

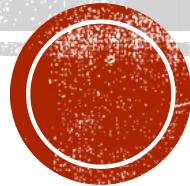
KERNELS

CS 412 Introduction to Machine Learning

Prof. Zheleva

February 22, 2017

Reading Assignment: CML: 11.1-11.4



LAST TIME: LINEAR MODELS

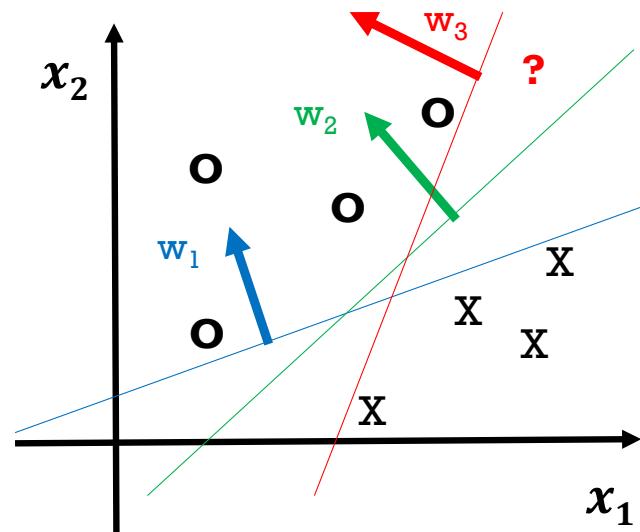
- Subgradients
 - When functions are not differentiable
- Closed form solutions to optimization problems
- (Linear) Support Vector Machines
 - Computing and optimizing the margin



LINEAR MODELS

- A linear classifier is a hyperplane that separates positive from negative examples (\mathbf{w}, b)
 - The prediction is a linear combination of feature values x
 - Examples: perceptron, linear SVM
- At test time, we check on what side of the hyperplane examples fall

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



LEARNING A LINEAR CLASSIFIER AS AN OPTIMIZATION PROBLEM

- Structural risk minimization framework
 - Tradeoff between low training error and a solution that is simple

$$\min_{w,b} L(w, b) = \min_{w,b} \left(\sum_{n=1}^N l(y_n, w \cdot x_n + b) + \lambda R(w, b) \right)$$

Variables over which we are optimizing the objective function

Objective function

Loss function
Measures how well classifier fits training data

Regularizing function
Helps with avoiding overfitting
Penalizes complex functions
 λ is a hyperparameter



APPROXIMATING THE 0/1 LOSS WITH CONVEX SURROGATE LOSS FUNCTIONS

- Examples (with $b=0$)

Zero/one: $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$

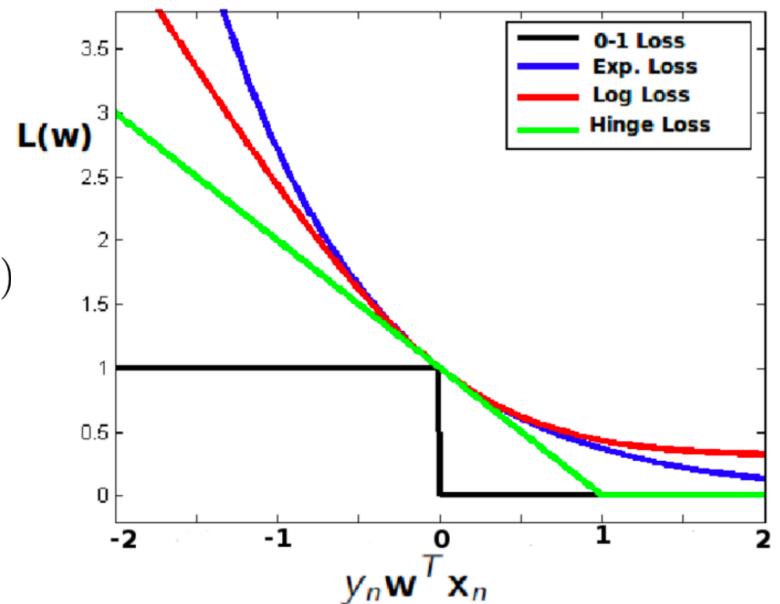
Hinge: $\ell^{(\text{hng})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$

Logistic: $\ell^{(\log)}(y, \hat{y}) = \frac{1}{\log 2} \log (1 + \exp[-y\hat{y}])$

Exponential: $\ell^{(\exp)}(y, \hat{y}) = \exp[-y\hat{y}]$

Squared: $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

- All are convex upper bounds on 0/1 loss
- Some smooth, some more sensitive to outliers



OPTIMIZATION AS MATRIX OPERATIONS

$$\hat{Y} = [Xw]_n = \sum_d X_{n,d} w_d$$

$$\underbrace{\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}}_w = \underbrace{\begin{bmatrix} \sum_d x_{1,d} w_d \\ \sum_d x_{2,d} w_d \\ \vdots \\ \sum_d x_{N,d} w_d \end{bmatrix}}_{\hat{Y}} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\hat{Y}}$$

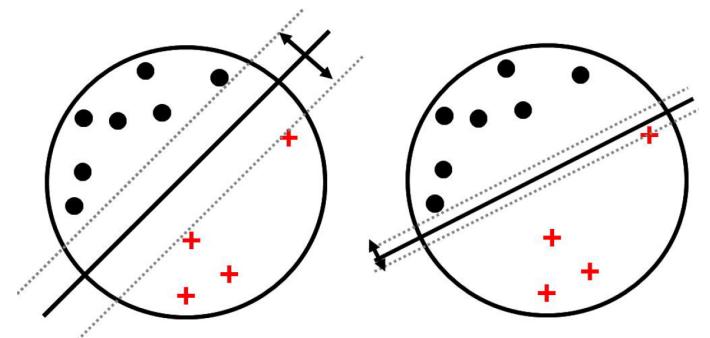


HARD-MARGIN VS SOFT-MARGIN SVM

- **Hard-margin SVM:** if *data not linearly separable*, feasible region is empty
 - No parameters can satisfy the constraints
- **Soft-margin SVM:** use slack parameters

$$\min_{w,b,\xi} \underbrace{\frac{1}{\gamma(w,b)}}_{\text{Large margin}} + C \sum_n \xi_n$$

Small slack:
Amount you need to pay, so a point is not misclassified



$$\begin{aligned} \text{Subject to } & y_n(w \cdot x_n + b) \geq 1 - \xi_n \quad (\forall n) \\ & \xi_n \geq 0 \end{aligned}$$



SVM OPTIMIZATION

- Recovering slacks when given \mathbf{w} and \mathbf{b}

$$\xi_n = \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases}$$

- SVM as unconstrained optimization problem

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n l^{(hinge)}(y_n(\mathbf{w} \cdot \mathbf{x}_n + b))$$



Large margin Small slack:

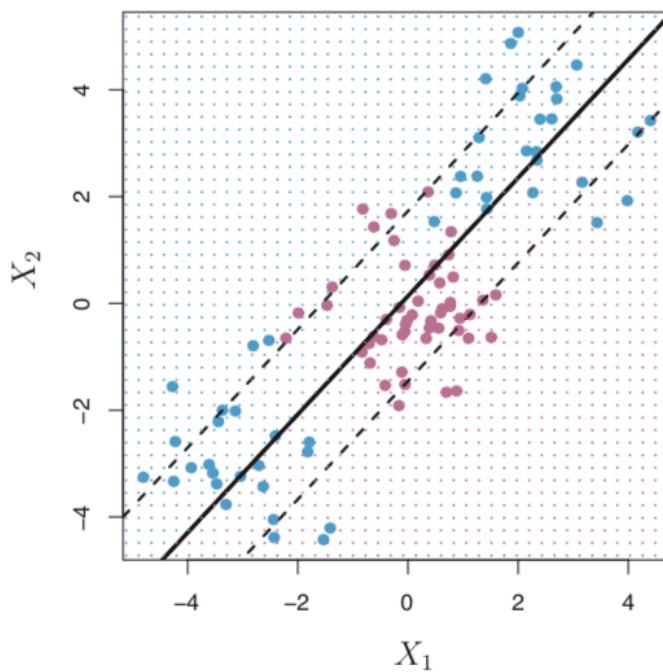
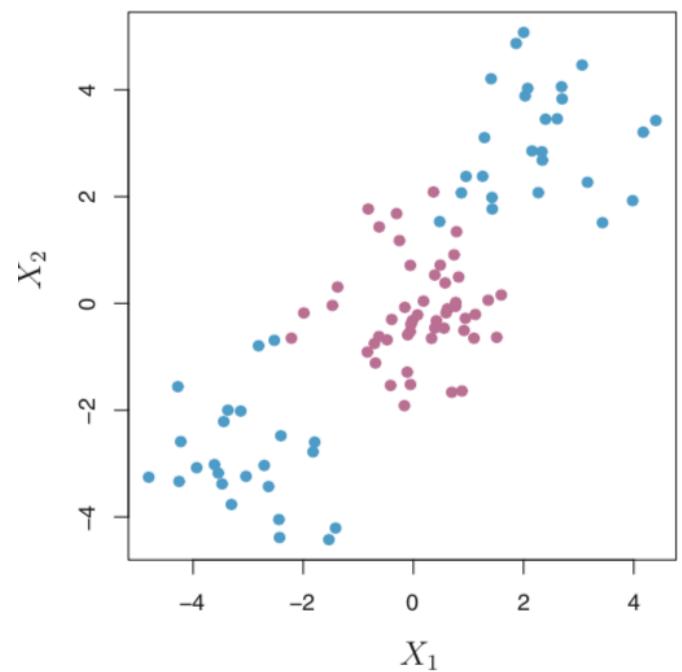


TODAY: KERNELS

- Problem: linear classifiers
 - Easy to implement and easy to optimize
 - But limited to linear decision boundaries
- What can we do about it?
 - Today's topic: Kernels



2D DATA AND A LINEAR SVM



KERNEL METHODS

- Goal: keep advantages of linear models, but make them capture non-linear patterns in data!
- How?
 - Use features of features
 - Map data to higher dimensions where it exhibits linear patterns

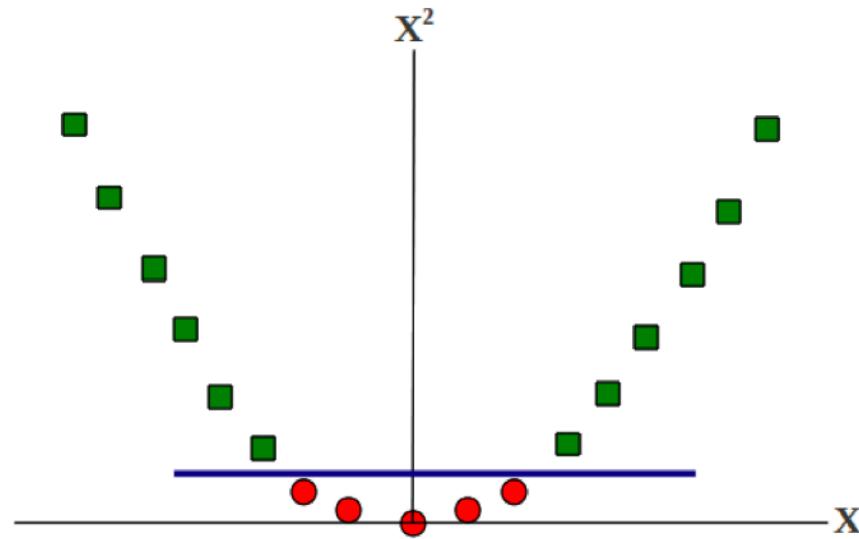


CLASSIFYING NON-LINEARLY SEPARABLE DATA WITH A LINEAR CLASSIFIER

- Example 1: Non-linearly separable data in 1D



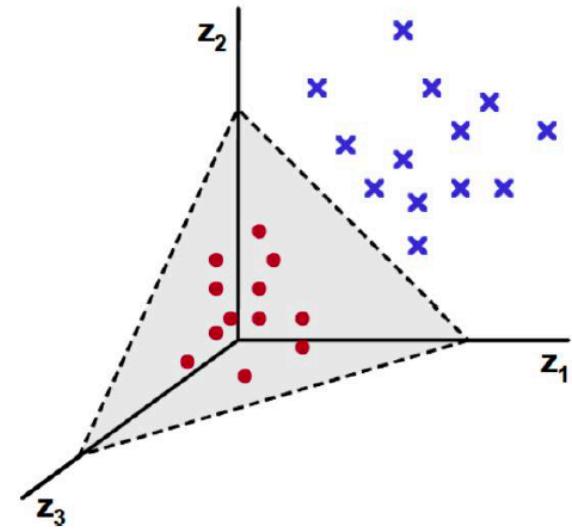
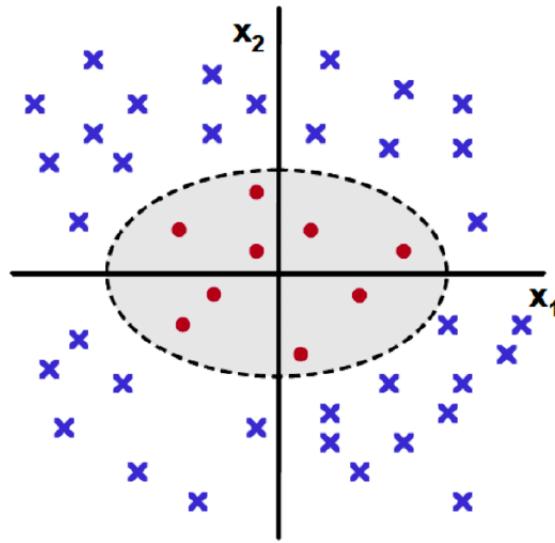
- Becomes linearly separable in new 2D space
- Defined by the following mapping: $(x) \rightarrow (x, x^2)$



CLASSIFYING NON-LINEARLY SEPARABLE DATA WITH A LINEAR CLASSIFIER

- Example 2: Non-linearly separable data in 2D
- Becomes linearly separable in 3D defined by following mapping:

$$x = (x_1, x_2) \rightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



DEFINING FEATURE MAPPINGS

- Map an original feature vector $x = (x_1, x_2, x_3, \dots, x_D)$ to an expanded version $\phi(x)$
- Example: quadratic feature mapping represents feature combinations

$$\begin{aligned}\phi(x) = & (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, \dots, \sqrt{2}x_D, \\& x_1^2, x_1x_2, x_1x_3, \dots, x_1x_D, \\& x_2x_1, x_2^2, x_2x_3, \dots, x_2x_D, \\& x_3x_1, x_3x_2, x_3^2, \dots, x_3x_D, \\& \dots, \\& x_Dx_1, x_Dx_2, x_Dx_3, \dots, x_D^2)\end{aligned}$$



FEATURE MAPPINGS

- **Pros:** can help turn non-linear classification problem into linear problem
- **Cons:** “feature explosion” creates issues when training linear classifier in new feature space
 - Computational: More computationally expensive to train
 - Statistical: More training examples needed to avoid overfitting



KERNEL METHODS

- Goal: keep advantages of linear models, but make them capture non-linear patterns in data!
- How?
 - By mapping data to higher dimensions where it exhibits linear patterns
 - **By rewriting linear models so that the mapping never needs to be explicitly computed**



THE KERNEL TRICK

1. Rewrite learning algorithms so they only depend on **dot products between two examples**
 2. Replace dot product $\phi(x) \cdot \phi(z)$
by kernel function $k(x,z)$
which computes the dot product **implicitly**
-
- **Intuition:** kernel function quantifies the similarity of two examples
 - For example, the linear kernel $k(x,z) = x \cdot z$ is equivalent to Pearson correlation



KERNEL FUNCTION EXAMPLE 1

- Consider two examples $x = (x_1, x_2)$ and $z = (z_1, z_2)$
- Let's assume we are given a **quadratic kernel k** that takes as inputs x and z

$$\begin{aligned} k(x, z) &= (x \cdot z)^2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2) \cdot (z_1^2, \sqrt{2}z_1 z_2, z_2^2) \\ &= \phi(x) \cdot \phi(z) \end{aligned}$$

- k implicitly defines a mapping ϕ to a higher dimensional space

$$\phi(x) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)$$



KERNEL FUNCTION EXAMPLE 2

$$\phi(x) = \langle 1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, \dots, \sqrt{2}x_D, \\ x_1^2, x_1x_2, x_1x_3, \dots, x_1x_D, \\ x_2x_1, x_2^2, x_2x_3, \dots, x_2x_D, \\ x_3x_1, x_3x_2, x_3^2, \dots, x_3x_D, \\ \dots, \\ x_Dx_1, x_Dx_2, x_Dx_3, \dots, x_D^2 \rangle$$

What is the function $k(x,z)$ that can implicitly compute the dot product $\phi(x) \cdot \phi(z)$?

$$\phi(x) \cdot \phi(z) = 1 + 2x_1z_1 + 2x_2z_2 + \dots + 2x_Dz_D + x_1^2z_1^2 + \dots + x_1x_Dz_1z_D + \\ \dots + x_Dx_1z_Dz_1 + x_Dx_2z_Dz_2 + \dots + x_D^2z_D^2 \quad (11.2)$$

$$= 1 + 2 \sum_d x_d z_d + \sum_d \sum_e x_d x_e z_d z_e \quad (11.3)$$

$$= 1 + 2x \cdot z + (x \cdot z)^2 \quad (11.4)$$

$$= (1 + x \cdot z)^2 \quad (11.5)$$



KERNELS: FORMALLY DEFINED

- Each kernel has an associated feature mapping ϕ
- ϕ takes input $x \in \mathcal{X}$ (input space) and maps it to \mathcal{F} (feature space)
- Kernel $k(\mathbf{x}, \mathbf{z})$ takes two inputs and gives their similarity in \mathcal{F} space

$$\begin{aligned}\phi: \mathcal{X} &\rightarrow \mathcal{F} \\ k: \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R}, \quad k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z})\end{aligned}$$

- \mathcal{F} needs to be a vector space with a dot product defined in it
 - Also called Hilbert space



MERCER'S CONDITION

- Can any function be used as a kernel function?
 - No. It must satisfy Mercer's condition
- For k to be a kernel function
 - There must exist a Hilbert Space \mathcal{F} for which k defines a dot product
 - The above is true if K is a positive semi-definite (psd) function
 - For all square integrable functions f , the following holds:

$$\iint f(x)k(x,z)f(z)dxdz$$

Non-trivial to show



CONSTRUCTING KERNEL COMBINATIONS

- If k_1 and k_2 are kernel functions, then the following are as well
 - Direct sum $k(x, z) = k_1(x, z) + k_2(x, z)$
 - Scalar product: $k(x, z) = \alpha k_1(x, z)$
 - Direct product: $k(x, z) = k_1(x, z)k_2(x, z)$
- Can use Mercer's condition to prove



COMMONLY USED KERNEL FUNCTIONS

Linear (trivial) Kernel:

$$k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z} \text{ (mapping function } \phi \text{ is identity - no mapping)}$$

Quadratic Kernel:

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z})^2 \quad \text{or} \quad (1 + \mathbf{x}^\top \mathbf{z})^2$$

Polynomial Kernel (of degree d):

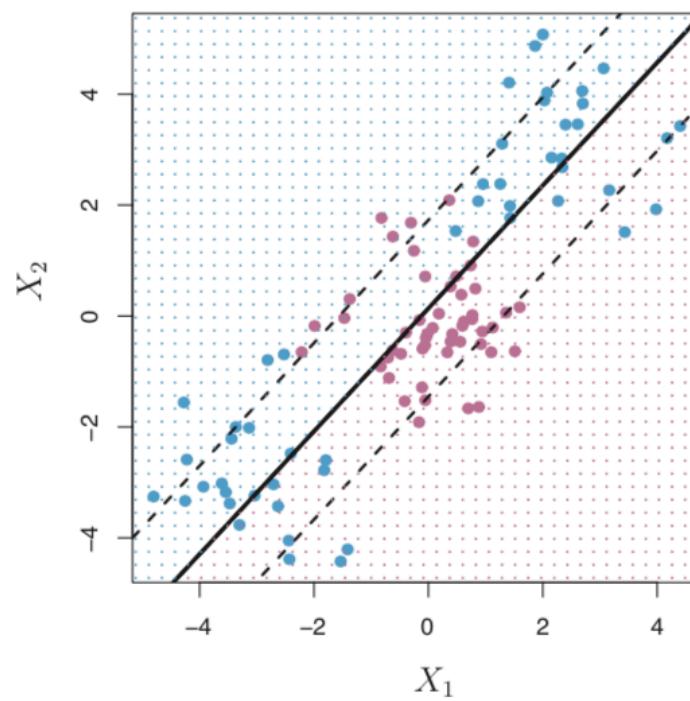
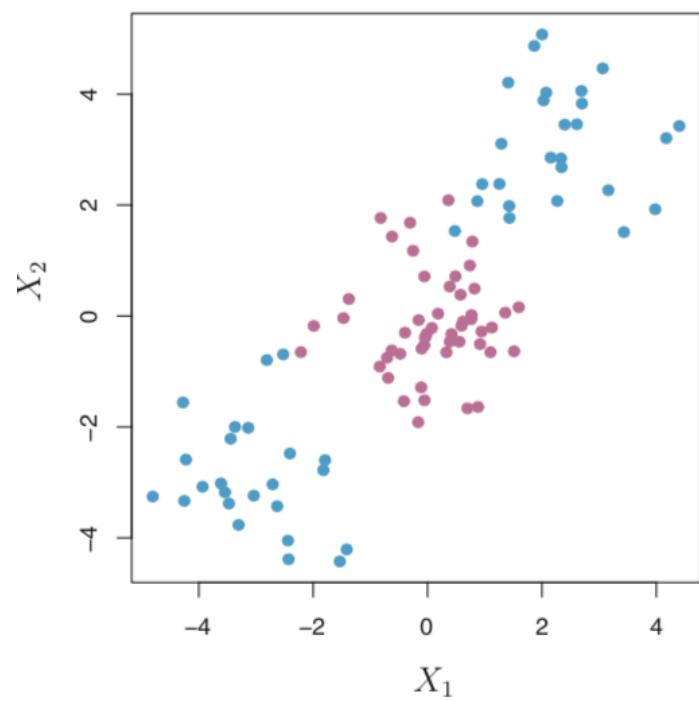
$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^\top \mathbf{z})^d \quad \text{or} \quad (1 + \mathbf{x}^\top \mathbf{z})^d$$

Radial Basis Function (RBF) Kernel:

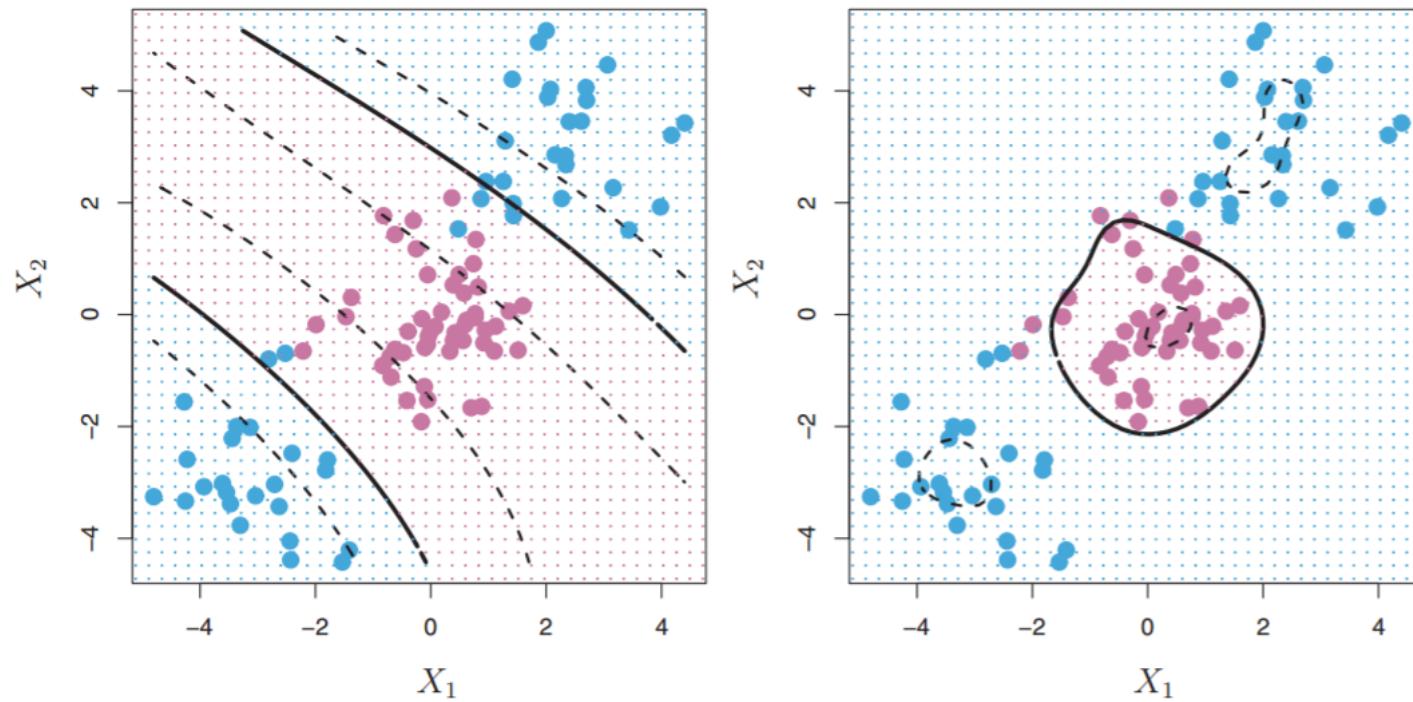
$$k(\mathbf{x}, \mathbf{z}) = \exp[-\gamma ||\mathbf{x} - \mathbf{z}||^2]$$



2D DATA AND A LINEAR KERNEL



DEGREE 3 POLYNOMIAL AND RADIAL KERNELS



$$k(\mathbf{x}, \mathbf{z}) = \exp[-\gamma ||\mathbf{x} - \mathbf{z}||^2]$$



How can you decide which kernel to use?



THE KERNEL TRICK

- Rewrite learning algorithms so they only depend on **dot products between two examples**
- Replace dot product $\phi(x)^T \phi(z)$
by kernel function $k(x,z)$
which computes the dot product **implicitly**



“KERNELIZING” THE PERCEPTRON

Algorithm 5 PERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$                                 // initialize weights
2:  $b \leftarrow 0$                                                        // initialize bias
3: for  $iter = 1 \dots MaxIter$  do
4:   for all  $(x, y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$                                      // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$                       // update weights
8:        $b \leftarrow b + y$                                                  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 
```



“KERNELIZING” THE PERCEPTRON

- Goal: Remove the explicit dependence of this algorithm on ϕ and w

Algorithm 29 PERCEPTRONTRAIN($D, MaxIter$)

```
1:  $w \leftarrow 0, b \leftarrow 0$  // initialize weights and bias
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x,y) \in D$  do
4:      $a \leftarrow w \cdot \phi(x) + b$  // compute activation for this example
5:     if  $ya \leq 0$  then
6:        $w \leftarrow w + y \phi(x)$  // update weights
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10:  end for
11:  return  $w, b$ 
```

Can we apply the kernel trick?
Not yet, we need to rewrite the algorithm
using dot products between examples



“KERNELIZING” THE PERCEPTRON

- **Perceptron Representer Theorem:** During a run of the perceptron algorithm, the weight vector \mathbf{w} can always be represented as a linear combination of the expanded training data $\phi(x_1), \phi(x_2), \dots, \phi(x_N)$ ”
- Proof by induction
 - Suppose it's true before the kth update, and the update is on example n

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

- The new weight vector is

$$\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i) + y_n \phi(\mathbf{x}_n) = \sum_i (\alpha_i + y_n[i = n]) \phi(\mathbf{x}_i)$$



“KERNELIZING” THE PERCEPTRON

- We can use the perceptron representer theorem to compute activations as a **dot product** between examples

$$\mathbf{w} \cdot \phi(\mathbf{x}) + b = \left(\sum_n \alpha_n \phi(\mathbf{x}_n) \right) \cdot \phi(\mathbf{x}) + b \quad \text{definition of } \mathbf{w}$$

(11.6)

$$= \sum_n \alpha_n [\phi(\mathbf{x}_n) \cdot \phi(\mathbf{x})] + b \quad \text{dot products are linear}$$

(11.7)



“KERNELIZING” THE PERCEPTRON

Algorithm 30 KERNELIZEDPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```
1:  $\alpha \leftarrow 0, b \leftarrow 0$  // initialize coefficients and bias
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x_n, y_n) \in \mathbf{D}$  do
4:      $a \leftarrow \sum_m \alpha_m \phi(x_m) \cdot \phi(x_n) + b$  // compute activation for this example
5:     if  $y_n a \leq 0$  then
6:        $\alpha_n \leftarrow \alpha_n + y_n$  // update coefficients
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10:  end for
11:  return  $\alpha, b$ 
```

Same training algorithm but no explicit reference
to weights
Only depends on dot products between examples

We can apply the kernel trick!



KERNEL METHODS

- Goal: keep advantages of linear models, but make them capture non-linear patterns in data
- How?
 - By mapping data to higher dimensions where it exhibits linear patterns
 - By rewriting linear models so that the (higher dimensional) mapping never needs to be explicitly computed
 - Making the representation sparser and the computation more efficient!



DISCUSSION

- Other algorithms can be “kernelized”:
 - We’ll talk about kernels and SVM next
- Do Kernels address all the downsides of “feature explosion”?
 - Helps reduce computation cost during training
 - But overfitting remains an issue



SUMMARY

- Kernel functions
 - What they are, why they are useful, how they relate to feature combinations
- Kernelized perceptron
 - You should be able to derive it and implement it



ANNOUNCEMENTS

- Homework 3
 - Released by the end of this week
- My office hours next week
 - Wednesday 2-4pm – cancelled (on travel)
 - Available Monday 3-5pm or by appointment



ACKNOWLEDGEMENTS

- These slides use materials by Marine Carpuat and Piyush Rai

