

CSCI567 Machine Learning (Spring 2019)

Prof. Victor Adamchik

U of Southern California

Feb. 19, 2019

()

February 19, 2019 1 / 52

Midterm Exam

Instructions:

- This is a closed-book exam.
- Questions should be answered concisely.
- Write legibly, avoid cursive writings.
- Partial credit may be available.
- Use only scratch paper provided in the exam packet.

()

February 19, 2019 3 / 52

Administration

Midterm Exam

- time: Feb. 27 (Wednesday) from 5 pm to 8 pm
- location: THH 101 and THH 201
- you can start preparing after this lecture
- practice exam will be posted
- TA's Review: next Tuesday during the lecture time

()

February 19, 2019 2 / 52

Outline

1 Kernel methods

2 Clustering

()

February 19, 2019 4 / 52

1 Kernel methods

- Dual formulation of linear regression
- Kernel Trick
- Kernelizing ML algorithms

2 Clustering

Case study: regularized linear regression

Recall the regularized least square solution:

$$\begin{aligned} \mathbf{w}^* &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ \mathbf{w}^* &= (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y} \end{aligned} \quad \left| \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}, \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \right.$$

Here $\mathbf{X}^T \mathbf{X}$ is $D \times D$ matrix, and $\Phi^T \Phi$ is $M \times M$ matrix,

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

We will rewrite the solution in a different form.

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

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

- neural network is one approach: learn ϕ from data
- **kernel method** is another one: sidestep the issue of choosing ϕ by using *kernel functions*

Another solution

Another minimizer is

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

where $\mathbf{K} = \Phi \Phi^T \in \mathbb{R}^{N \times N}$ is the Gram/Kernel matrix and $\boldsymbol{\alpha}$ is a new vector.

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

Gram matrix

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

$$\phi(x_i)^T \phi(x_j)$$

Therefore,

$$K = \Phi\Phi^T = \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \cdots & \phi(x_1)^T \phi(x_N) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \cdots & \phi(x_2)^T \phi(x_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_N)^T \phi(x_1) & \phi(x_N)^T \phi(x_2) & \cdots & \phi(x_N)^T \phi(x_N) \end{pmatrix} \in \mathbb{R}^{N \times N}$$

()

February 19, 2019 9 / 52

Another solution

Here we prove that two solutions

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

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

are the same.

$$\begin{aligned} & (\Phi^T \Phi + \lambda I_1)^{-1} \Phi^T y \\ &= (\Phi^T \Phi + \lambda I_1)^{-1} \Phi^T (\Phi \Phi^T + \lambda I_2) (\Phi \Phi^T + \lambda I_2)^{-1} y \\ &= (\Phi^T \Phi + \lambda I_1)^{-1} (\Phi^T \Phi \Phi^T + \lambda \Phi^T) (\Phi \Phi^T + \lambda I_2)^{-1} y \\ &= (\Phi^T \Phi + \lambda I_1)^{-1} (\Phi^T \Phi + \lambda I_1) \Phi^T (\Phi \Phi^T + \lambda I_2)^{-1} y \\ &= \Phi^T (\Phi \Phi^T + \lambda I_2)^{-1} y \end{aligned}$$

()

February 19, 2019 10 / 52

Then what is the difference?

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

More importantly, computing $(K + \lambda I)^{-1}$ *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 $\phi(x)^T \phi(x')$.*

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

()

February 19, 2019 11 / 52

Why is this helpful?

The prediction of w^* on a new example x is

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

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

Kernel methods are exactly about computing inner products *without knowing ϕ* .

()

February 19, 2019 12 / 52

Examples of kernel matrix

3 data points in \mathbb{R}

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

ϕ 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

$$\begin{aligned} \mathbf{K} &= \begin{pmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \phi(x_1)^T \phi(x_3) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \phi(x_2)^T \phi(x_3) \\ \phi(x_3)^T \phi(x_1) & \phi(x_3)^T \phi(x_2) & \phi(x_3)^T \phi(x_3) \end{pmatrix} \\ &= \begin{pmatrix} 4 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 4 \end{pmatrix} \end{aligned}$$

()

February 19, 2019 13 / 52

()

February 19, 2019 14 / 52

Example of the kernel trick

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

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

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

$$\begin{aligned} \phi(\mathbf{x})^T \phi(\mathbf{x}') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x}^T \mathbf{x}')^2 = k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

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

$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{2D}$ is parameterized by θ :

$$\phi_\theta(\mathbf{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \vdots \\ \cos(\theta x_D) \\ \sin(\theta x_D) \end{pmatrix}$$

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

$$\begin{aligned} \phi_\theta(\mathbf{x})^T \phi_\theta(\mathbf{x}') &= \sum_{d=1}^D \cos(\theta x_d) \cos(\theta x'_d) + \sin(\theta x_d) \sin(\theta x'_d) \\ &= \sum_{d=1}^D \cos(\theta(x_d - x'_d)) = k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

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

()

February 19, 2019 15 / 52

()

February 19, 2019 16 / 52

More complicated example

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

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

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

$$\begin{aligned} \phi_L(\mathbf{x})^T \phi_L(\mathbf{x}') &= \sum_{\ell=0}^L \phi_{\frac{2\pi\ell}{L}}(\mathbf{x})^T \phi_{\frac{2\pi\ell}{L}}(\mathbf{x}') \\ &= \sum_{\ell=0}^L \sum_{d=1}^D \cos\left(\frac{2\pi\ell}{L}(x_d - x'_d)\right) \end{aligned}$$

()

February 19, 2019 17 / 52

Infinite dimensional mapping

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

product:

$$\begin{aligned} \phi_\infty(\mathbf{x})^T \phi_\infty(\mathbf{x}') &= \int_0^{2\pi} \sum_{d=1}^D \cos(\theta(x_d - x'_d)) d\theta \\ &= \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d} \end{aligned}$$

Again, a simple function of the original features.

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

()

February 19, 2019 18 / 52

Kernel functions

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

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

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

Examples we have seen

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$$

$$k(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \frac{\sin(2\pi(x_d - x'_d))}{x_d - x'_d}$$

()

February 19, 2019 19 / 52

Using kernel functions

Choosing a nonlinear basis ϕ 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:

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

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

()

February 19, 2019 20 / 52

Using kernel functions

The prediction on a new example \mathbf{x} is

$$\begin{aligned} \mathbf{w}^{*\top} \phi(\mathbf{x}) &= \left(\sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^\top \right) \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n (\phi(\mathbf{x}_n)^\top \phi(\mathbf{x})) \\ &= \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}) \end{aligned}$$

()

February 19, 2019 21 / 52

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

Verify using the definition of kernel!

()

February 19, 2019 23 / 52

More examples of kernel functions

Two most commonly used kernel functions in practice:

Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^d$$

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

Gaussian kernel or Radial basis function (RBF) kernel

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

for some $\sigma > 0$.

Think about *what the corresponding ϕ is* for each kernel.

()

February 19, 2019 22 / 52

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
- ...

()

February 19, 2019 24 / 52

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

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

()

February 19, 2019 25 / 52

Kernelizing NNC

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

But we can convert the norm

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

into the following kernel function

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

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

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

()

February 19, 2019 27 / 52

Kernelizing NNC

Kernel Binary Classification Algorithm.

Given

- training data $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$, where $y_n \in \{-1, 1\}$
- kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$
- input \mathbf{x} to classify

Return the class given by

$$\text{sign} \left(\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) y_n \right)$$

()

February 19, 2019 26 / 52

Kernelizing Perceptron

Perceptron Algorithm:

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

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

()

February 19, 2019 28 / 52

Kernelizing Perceptron

Kernelized Perceptron Algorithm:

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

()

February 19, 2019 29 / 52

Kernelizing Perceptron

XOR example

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

Four 2-dimensional training points:

Gram matrix \mathbf{K} :

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

$$k(\mathbf{x}_1, \mathbf{x}_1) = ((1, 1)^T (1, 1))^2 = 4 \quad k(\mathbf{x}_1, \mathbf{x}_2) = ((1, 1)^T (-1, 1))^2 = 0$$

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

()

February 19, 2019 30 / 52

Kernelizing Perceptron

XOR example

Initialization: $\alpha = (0, 0, 0, 0)$.

First round

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

Thus, $\alpha_1 = y_1 = 1$.

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

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

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

()

February 19, 2019 31 / 52

Kernelizing Perceptron

XOR example

$\alpha = (1, 0, 0, 0)$.

Second round.

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

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

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

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

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

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

()

February 19, 2019 32 / 52

1 Kernel methods

2 Clustering

- Problem setup
- K-means algorithm

There are different types of machine learning problems

- **supervised learning** (what we have discussed by now)
All data is **labeled**
Aim to **predict**, e.g. classification and regression
- **unsupervised learning**
All data is **unlabeled**
Aim to **discover hidden and latent patterns and explore data**

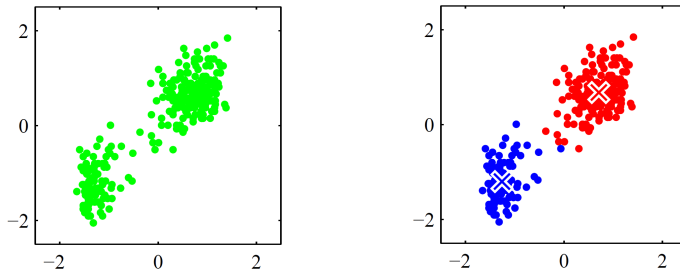
Today's focus: one important unsupervised learning problem: **clustering**

Clustering: informal definition

Given: a set of data points (feature vectors), *without labels*

Output: group the data into some clusters, which means

- **assign** each point to a specific cluster
- find the **center** (representative/prototype/...) of each cluster

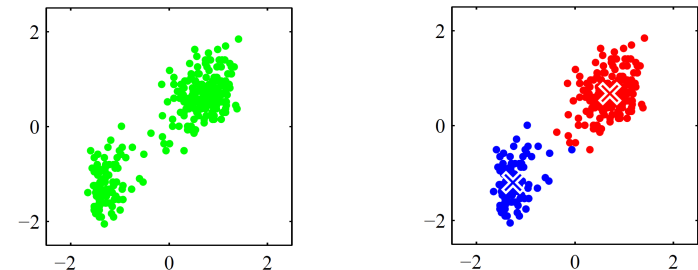


Clustering: formal definition

Given: data points $x_1, \dots, x_N \in \mathbb{R}^D$ and #clusters K we want to find

Output: group the data into K clusters, which means

- find an **assignment** $\gamma_{nk} \in \{0, 1\}$ s.t. if a data point $n \in [N]$ belongs to a cluster $k \in [K]$ then $\gamma_{nk} = 1$ and $\sum_{k \in [K]} \gamma_{nk} = 1$.
- find the cluster **centers** $\mu_1, \dots, \mu_K \in \mathbb{R}^D$



Many applications

One example: **image compression** (vector quantization)

- each pixel is a point
- perform clustering over these points
- **replace each point by the center** of the cluster it belongs to



Original image

Large $K \rightarrow$ Small K

()

February 19, 2019 37 / 52

Alternating minimization

Instead, use a heuristic that **alternatively minimizes over $\{\gamma_{nk}\}$ and $\{\mu_k\}$** :

Initialize $\{\gamma_{nk}^{(1)}\}$ and $\{\mu_k^{(1)}\}$

For $t = 1, 2, \dots$

- fix centers $\{\mu_k^{(t)}\}$, find assignments $\{\gamma_{nk}^{(t+1)}\}$

$$\{\gamma_{nk}^{(t+1)}\} = \operatorname{argmin}_{\{\gamma_{nk}\}} F(\{\gamma_{nk}\}, \{\mu_k^{(t)}\})$$

- fix assignments $\{\gamma_{nk}^{(t+1)}\}$, find new centers $\{\mu_k^{(t+1)}\}$

$$\{\mu_k^{(t+1)}\} = \operatorname{argmin}_{\{\mu_k\}} F(\{\gamma_{nk}^{(t+1)}\}, \{\mu_k\})$$

()

February 19, 2019 39 / 52

Formal Objective

Key difference from supervised learning problems: no labels given, which means *no ground-truth to even measure the quality of your answer!*

Still, we can turn it into an optimization problem, e.g. through the popular **“K-means” objective**: find γ_{nk} and μ_k to minimize

$$F(\{\gamma_{nk}\}, \{\mu_k\}) = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

i.e. the **sum of distances of each point to its center**.

Unfortunately, finding the exact minimizer is **NP-hard!**

()

February 19, 2019 38 / 52

A closer look

The first step (fixed centers, find assignments)

$$\operatorname{argmin}_{\{\gamma_{nk}\}} F(\{\gamma_{nk}\}, \{\mu_k\}) = \operatorname{argmin}_{\{\gamma_{nk}\}} \sum_n \sum_k \gamma_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

is simply to **assign each \mathbf{x}_n to the closest μ_k** , i.e.

$$\gamma_{nk} = \mathbb{I} \left[k = \operatorname{argmin}_c \|\mathbf{x}_n - \mu_c\|_2^2 \right]$$

for all $k \in [K]$ and $n \in [N]$.

()

February 19, 2019 40 / 52

A closer look

The second step (fixed assignments, find centers)

$$\operatorname{argmin}_{\{\mu_k\}} F(\{\gamma_{nk}\}, \{\mu_k\}) = \operatorname{argmin}_{\{\gamma_{nk}\}} \sum_n \sum_k \gamma_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

We will do it for each cluster.

The center is simply **an average of the points in that cluster** (hence the name)

$$\mu_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$

for each $k \in [K]$.

()

February 19, 2019 41 / 52

The K-means algorithm, S. Lloyd (1957)

Step 0 Initialization (choose K centers)

Step 1 Fix the centers μ_1, \dots, μ_K , **assign each point to the closest center**:

$$\gamma_{nk} = \mathbb{I} \left[k == \operatorname{argmin}_c \|\mathbf{x}_n - \mu_c\|_2^2 \right]$$

Step 2 Fix the assignment $\{\gamma_{nk}\}$, **update the centers**

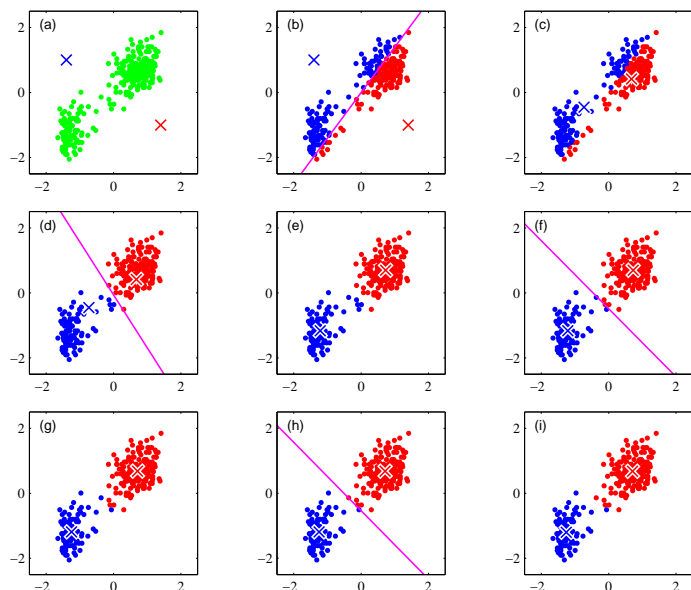
$$\mu_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$

Step 3 Repeat Steps 1 and 2 until the centers no longer change.

()

February 19, 2019 42 / 52

An example

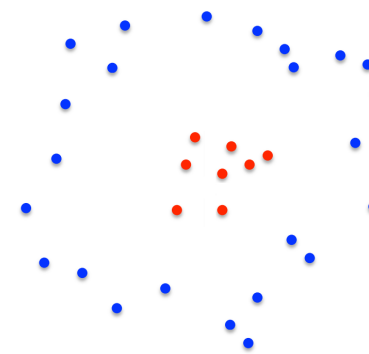


()

February 19, 2019 43 / 52

K-means algorithm is a heuristic!

K-means is not always able to properly cluster:



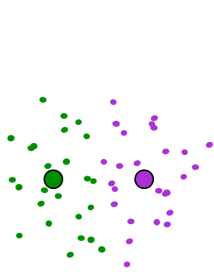
()

February 19, 2019 44 / 52

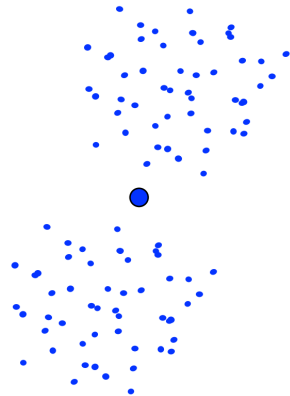
K-means algorithm is a heuristic!

It does matter how you initialize the centers!

In the following example $K = 3$:



Would be better to have one cluster here



... and two clusters here

()

February 19, 2019 45 / 52

How to initialize?

A bad selection for the initial centers can lead to a very poor clustering of data.

It also may lead a very long to converge.

There are [different ways to initialize](#):

- randomly pick K points as initial centers
- as it turns out, good initial centers are ones that aren't close to each other. (e.g. [K-means++](#), 2007)

()

February 19, 2019 46 / 52

How to initialize?

The K-means++ algorithm.

The algorithm selects initial centers that aren't close to each other, then uses K-means algorithm for clustering.

The high-level pseudo-code for the K-means++:

- select a data point at random as the first center
- loop $K-1$ times
 - ▶ compute distance squared $d(x)^2$ from each point to the nearest cluster center
 - ▶ select a point that has largest probability $\frac{d(x)^2}{\sum_x d(x)^2}$ as the next center

()

February 19, 2019 47 / 52

Convergence

It will [converge in a finite number of iterations](#) to a **local** minimum.

- objective decreases at each step
- objective is lower bounded by 0
- #possible_assignments is finite (K^N , exponentially large though)
- it may take **exponentially** many iterations to converge
- it might not converge to the **global** minimum

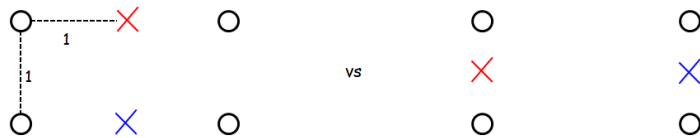
()

February 19, 2019 48 / 52

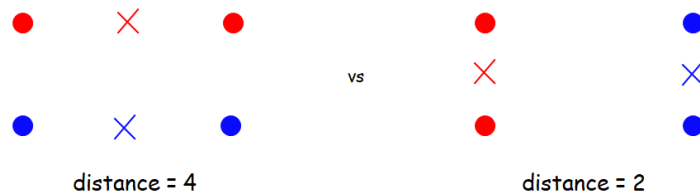
Local minimum v.s global minimum

Simple example: 4 data points, 2 clusters, 2 different initializations.

We initialize the centers by the mean of two points.



K-means converges immediately in both cases.



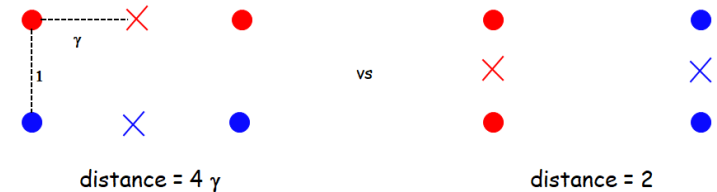
()

February 19, 2019 49 / 52

Local minimum v.s global minimum

In the left picture we get a local minimum, but in the right - a global minimum!

Moreover, local minimum can be *arbitrarily worse* if we increase the width of this "rectangle" to 2γ .



So, we get stuck at a local minimum.

Initialization matters a lot!

()

February 19, 2019 50 / 52

Cluster Quality Measures

We need to define a measure of cluster quality Q and then try different values of K until we get an optimal value for Q

There are different metrics for evaluating clustering algorithms, depending on what types of clusters we want

K-means emphasizes similarity of data within clusters:

$$Q = \sum_{k=1}^K \frac{1}{C_k} \sum_{x \in C_k} \|x - \mu_k\|_2^2$$

where C_k is the number of data points in cluster k .

()

February 19, 2019 51 / 52

Cluster Quality Measures

Other Quality measures:

The aim is to identify sets of clusters that are compact and at the same time are well separated

- Dunn Index
- Davies-Bouldin Index
- Silhouette Index

()

February 19, 2019 52 / 52