

Learning theory, Kernel Methods

Sriram Sankararaman

The instructor gratefully acknowledges 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.

Outline

- 1 Learning theory
- 2 Kernel methods
- 3 Example
- 4 Kernels
- 5 Another example

What concept is learnable?

Learning Conjunctions

Learning Conjunctions -- Algorithm

Training data

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

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

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

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

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

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

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

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

A simple learning algorithm (*Elimination*)

- Discard all negative examples
- Build a conjunction using the features that are common to all positive conjunctions

$$h = x_1 \wedge x_3 \wedge x_4 \wedge x_5 \wedge x_{100}$$

$x_1, x_3, x_4, x_5, x_{100}$
0, 1, 1, 1, 1

Positive examples *eliminate* irrelevant features

Learning Conjunctions: Analysis

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

$$O(n \log n) \leq O(n^2)$$

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

Poly in n , $1/\delta$, $1/\epsilon$

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

n : # literals

If we see these many training examples, then the algorithm will produce a conjunction that, with high probability, will make few errors

Requirements of Learning

- ❖ Cannot expect a learner to learn a concept *exactly*
- ❖ There will generally be multiple concepts consistent with the available data
- ❖ Unseen examples could *potentially* have any label
- ❖ We may misclassify *uncommon* examples that do not show up in the training set

Requirements of Learning

- ❖ Cannot expect a learner to learn a concept exactly
 - ❖ There will generally be multiple concepts consistent with the available data
 - ❖ Unseen examples could *potentially* have any label
 - ❖ We may misclassify *uncommon* examples that do not show up in the training set
- ❖ Cannot always expect to learn a *close approximation* to the target concept
 - ❖ Sometimes the training set will not be representative

Probably approximately correctness

- ❖ The only realistic expectation of a good learner is that with high probability it will learn a close approximation to the target concept

$1 - \delta$ ϵ is small

- ❖ In Probably Approximately Correct (PAC) learning, one requires that

- ❖ given small parameters ϵ and δ ,

- ❖ With probability at least $1 - \delta$, a learner produces a hypothesis with error at most ϵ

- ❖ The reason we can hope for this is the consistent distribution assumption

Training & Test data have same distribution ✓

PAC Learnability

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 $\text{size}(H)$.

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

efficiently learnability

- ❖ The concept class C is *efficiently learnable* if L can produce the hypothesis in time that is polynomial in $1/\epsilon$, $1/\delta$, n and $\text{size}(H)$

PAC Learnability

- ❖ We impose two limitations
 - ❖ Polynomial sample complexity (information theoretic constraint)
 - ❖ Is there enough information in the sample to distinguish a hypothesis h that approximate f ?
 - ❖ Polynomial time complexity (computational complexity)
 - ❖ Is there an efficient algorithm that can process the sample and produce a good hypothesis h ?
- Worst Case definition: the algorithm must meet its accuracy
- ❖ for every distribution (The distribution free assumption)
 - ❖ for every target function f in the class C

Example: Learning Conjunctions

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

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

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

m is *polynomial* in $\underline{1/\epsilon}$, $\underline{1/\delta}$, \underline{n} and $\underline{\text{size}(H)}$

A general result

$|H|$ \leftarrow small \Leftrightarrow simple

Training error
 $err_S(H) = 0$

Let H be any hypothesis space.

With probability $1 - \delta$ a hypothesis h that is consistent with a training set of size m will have an error $< \epsilon$ 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.

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 $< \epsilon$ on future examples if

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

It expresses a preference towards smaller hypothesis spaces

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

Complicated/larger hypothesis spaces are not necessarily bad. But simpler ones are unlikely to fool us by being consistent with many examples!

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

$$h(x) = f(\underline{w^T x + b})$$

$$\text{sign}(w^T x + b)$$

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)

$$\hat{y} = \text{Sign}(w^T x + b)$$

$$h = g(w^T x + b)$$

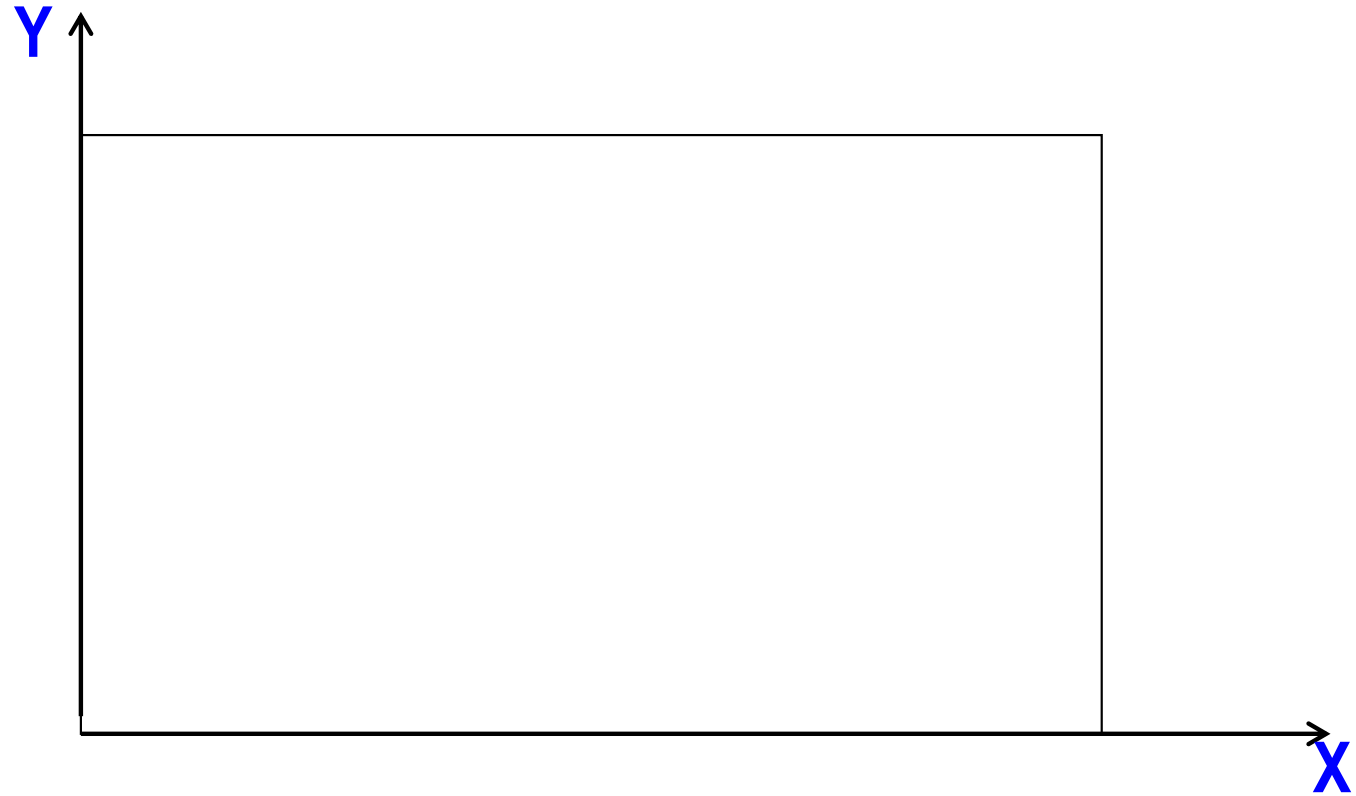
$$\hat{y} = g(v^T h + c)$$

Not expressed
← XOR

↙ Expressed

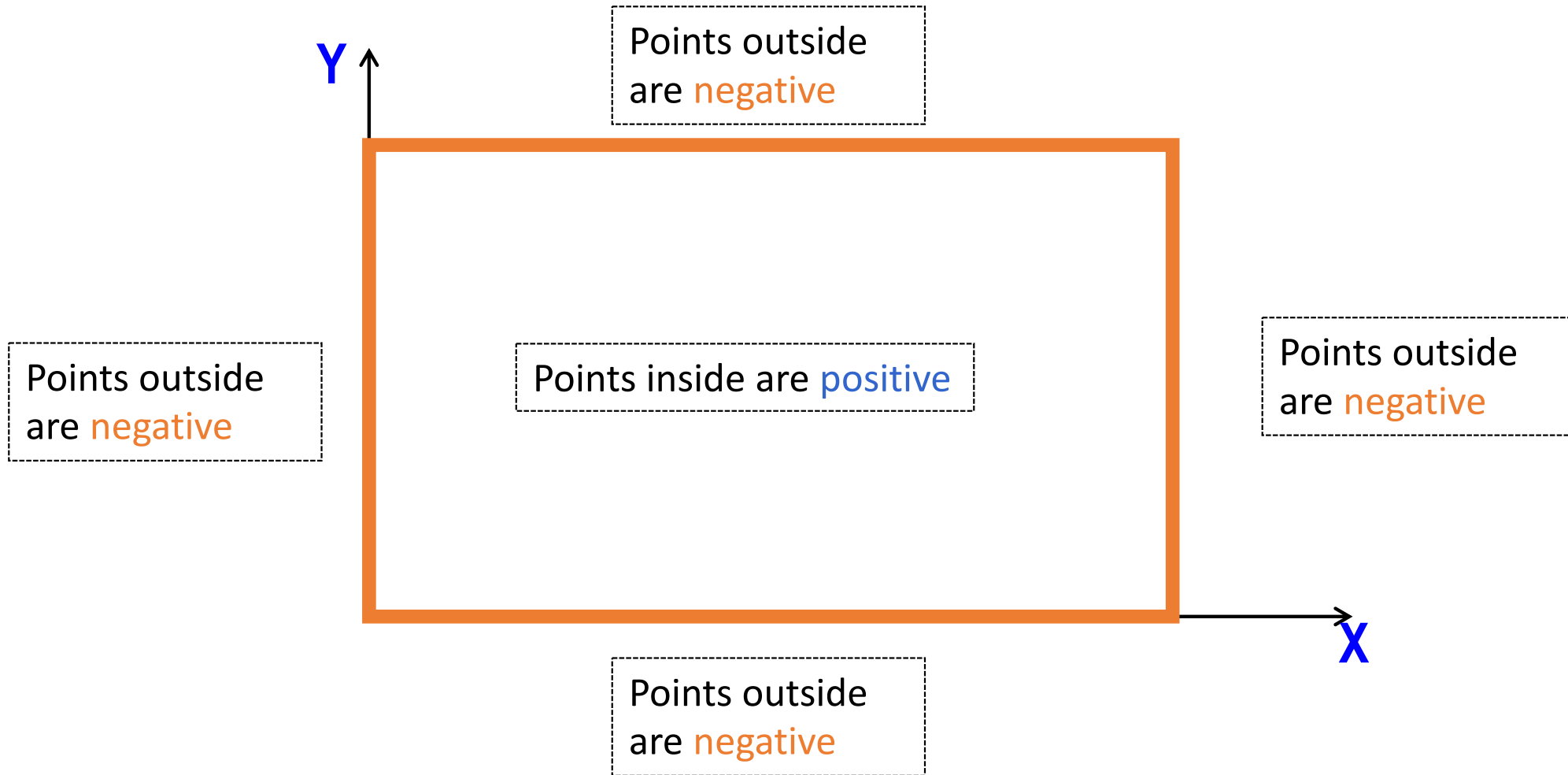
Learning Rectangles

Assume the target concept is an axis parallel rectangle



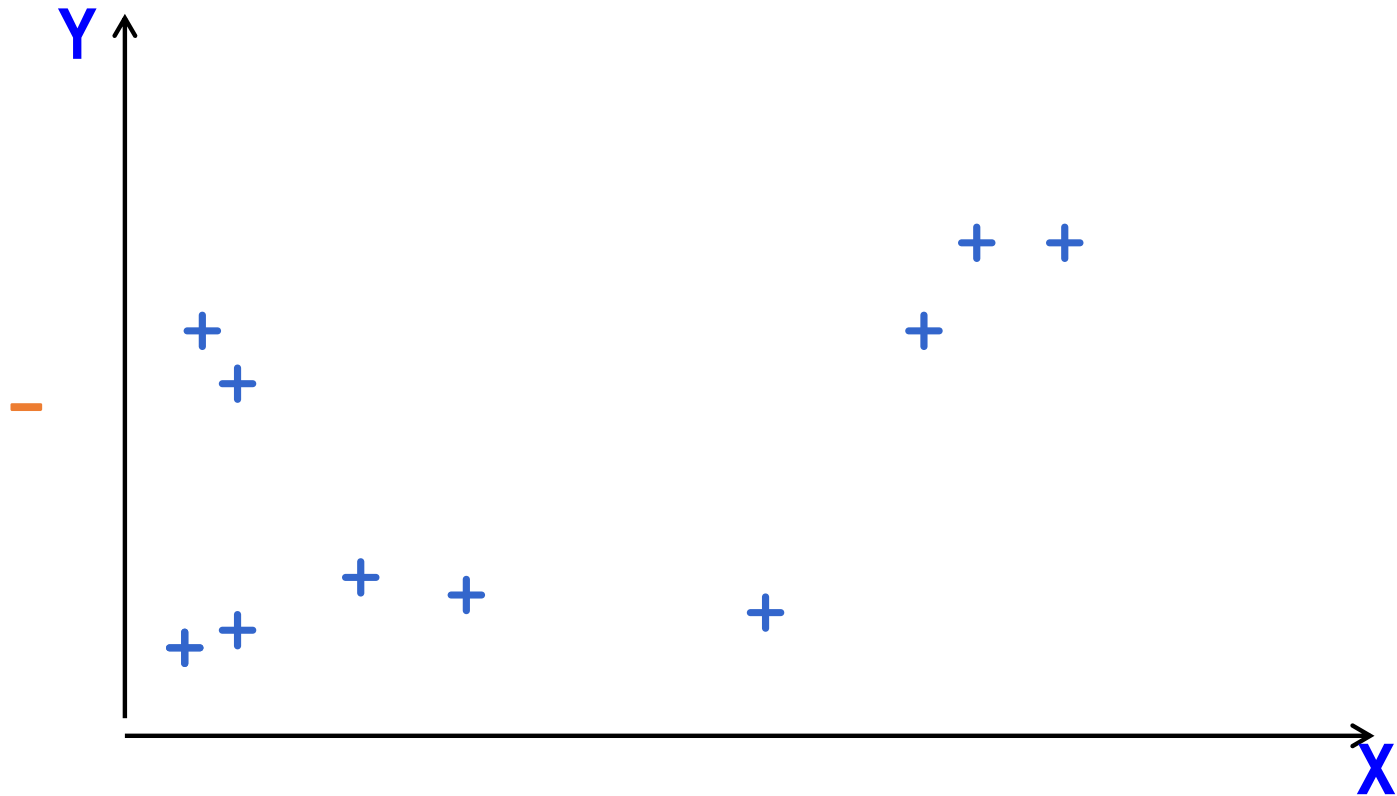
Learning Rectangles

Assume the target concept is an axis parallel rectangle



Learning Rectangles

Assume the target concept is an axis parallel rectangle

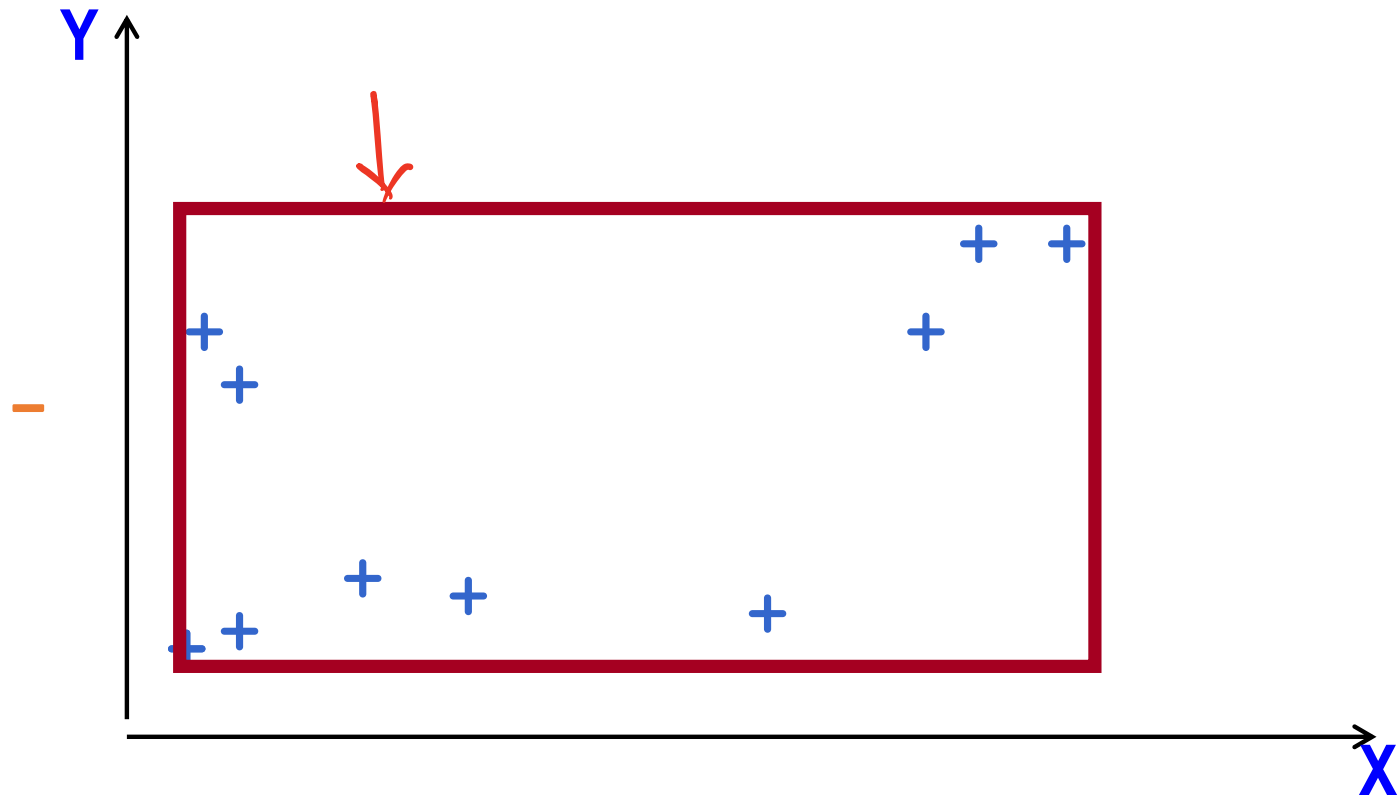


Will we be able to learn the target rectangle?

Can we come close?

Learning Rectangles

Assume the target concept is an axis parallel rectangle

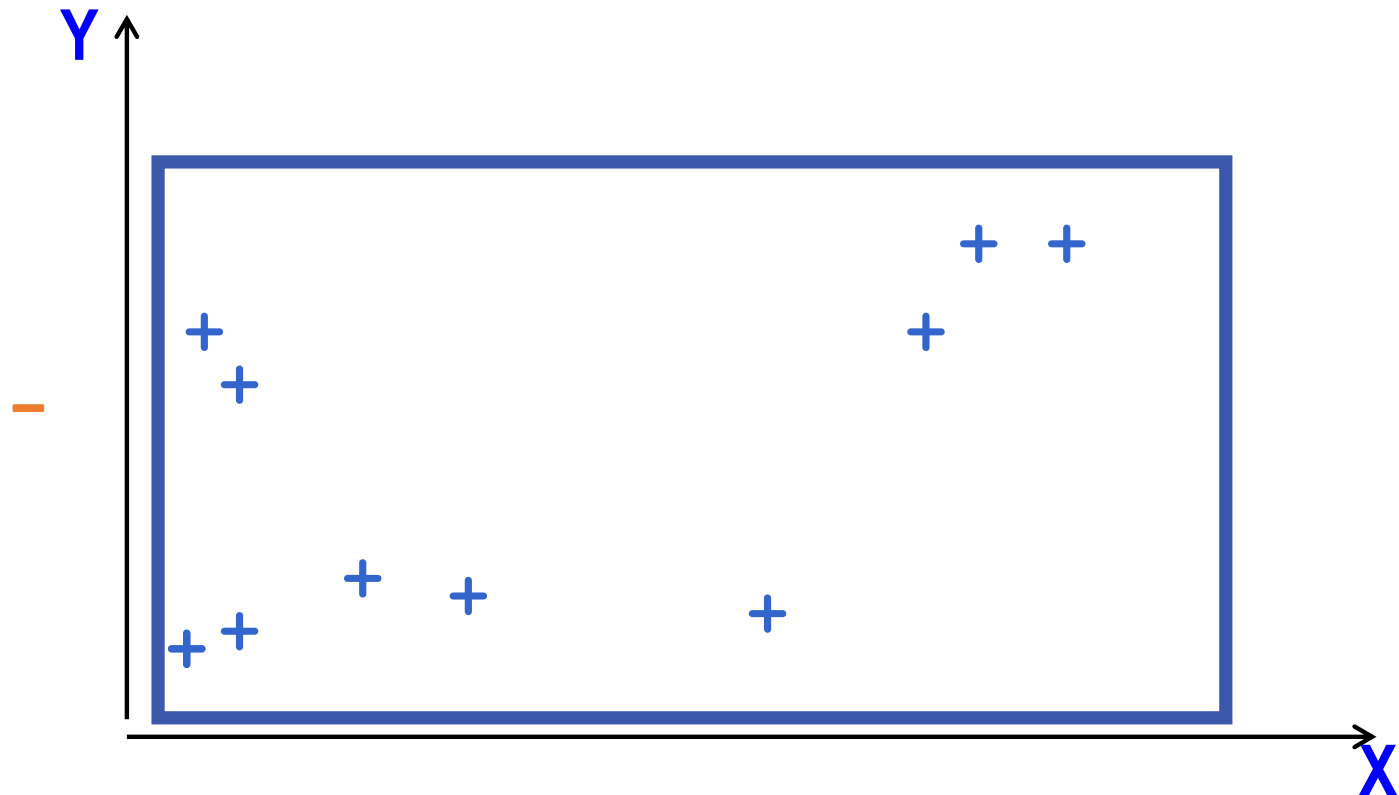


Key observation: Despite there are infinite # hypothesis
The blue & red rectangles have the same predictions

Can we come close?

Learning Rectangles

Assume the target concept is an axis parallel rectangle

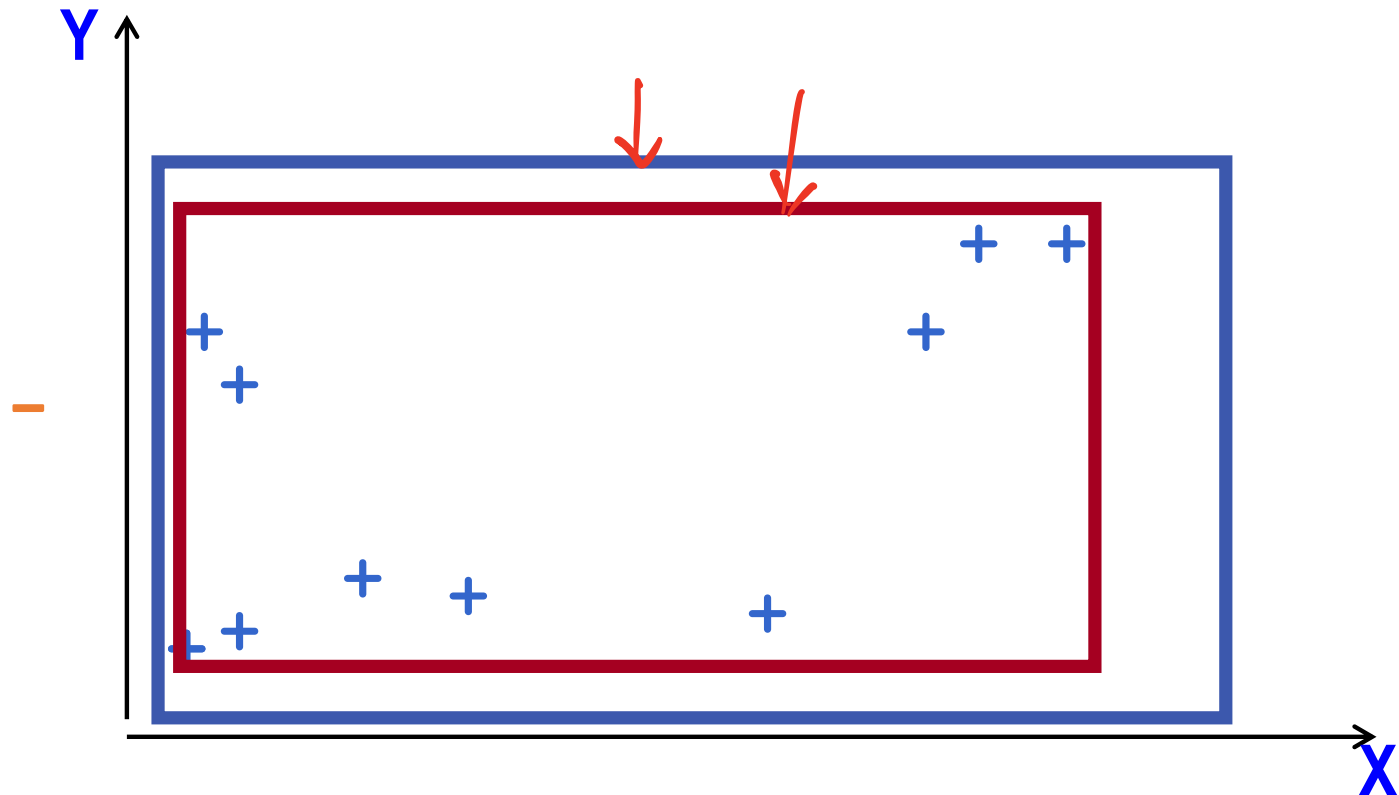


Key observation: Despite there are infinite # hypothesis
The blue & red rectangles have the same predictions

Can we come close?

Learning Rectangles

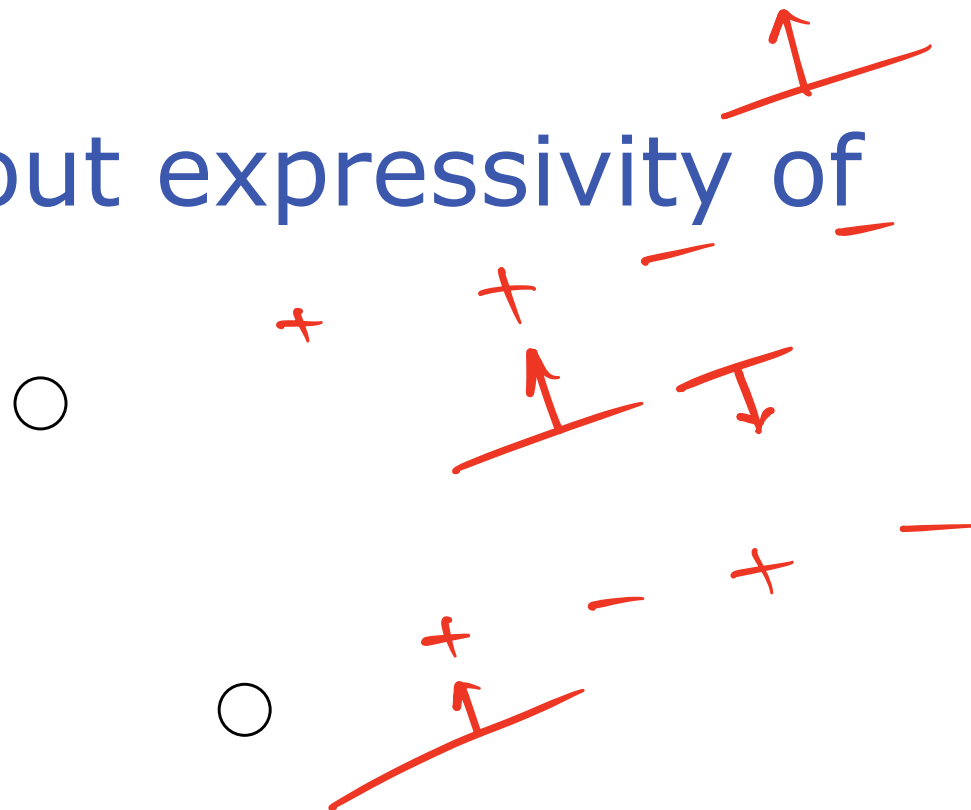
Assume the target concept is an axis parallel rectangle



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

Can we come close?

Let's think about expressivity of functions

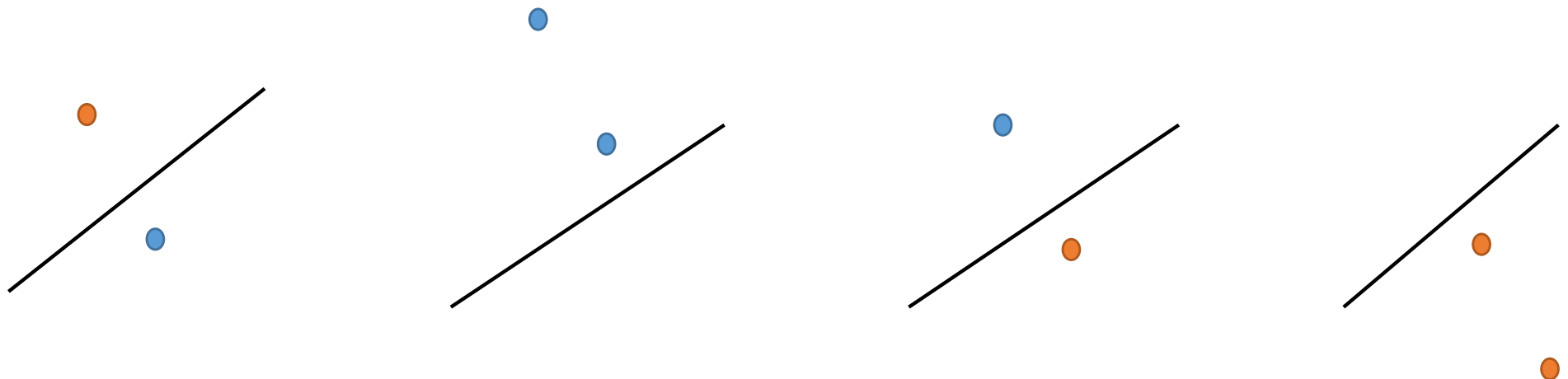


Suppose we have two points.

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

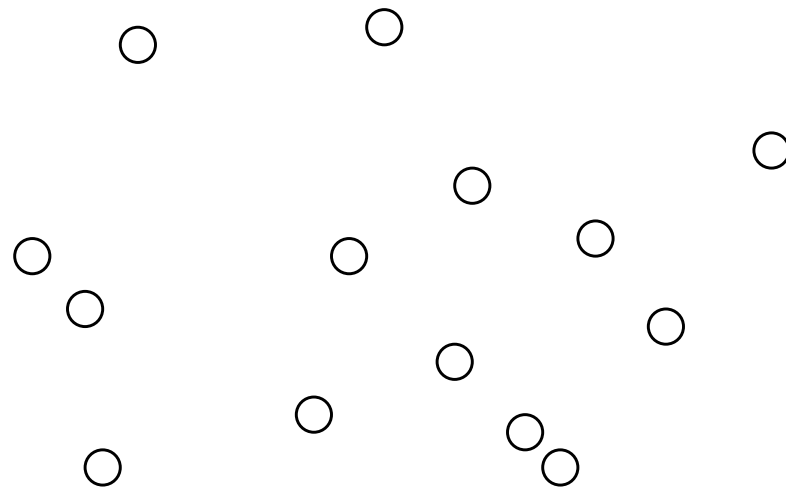


Suppose we have two points.

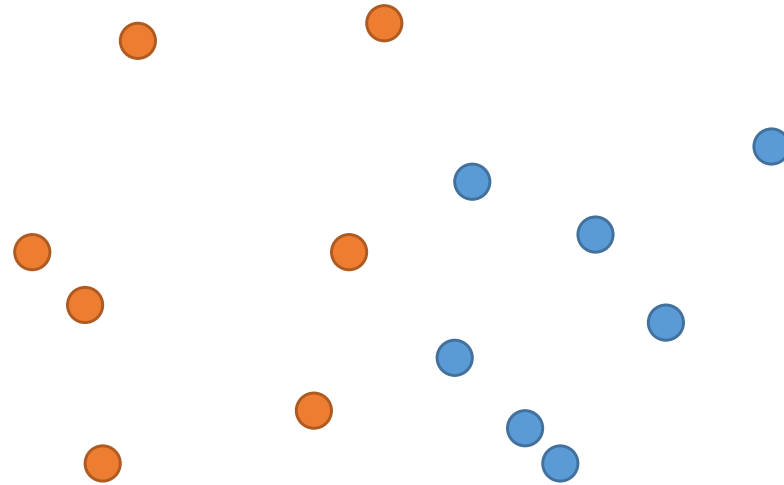
Can linear classifiers correctly classify any labeling of these points?

Linear functions are expressive enough to *shatter* 2 points

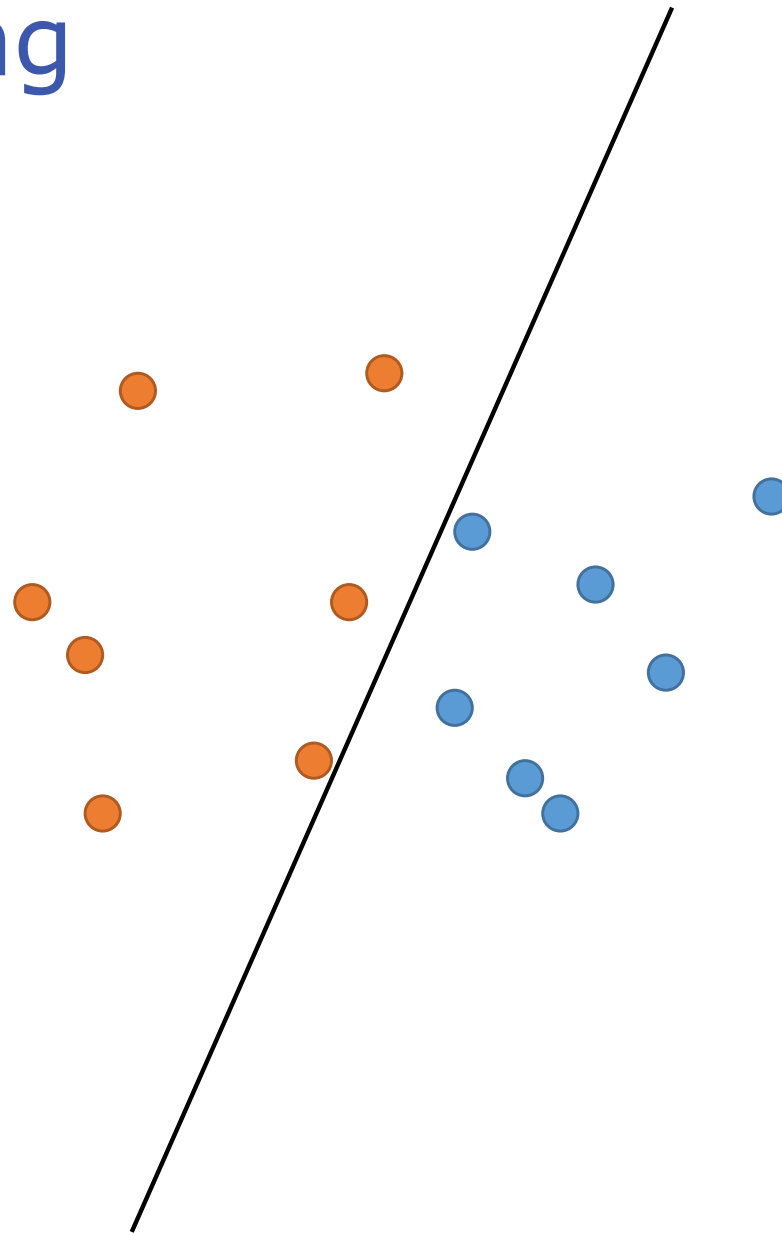
Shattering



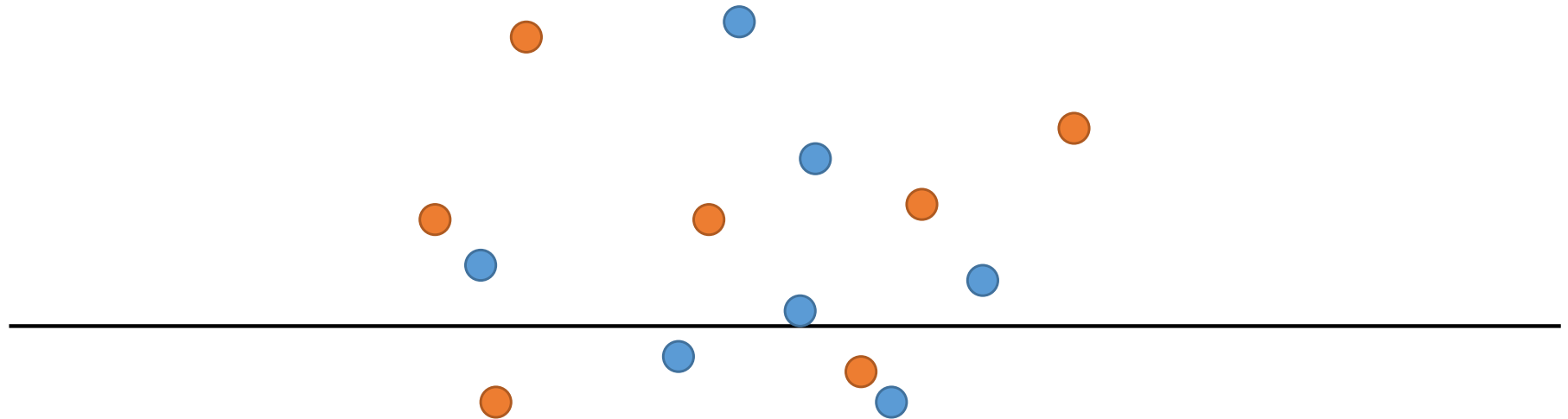
Shattering



Shattering

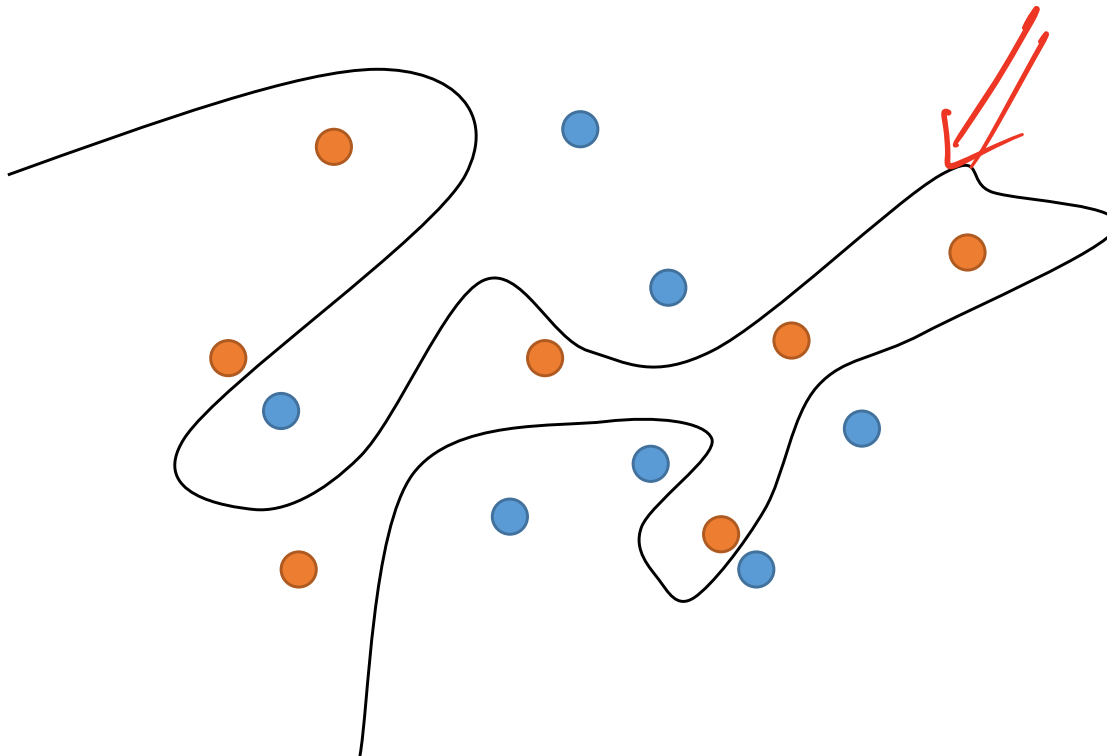


Shattering



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

Shattering



Linear functions are not expressive to shatter fourteen points

Because there is a labeling that can not be separated by them

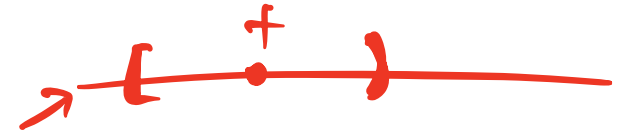
Of course, a more complex function could separate them

Shattering

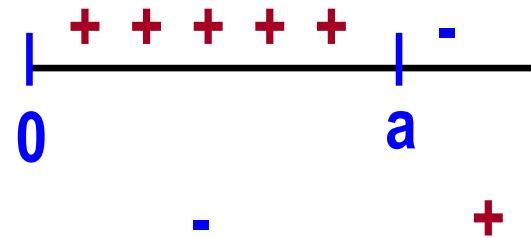
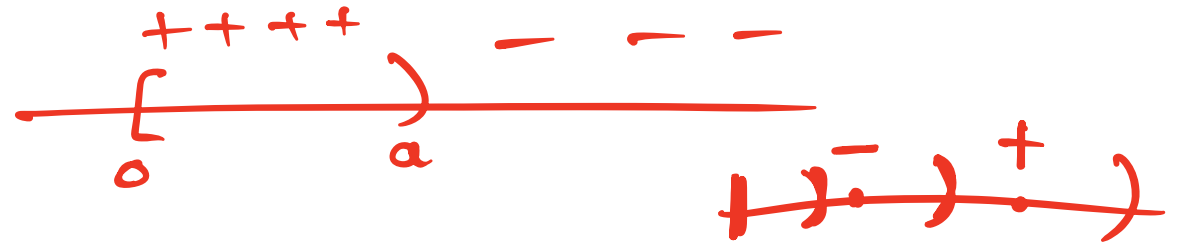
Definition: A set S of examples is shattered by a set of functions H if for every partition of the examples in S into positive and negative examples there is a function in H that gives exactly these labels to the examples

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

Left bounded intervals



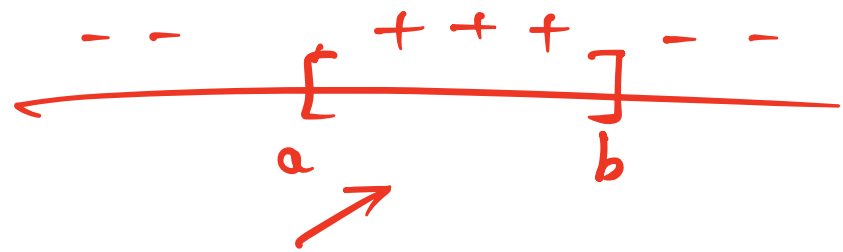
Example 1: Hypothesis class of left bounded intervals on the real axis: $[0, a)$ for some real number $a > 0$



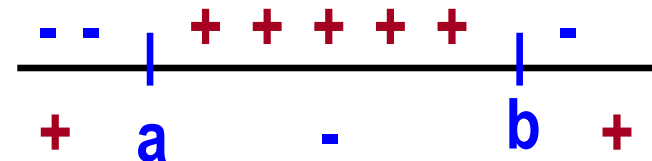
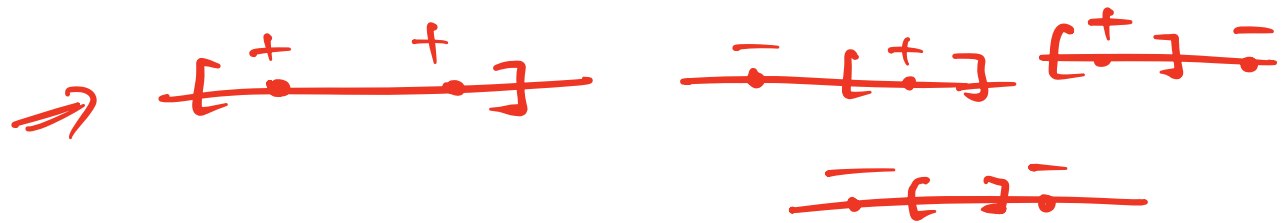
Sets of **two** points **cannot** be shattered

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

Shattering: The adversarial game

You



You: Hypothesis class H can shatter these d points

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. I will find another one

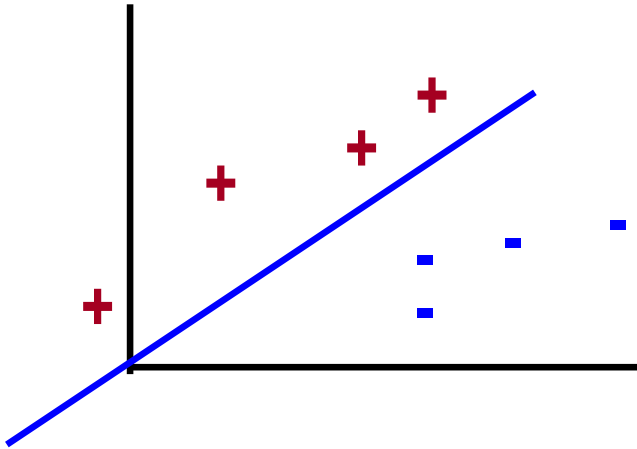
Vapnik-Chervonenkis Dimension

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, $VC(H) \geq \underline{d}$
- ❖ Even one subset will do
- ❖ If **no subset** of size d can be shattered, then $VC(H) < d$

Example 3: 2-D Half spaces in a plane

↑
Lines



Can one point be shattered?

Is there any two points can be shattered?

Is there any three points?

Can any three points be shattered?

$$VC(\text{Lines}) \geq 3$$



Example Half spaces in a plane

- ❖ Prove $VC \geq 1$

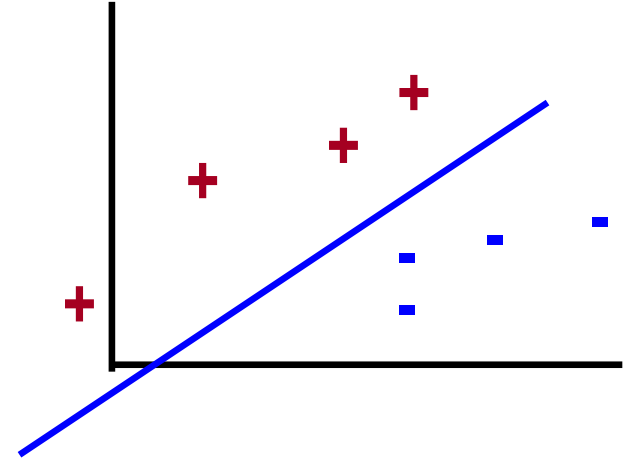
- ❖ Show any point can be shattered

- ❖ Prove $VC \geq 2$

- ❖ Show there exists 2 points can be shattered

- ❖ Prove $VC \geq 3$

- ❖ Show there exists 3 points can be shattered



+	-	+	-
x	x	x	x

+

-

-

+

+

-

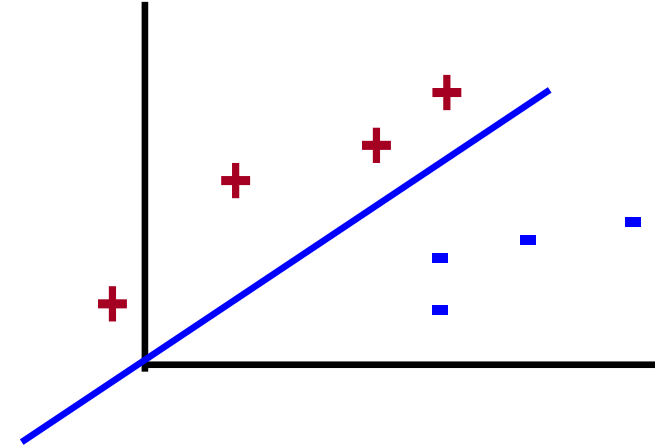
+

+

VC (Line) < 4

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 –
- ❖ Otherwise, make a convex hull. Clockwisely, label them + - + - ;
- ❖ Otherwise, one point inside the convex hull of the other three points; label 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 $n+1$

This term may decrease

This term will decrease

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

Generalization error

Training error

m

This formulation will be provided in the exam

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

Outline

- 1 Learning theory
- 2 Kernel methods
 - Motivation
- 3 Example
- 4 Kernels
- 5 Another example

Motivation

- Linear models are convenient.
 - ▶ Computationally efficient for learning (training) and prediction.
- We would like our models to be “expressive”.
 - ▶ If it is not expressive enough, it will underfit.
 - ▶ Too expressive models can overfit.

Motivation

How to increase the expressive power of linear models?

- Map the feature vector \mathbf{x} to an expanded version
 $\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$.
- Use a non-linear basis function $\phi(\mathbf{x})$ as input to linear model.

Motivation

How to increase the expressive power of linear models?

- Map the feature vector \mathbf{x} to an expanded version $\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D \rightarrow \mathbf{z} \in \mathbb{R}^M$.
- Use a non-linear basis function $\phi(\mathbf{x})$ as input to linear model.
- **Difficulty:** When M is large, computational difficulty.

Motivation

How to choose nonlinear basis function for regression?

$$w^T \underline{\phi(x)}$$

where $\phi(\cdot)$ maps the original feature vector x to a M -dimensional *new* feature vector. In the following, we will show that we can sidestep the issue of choosing which $\phi(\cdot)$ to use — instead, we will choose *equivalently* a *kernel function* that are often computationally efficient to work with.

Example

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_D \\ x_1^2 \\ x_1x_2 \\ \vdots \\ x_1x_D \\ x_2x_1 \\ x_2^2 \\ \vdots \\ x_2x_D \\ \vdots \\ x_Dx_1 \\ \vdots \\ x_D^2 \end{pmatrix}$$

Handwritten red annotations:

- A red arrow points to the constant term 1.
- A red curly brace groups the terms $\sqrt{2}x_1, \sqrt{2}x_2, \dots, \sqrt{2}x_D$.
- A red curly brace groups the quadratic terms $x_1^2, x_1x_2, \dots, x_2x_D, \dots, x_Dx_1, x_D^2$.
- Handwritten red text $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix}$ is shown to the right.
- Handwritten red text $O(D^2)$ is shown below the quadratic terms.

Outline

- 1 Learning theory
- 2 Kernel methods
- 3 Example**
- 4 Kernels
- 5 Another example

Kernelized nearest neighbors

$$\mathbf{x}_m \in \mathbb{R}^D \quad \mathbf{x}_n \in \mathbb{R}^D$$

$$(\mathbf{x}_m - \mathbf{x}_n)^T (\mathbf{x}_m - \mathbf{x}_n)$$

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

To perform classification in a transformed feature space, we only need to compute distances in this space:

$$d(\phi(\mathbf{x}_m), \phi(\mathbf{x}_n)) = \|\phi(\mathbf{x}_m) - \phi(\mathbf{x}_n)\|_2^2 = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_m) + \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_n) - 2\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

Computing inner products

Many learning algorithms can be rewritten to depend on the instances $\underline{x_i}, \underline{x_j}$ only through inner products $\phi(\underline{x_i})^T \phi(\underline{x_j})$.

Why is this helpful ?

Computing inner products

$$\rightarrow \phi(x) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \vdots \\ \sqrt{2}x_D \\ x_1x_1 \\ x_1x_2 \\ \vdots \\ x_Dx_D \end{pmatrix} \quad O(D^2)$$

Many learning algorithms can be rewritten to depend on the instances x_i, x_j only through inner products $\phi(x_i)^T \phi(x_j)$.

Why is this helpful? Computing $\phi(x)$ is $O(D^2)$.

$$\begin{aligned} \phi(x_i)^T \phi(x_j) &= \underbrace{(1 + x_i^T x_j)}_{\substack{x_i \in \mathbb{R}^D \\ x_j \in \mathbb{R}^D}}^2 \quad O(D) \end{aligned}$$

$$\phi(x_i)^T \phi(x_j) = (1 + x_i^T x_j)^2$$

However, inner product can be computed in $O(D)$.

Even bigger gains when the dimensionality of the feature space (M) increases.

$$\begin{aligned} \phi(x_i)^T \phi(x_j) &= \begin{pmatrix} 1 & \sqrt{2}x_{i1} & \sqrt{2}x_{i2} & x_{i1}^2 & x_{i1}x_{i2} & x_{i2}^2 \end{pmatrix} \begin{pmatrix} 1 \\ \sqrt{2}x_{j1} \\ \sqrt{2}x_{j2} \\ x_{j1}^2 \\ x_{j1}x_{j2} \\ x_{j2}^2 \end{pmatrix} \\ x_i &= \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \quad x_j = \begin{pmatrix} x_{j1} \\ x_{j2} \end{pmatrix} \end{aligned}$$

Computing inner products

Let us define a function that computes inner products in the transformed feature space

$$\underline{k(\mathbf{x}_i, \mathbf{x}_j)} = \underline{\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)}$$

If we can compute the value of k without explicitly computing ϕ , we will have a computational advantage.

$$\begin{aligned} \Rightarrow \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) &= 1 + \underbrace{(2x_{i1}x_{j1})}_{\text{circled}} + \underbrace{(2x_{i2}x_{j2})}_{\text{circled}} \\ &\quad + x_{i1}^2 x_{j1}^2 + x_{i1} x_{i2} x_{j1} x_{j2} + x_{i2}^2 x_{j2}^2 \\ &\quad + x_{i2} x_{j1} x_{j2} x_{i1} + x_{i2}^2 x_{j2}^2 \\ \Rightarrow \underbrace{(1 + x_{i1}^T x_{j1})}_{\text{underlined}} &= \underbrace{(1 + x_{i1} x_{j1} + x_{i2} x_{j2})}_{\text{underlined}} \\ &\quad + (x_{i1} x_{j1})^2 + (x_{i2} x_{j2})^2 + 2(x_{i1} x_{j1}) + 2(x_{i2} x_{j2}) \\ &\quad + 2(x_{i1} x_{j1} x_{i2} x_{j2}) \end{aligned}$$

Kernelized nearest neighbors

$$\cancel{d(\phi(x_m), \phi(x_n)) = \phi(x_m)^T \phi(x_m) + \phi(x_n)^T \phi(x_n) - 2 \phi(x_m)^T \phi(x_n)}$$

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

$$\underline{d^{\text{KERNEL}}}(\underline{x_m}, \underline{x_n}) = k(\underline{x_m}, \underline{x_m}) + k(\underline{x_n}, \underline{x_n}) - 2k(\underline{x_m}, \underline{x_n})$$

The distance is equivalent to compute the distance between $\phi(x_m)$ and $\phi(x_n)$

$$d^{\text{KERNEL}}(\underline{x_m}, \underline{x_n}) = d(\underline{\phi(x_m)}, \underline{\phi(x_n)})$$

where the $\phi(\cdot)$ is the nonlinear mapping function implied by the kernel function. The nearest neighbor of a point \underline{x} is thus found with

$$\arg \min_n d^{\text{KERNEL}}(\underline{x}, \underline{x_n})$$

Outline

- 1 Learning theory
- 2 Kernel methods
- 3 Example
- 4 Kernels**
 - Kernel matrix and kernel functions
 - Kernelized machine learning methods
- 5 Another example

Inner products between features

Let us examine more closely the inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

Inner products between features

Let us examine more closely the inner products $\phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

This gives rise to an inner product in a special form,

$$\begin{aligned} \phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n) &= x_{m1}^2 x_{n1}^2 + 2x_{m1}x_{m2}x_{n1}x_{n2} + x_{m2}^2 x_{n2}^2 \\ &= (x_{m1}x_{n1} + x_{m2}x_{n2})^2 = \underline{\underline{(\mathbf{x}_m^\top \mathbf{x}_n)^2}} \end{aligned}$$

Namely, the inner product can be computed by a function $(\mathbf{x}_m^\top \mathbf{x}_n)^2$ defined in terms of the original features, *without computing $\phi(\cdot)$* .

Kernel functions

Some intuition

- To use kernelized nearest neighbors, we need to be able to compute an inner product between a test point and any training point.
- Since we need this for any possible pair of points, we need a function that takes a pair of points and computes an inner product.
- This is the kernel function $k(\cdot, \cdot)$.
- Given two inputs, $k(\cdot, \cdot)$ tells us how “similar” or “close” these inputs are in the space defined by the function ϕ .

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Only depends on difference between two inputs
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)!

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $\mathbf{x}_m, \mathbf{x}_n \in \mathbb{R}^D$, $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Only depends on difference between two inputs
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)!

These kernels have hyperparameters to be tuned: d , c , σ^2

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Kernel functions

Definition: a kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Not very useful though

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Example that is not a kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2$$

Conditions for being a positive semidefinite kernel function

Mercer theorem (loosely), a bivariate function $k(\cdot, \cdot)$ is a kernel function, if and only if, for *any N and any $\mathbf{x}_1, \mathbf{x}_2, \dots, \text{and } \mathbf{x}_N$* , the matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is positive semidefinite.

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

This matrix cannot be positive semidefinite as it has both *negative* and positive eigenvalues.

Recap: why use kernel functions?

Without specifying $\phi(\cdot)$, the kernel matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is exactly the same as

$$\begin{aligned} \mathbf{K} &= \Phi \Phi^T \\ &= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

‘Kernel trick’

Many learning methods depend on computing *inner products* between features — we have seen the example of regularized least squares. For those methods, we can use a kernel function in the place of the inner products, i.e., “*kernelizing*” the methods, thus, introducing nonlinear features.

When we talk about support vector machines, we will see the trick one more time.

There are infinite numbers of kernels to use!

Rules of composing kernels (this is just a partial list)

- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $ck(\mathbf{x}_m, \mathbf{x}_n)$ is also if $c > 0$.
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $\alpha k_1(\mathbf{x}_m, \mathbf{x}_n) + \beta k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also if $\alpha, \beta \geq 0$
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $k_1(\mathbf{x}_m, \mathbf{x}_n)k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also.
- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $e^{k(\mathbf{x}_m, \mathbf{x}_n)}$ is also.
- ...

In practice, choosing an appropriate kernel is an “art”

People typically start with polynomial and Gaussian RBF kernels or incorporate domain knowledge.

Outline

- 1 Learning theory
- 2 Kernel methods
- 3 Example
- 4 Kernels
- 5 Another example
 - Kernelized perceptron

Kernelized perceptron

Algorithm 1 PerceptronTrain ($\mathcal{D}, MaxIter$)

```
1:  $w \leftarrow 0$ 
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x, y) \in \mathcal{D}$  do
4:      $a \leftarrow w^T \phi(x)$ 
5:     if  $ya \leq 0$  then
6:        $w \leftarrow w + y\phi(x)$ 
7:     end if
8:   end for
9: end for
10: return  $w$ 
```

Prediction: $\hat{y} = \text{sign}(w^T \phi(x))$.

How to kernelize?

- At any iteration, $\mathbf{w} = \sum_n \alpha_n \phi(\mathbf{x}_n)$
- Here α_n is the number of mistakes made on example n .
- What will the prediction on a new sample \mathbf{x} be ?
- The activation a can then be computed as

$$\begin{aligned} a &= \mathbf{w}^T \phi(\mathbf{x}) \\ &= \left(\sum_n \alpha_n \phi(\mathbf{x}_n) \right)^T \phi(\mathbf{x}) \\ &= \sum_n \alpha_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) \end{aligned}$$

We can kernelize the prediction so that:

$$\hat{y} = \sum_n \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

Kernelized perceptron

Algorithm 2 KernelizedPerceptronTrain ($\mathcal{D}, \text{MaxIter}$)

```
1:  $\alpha \leftarrow 0$ 
2: for  $iter = 1 \dots \text{MaxIter}$  do
3:   for all  $(\mathbf{x}, y) \in \mathcal{D}$  do
4:      $a \leftarrow \sum_n \alpha_n k(\mathbf{x}_n, \mathbf{x})$ 
5:     if  $ya \leq 0$  then
6:        $\alpha \leftarrow \alpha + y$ 
7:     end if
8:   end for
9: end for
10: return  $(\alpha)$ 
```

Summary

- Kernels allow us to design algorithms that use rich set of features while being computationally efficient.
- Many machine learning algorithms can be “kernelized”.
- Picking kernels is an art.
- We still need to tune hyperparameters.