

Speeding Up Kernel Methods, and Intro to Unsupervised Learning

Piyush Rai

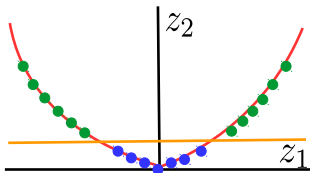
Introduction to Machine Learning (CS771A)

September 11, 2018



Recap: Nonlinear Mappings

- Idea: Use a nonlinear mapping $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ to map original data to a high-dim space, e.g.,



- Learn a linear model in the new space using the mapped inputs $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)$
- Equivalent to learning a nonlinear model on the original data $\mathbf{x}_1, \dots, \mathbf{x}_N$
- The mappings can be explicitly defined, or implicitly defined via a **kernel function** k , s.t.

$$k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

- Benefit of using kernels: Don't need to explicitly compute the mappings (M can be very large)
- Many ML algos only have data appearing as inner products. Can **kernelize** such algos

Recap: Nonlinear Mappings and Kernels

- A kernel function $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$ defines inner-product similarity between two inputs
- This is a Euclidean similarity in ϕ space but a “nonlinear” similarity in original space
- Some popular examples of kernel functions

$$k(\mathbf{x}_n, \mathbf{x}_m) = \mathbf{x}_n^\top \mathbf{x}_m \quad (\text{Linear kernel})$$

$$k(\mathbf{x}_n, \mathbf{x}_m) = (1 + \mathbf{x}_n^\top \mathbf{x}_m)^2 \quad (\text{Quadratic kernel})$$

$$k(\mathbf{x}_n, \mathbf{x}_m) = (1 + \mathbf{x}_n^\top \mathbf{x}_m)^d \quad (\text{Polynomial kernel of degree } d)$$

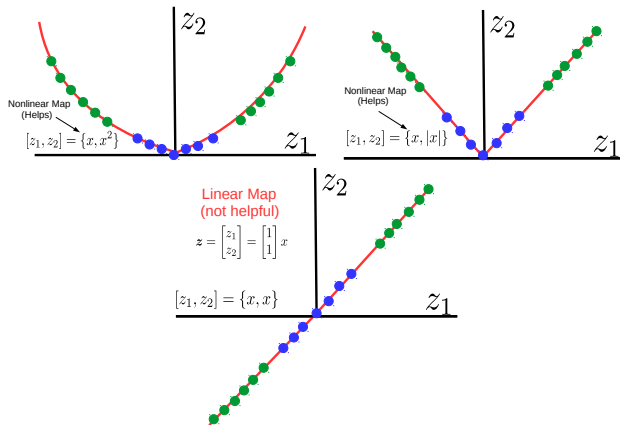
$$k(\mathbf{x}_n, \mathbf{x}_m) = \exp[-\gamma \|\mathbf{x}_n - \mathbf{x}_m\|^2] \quad (\text{RBF/Gaussian kernel})$$

- Each of these kernels have an associated feature mapping $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$, e.g.,
 - Quadratic kernel: $\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \dots, \sqrt{2}x_1x_2, \dots, x_1^2, \dots]$. Here $M = O(D + D^2)$ dimensional
 - RBF kernel: $\phi(\mathbf{x})$ is infinite dimensional (saw in the last class). Here $M = \infty$
- Again, remember that when using kernels, we don't have to compute ϕ explicitly



Recap: Nonlinear Mappings and Kernels

- Not every high-dim mapping is helpful. The mapping ϕ must be nonlinear



Kernel Methods can be Slow

Training phase can be slow (if N is very large)

Soft-Margin SVM: $\max_{\alpha \leq C} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha$

\downarrow
N x N size

Storing the learned model may be expensive

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \phi(\mathbf{x}_n)$$

\downarrow
Possibly very high-dimensional

Testing (prediction) phase can be slow (scales in N or at least the number of support vectors)

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



Kernel Methods can be Slow

Training phase can be slow (if N is very large)

Dual form of Ridge Regression: $\mathbf{w} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \mathbf{X}^\top \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \mathbf{x}_n$

Kernelized Ridge Regression: $\mathbf{w} = \phi(\mathbf{X})^\top (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \phi(\mathbf{X})^\top \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$

N x N size (also need to invert it)

Storing the learned model may be expensive

$$\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$$

Possibly very high-dimensional

Testing (prediction) phase can be slow (scales in N)

$$y = \mathbf{w}^\top \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$



Speeding Up Kernel Methods

- Kernel methods are slow at training and test time
- Would have been nice if we could easily compute the mapping $\phi(\mathbf{x})$ associated with kernel k
- Then we could apply linear models **directly on $\phi(\mathbf{x})$** without having to kernelize
- But this is in general not possible since $\phi(\mathbf{x})$ is very high dimensional
- Instead of a high-dim $\phi(\mathbf{x})$, can we get a good set of **low-dim features** $\psi(\mathbf{x}) \in \mathbb{R}^L$ using the kernel?
- If $\psi(\mathbf{x})$ is a good approximation of $\phi(\mathbf{x})$, then we can just use $\psi(\mathbf{x})$ in a linear model

$$\text{"Goodness" Criterion: } \psi(\mathbf{x}_n)^\top \psi(\mathbf{x}_m) \approx \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

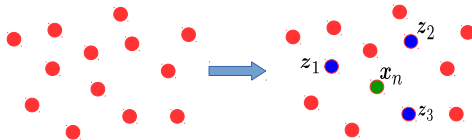
i.e., we want $\psi(\mathbf{x}_n)^\top \psi(\mathbf{x}_m) \approx k(\mathbf{x}_n, \mathbf{x}_m)$

- We will see two popular approaches: **Landmarks** and **Random Features**



Using Kernels to “Extract” Good Features: Landmarks

- Suppose we choose a small set of L “landmark” inputs $\mathbf{z}_1, \dots, \mathbf{z}_L$ in the training data



$$\psi(\mathbf{x}_n) = [k(\mathbf{z}_1, \mathbf{x}_n), k(\mathbf{z}_2, \mathbf{x}_n), k(\mathbf{z}_3, \mathbf{x}_n)] \in \mathbb{R}^3$$

- For each input \mathbf{x}_n , using a kernel k , define an L -dimensional feature vector as follows

$$\psi(\mathbf{x}_n) = [k(\mathbf{z}_1, \mathbf{x}_n), k(\mathbf{z}_2, \mathbf{x}_n), \dots, k(\mathbf{z}_L, \mathbf{x}_n)]$$

- $\psi(\mathbf{x}_n) \in \mathbb{R}^L$ is such that $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) \approx \psi(\mathbf{x}_n)^\top \psi(\mathbf{x}_m)$
- Can now apply a linear model on the ψ representation (L -dimensional now) of the inputs
- This will be fast both at training as well as test time if L is small
- No need to kernelize the linear model and work with kernels (but still reap their benefits :-))
- Note: The landmarks need not be actual inputs. Can even be learned from data.



Using Kernels to “Extract” Good Features: Random Features

- Many kernel functions can be written as[†]

$$k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) = \mathbb{E}_{\mathbf{w} \sim p(\mathbf{w})}[t_{\mathbf{w}}(\mathbf{x}_n)t_{\mathbf{w}}(\mathbf{x}_m)]$$

where $t_{\mathbf{w}}(\cdot)$ is a **scalar-valued function** with parameters $\mathbf{w} \in \mathbb{R}^D$ from some distribution $p(\mathbf{w})$

- Example: For the **RBF kernel**, $t_{\mathbf{w}}(\cdot)$ is **cosine function** and $p(\mathbf{w})$ is **zero mean Gaussian**

$$k(\mathbf{x}_n, \mathbf{x}_m) = \mathbb{E}_{\mathbf{w} \sim p(\mathbf{w})}[\cos(\mathbf{w}^\top \mathbf{x}_n) \cos(\mathbf{w}^\top \mathbf{x}_m)]$$

- Given $\mathbf{w}_1 \dots, \mathbf{w}_L$ drawn from $p(\mathbf{w})$, using Monte-Carlo approximation of expectation above

$$k(\mathbf{x}_n, \mathbf{x}_m) \approx \frac{1}{L} \sum_{\ell=1}^L \cos(\mathbf{w}_\ell^\top \mathbf{x}_n) \cos(\mathbf{w}_\ell^\top \mathbf{x}_m) = \psi(\mathbf{x}_n)^\top \psi(\mathbf{x}_m)$$

where $\psi(\mathbf{x}_n) = \frac{1}{\sqrt{L}}[\cos(\mathbf{w}_1^\top \mathbf{x}_n), \dots, \cos(\mathbf{w}_L^\top \mathbf{x}_n)]$ is an L -dim. feature vector (L needs to be set)

- Can apply a linear model on this L -dim representation of the data (no need to kernelize)
- Such techniques exist for several kernels (RBF, polynomial, etc)

[†] **Random Features for Large-Scale Kernel Machines** (Ben and Retch, NIPS 2007. Note: This paper actually won the test-of-time award at NIPS 2017)



Other Techniques for Speeding Up Kernel Methods

- Reducing the number of support vectors (for SVM based models), For example,
 - Learn the kernelized SV. Identify the support vectors.
 - Cluster the support vectors
 - Pick one SV from each cluster, retrain SVM using the chosen SVs
- Low-rank approximations of kernel matrix (Nyström approximation)



Kernel Methods: Some Final Comments

- Sometimes, even linear models can be trained via kernelization (but with [linear kernel](#))
- Benefit? Well, this may be beneficial sometimes due to computational reasons
- For example, ridge regression requires solving

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

.. where we learn \mathbf{w} by inverting a $D \times D$ matrix

- Instead, the dual version of Ridge Regression, as we saw earlier, requires solving

$$\mathbf{w} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \mathbf{X}^\top \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \mathbf{x}_n$$

.. where we learn \mathbf{w} in terms of $\boldsymbol{\alpha}$ by inverting $N \times N$ matrix

- Even when working with linear model, if $D > N$, the latter way may be preferable
- Similar considerations apply to other kernelizable models too (e.g., SVM)
- If linear model is what you want, still makes sense to look at the relative values of N and D to decide whether to go for the dual (kernelized) formulation of the problem with a linear kernel



Kernel Methods: Some Final Comments

- Kernel methods give us good features to make learning easier
- However, these features are pre-defined (due to the choice of kernel)
- Example: Consider the quadratic kernel applied to input $\mathbf{x} = [x_1, x_2]$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, x_1^2, \sqrt{2}x_1x_2, x_2, x_2^2, \sqrt{2}x_2] \quad (\text{fixed definition for } \phi)$$

- Another alternative is to **learn** good features from data
- We will revisit this when we talk about **deep neural networks**

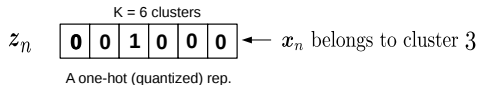


Unsupervised Learning



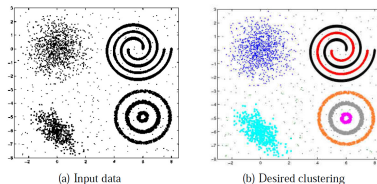
Unsupervised Learning

- Roughly speaking, it is about learning interesting structures in the data (unsupervisedly!)
- There is no supervision (no labels/responses), only inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$
- Some examples of unsupervised learning
 - **Clustering**: Grouping similar inputs together (and dissimilar ones far apart)
 - **Dimensionality Reduction**: Reducing the data dimensionality
 - **Estimating the probability density** of data (which distribution “generated” the data)
- Most unsupervised learning algos can also be seen as learning a **new representation** of data
 - Typically a **compressed** representation, e.g., clustering can be used to get a one-hot representation



Clustering

- Given: N unlabeled examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; no. of desired partitions K
- Goal: Group the examples into K “homogeneous” partitions



Picture courtesy: “Data Clustering: 50 Years Beyond K-Means”, A.K. Jain (2008)

- Loosely speaking, it is classification without ground truth labels
- A good clustering is one that achieves:
 - High within-cluster similarity
 - Low inter-cluster similarity



Similarity can be Subjective

- Clustering only looks at similarities, no labels are given
- Without labels, similarity can be hard to define



- Thus using the right distance/similarity is very important in clustering
- Also important to define/ask: “Clustering based on what”?

Picture courtesy: http://www.guy-sports.com/humor/videos/powerpoint_presentation_dogs.htm



Clustering: Some Examples

- Document/Image/Webpage Clustering
- Image Segmentation (clustering pixels)



- Clustering web-search results
- Clustering (people) nodes in (social) networks/graphs
- .. and many more..

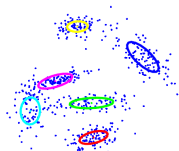
Picture courtesy: <http://people.cs.uchicago.edu/~pff/segment/>



Types of Clustering

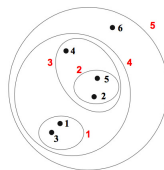
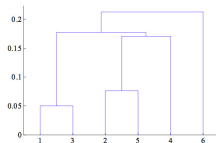
① Flat or Partitional clustering

- Partitions are **independent of each other**



② Hierarchical clustering

- Partitions can be visualized using a tree structure (a **dendrogram**)



- Possible to view partitions at **different levels of granularities** by “cutting” the tree at some level

Flat Clustering: K -means algorithm (Lloyd, 1957)

- **Input:** N examples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$; $\mathbf{x}_n \in \mathbb{R}^D$; the number of partitions K
- **Desired Output:** Cluster assignments of these N examples and K cluster means μ_1, \dots, μ_K
- **Initialize:** K cluster means μ_1, \dots, μ_K , each $\mu_k \in \mathbb{R}^D$
 - Usually initialized randomly, but good initialization is crucial; many smarter initialization heuristics exist (e.g., K -means++, Arthur & Vassilvitskii, 2007)
- **Iterate:**
 - (Re)-Assign each example \mathbf{x}_n to its closest cluster center (based on the smallest Euclidean distance)

$$\mathcal{C}_k = \{n : k = \arg \min_k \|\mathbf{x}_n - \mu_k\|^2\}$$

(\mathcal{C}_k is the set of examples assigned to cluster k with center μ_k)

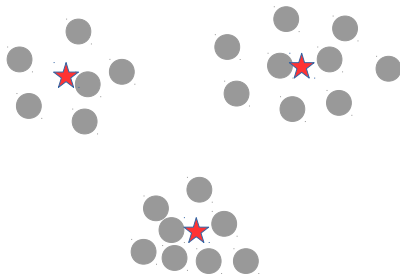
- Update the cluster means

$$\mu_k = \text{mean}(\mathcal{C}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$

- Repeat while not converged
- Stop when cluster means or the “loss” does not change by much



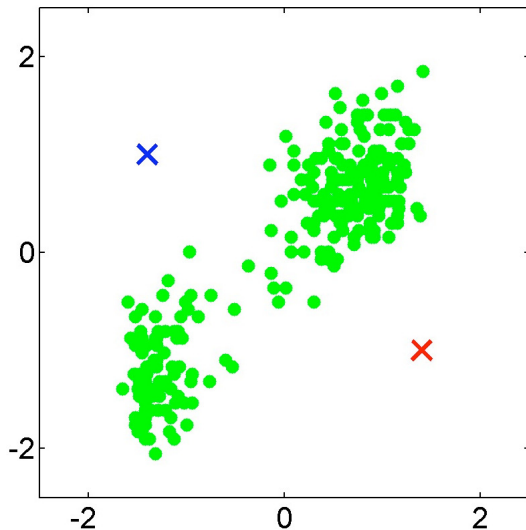
K -means = Prototype Classification (with unknown labels)



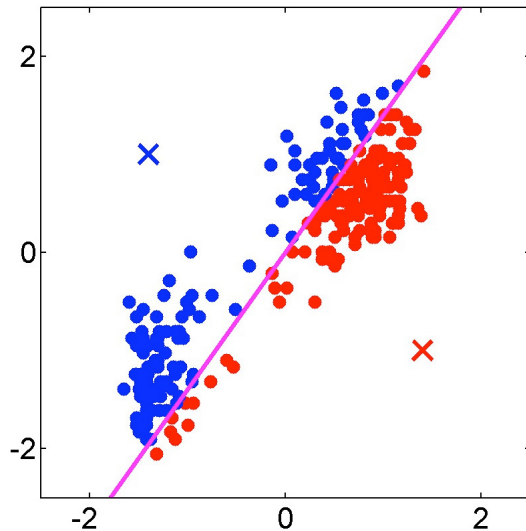
- Guess the means
- Predict the labels
- Recompute the means
- Repeat



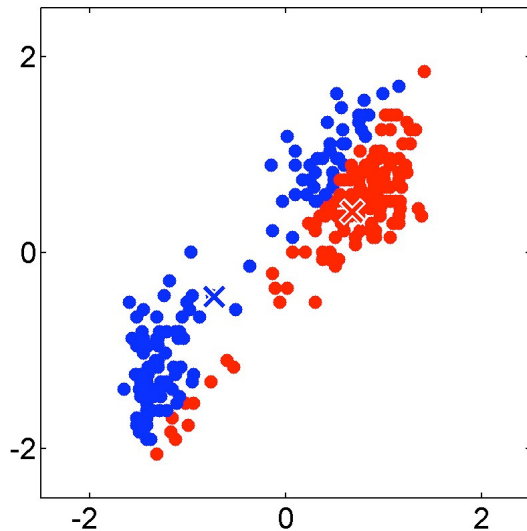
K -means: Initialization (assume $K = 2$)



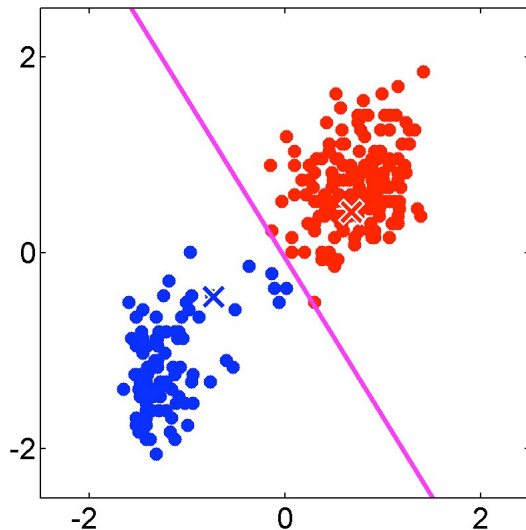
K -means iteration 1: Assigning points



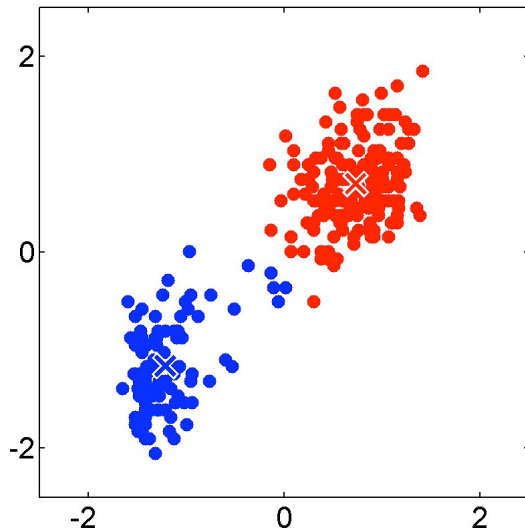
K -means iteration 1: Recomputing the centers



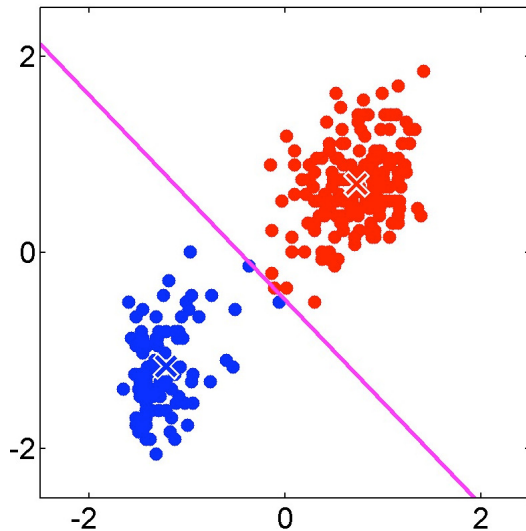
K -means iteration 2: Assigning points



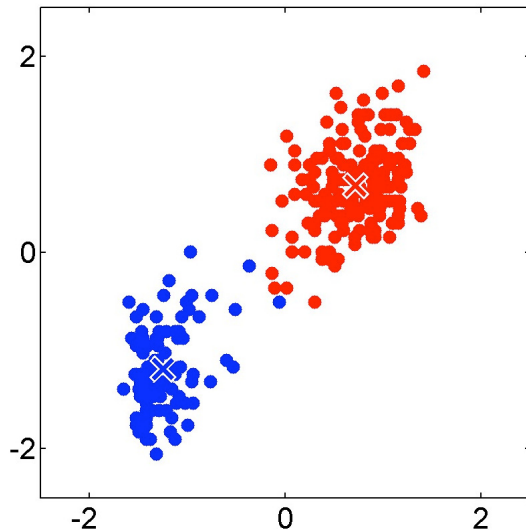
K -means iteration 2: Recomputing the centers



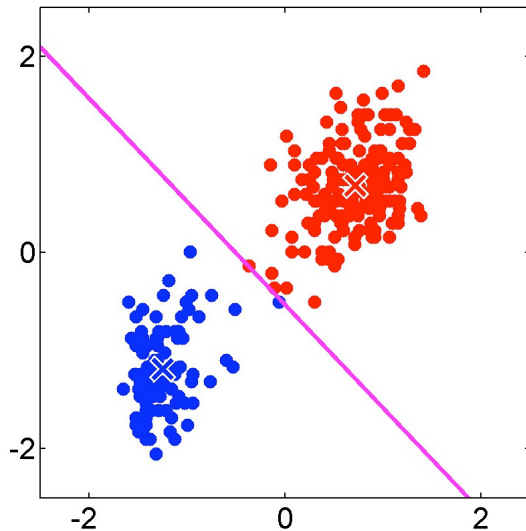
K-means iteration 3: Assigning points



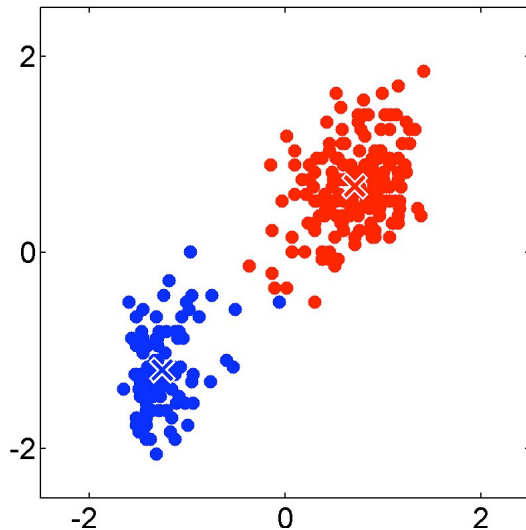
K -means iteration 3: Recomputing the centers



K -means iteration 4: Assigning points



K -means iteration 4: Recomputing the centers



What Loss Function is K -means Optimizing?

- Let μ_1, \dots, μ_K be the K cluster centroids (means)
- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if \mathbf{x}_n belongs to cluster k , and 0 otherwise
 - Note: $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$ represents a length K **one-hot encoding** of \mathbf{x}_n
- Define the **distortion** or **“loss”** for the cluster assignment of \mathbf{x}_n

$$\ell(\mu, \mathbf{x}_n, \mathbf{z}_n) = \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

- Total distortion over all points defines the K -means “loss function”

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2 = \underbrace{\|\mathbf{X} - \mathbf{Z}\mu\|_F^2}_{\text{matrix factorization view}}$$

where \mathbf{Z} is $N \times K$ (row n is \mathbf{z}_n) and μ is $K \times D$ (row k is μ_k)

- The K -means **problem** is to minimize this objective w.r.t. μ and \mathbf{Z}
 - Alternating optimization would give the K -means (Lloyd’s) algorithm we saw earlier!



Next Class: Clustering (Contd.)

