# CSCI-567: Machine Learning (Fall 2019)

Prof. Victor Adamchik

U of Southern California

Oct. 01, 2019

# Outline

# Outline

# Math formulation

An L-layer neural net can be written as

$$f(x) = h_L \left( W^L h_{L-1} \left( W^{L-1} \cdots h_1 \left( W^1 x \right) \right) \right)$$



To ease notation, for a given input $x$, define recursively

$$o^0 = x, \qquad a^\ell = W^\ell o^{\ell-1}, \qquad o^\ell = h_\ell(a^\ell) \qquad (\ell = 1, \ldots, L)$$

where

- $D_0 = D, D_1, \ldots, D_L$ are numbers of neurons at each layer
- $W^\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$ is the weights for layer $\ell$
- $a^\ell \in \mathbb{R}^{D_\ell}$ is input to layer $\ell$
- $o^\ell \in \mathbb{R}^{D_\ell}$ is output to layer $\ell$
- $h_\ell : \mathbb{R}^{D_\ell} \to \mathbb{R}^{D_\ell}$ is activation functions at layer $\ell$

# Putting everything into SGD

The backpropagation algorithm (1986) computes the gradient of the cost function for a single training example.

Initialize $\boldsymbol{W}^1, \ldots, \boldsymbol{W}^{\mathsf{L}}$ (all $\boldsymbol{0}$ or randomly). Repeat:

1. randomly pick one data point $n \in [\mathsf{N}]$
2. **forward propagation**: for each layer $\ell = 1, \ldots, \mathsf{L}$
   - compute $\boldsymbol{a}^\ell = \boldsymbol{W}^\ell \boldsymbol{o}^{\ell-1}$ and $\boldsymbol{o}^\ell = \boldsymbol{h}_\ell(\boldsymbol{a}^\ell)$ $\hspace{2em}$ $(\boldsymbol{o}^0 = \boldsymbol{x}_n)$
3. **backward propagation**: for each $\ell = L, \ldots, 1$

   - compute
   $$\frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}^\ell}$$

   - update weights
   $$\boldsymbol{W}^\ell \leftarrow \boldsymbol{W}^\ell - \lambda \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{W}^\ell} = \boldsymbol{W}^\ell - \lambda \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}^\ell} (\boldsymbol{o}^{\ell-1})^{\mathrm{T}}$$

# Overfitting

**Overfitting is very likely** since the models are too powerful.

Methods to overcome overfitting:
- data augmentation
- regularization
- dropout
- early stopping
- $\cdots$

# Conclusions for neural nets

Deep neural networks
- are hugely popular, achieving *best performance* on many problems
- do need *a lot of data* to work well
- take *a lot of time* to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters
- are still not well understood in theory

# Outline

1. Review of last lecture

2. NN discussion problems

3. Kernel methods

4. Kernel discussion problems

## Outline

## Motivation

Recall the question: *how to choose nonlinear basis $\phi : \mathbb{R}^D \to \mathbb{R}^M$?*

$$\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x})$$

- neural network is one approach: learn $\phi$ from data

- **kernel method** is another one: sidestep the issue of choosing $\phi$ by using *kernel functions*

## Case study: regularized linear regression

Recall the regularized least square solution:

$$\boldsymbol{w}^* = \left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X} + \lambda\boldsymbol{I}\right)^{-1}\boldsymbol{X}^{\mathrm{T}}\boldsymbol{y}$$
$$\boldsymbol{w}^* = \left(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}\right)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{y}$$

$$\boldsymbol{X} = \begin{pmatrix} \boldsymbol{x}_1^{\mathrm{T}} \\ \boldsymbol{x}_2^{\mathrm{T}} \\ \vdots \\ \boldsymbol{x}_{\mathsf{N}}^{\mathrm{T}} \end{pmatrix}, \boldsymbol{\Phi} = \begin{pmatrix} \boldsymbol{\phi}(\boldsymbol{x}_1)^{\mathrm{T}} \\ \boldsymbol{\phi}(\boldsymbol{x}_2)^{\mathrm{T}} \\ \vdots \\ \boldsymbol{\phi}(\boldsymbol{x}_{\mathsf{N}})^{\mathrm{T}} \end{pmatrix}$$

Here $\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}$ is D × D matrix, and $\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}$ is M × M matrix,

Issue: *M could be huge or even infinity!*

We will rewrite the solution in a different form.

## Another solution

**Another minimizer is** (we will prove it later on slide 16)

$$\boldsymbol{w}^* = \boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$$
$$\boldsymbol{w}^* = \boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{K} + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$$
$$\boldsymbol{w}^* = \boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\alpha} = \sum_{n=1}^{N} \alpha_n \boldsymbol{\phi}(\boldsymbol{x}_n)$$

where $\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} \in \mathbb{R}^{\mathsf{N}\times\mathsf{N}}$ is the Gram/Kernel matrix

and $\boldsymbol{\alpha}$ is a new vector.

Solution $\boldsymbol{w}^*$ is a linear combination of features!

## Gram matrix

We call $\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}}$ **Gram matrix** or **kernel matrix** where the $(i,j)$ entry is

$$\phi(\boldsymbol{x}_i)^{\mathrm{T}}\phi(\boldsymbol{x}_j)$$

Therefore,

$$\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}}$$

$$= \begin{pmatrix} \phi(\boldsymbol{x}_1)^{\mathrm{T}}\phi(\boldsymbol{x}_1) & \phi(\boldsymbol{x}_1)^{\mathrm{T}}\phi(\boldsymbol{x}_2) & \cdots & \phi(\boldsymbol{x}_1)^{\mathrm{T}}\phi(\boldsymbol{x}_N) \\ \phi(\boldsymbol{x}_2)^{\mathrm{T}}\phi(\boldsymbol{x}_1) & \phi(\boldsymbol{x}_2)^{\mathrm{T}}\phi(\boldsymbol{x}_2) & \cdots & \phi(\boldsymbol{x}_2)^{\mathrm{T}}\phi(\boldsymbol{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\boldsymbol{x}_N)^{\mathrm{T}}\phi(\boldsymbol{x}_1) & \phi(\boldsymbol{x}_N)^{\mathrm{T}}\phi(\boldsymbol{x}_2) & \cdots & \phi(\boldsymbol{x}_N)^{\mathrm{T}}\phi(\boldsymbol{x}_N) \end{pmatrix} \in \mathbb{R}^{\mathsf{N}\times\mathsf{N}}$$

## Examples of kernel matrix

3 data points in $\mathbb{R}$

$$x_1 = -1, x_2 = 0, x_3 = 1$$

$\phi$ is polynomial basis with degree 4:

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ x^3 \end{pmatrix}$$

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

## Calculation of the Gram matrix

$$\phi(x_1) = \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \quad \phi(x_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad \phi(x_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

**Gram/Kernel matrix**

$$\boldsymbol{K} = \begin{pmatrix} \phi(x_1)^{\mathrm{T}}\phi(x_1) & \phi(x_1)^{\mathrm{T}}\phi(x_2) & \phi(x_1)^{\mathrm{T}}\phi(x_3) \\ \phi(x_2)^{\mathrm{T}}\phi(x_1) & \phi(x_2)^{\mathrm{T}}\phi(x_2) & \phi(x_2)^{\mathrm{T}}\phi(x_3) \\ \phi(x_3)^{\mathrm{T}}\phi(x_1) & \phi(x_3)^{\mathrm{T}}\phi(x_2) & \phi(x_3)^{\mathrm{T}}\phi(x_3) \end{pmatrix}$$

$$= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix}$$

## Another solution

Here we prove that two solutions

$$\boldsymbol{w}^* = (\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I})^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{y}$$
$$\boldsymbol{w}^* = \boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$$

are the same.

$$\begin{aligned} &(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}_1)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{y} \\ &= (\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}_1)^{-1}\boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I}_2)(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I}_2)^{-1}\boldsymbol{y} \\ &= (\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}_1)^{-1}(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{\Phi}^{\mathrm{T}})(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I}_2)^{-1}\boldsymbol{y} \\ &= (\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}_1)^{-1}(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi} + \lambda\boldsymbol{I}_1)\boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I}_2)^{-1}\boldsymbol{y} \\ &= \boldsymbol{\Phi}^{\mathrm{T}}(\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} + \lambda\boldsymbol{I}_2)^{-1}\boldsymbol{y} \end{aligned}$$

## Then what is the difference?

First, computing $(\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}} + \lambda \mathbf{I})^{-1} = (\mathbf{K} + \lambda \mathbf{I})^{-1} = \boldsymbol{\alpha}$ can be more efficient than computing $(\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi} + \lambda \mathbf{I})^{-1}$ when N ≤ M.

More importantly, computing $(\mathbf{K} + \lambda \mathbf{I})^{-1}$ *only requires computing inner products in the new feature space!*

Now we can conclude that the exact form of $\phi(\cdot)$ is not essential; *all we need is computing inner products $\phi(\boldsymbol{x})^{T}\phi(\boldsymbol{x}')$.*

For some $\phi$ it is indeed possible to compute $\phi(\boldsymbol{x})^{\mathrm{T}}\phi(\boldsymbol{x}')$ without computing/knowing $\phi$. This is the *kernel trick*.

## Why is this helpful?

The prediction of $\boldsymbol{w}^*$ on a new example $\boldsymbol{x}$ is

$$\boldsymbol{w}^{*\mathrm{T}}\phi(\boldsymbol{x}) = \left(\sum_{n=1}^{N} \alpha_n \phi(\boldsymbol{x}_n)^{\mathrm{T}}\right)\phi(\boldsymbol{x}) = \sum_{n=1}^{N} \alpha_n \left(\phi(\boldsymbol{x}_n)^{\mathrm{T}}\phi(\boldsymbol{x})\right)$$

Therefore we do not really need to know a nonlinear mapping $\phi$, only inner products in the new feature space matter!

*Kernel methods* are exactly about computing inner products *without knowing $\phi$.*

## Example of the kernel trick

Consider the following polynomial basis $\phi : \mathbb{R}^2 \to \mathbb{R}^3$:

$$\phi(\boldsymbol{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \end{pmatrix}$$

What is the inner product between $\phi(\boldsymbol{x})$ and $\phi(\boldsymbol{x}')$?

$$\phi(\boldsymbol{x})^{\mathrm{T}}\phi(\boldsymbol{x}') = x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2$$
$$= (x_1 x_1' + x_2 x_2')^2 = (\boldsymbol{x}^{\mathrm{T}}\boldsymbol{x}')^2 = k(\boldsymbol{x}, \boldsymbol{x}')$$

Therefore, the inner product in the new space is simply a function of the inner product in the original space.

Count the number of multiplications.

## Another example

$\phi : \mathbb{R}^{\mathrm{D}} \to \mathbb{R}^{2\mathrm{D}}$ is parameterized by $\theta$:

$$\phi_\theta(\boldsymbol{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_{\mathrm{D}}) \\ \sin(\theta x_{\mathrm{D}}) \end{pmatrix}$$

What is the inner product between $\phi_\theta(\boldsymbol{x})$ and $\phi_\theta(\boldsymbol{x}')$?

$$\phi_\theta(\boldsymbol{x})^{\mathrm{T}}\phi_\theta(\boldsymbol{x}') = \sum_{d=1}^{\mathrm{D}} \cos(\theta x_d)\cos(\theta x_d') + \sin(\theta x_d)\sin(\theta x_d')$$
$$= \sum_{d=1}^{\mathrm{D}} \cos(\theta(x_d - x_d')) = k(\boldsymbol{x}, \boldsymbol{x}')$$

Once again, *the inner product in the new space is a simple function of the features in the original space.*

## More complicated example

Based on the previous example mapping $\phi_\theta$, we define a new one
$\phi_L : \mathbb{R}^D \to \mathbb{R}^{2D(L+1)}$ as follows:

$$\phi_L(\boldsymbol{x}) = \begin{pmatrix} \phi_0(\boldsymbol{x}) \\ \phi_{\frac{2\pi}{L}}(\boldsymbol{x}) \\ \phi_{2\frac{2\pi}{L}}(\boldsymbol{x}) \\ \vdots \\ \phi_{L\frac{2\pi}{L}}(\boldsymbol{x}) \end{pmatrix}$$

What is the inner product between $\phi_L(\boldsymbol{x})$ and $\phi_L(\boldsymbol{x}')$?

$$\phi_L(\boldsymbol{x})^{\mathrm{T}}\phi_L(\boldsymbol{x}') = \sum_{\ell=0}^{L} \phi_{\frac{2\pi\ell}{L}}(\boldsymbol{x})^{\mathrm{T}}\phi_{\frac{2\pi\ell}{L}}(\boldsymbol{x}')$$

$$= \sum_{\ell=0}^{L}\sum_{d=1}^{D} \cos\left(\frac{2\pi\ell}{L}(x_d - x_d')\right)$$

## Infinite dimensional mapping

Let us set $L \to \infty$. This means that $\phi_L(\boldsymbol{x})$ vector has infinite dimension. Clearly we cannot compute $\phi_L(\boldsymbol{x})$, but we can still compute the inner

product:

$$\phi_\infty(\boldsymbol{x})^{\mathrm{T}}\phi_\infty(\boldsymbol{x}') = \int_0^{2\pi} \sum_{d=1}^{D} \cos(\theta(x_d - x_d'))\, d\theta$$

$$= \sum_{d=1}^{D} \frac{\sin(2\pi(x_d - x_d'))}{x_d - x_d'}$$

Again, a simple function of the original features.

Note that using this mapping in linear regression, we are *learning a weight* $\boldsymbol{w}^*$ *with infinite dimension!*

## Kernel functions

**Definition**: a function $k : \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}$ is called a *(positive semidefinite) kernel function* if there exists a function $\phi : \mathbb{R}^D \to \mathbb{R}^M$ so that for any $\boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^D$,

$$k(\boldsymbol{x}, \boldsymbol{x}') = \phi(\boldsymbol{x})^{\mathrm{T}}\phi(\boldsymbol{x}')$$

Kernel functions are used to quantify similarity between a pair of points $\boldsymbol{x}$ and $\boldsymbol{x}'$.

Examples we have seen

$$k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^{\mathrm{T}}\boldsymbol{x}'$$
$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\mathrm{T}}\boldsymbol{x}')^2$$
$$k(\boldsymbol{x}, \boldsymbol{x}') = \sum_{d=1}^{D} \frac{\sin(2\pi(x_d - x_d'))}{x_d - x_d'}$$

## Using kernel functions

Choosing a nonlinear basis $\phi$ becomes choosing a kernel function.

As long as computing the kernel function is more efficient, we should apply the kernel trick.

**Gram/kernel matrix** becomes:

$$\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} = \begin{pmatrix} k(\boldsymbol{x}_1, \boldsymbol{x}_1) & k(\boldsymbol{x}_1, \boldsymbol{x}_2) & \cdots & k(\boldsymbol{x}_1, \boldsymbol{x}_N) \\ k(\boldsymbol{x}_2, \boldsymbol{x}_1) & k(\boldsymbol{x}_2, \boldsymbol{x}_2) & \cdots & k(\boldsymbol{x}_2, \boldsymbol{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\boldsymbol{x}_N, \boldsymbol{x}_1) & k(\boldsymbol{x}_N, \boldsymbol{x}_2) & \cdots & k(\boldsymbol{x}_N, \boldsymbol{x}_N) \end{pmatrix}$$

In fact, $k$ is a kernel if and only if $\boldsymbol{K}$ is positive semidefinite for *any $\boldsymbol{x}_1$, $\boldsymbol{x}_2$, ..., $\boldsymbol{x}_N$* (**Mercer theorem**).

## Using kernel functions

The prediction on a new example $\boldsymbol{x}$ is

$$\boldsymbol{w}^{*\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}) = \left(\sum_{n=1}^{N}\alpha_n\boldsymbol{\phi}(\boldsymbol{x}_n)^{\mathrm{T}}\right)\boldsymbol{\phi}(\boldsymbol{x}) = \sum_{n=1}^{N}\alpha_n\left(\boldsymbol{\phi}(\boldsymbol{x}_n)^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x})\right)$$
$$= \sum_{n=1}^{N}\alpha_n k(\boldsymbol{x}_n, \boldsymbol{x})$$

## More examples of kernel functions

Two most commonly used kernel functions in practice:

**Polynomial kernel**

$$k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\mathrm{T}}\boldsymbol{x}' + c)^d$$

for $c \geq 0$ and $d$ is a positive integer.

**Gaussian kernel or Radial basis function (RBF) kernel**

$$k(\boldsymbol{x}, \boldsymbol{x}') = e^{-\frac{\|\boldsymbol{x}-\boldsymbol{x}'\|_2^2}{2\sigma^2}}$$

for some $\sigma > 0$.

Think about *what the corresponding $\phi$ is* for each kernel.

## Composing kernels

Creating more kernel functions using the following rules:

If $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are kernels, the followings are kernels too

- **linear combination**: $\alpha k_1(\cdot, \cdot) + \beta k_2(\cdot, \cdot)$ if $\alpha, \beta \geq 0$

- **product**: $k_1(\cdot, \cdot)k_2(\cdot, \cdot)$

- **exponential**: $e^{k(\cdot, \cdot)}$

- $\cdots$

Verify using the definition of kernel!

## Kernelizing ML algorithms

There are two **main aspects** of kernelized algorithms:

- the solution is expressed as a linear combination of training examples
- algorithm relies only on inner products between data points

Kernel trick is applicable to **many** ML algorithms:

- nearest neighbor classifier
- perceptron
- logistic regression
- $\cdots$

## Kernelizing NNC

For NNC with **L2 distance**, $\|x - x'\|_2^2$ is not a valid kernel.

But we can covert the norm

$$d(x, x') = \|x - x'\|_2^2 = x^{\mathrm{T}}x + x'^{\mathrm{T}}x' - 2x^{\mathrm{T}}x'$$

into the following kernel function

$$d^{\mathrm{KERNEL}}(x, x') = k(x, x) + k(x', x') - 2k(x, x')$$

which by definition is the **L2 distance in a new feature space**

$$d^{\mathrm{KERNEL}}(x, x') = \|\phi(x) - \phi(x')\|_2^2$$

## Kernelizing NNC

Regular KNN algorithm has two shortcomings.

- all neighbors receive equal weight
- the number of neighbors must be chosen globally.

Kernel addresses these issues.

Instead of selected nearest neighbors, all neighbors are used, but with different weights.

Closer neighbors receive higher weight.

The weighting function is a kernel.

One of the most common way is the Gaussian kernel

$$k(x, x') = e^{-\frac{\|x - x'\|_2^2}{2\sigma^2}}$$

## Kernelizing NNC

**Kernel Binary Classification Algorithm.**

Given

- training data $\mathcal{D} = \{(x_n, y_n), n = 1, 2, \dots, \mathsf{N}\}$, where $y_n \in \{-1, 1\}$
- kernel function $k(x_i, x_j)$
- input $x$ to classify

Return the class given by

$$\mathrm{sign}\left(\sum_{n=1}^{\mathsf{N}} k(x, x_n)y_n\right)$$

## Kernelizing Perceptron

**Perceptron Algorithm:**

- Pick $(x_n, y_n)$ randomly
- Compute $y^* = \mathrm{sign}(w^{\mathrm{T}}x_n)$
- If $y^* \neq y_n$ then
  - $w = w + y_n x_n$
- Solution $w$ is a linear combination of features.

- Pick $(x_n, y_n)$ randomly
- Compute $y^* = \mathrm{sign}(w^{\mathrm{T}}\phi(x_n))$
- If $y^* \neq y_n$ then
  - $w = w + y_n\phi(x_n)$
  - $\alpha_n = \alpha_n + y_n$
- Solution $w = \sum_n \alpha_n\phi(x_n)$ is a linear combination of features.

# Kernelizing Perceptron

## Kernelized Perceptron Algorithm:

- Pick $(\boldsymbol{x}_n, y_n)$ randomly
- Compute $y^* = \text{sign}(\boldsymbol{w}^{\text{T}}\boldsymbol{\phi}(\boldsymbol{x}_n))$
- If $y^* \neq y_n$ then
  - $\boldsymbol{w} = \boldsymbol{w} + y_n\boldsymbol{\phi}(\boldsymbol{x}_n)$
  - $\alpha_n = \alpha_n + y_n$
- Solution $\boldsymbol{w} = \sum_n \alpha_n\boldsymbol{\phi}(\boldsymbol{x}_n)$ is a linear combination of features.

- Pick $(\boldsymbol{x}_n, y_n)$ randomly
- Compute
  $y^* = \text{sign}\left(\sum_i \alpha_i k(\boldsymbol{x}_i, \boldsymbol{x}_n)\right)$
- If $y^* \neq y_n$ then
  - $\alpha_n = \alpha_n + y_n$
- Prediction $\sum_n \alpha_n k(\boldsymbol{x}_n, \boldsymbol{x})$.

Note, in the kernelized algorithm we do not need to compute weights $\boldsymbol{w}$.

---

# Kernelizing Perceptron

## XOR example

Kernel function: $k(\boldsymbol{x}, \boldsymbol{x}') = (\boldsymbol{x}^{\text{T}}\boldsymbol{x}')^2$

Four 2-dimensional training points:       Gram matrix $\boldsymbol{K}$:

$$\begin{pmatrix} x_1 & x_2 & y \\ 1 & 1 & 1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \\ 1 & -1 & -1 \end{pmatrix} \qquad \begin{pmatrix} 4 & 0 & 4 & 0 \\ 0 & 4 & 0 & 4 \\ 4 & 0 & 4 & 0 \\ 0 & 4 & 0 & 4 \end{pmatrix}$$

$k(\boldsymbol{x}_1, \boldsymbol{x}_1) = \left((1,1)^{\text{T}}(1,1)\right)^2 = 4 \qquad k(\boldsymbol{x}_1, \boldsymbol{x}_2) = \left((1,1)^{\text{T}}(-1,1)\right)^2 = 0$

$k(\boldsymbol{x}_1, \boldsymbol{x}_3) = \left((1,1)^{\text{T}}(-1,-1)\right)^2 = 4 \quad k(\boldsymbol{x}_1, \boldsymbol{x}_4) = \left((1,1)^{\text{T}}(1,-1)\right)^2 = 0$

---

# Kernelizing Perceptron

## XOR example

Initialization: $\boldsymbol{\alpha} = (0, 0, 0, 0)$.

First round

$\boldsymbol{x}_1$: compute $y^* = \text{sign}\left(\sum_{i=1}^4 \alpha_i k(\boldsymbol{x}_i, \boldsymbol{x}_1)\right) = \text{sign}(0) = -1 \neq y_1$
Thus, $\alpha_1 = y_1 = 1$.

$\boldsymbol{x}_2$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(0,4,0,4)) = \text{sign}(0) = -1 = y_2$

$\boldsymbol{x}_3$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(4,0,4,0)) = \text{sign}(4) = 1 = y_3$

$\boldsymbol{x}_4$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(0,4,0,4)) = \text{sign}(0) = -1 = y_4$

---

# Kernelizing Perceptron

## XOR example

$\boldsymbol{\alpha} = (1, 0, 0, 0)$.

Second round.

$\boldsymbol{x}_1$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(4,0,4,0)) = \text{sign}(4) = 1 = y_1$

$\boldsymbol{x}_2$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(0,4,0,4)) = \text{sign}(0) = -1 = y_2$

$\boldsymbol{x}_3$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(4,0,4,0)) = \text{sign}(4) = 1 = y_3$

$\boldsymbol{x}_4$: compute $y^* = \text{sign}((1,0,0,0)^{\text{T}}(0,4,0,4)) = \text{sign}(0) = -1 = y_4$

Converged! The prediction on a new example $\boldsymbol{x}$ is

$$\sum_{i=1}^4 \alpha_i k(\boldsymbol{x}_i, \boldsymbol{x}) = k(\boldsymbol{x}_1, \boldsymbol{x}) = \left((1,1)^{\text{T}}\boldsymbol{x}\right)^2 = (x_1 + x_2)^2$$

# Outline