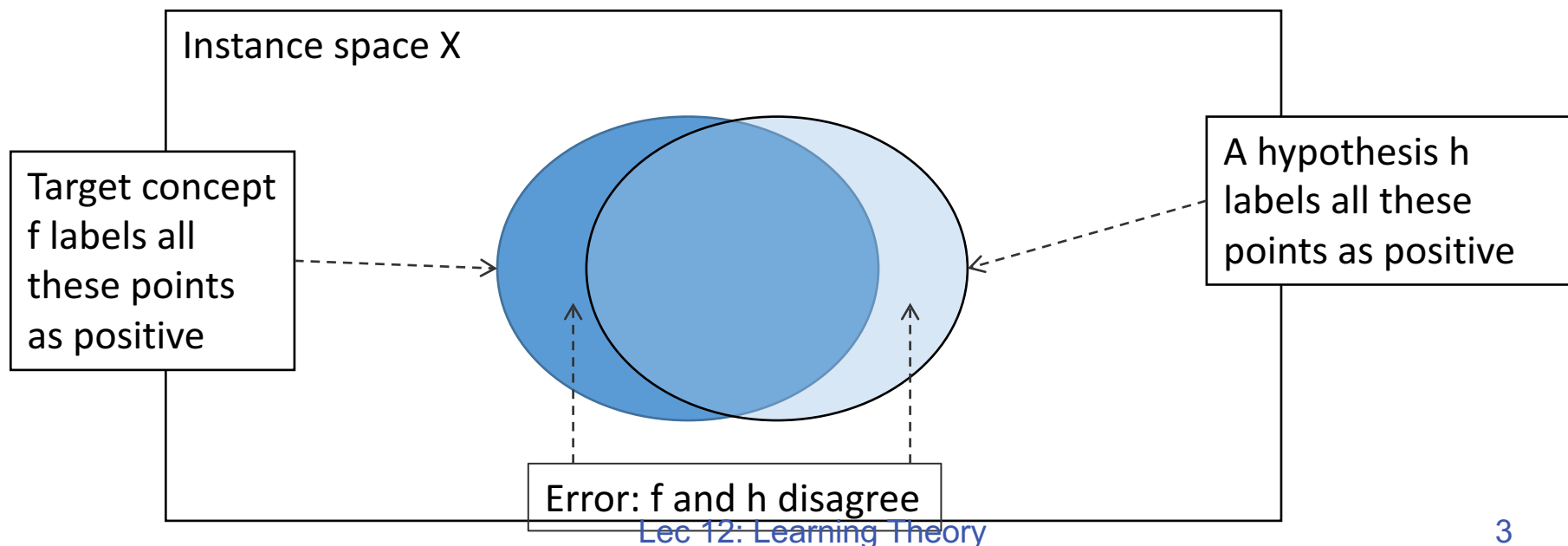
- ❖ Shattering and the VC dimension

Recap: Error of a hypothesis

Definition

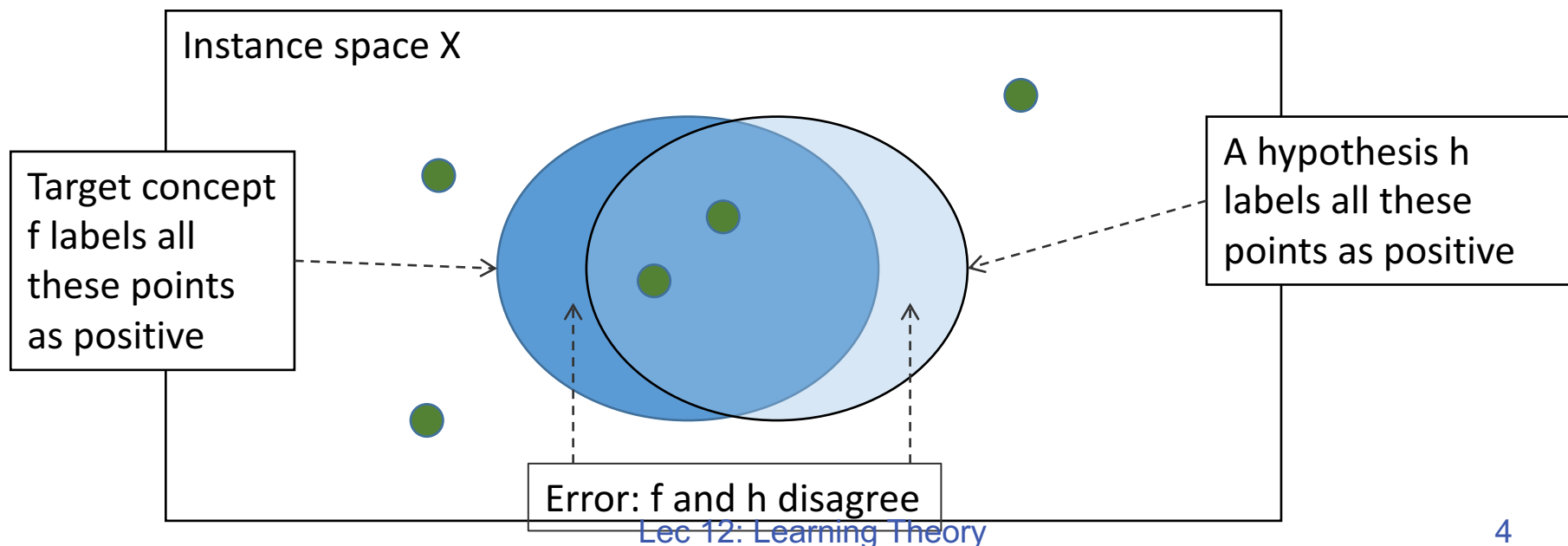
Given a distribution D over examples, the *error* of a hypothesis h with respect to a target concept f is

$$\text{err}_D(h) = \Pr_{x \sim D}[h(x) \neq f(x)]$$



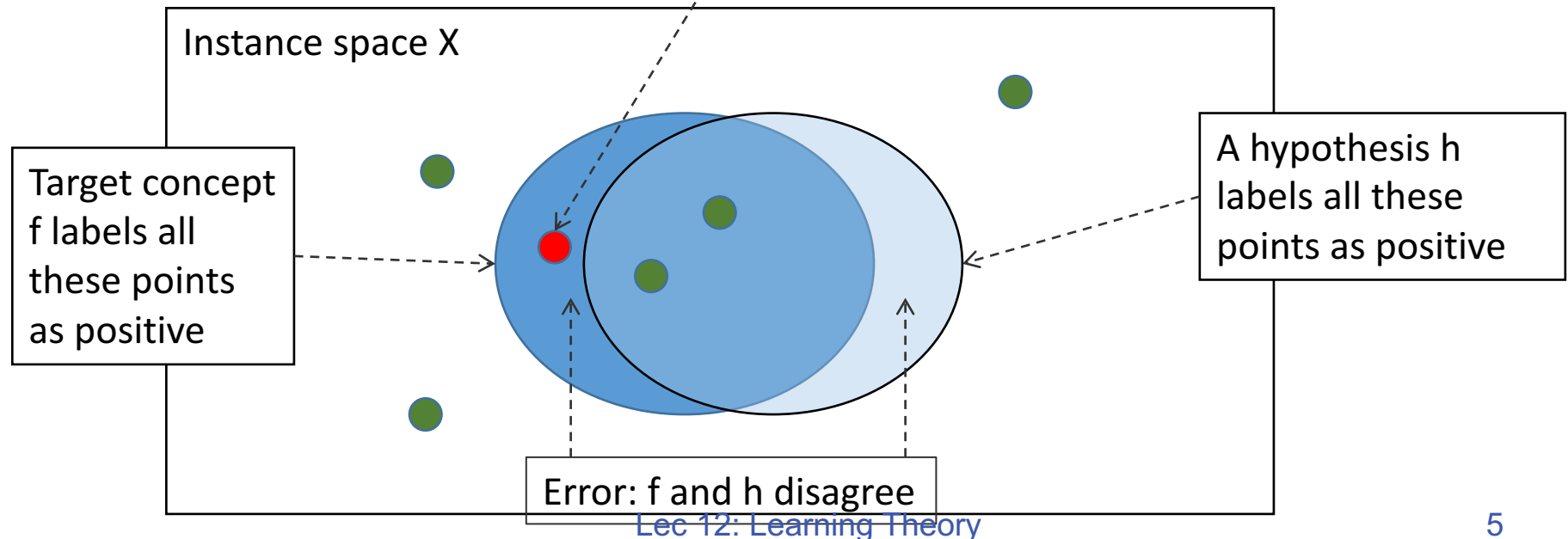
Recap: Error of a hypothesis

Overfitting: You may have a learned model that is consistent with the training data but still makes mistakes.



Recap: Error of a hypothesis

With the IID sampling assumption, we either have seen this example in the training phase, or it is unlikely to see it in the test time.



Requirements of Learning

- ❖ Cannot expect a learner to learn a concept *exactly*
- ❖ Instead, we “agree” to misclassify *uncommon* examples that do not show up in the training set

PAC Learnability

Turing Award: [Leslie Valiant](#).

Consider a concept class C defined over an instance space X (containing instances of length n), and a learner L using a hypothesis space H

The concept class C is **PAC learnable** by L using H if for all $f \in C$, for all distribution D over X , and fixed $\epsilon > 0$, $\delta < 1$, given m examples sampled i.i.d. according to D , the algorithm L produces, with probability at least $(1 - \delta)$, a hypothesis $h \in H$ that has error at most ϵ , where m is *polynomial* in $1/\epsilon$, $1/\delta$, n and $\text{size}(H)$

Intuition of PAC Learnability

With the IID sampling assumption, if a concept is reasonable. After, we saw enough samples, it is unlikely to have many these red points

Instance space X

Target concept f labels all these points as positive

A hypothesis h labels all these points as positive

Error: f and h disagree

$$f = x_2 \wedge x_3 \wedge x_4 \wedge x_5 \wedge x_{100}$$

Recap: Learning Conjunctions

❖ Protocol 1:

Teacher provides a set of example $(x, f(x))$

❖ $\langle (1, 1, 1, 1, 1, 1, \dots, 1, 1), 1 \rangle$

❖ $\langle (1, 1, 1, 0, 0, 0, \dots, 0, 0), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 1, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 1, 1, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 0, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 0, 0, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 1, \dots, 0, 1), 1 \rangle$

❖ $\langle (0, 1, 0, 1, 0, 0, \dots, 0, 1, 1), 0 \rangle$

What would f look like?

Whenever the output is 1, x_1 is present

With the given data, we only learned an *approximation* to the true concept.
Is it good enough?

Recap: Learning Conjunctions: Analysis

Theorem: Suppose we are learning a conjunctive concept with n dimensional Boolean features using m training examples. If

$$m > \frac{n}{\epsilon} \left(\log(n) + \log \left(\frac{1}{\delta} \right) \right)$$

then, with probability $> 1 - \delta$, the error of the learned hypothesis $\text{err}_D(h)$ will be less than ϵ .

A general result

Let H be any hypothesis space.

With probability $1 - \delta$ a hypothesis $h \rightarrow H$ that is **consistent** with a training set of size m will have an error $\leq \epsilon$ on future examples if

$$m > \frac{1}{\epsilon} \left(\ln(|H|) + \ln \frac{1}{\delta} \right)$$

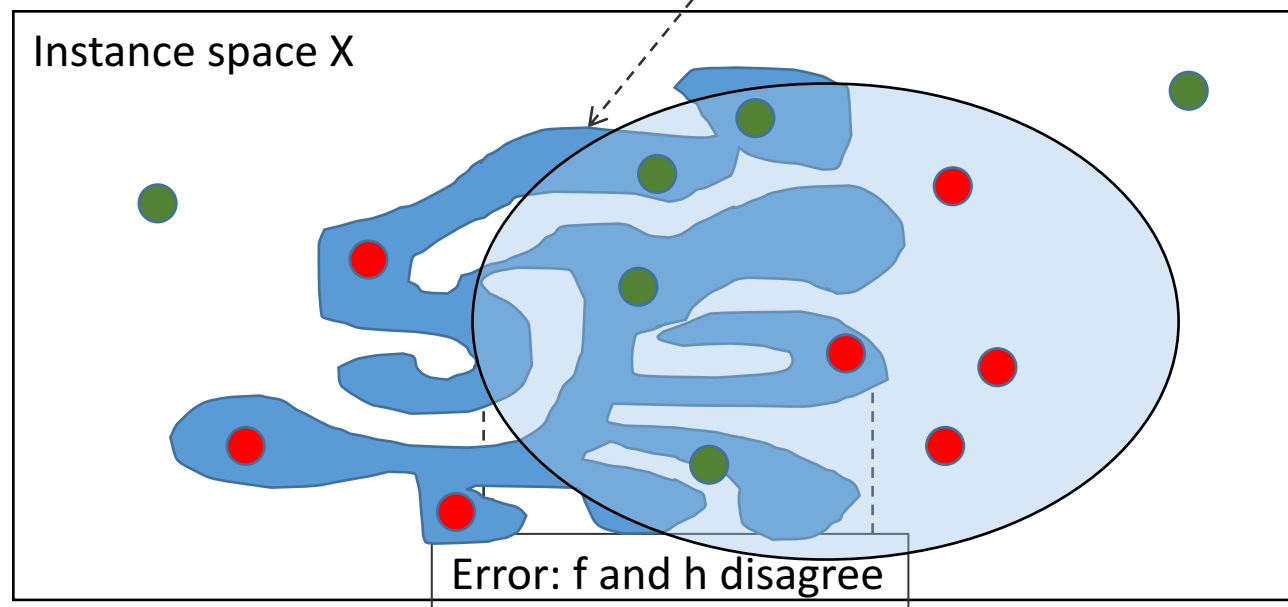
1. Expecting lower error increases sample complexity (i.e more examples needed for the guarantee)

2. If we have a larger hypothesis space, then we will make learning harder (i.e higher sample complexity)

3. If we want a higher confidence in the classifier we will produce, sample complexity will be higher.

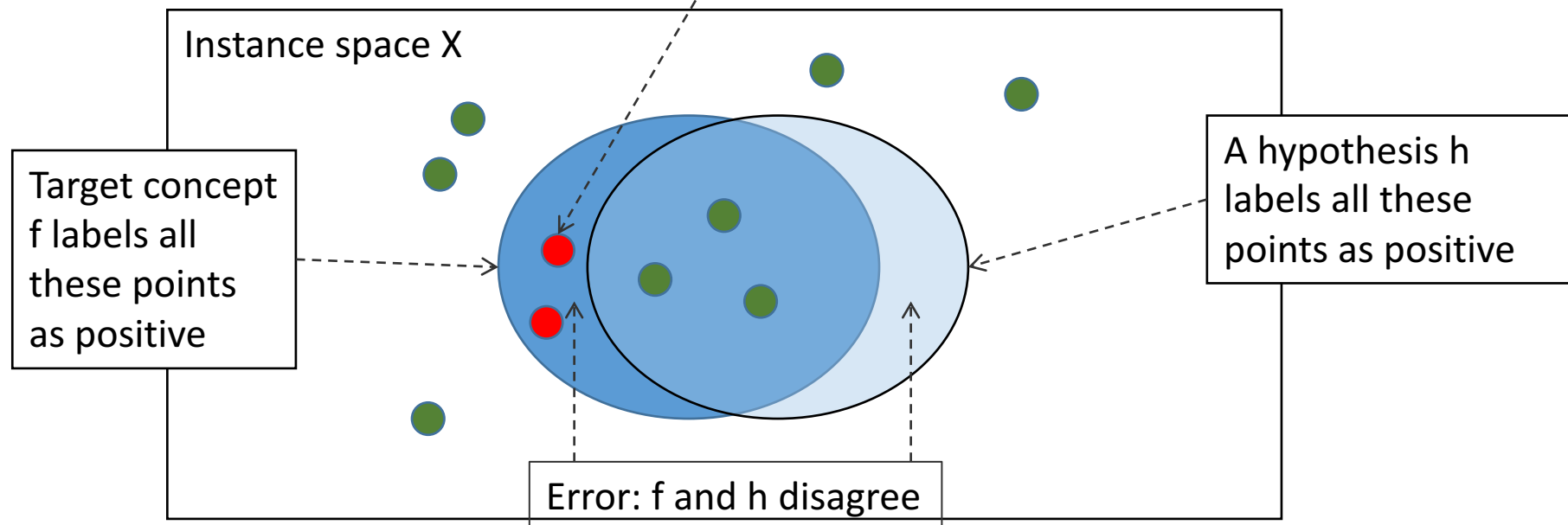
Intuition of PAC Learnability

With the IID sampling assumption, if a concept is too complicated. We need to see exponential number of samples, such that we can rule out those red points

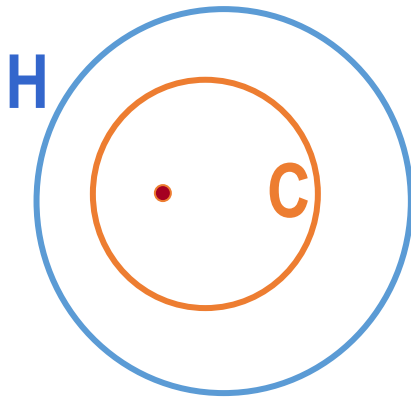


Intuition of PAC Learnability

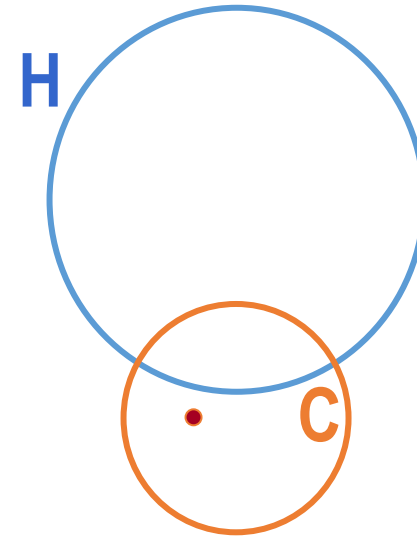
If a concept is simple:



What if the concept space is different from the hypothesis space?
[not in exam]



It is fine, we can still find the right function



The training error will not be zero

Generalization bound [not in exam]

A bound on how much the true error will deviate from the training error. If we have more than m examples, then with high probability $1 - \delta$

$$\underset{\substack{\text{Generalization error} \nearrow}}{err_D(h)} - \underset{\substack{\nearrow \text{Training error}}}{err_S(h)} \leq \sqrt{\frac{\ln |H| + \ln(1/\delta)}{2m}}$$

Generalization bound

A bound on how much the true error will deviate from the

Now, we know if $\text{size}(H)$ is finite, we can define what is learnable. This works for Boolean functions.

Next question: What if $\text{size}(H)$ is infinity?

This lecture: Computational Learning Theory

- ❖ The Theory of Generalization
- ❖ Probably Approximately Correct (PAC) learning
- ❖ Shattering and the VC dimension

Infinite Hypothesis Space

- ❖ The previous analysis was restricted to finite hypothesis spaces
- ❖ Some infinite hypothesis spaces are more expressive than others
 - ❖ Linear threshold function vs. a combination of LTUs
- ❖ Need a measure of the expressiveness of an infinite hypothesis space other than its size

A general result

If $|H|$ is infinite, m is always infinite as well.

$$m > \frac{1}{\epsilon} \left(\ln(|H|) + \ln \frac{1}{\delta} \right)$$

1. Expecting lower error increases sample complexity (i.e more examples needed for the guarantee)

2. If we have a larger hypothesis space, then we will make learning harder (i.e higher sample complexity)

3. If we want a higher confidence in the classifier we will produce, sample complexity will be higher.

Vapnik-Chervonenkis dimension

- ❖ The Vapnik-Chervonenkis dimension (**VC dimension**) provides such a measure
- ❖ “What is the expressive *capacity* of a set of functions?”
- ❖ Analogous to $|H|$, there are bounds for sample complexity using $VC(H)$

VC dimension and consistent learners [not in exam]

- ❖ Using $VC(H)$ as a measure of expressiveness we have a sample complexity bound for infinite hypothesis spaces
- ❖ Given a sample D with m examples, find some $h \rightarrow H$ is consistent with all m examples. If

$$m > \frac{1}{\epsilon} \left(8VC(H) \log \frac{13}{\epsilon} + 4 \log \frac{2}{\delta} \right)$$

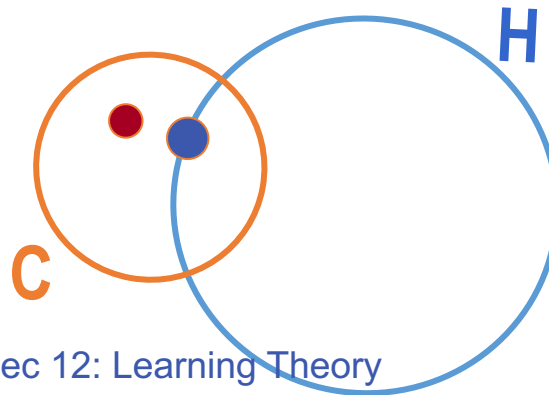
Then with probability at least $(1-\delta)$, h has error less than ϵ .

You don't need to remember this equation
but just need to understand the meaning

Generation bound for agnostic learner [not in exam]

If we have m examples, then with probability $1 - \delta$, the true error of a hypothesis h with training error $\text{err}_S(h)$ is bounded by

$$\text{err}_D(h) \leq \text{err}_S(h) + \sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$



Intuition of VC dimension

- ❖ Although there are infinitely many hypotheses, many of them are similar
- ❖ The idea of learning is by eliminating incorrect hypotheses
 - ❖ We can eliminate infinite # hypotheses for each training sample

Recap: Learning Conjunctions

❖ Protocol 1:

Teacher provides a set of example $(x, f(x))$

❖ $\langle (1, 1, 1, 1, 1, 1, \dots, 1, 1), 1 \rangle$

❖ $\langle (1, 1, 1, 0, 0, 0, \dots, 0, 0), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 1, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 1, 1, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 0, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 0, 0, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 1, \dots, 0, 1), 1 \rangle$

❖ $\langle (0, 1, 0, 1, 0, 0, \dots, 0, 1, 1), 0 \rangle$

$$x_1 \wedge x_2 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_1 \wedge x_2 \wedge x_3 \dots \dots x_{99}$$

$$x_1 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_2 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_2 \wedge x_3 \dots \dots x_{99}$$

$$x_3 \dots \dots x_{99} \wedge x_{100}$$

Recap: Learning Conjunctions

❖ Protocol 1:

Teacher provides a set of example $(x, f(x))$

❖ $\langle (0, 1, 1, 1, 1, 1, \dots, 1, 1), 1 \rangle$

❖ $\langle (1, 1, 1, 0, 0, 0, \dots, 0, 0), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 1, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 1, 1, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 0, \dots, 0, 0, 1), 1 \rangle$

❖ $\langle (1, 0, 1, 0, 0, 0, \dots, 0, 1, 1), 0 \rangle$

❖ $\langle (1, 1, 1, 1, 1, 1, \dots, 0, 1), 1 \rangle$

❖ $\langle (0, 1, 0, 1, 0, 0, \dots, 0, 1, 1), 0 \rangle$

❖ One instance can eliminate many hypothesis


$$x_1 \wedge x_2 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_1 \wedge x_2 \wedge x_3 \dots \dots x_{99}$$

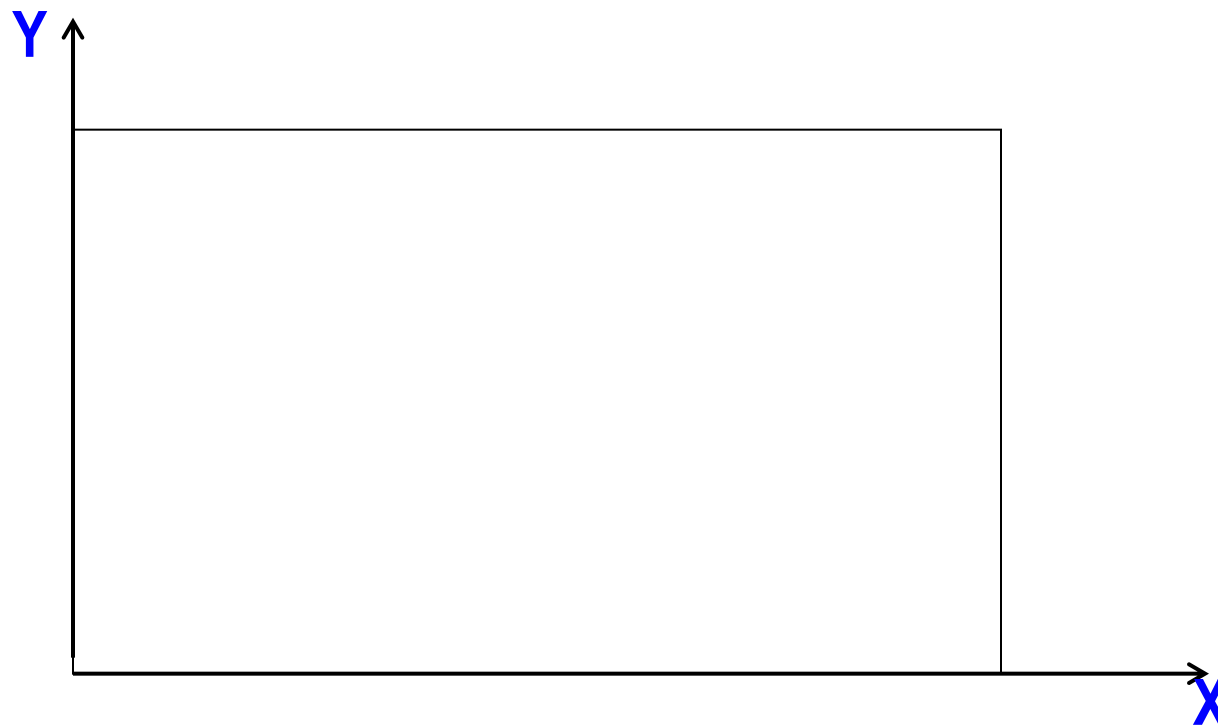
$$x_1 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_2 \wedge x_3 \dots \dots x_{99} \wedge x_{100}$$

$$x_2 \wedge x_3 \dots \dots x_{99}$$

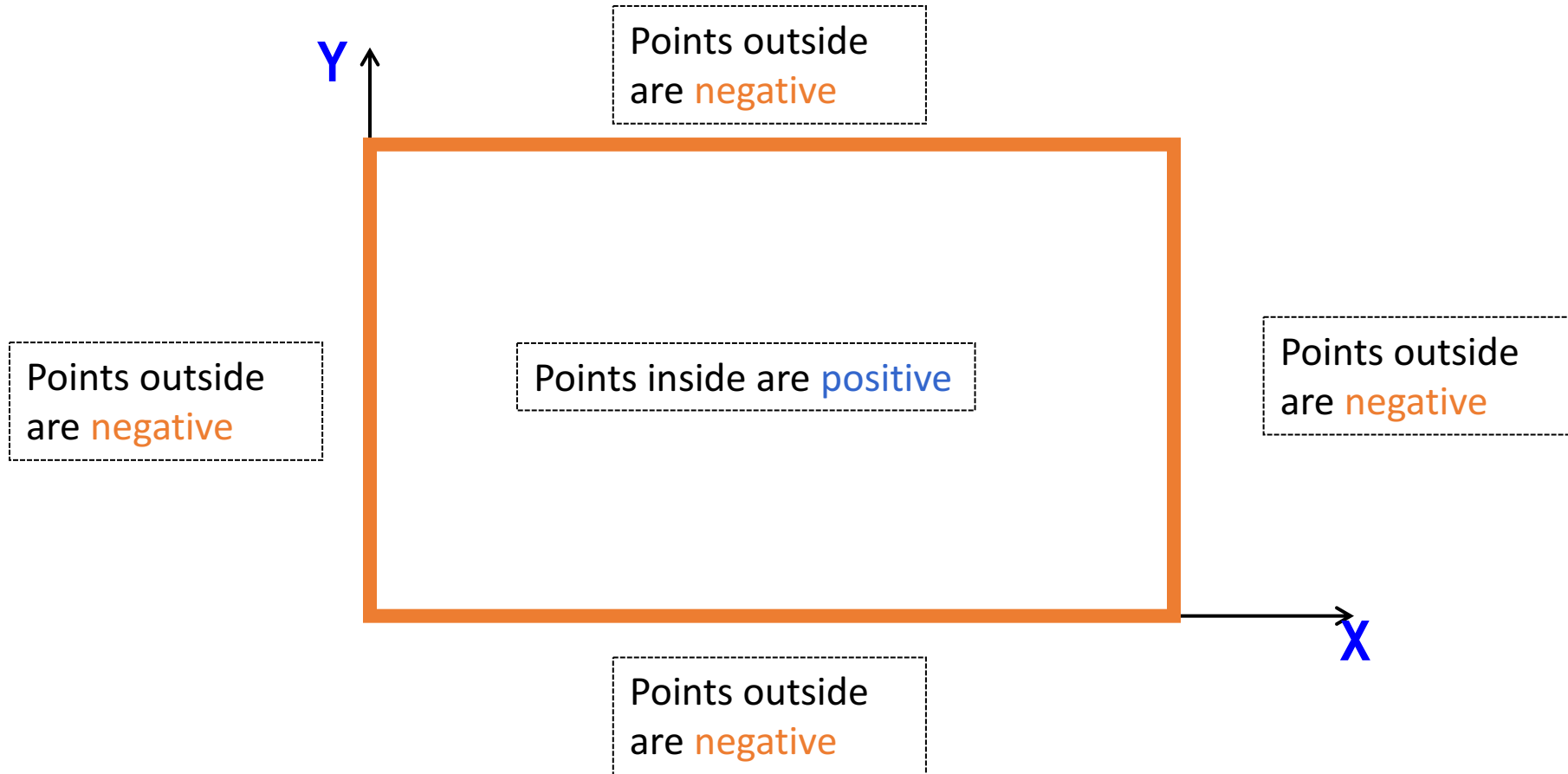
$$x_3 \dots \dots x_{99} \wedge x_{100}$$

Intuition of VC dimension: Learning Rectangles



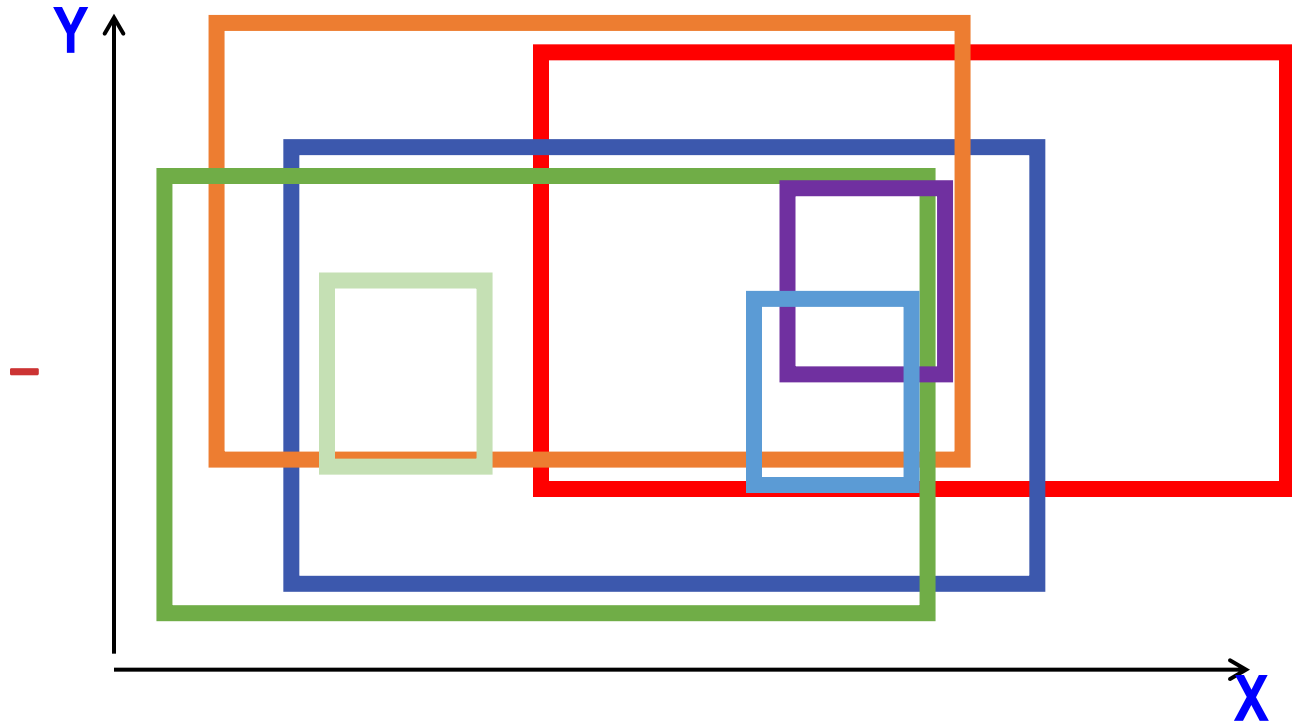
Assume the target concept is an axis parallel rectangle

Learning Rectangles



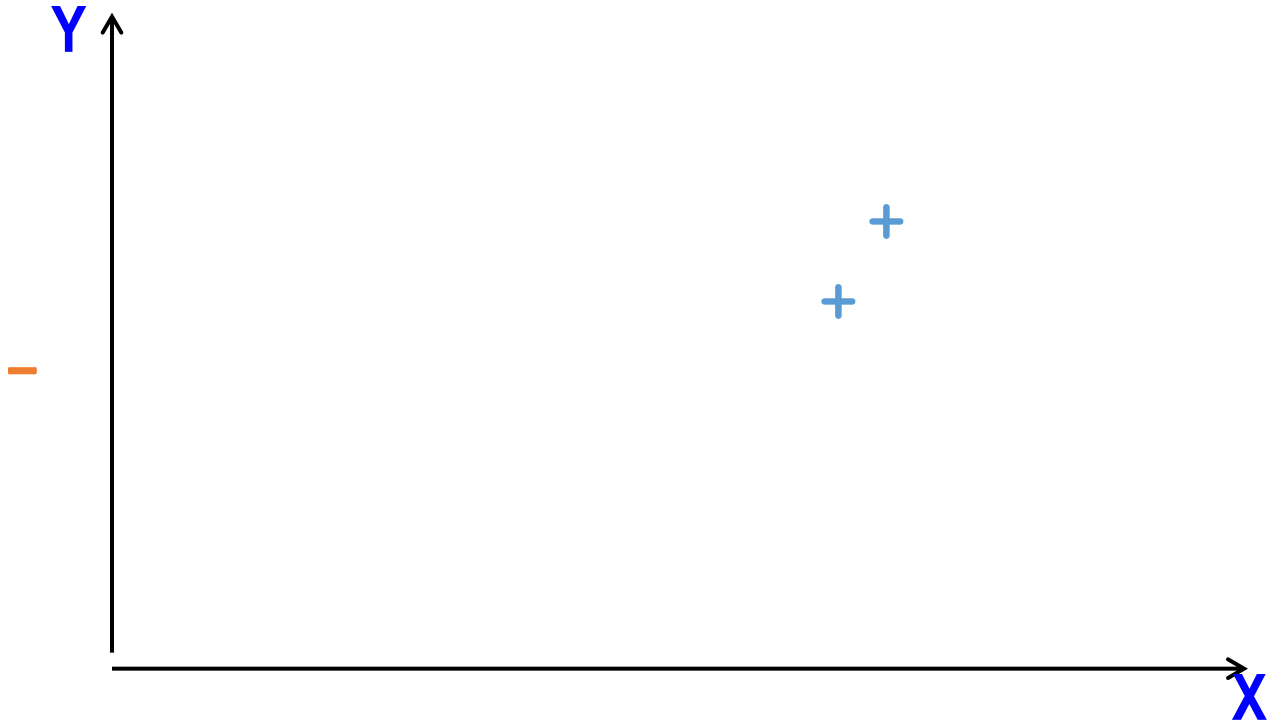
Learning Rectangles

Assume the target concept is an axis parallel rectangle



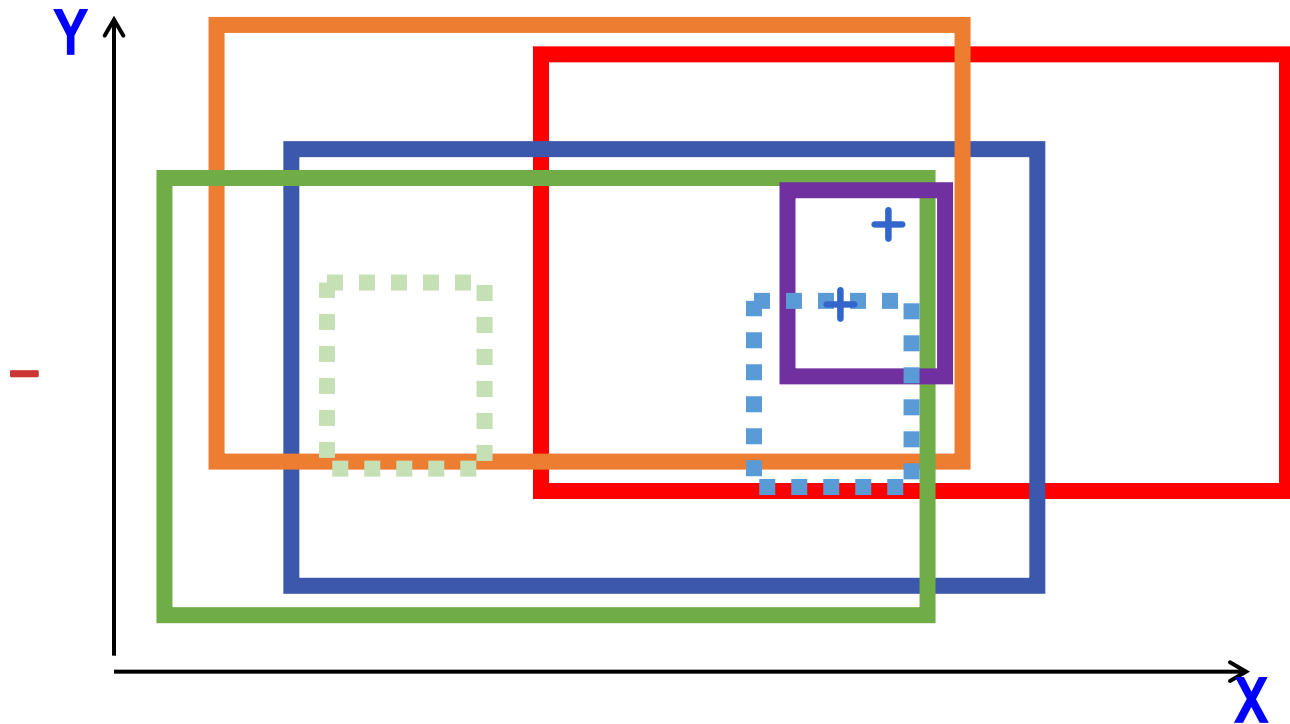
Learning Rectangles

Assume the target concept is an axis parallel rectangle



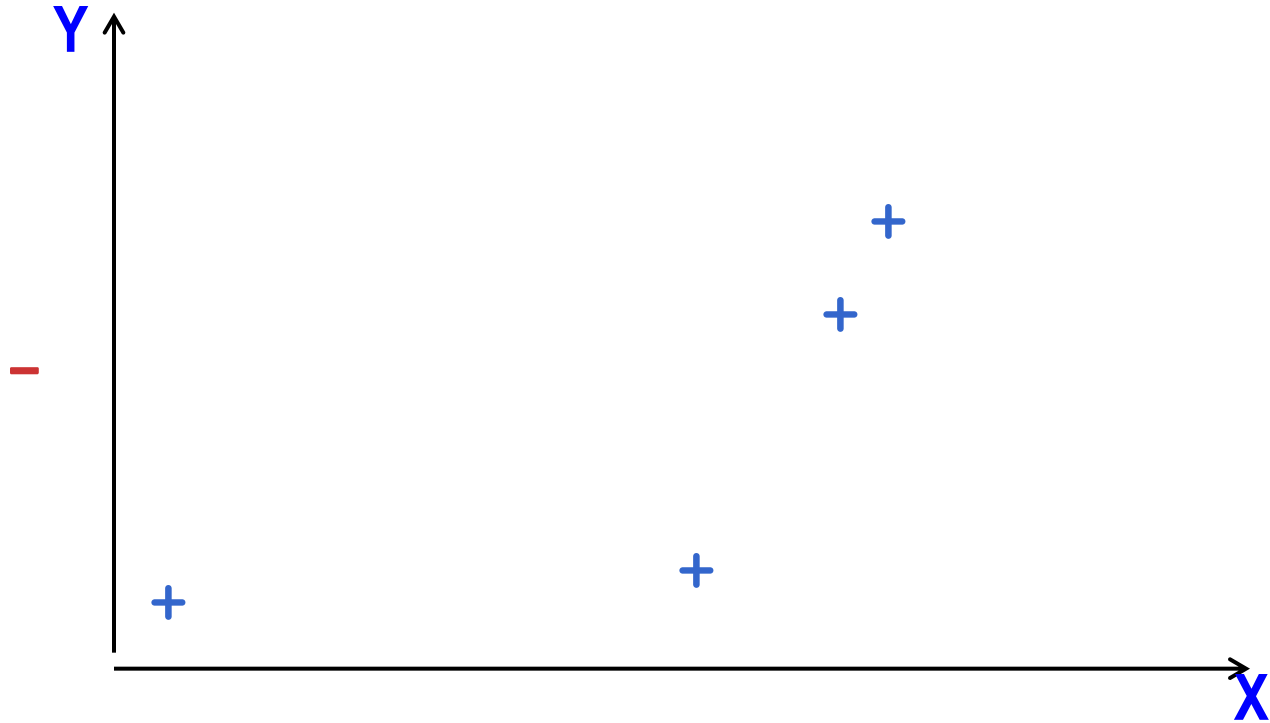
Learning Rectangles

Assume the target concept is an axis parallel rectangle



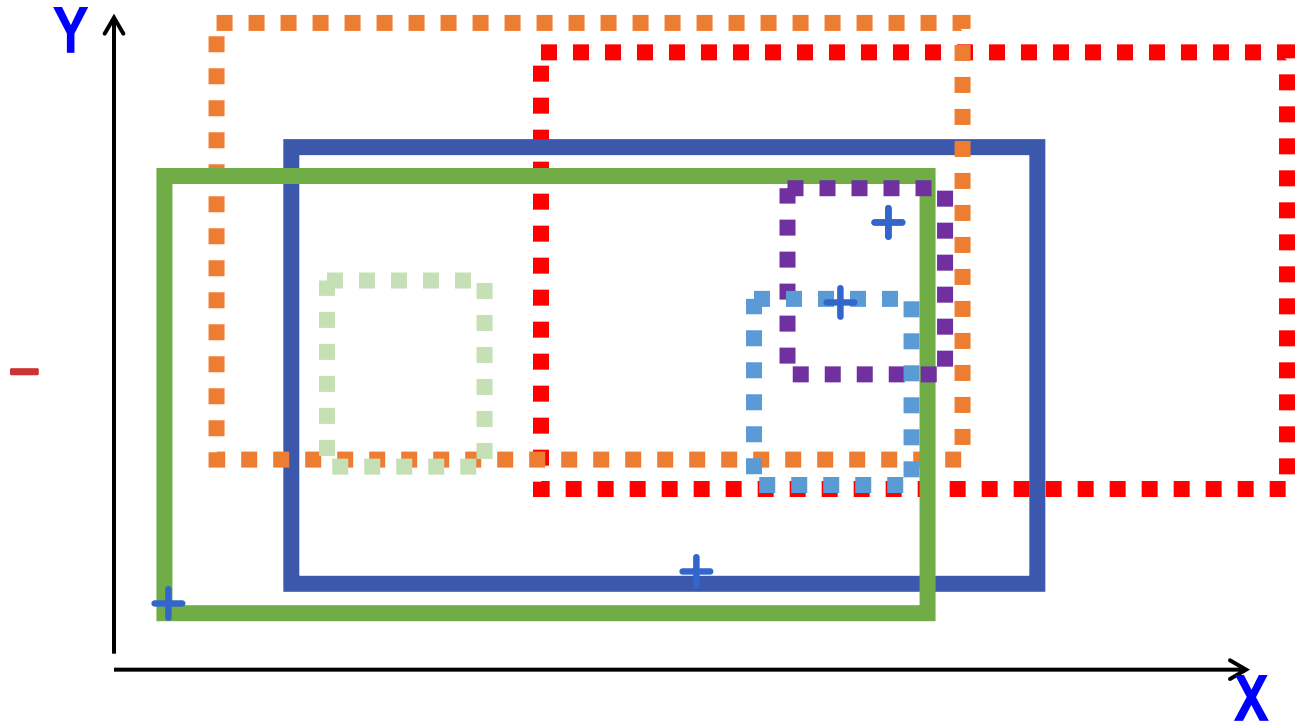
Learning Rectangles

Assume the target concept is an axis parallel rectangle



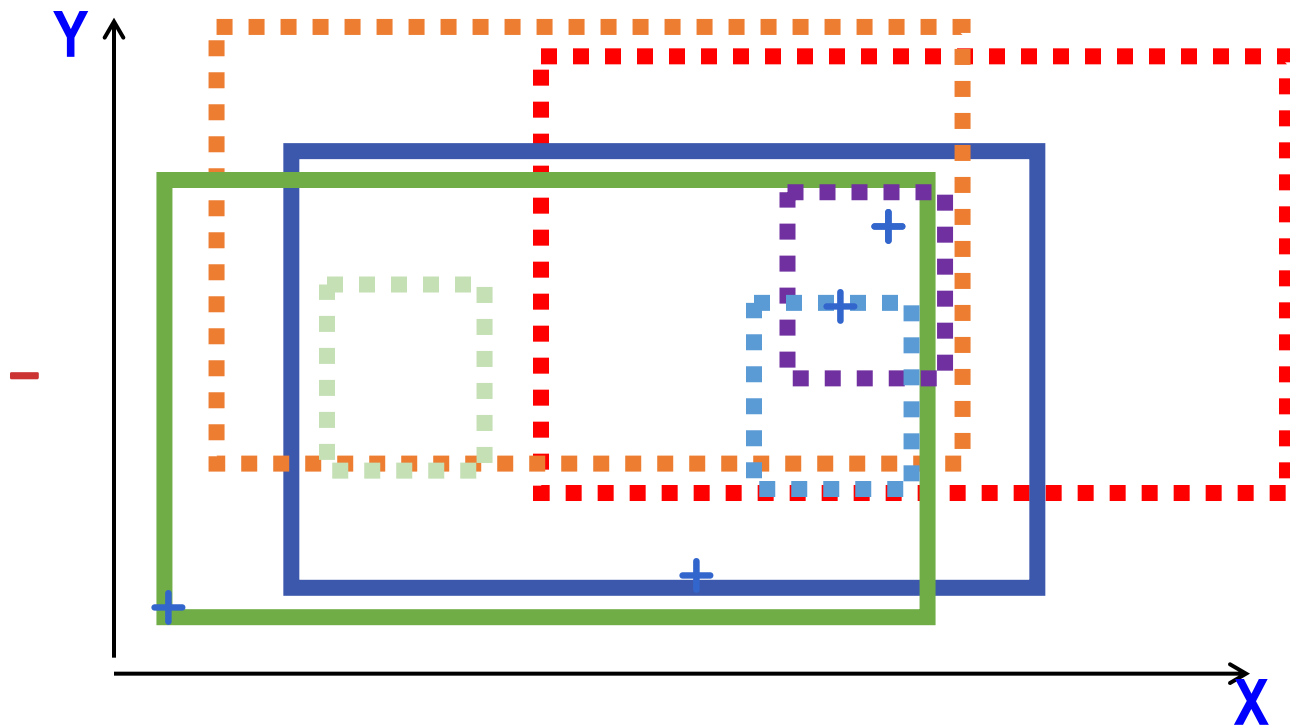
Learning Rectangles

Assume the target concept is an axis parallel rectangle



Learning Rectangles

Assume the target concept is an axis parallel rectangle



Key observation: Despite there are infinite # hypothesis
The blue & red rectangles have the same predictions

Let's think about expressivity of functions

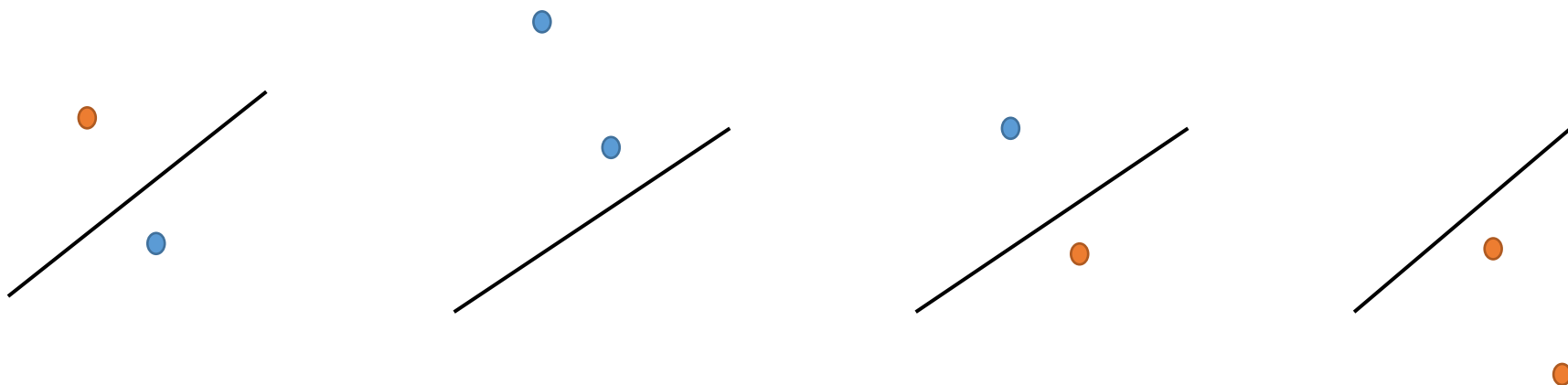


Suppose we have two points.

Can linear classifiers correctly classify any labeling of these points?

Linear functions are expressive enough to *shatter* 2 points

Let's think about expressivity of functions

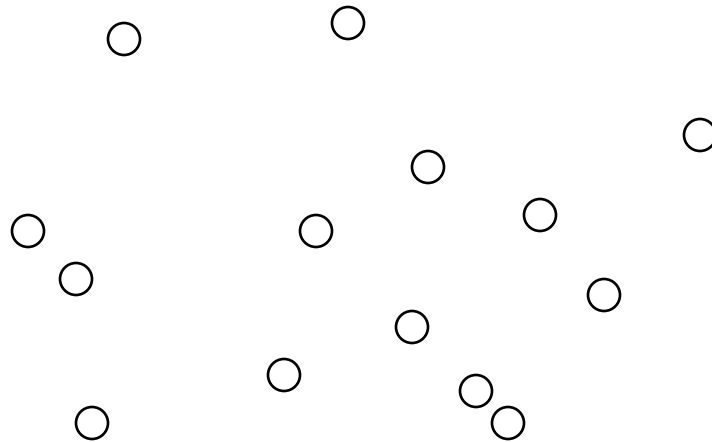


Suppose we have two points.

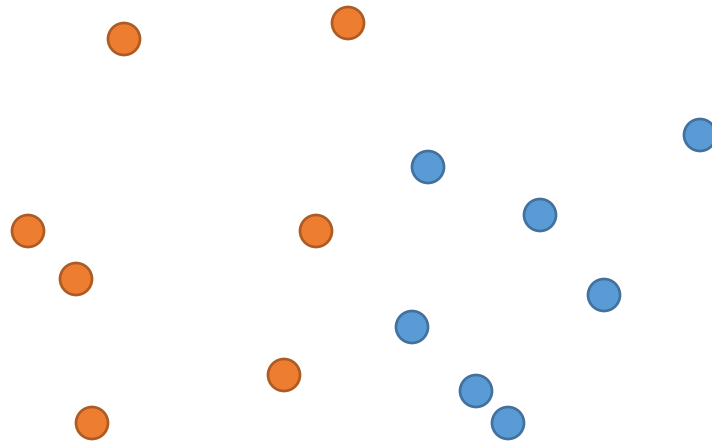
Can linear classifiers correctly classify any labeling of these points?

Linear functions are expressive enough to *shatter* 2 points

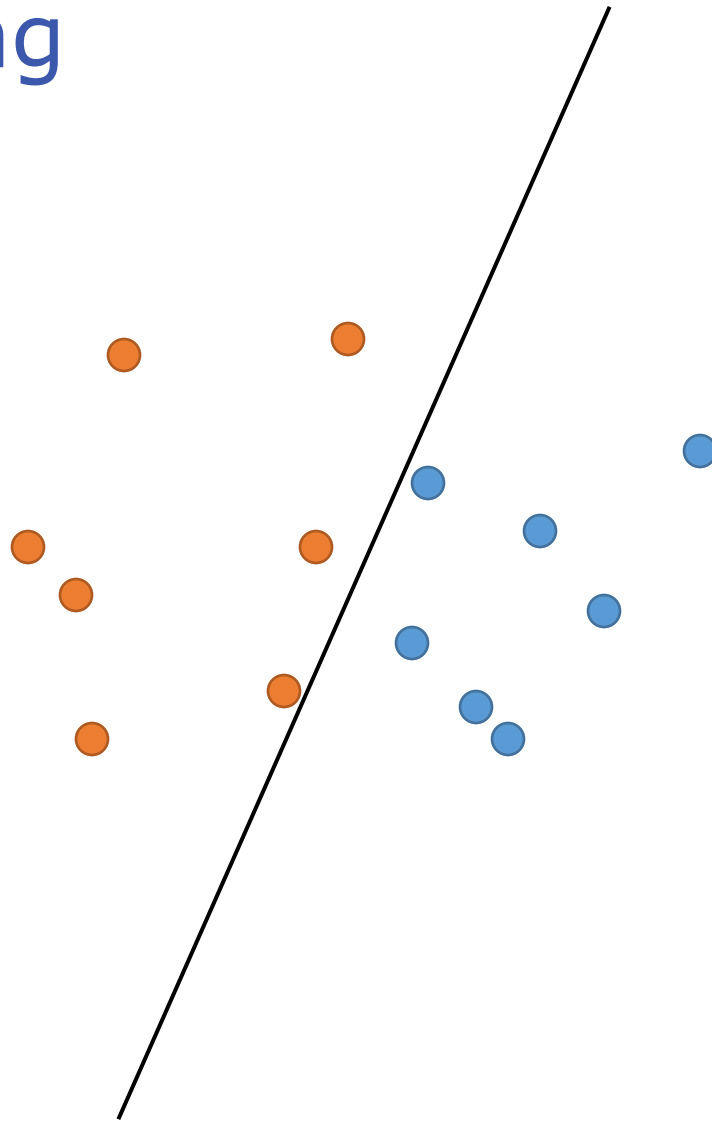
Shattering



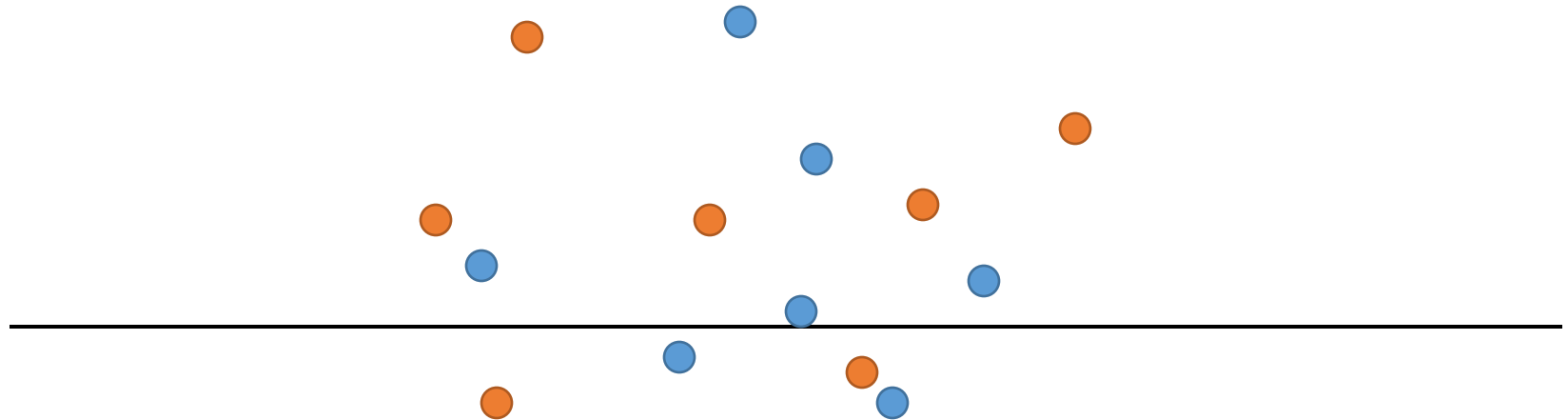
Shattering



Shattering

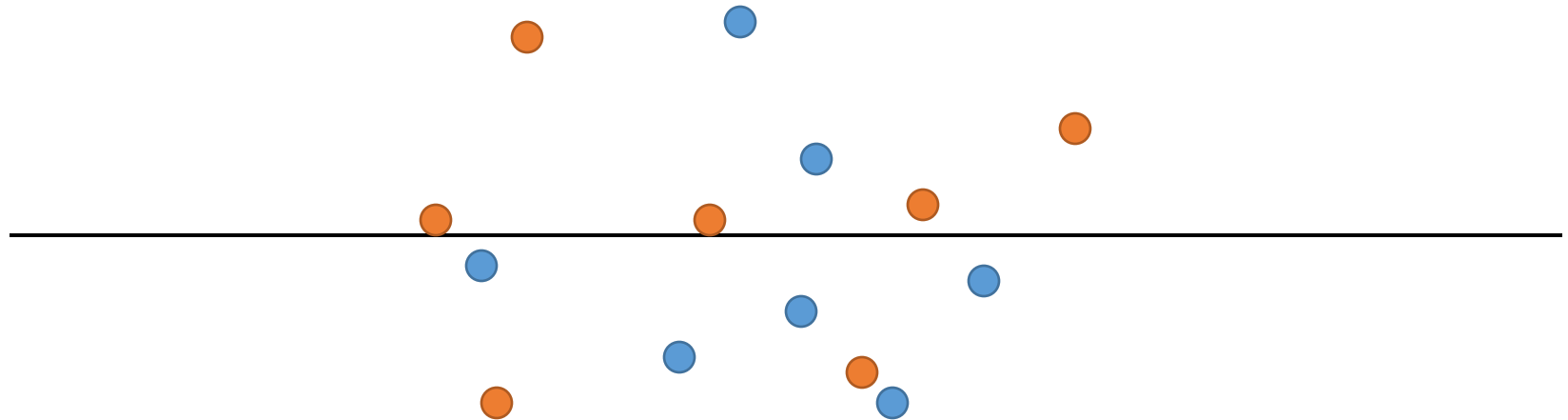


Shattering



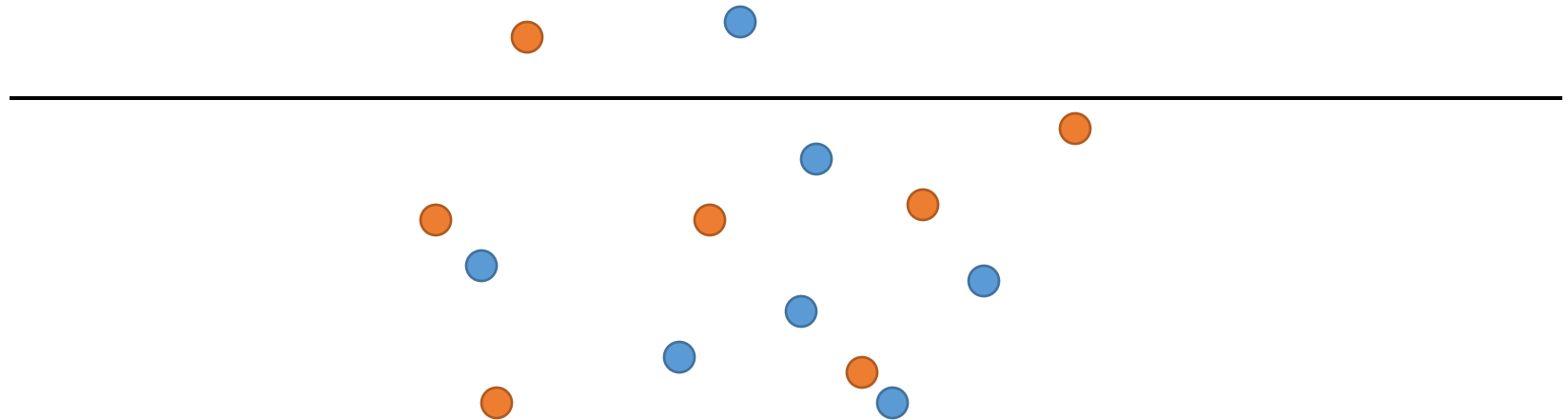
This particular labeling of the points can not be separated by *any* line

Shattering



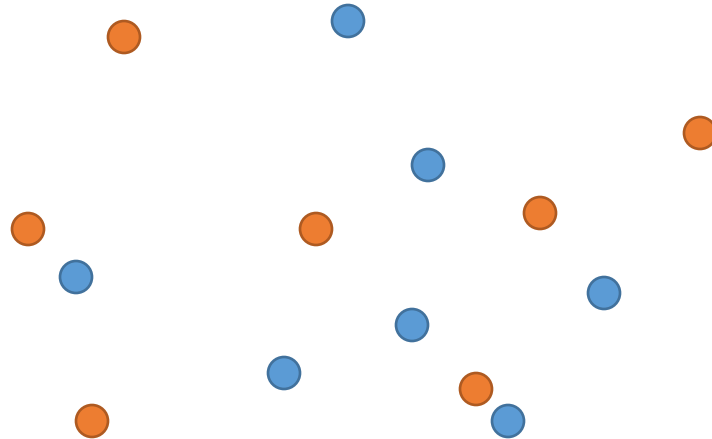
This particular labeling of the points can not be separated by *any* line

Shattering



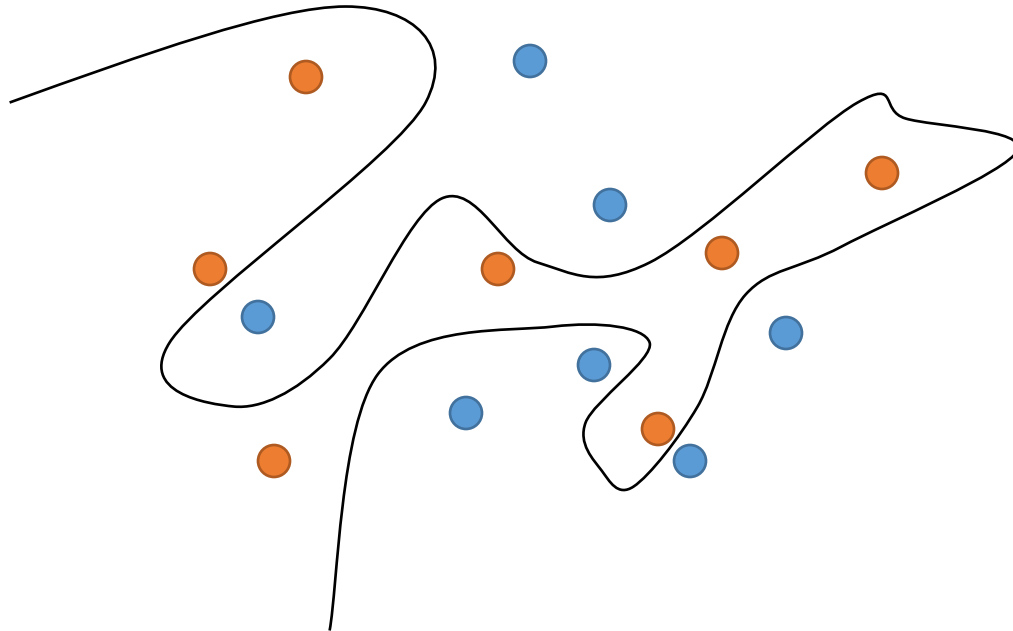
This particular labeling of the points can not be separated by *any* line

Shattering



This particular labeling of the points can not be separated by *any* line

Shattering



Linear functions are not expressive to shatter fourteen points
Because there is a labeling that can not be separated by them
Of course, a more complex function could separate them

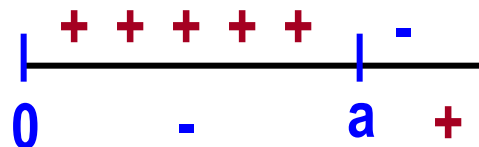
Shattering

Definition: A set S of examples is **shattered** by a set of functions H if **for every** partition of the examples in S into positive and negative examples **there is** a function in H that gives exactly these labels to the examples

Intuition: A rich set of functions shatters large sets of points

Left bounded intervals

Example 1: Hypothesis class of left bounded intervals on the real axis: $[0, a)$ for some real number $a > 0$

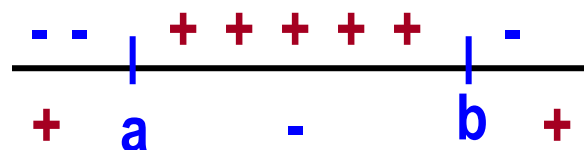


Sets of **two** points **cannot** be shattered

That is: given two points, you can label them in such a way that no concept in this class will be consistent with their labeling

Real intervals

Example 2: Hypothesis class is the set of intervals on the real axis: $[a,b]$, for some real numbers $b > a$



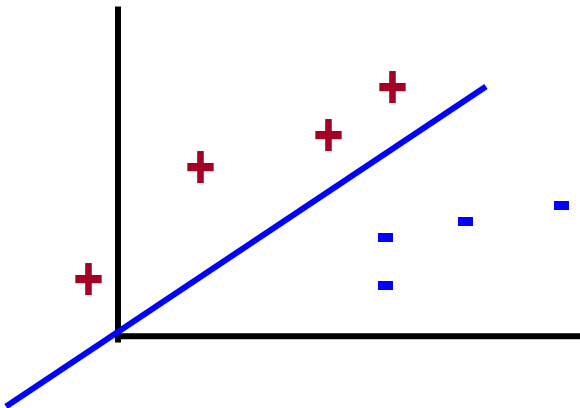
All sets of one or two points can be shattered

But some sets of **three** points **cannot** be shattered

Shattering

Definition: A set S of examples is **shattered** by a set of functions H if for every partition of the examples in S into positive and negative examples there is a function in H that gives exactly these labels to the examples

Example 3: 2-D Half spaces in a plane



Can one point be shattered?

Is there any two points can be shattered?

Is there any three points?

Can any three points be shattered?

Vapnik-Chervonenkis Dimension

Definition: The **VC dimension** of hypothesis space H over instance space X is the size of the largest finite subset of X that is shattered by H

- ❖ If there **exists** any subset of size d that can be shattered, $VC(H) \geq d$
 - ❖ Even one subset will do
- ❖ If **no subset** of size d can be shattered, then $VC(H) < d$

Shattering: The adversarial game

You



You: Hypothesis class H can shatter these d points I provide

You: Aha! There is a function $h \in H$ that correctly predicts your evil labeling

An adversary



Adversary: That's what you think! Here is a labeling that will defeat you.

Adversary: Argh! You win this round. But I'll be back.....

Example Half spaces in a plane

- ❖ Prove $VC \geq 1$

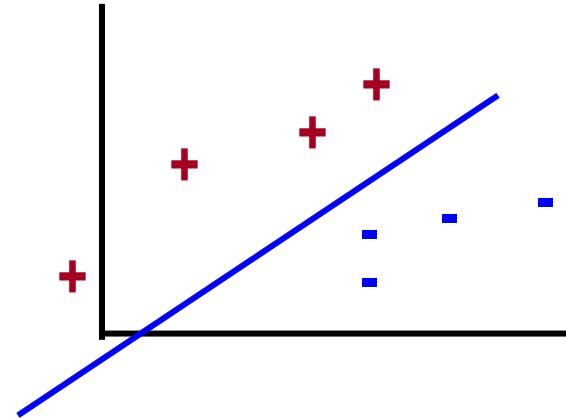
- ❖ Show any point can be shattered

- ❖ Prove $VC \geq 2$

- ❖ Show there exists 2 points can be shattered

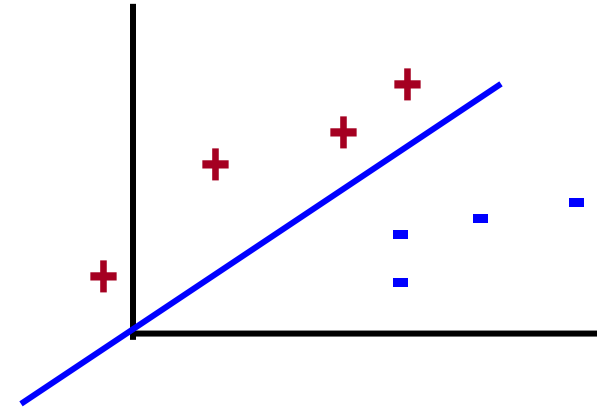
- ❖ Prove $VC \geq 3$

- ❖ Show there exists 3 points can be shattered



Example Half spaces in a plane

- ❖ Prove $VC < 4$
 - ❖ Show **no** 4 points can be shattered
- ❖ Therefore, $VC = 3$



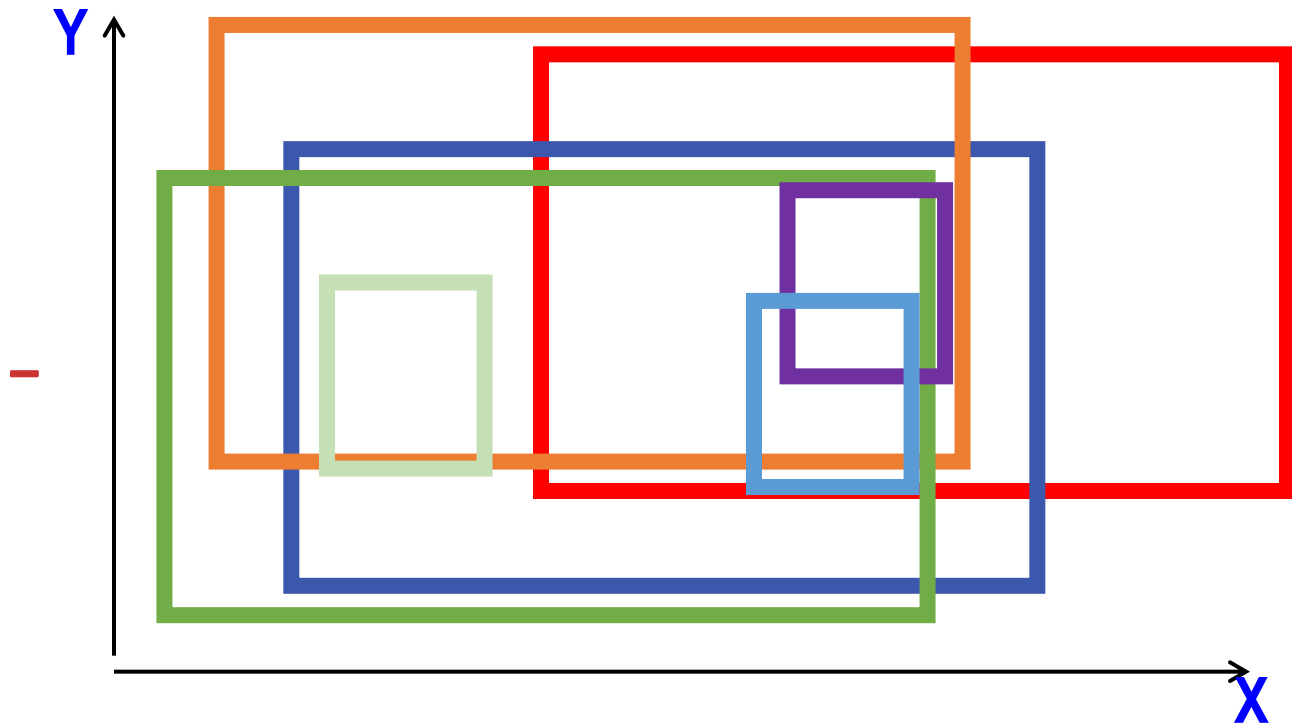
- ❖ Suppose three of them lie on the same line, label the outside points + and the inner one –
- ❖ Other wise, make a convex hull. Label points outside + and the inner one –
- ❖ **Four** points **cannot** be shattered!

VC dimension of Half spaces

- ❖ In general, the VC dimension of an n -dimensional linear function is $n+1$
- ❖ Give the same δ and m

$$err_D(h) \leq err_S(h) + \sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

Exercise



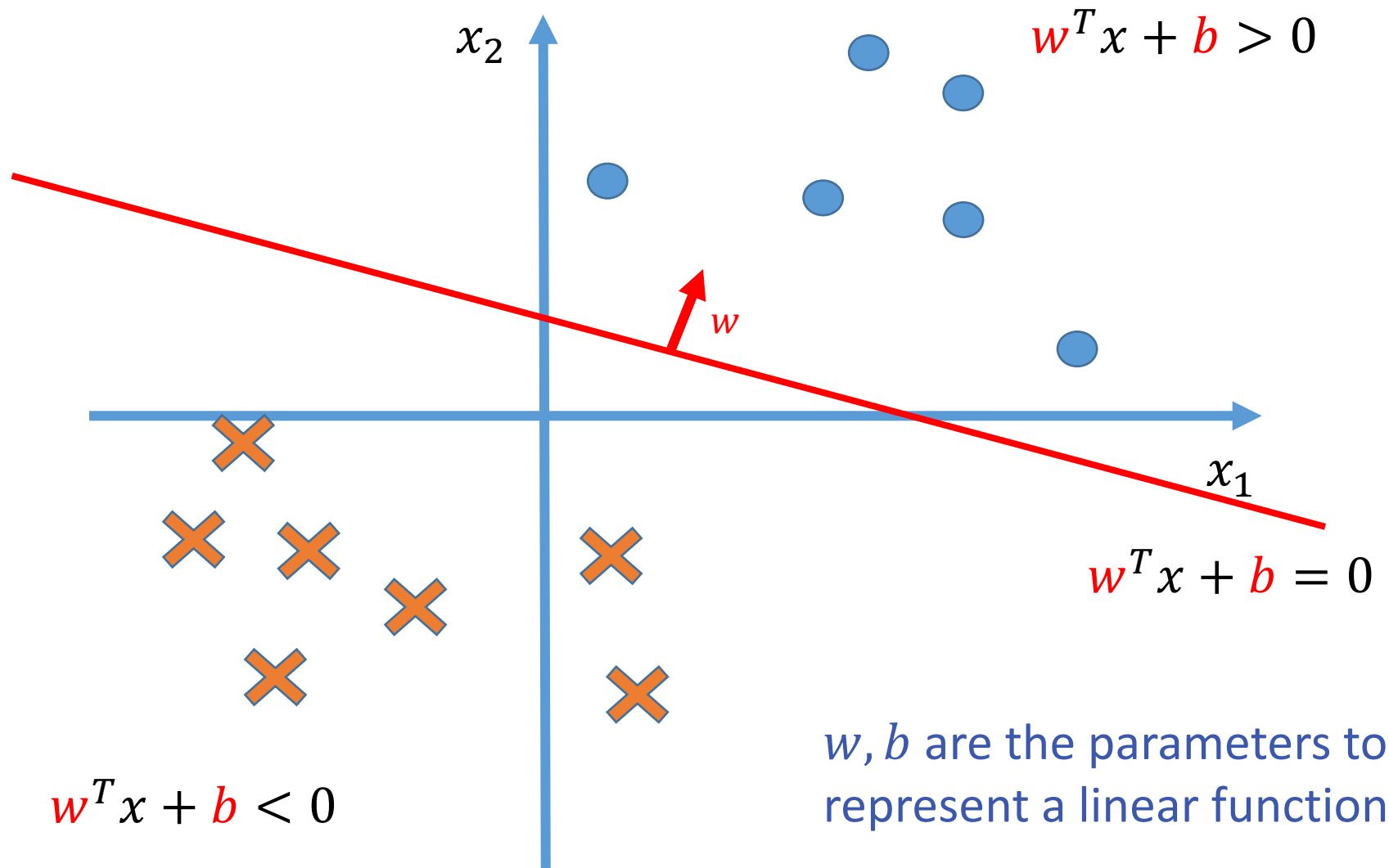
What is the VC dimension for the rectangle concept space?

Computational Learning Theory

- ❖ The Theory of Generalization
 - ❖ Using training instance to rule out incorrect hypotheses
- ❖ Probably Approximately Correct (PAC) learning
 - ❖ How many examples you need to see to obtain a learned function with error $\leq \epsilon$
- ❖ Shattering and the VC dimension

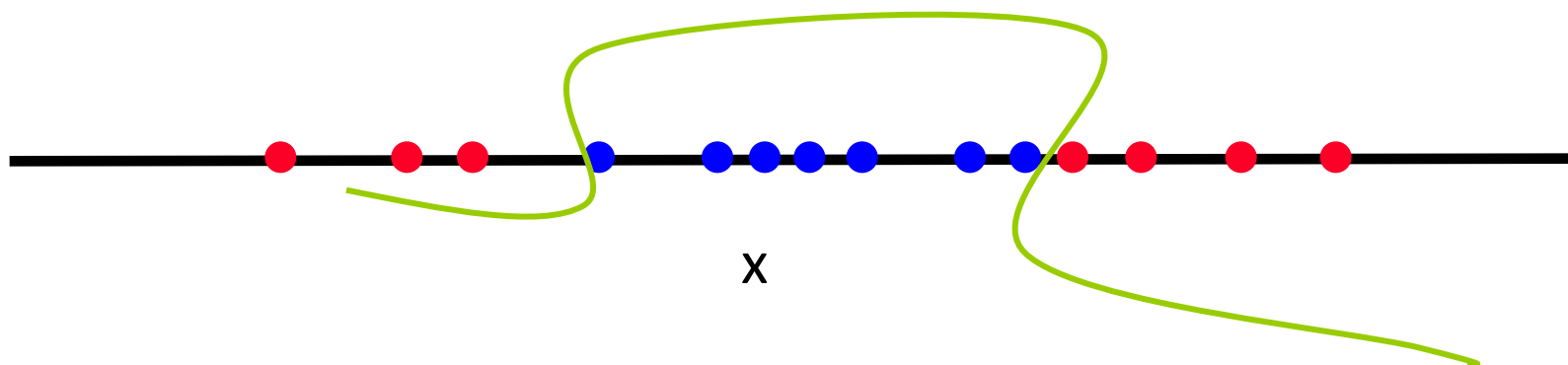
Kernel and Kernel methods

Hypothesis space: linear model



Functions Can be Made Linear

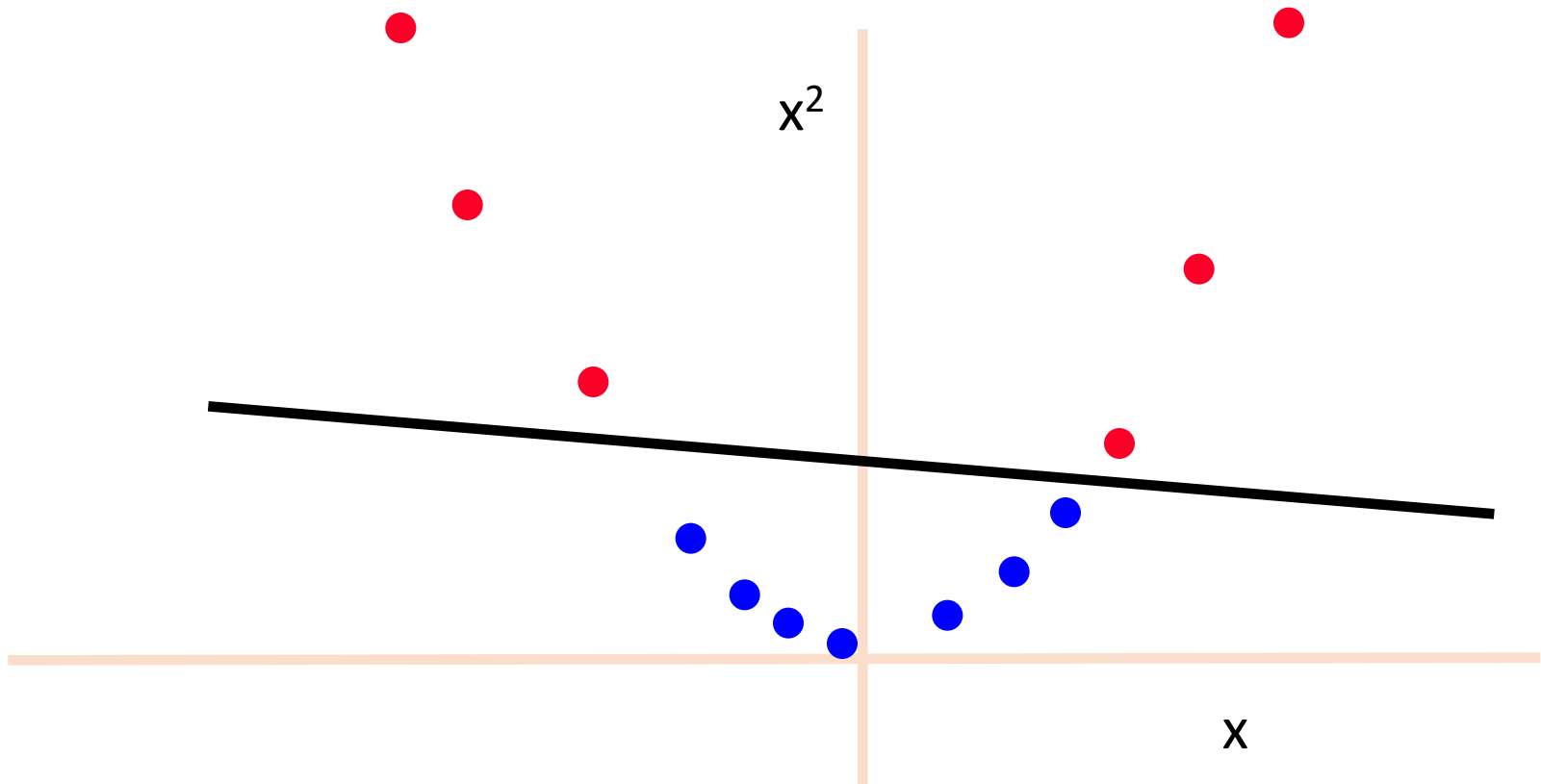
- ❖ Data are not linearly separable in one dimension
- ❖ Not separable if you insist on using a specific class of functions



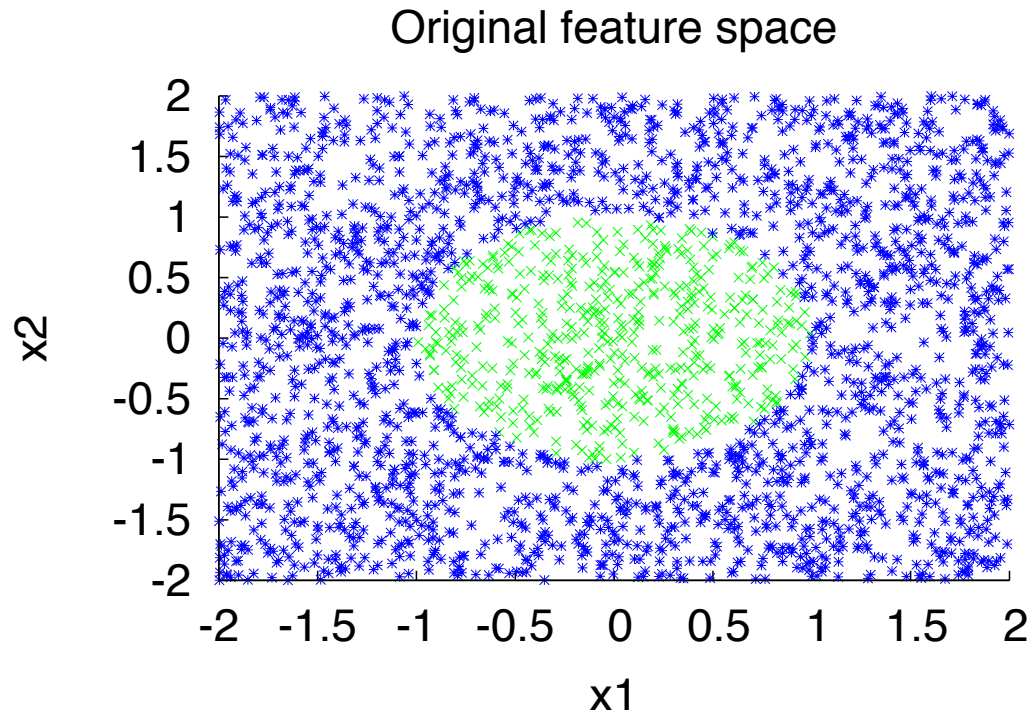
Can we do some mapping to make it linear spreadable?

Blown Up Feature Space

❖ Data are separable in $\langle x, x^2 \rangle$ space

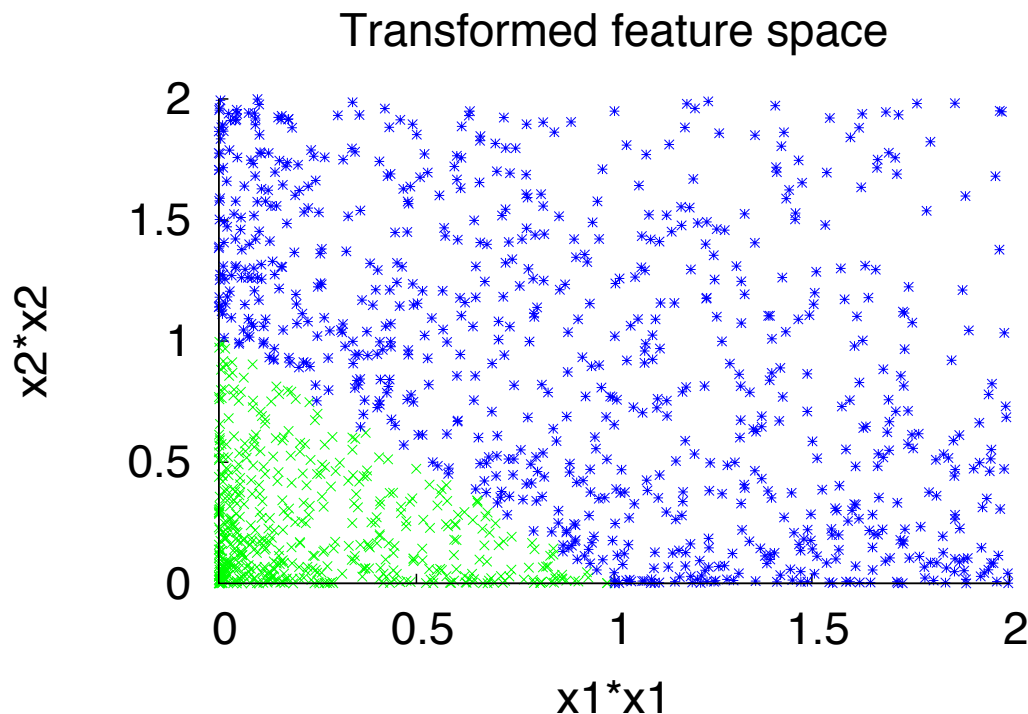


Making data linearly separable



$$f(\mathbf{x}) = 1 \text{ iff } x_1^2 + x_2^2 \leq 1$$

Making data linearly separable



Transform data: $\mathbf{x} = (x_1, x_2) \Rightarrow \mathbf{x}' = (x_1^2, x_2^2)$
 $f(\mathbf{x}') = 1$ iff $x'_1 + x'_2 \leq 1$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$

2. For (\mathbf{x}, y) in \mathcal{D} :

3. if $y(\mathbf{w}^\top \mathbf{x}) \leq 0$

Assume $y \in \{1, -1\}$

4. $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$

5.

6. Return \mathbf{w}

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(\mathbf{w}^\top \mathbf{x}^{\text{test}})$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^{2n}$

2. For (\mathbf{x}, y) in \mathcal{D} :

3. if $y \mathbf{w}^T \begin{bmatrix} x \\ x^2 \end{bmatrix} \leq 0$

Assume $y \in \{1, -1\}$

4. $\mathbf{w} \leftarrow \mathbf{w} + y \begin{bmatrix} x \\ x^2 \end{bmatrix}$

5.

6.

What if our mapping function is more complex?

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(\mathbf{w}^T \begin{bmatrix} x \\ x^2 \end{bmatrix})$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$

2. For (\mathbf{x}, y) in \mathcal{D} :

3. if $y(\mathbf{w}^\top \mathbf{x}) \leq 0$

Assume $y \in \{1, -1\}$

4. $\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$

5.

6. Return \mathbf{w}

Observation: \mathbf{w} is a combination of the input instances!!

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(\mathbf{w}^\top \mathbf{x}^{\text{test}})$

Dual Representation

$$\text{if } y(\mathbf{w}^\top \mathbf{x}) \leq 0$$

$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$

- ❖ Let \mathbf{w} be an initial weight vector for perceptron. Let $(x_1, +)$, $(x_2, +)$, $(x_3, -)$, $(x_4, -)$ be examples and assume mistakes are made on x_1 , x_2 and x_4 .
- ❖ What is the resulting weight vector?

$$\mathbf{w} = \mathbf{w} + x_1 + x_2 - x_4$$

- ❖ In general, the weight vector \mathbf{w} can be written as a linear combination of examples:

$$\mathbf{w} = \sum_{1..m} \alpha_i y_i x_i$$

- ❖ Where α_i is the **number of mistakes** made on x_i .

Predicting with linear classifiers

- ❖ Prediction = $\text{sgn}(\mathbf{w}^T \mathbf{x})$ and $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$
- ❖ That is, we just showed that

$$\mathbf{w}^T \mathbf{x} = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}$$

- ❖ We only need to compute dot products between training examples and the new example \mathbf{x}
- ❖ This is true even if we map examples to a high dimensional space

$$\mathbf{w}^T \phi(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$$

Many learning algorithm require to compute inner products

❖ Perceptron:

$$y(\mathbf{w}^T \mathbf{x}) \leq 0$$

❖ K-NN:

$$\text{similarity}(\mathbf{x}, \mathbf{x}^{\text{neighbor}}) = \mathbf{x}^T \mathbf{x}^{\text{neighbor}}$$

$$\text{dist}(\mathbf{x}, \mathbf{x}^{\text{neighbor}}) = \|\mathbf{x} - \mathbf{x}^{\text{neighbor}}\|^2$$

$$\text{dist}(\mathbf{x}, \mathbf{x}^{\text{neighbor}}) = \|\mathbf{x}\|^2 + \|\mathbf{x}^{\text{neighbor}}\|^2 - 2\mathbf{x}^T \mathbf{x}^{\text{neighbor}}$$

Is there a smarter way to compute the inner product?

Dot products in high dimensional spaces

Let us define a dot product in the high dimensional space

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

Dot products in high dimensional spaces

Let us define a dot product in the high dimensional space

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

So prediction with this *high dimensional lifting map* is

$$\text{sgn}(\mathbf{w}^T \phi(\mathbf{x})) = \text{sgn} \left(\sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) \right)$$

because $\mathbf{w}^T \phi(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$

Dot products in high dimensional spaces

Let us define a dot product in the high dimensional space

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$$

So prediction with this *high dimensional lifting map* is

If we can compute the value of K *without explicitly writing the blown up representation*, then we will have a computational advantage.

because $\mathbf{w}^T \phi(\mathbf{x}) = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$

Example: Polynomial Kernel


- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

 All degree zero terms

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1 x_2, \dots, x_{n-1} x_n]^T$$

All degree zero terms

All degree one terms

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

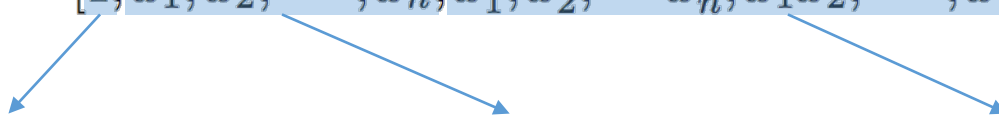
$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

The diagram illustrates the structure of the polynomial kernel mapping ϕ . The vector $\phi(x_1, x_2, \dots, x_n)$ is composed of three distinct groups of terms, each highlighted with a light blue background in the original image. Three blue arrows point from these groups to descriptive labels below the equation:

- The first group, containing the constant term 1, is labeled "All degree zero terms".
- The second group, containing the linear terms x_1, x_2, \dots, x_n , is labeled "All degree one terms".
- The third group, containing the quadratic terms $x_1^2, x_2^2, \dots, x_n^2$ and the cross terms $x_1x_2, \dots, x_{n-1}x_n$, is labeled "All degree two terms".

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$


The diagram shows three blue arrows pointing from the vector components to their respective degree categories: the first arrow points from the constant term '1' to 'All degree zero terms'; the second arrow points from the linear terms x_1, x_2, \dots, x_n to 'All degree one terms'; the third arrow points from the quadratic and cross terms $x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n$ to 'All degree two terms'.

All degree zero terms All degree one terms All degree two terms

and compute the dot product $A = \phi(\mathbf{x})^T \phi(\mathbf{z})$
[takes time]

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space**

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

and compute the dot product $\mathbf{A} = \phi(\mathbf{x})^T \phi(\mathbf{z})$
[takes time]

- ❖ Instead, in the original space, compute

$$\mathbf{B} = K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$$

Theorem: $\mathbf{A} = \mathbf{B}$ (Coefficients do not really matter)

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space** [for example, quadratic]

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

and compute the dot product $A = \phi(\mathbf{x})^T \phi(\mathbf{z})$ [takes time]

- ❖ Instead, in the original space, compute

$$B = K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$$

Example: Polynomial Kernel

- ❖ Given two examples \mathbf{x} and \mathbf{z} we want to map them to a **high dimensional space** [for example, quadratic]

$$\phi(x_1, x_2, \dots, x_n) = [1, x_1, x_2, \dots, x_n, x_1^2, x_2^2, \dots, x_n^2, x_1x_2, \dots, x_{n-1}x_n]^T$$

and compute the dot product $\mathbf{A} = \phi(\mathbf{x})^T \phi(\mathbf{z})$ [takes time]

- ❖ Instead, in the original space, compute

$$\mathbf{B} = K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$$

Claim: Compute \mathbf{B} instead of \mathbf{A} (Coefficients do not really matter)

The Kernel Trick

Suppose we wish to compute

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z})$$

Here ϕ maps \mathbf{x} and \mathbf{z} to a high dimensional space

The Kernel Trick: Save time/space by computing the value of $K(\mathbf{x}, \mathbf{z})$ by performing operations in the original space (without a feature transformation!)