Regularization in Unsupervised Learning

Guille Canas, Lorenzo Rosasco

MIT, 9.520 Class 22

May 2, 2012

L. Rosasco Regularization in Unsupervised Learning

A ▶

< ∃ > <

ъ

- ⊒ →

Goal To discuss some problems in unsupervised learning and in particular introduce a statistical learning framework for learning data representation/reconstruction under constraints, discussing the role played by regularization theory and regularized algorithms.

・ 同 ト ・ ヨ ト ・ ヨ ト …

- From supervised learning to learning data representation
- Unsupervised Learning Algorithms
- Computations
- Final Comments

・ 同 ト ・ ヨ ト ・ ヨ ト

3

- we are given a training set $S = ((x_1, y_1), \dots, (x_n, y_n))$ sampled i.i.d. with respect to p(x, y).
- the problem is: to predict the best possible output for **new** input
- find

$$f(x_{new}) = y_{new}.$$

・ 同 ト ・ ヨ ト ・ ヨ ト

The problem can be formalized fixing some concepts.

- loss function V(y, f(x))
- hypotheses space $\mathcal{H} : \{f \mid f : X \to \mathbb{R}\}$
- a learning algorithm maps the data into a function, i.e. $A(S) = f_S \in \mathcal{H}$

・ 同 ト ・ ヨ ト ・ ヨ ト

The problem can be formalized fixing some concepts.

- loss function V(y, f(x))
- hypotheses space $\mathcal{H} : \{f \mid f : X \to \mathbb{R}\}$
- a learning algorithm maps the data into a function, i.e. $A(S) = f_S \in \mathcal{H}$

・ 同 ト ・ ヨ ト ・ ヨ ト

A good algorithm is such that

$$\mathbb{P}(|I_{\mathcal{S}}[f_{\mathcal{S}}] - I[f_{\mathcal{S}}]| \ge \epsilon)$$

is small,or

$$\mathbb{P}(I[f_{\mathcal{S}}] - \inf_{f \in \mathcal{H}} I[f] \ge \epsilon)$$

is small, where $I[f] = \mathbb{E}[V(y, f(x))]$.

ъ

프 🖌 🖌 프 🕨

Typically we use algorithm based on the (regularized) empirical risk minimization.

Constrained minimization

 $\min_{f \in \mathcal{H}_{\lambda}} I_{\mathcal{S}}[f_{\mathcal{S}}]$

or penalized minimization

 $\min_{f\in\mathcal{H}}\{I_{\mathcal{S}}[f_{\mathcal{S}}]+\lambda R(f)\}.$

・ 同 ト ・ ヨ ト ・ ヨ ト …

3

- we are given a training set S = (x₁,..., x_n) sampled i.i.d. with respect to p(x).
- the problem is to extract from data "useful" information about p.....

・ 同 ト ・ ヨ ト ・ ヨ ト …

Unsupervised Learning is not One but Many Problems

- density estimation
- dimensionality reduction
- clustering
- manifold learning
- o correlation analysis
- association rules
- networks/graph analysis
- dictionary learning, vector quantizaton, data representation/reconstruction

Unsupervised learning is ubiquitous: many algorithms, no unified framework.

個人 くほん くほん

Data representation is the preliminary step to supervised learning. Representation are often designed/engineered, but ideally should be *learned*.

We have discussed how a *good representation* should be discriminant and yet invariant to *task irrelevant* transformations.

Ideally such a good representation should reduce the sample complexity of subsequent supervised tasks.

Some questions before we even start..

- How can we know what is "task relevant" if we don't have labels (we don't know what's the task)?
- What do we mean with learning a representation?
- What guarantees we have? (stability, overfitting...)

- From supervised learning to learning data representation
- Unsupervised Learning Algorithms
- Computations
- Final Comments

・ 同 ト ・ ヨ ト ・ ヨ ト

3

We are going to replace

discrimination \mapsto reconstruction+constraints

for the time being we are not going to consider the problem of invariance (see last lectures to see how this can be done).

Given a training set S, we are going to consider algorithms that minimize the following empirical (reconstruction) error,

$$I_{S}[C] = \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, C)$$

where:

- C is a subset of X
- d²(x_i, C) = min_{v∈C} ||x v||² is the square distance of a point x to the set C

We are going to consider (constrained) algorithms based on

 $\min_{\mathcal{C}\in\mathcal{H}_k}I_{\mathcal{S}}[\mathcal{C}]$

where $\mathcal{H}_k \subset \{C \mid C \subset X\}$ is a hypotheses space of suitable subsets of *X*.

The above algorithm returns an estimator $A(S) = C_S \in \mathcal{H}_k$.

- the above error measure depends on the distribution, points that are more likely to be sampled contribute more to the error
- we can alternatively consider penalized algorithms

$$\min_{C\in\mathcal{H}}I_{\mathcal{S}}[C]+R(C)$$

where $\mathcal{H} \subset \{C \mid C \subset X\}$ is an essentially unconstrained hypotheses space and R(C) a regularizer that penalize *complex/large* sets.

The above framework is general enough to encompass a variety of algorithms.

- PCA
- Sparse coding
- K-Means
- K-Flats
- Non Negative Matrix Factorization
- ...

Each algorithm is defined by a suitable choice of \mathcal{H}_k .

・ 同 ト ・ ヨ ト ・ ヨ ト …

Hypothesis space

$$\mathcal{H} = \{ C \subseteq \mathcal{X} : C = \{ x : x = Tb = \sum_{j=1}^{k} b^{j} t_{j}, \text{ with } T \in \mathcal{T}, b \in \mathcal{B} \} \}$$

- \mathcal{T} linear transformations $\mathcal{B} \mapsto \mathcal{X}$, $Tb = \sum_{j=1}^{k} b^{j} t_{j}$, where t_{j} are the **code vectors** defining a **dictionary**.
- The choice of B determines the encoding x → (b¹,..., b^k) and the algorithm.

・ 同 ト ・ ヨ ト ・ ヨ ト ・

Hypothesis space

$$\mathcal{H} = \{ C \subseteq \mathcal{X} : C = \{ x : x = Tb = \sum_{j=1}^{k} b^{j} t_{j}, \text{ with } T \in \mathcal{T}, b \in \mathcal{B} \} \}$$

- \mathcal{T} linear transformations $\mathcal{B} \mapsto \mathcal{X}$, $Tb = \sum_{j=1}^{k} b^{j} t_{j}$, where t_{j} are the **code vectors** defining a **dictionary**.
- The choice of B determines the encoding x → (b¹,...,b^k) and the algorithm.

ヘロト ヘアト ヘビト ヘビト

Hypothesis space

$$\mathcal{H} = \{ C \subseteq \mathcal{X} : C = \{ x : x = Tb = \sum_{j=1}^{k} b^{j} t_{j}, \text{ with } T \in \mathcal{T}, b \in \mathcal{B} \} \}$$

- \mathcal{T} linear transformations $\mathcal{B} \mapsto \mathcal{X}$, $Tb = \sum_{j=1}^{k} b^{j} t_{j}$, where t_{j} are the **code vectors** defining a **dictionary**.
- The choice of B determines the encoding x → (b¹,..., b^k) and the algorithm.

・ 同 ト ・ ヨ ト ・ ヨ ト …

- Input: samples S.
- **Output**: mapping T_S (determines set $C_S = T_S(B)$).
- **Dictionary**: *T* in coordinates (t_1, \ldots, t_k) .
- **Encoding**: given x, encoding is $\operatorname{argmin}_{b \in \mathcal{B}} d^2(x, T(b))$
- C = T(B) are encoded with no error

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

- Input: samples S.
- **Output**: mapping T_S (determines set $C_S = T_S(B)$).
- **Dictionary**: *T* in coordinates (t_1, \ldots, t_k) .
- **Encoding**: given x, encoding is $\operatorname{argmin}_{b \in \mathcal{B}} d^2(x, T(b))$
- C = T(B) are encoded with no error

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

PCA:

- Let $\mathcal{B} = \mathbb{R}^k$
- Sets are of the form $C = T(\mathbb{R}^k)$, with T linear
 - C is a k-dimensional linear subspaces of \mathcal{X}

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\mathbb{R}^k))$$
$$= \min_{\text{rank-}k \text{ linear projection } \pi} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, \pi x_i)$$

• Free parameter: k (maximum dimension of C)

ヘロト ヘワト ヘビト ヘビト

PCA:

- Let $\mathcal{B} = \mathbb{R}^k$
- Sets are of the form $C = T(\mathbb{R}^k)$, with T linear
 - C is a k-dimensional linear subspaces of \mathcal{X}

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\mathbb{R}^k))$$
$$= \min_{\text{rank-}k \text{ linear projection } \pi} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, \pi x_i)$$

• Free parameter: k (maximum dimension of C)

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ .

PCA:

- Let $\mathcal{B} = \mathbb{R}^k$
- Sets are of the form $C = T(\mathbb{R}^k)$, with T linear
 - C is a k-dimensional linear subspaces of \mathcal{X}

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, T(\mathbb{R}^{k}))$$
$$= \min_{\text{rank-}k \text{ linear projection } \pi} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, \pi x_{i})$$

• Free parameter: k (maximum dimension of C)

ヘロト 人間 とくほとくほとう

э.

PCA:

- Let $\mathcal{B} = \mathbb{R}^k$
- Sets are of the form $C = T(\mathbb{R}^k)$, with T linear
 - C is a k-dimensional linear subspaces of \mathcal{X}

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, T(\mathbb{R}^{k}))$$
$$= \min_{\text{rank-}k \text{ linear projection } \pi} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, \pi x_{i})$$

• Free parameter: k (maximum dimension of C)

・ロト ・同ト ・ヨト ・ヨトー

3

- Let $\mathcal{B} = \{e_1, \ldots, e_k\}$
- Sets of the form $C = T(B) \subseteq X$, with T linear mapping

• Arbitrary sets $C \subset \mathcal{X}$ of size |C| = k

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\{e_1, \dots, e_k\})))$$
$$= \min_{C = \{m_1, \dots, m_k\}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^2(x_i, m_j)$$

• Free parameter: k (size of C)

<ロト <回 > < 注 > < 注 > 、

- Let $\mathcal{B} = \{e_1, \ldots, e_k\}$
- Sets of the form $C = T(B) \subseteq X$, with T linear mapping
 - Arbitrary sets $C \subset \mathcal{X}$ of size |C| = k

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\{e_1, \dots, e_k\})))$$
$$= \min_{C = \{m_1, \dots, m_k\}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^2(x_i, m_j)$$

• Free parameter: k (size of C)

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ …

- Let $\mathcal{B} = \{e_1, \ldots, e_k\}$
- Sets of the form $C = T(B) \subseteq X$, with T linear mapping
 - Arbitrary sets $C \subset \mathcal{X}$ of size |C| = k

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, T(\{e_{1}, \dots, e_{k}\})))$$
$$= \min_{C = \{m_{1}, \dots, m_{k}\}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

• Free parameter: k (size of C)

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ …

- Let $\mathcal{B} = \{e_1, \ldots, e_k\}$
- Sets of the form $C = T(B) \subseteq X$, with T linear mapping
 - Arbitrary sets $C \subset \mathcal{X}$ of size |C| = k

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, T(\{e_{1}, \dots, e_{k}\})))$$
$$= \min_{C = \{m_{1}, \dots, m_{k}\}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

• Free parameter: k (size of C)

<四> < 回 > < 回 > < 回 > -

э.

K-Flats

K-Flats:

- Let $\mathcal{B} \subset \mathbb{R}^{m \times k}$
- $\mathcal{B} = \cup_{j=1}^{k} \operatorname{span} \left(\boldsymbol{e}_{j,1}, \dots, \boldsymbol{e}_{j,m} \right) \neq \mathbb{R}^{m \times k}$
- Sets are collections of k m-dimensional linear subspaces of X (m-flats)

$$I_{S,k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\mathcal{B}))$$
$$= \min_{\substack{C = \{\pi_1, \dots, \pi_k\} \\ \pi_j \text{ rank-}k \text{ linear projection}}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^2(x_i, \pi_j x_i)$$

• Free parameters: *k*, *m* (number and dimension of flats)

・ロト ・ 理 ト ・ ヨ ト ・

ъ

K-Flats

K-Flats:

- Let $\mathcal{B} \subset \mathbb{R}^{m \times k}$
- $\mathcal{B} = \cup_{j=1}^{k} \operatorname{span} \left(\boldsymbol{e}_{j,1}, \dots, \boldsymbol{e}_{j,m} \right) \neq \mathbb{R}^{m \times k}$
- Sets are collections of k m-dimensional linear subspaces of X (m-flats)

$$I_{\mathcal{S},k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^{2}(x_{i}, T(\mathcal{B}))$$
$$= \min_{\substack{C = \{\pi_{1}, \dots, \pi_{k}\} \\ \pi_{j} \text{ rank-}k \text{ linear projection}}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, \pi_{j}x_{i})$$

• Free parameters: *k*, *m* (number and dimension of flats)

・ロト ・ 理 ト ・ ヨ ト ・

ъ

K-Flats

K-Flats:

- Let $\mathcal{B} \subset \mathbb{R}^{m \times k}$
- $\mathcal{B} = \cup_{j=1}^{k} \operatorname{span} \left(\boldsymbol{e}_{j,1}, \dots, \boldsymbol{e}_{j,m} \right) \neq \mathbb{R}^{m \times k}$
- Sets are collections of k m-dimensional linear subspaces of X (m-flats)

$$I_{\mathcal{S},k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(\mathcal{B}))$$
$$= \min_{\substack{C = \{\pi_1, \dots, \pi_k\} \\ \pi_j \text{ rank-}k \text{ linear projection}}} \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^2(x_i, \pi_j x_i)$$

• Free parameters: *k*, *m* (number and dimension of flats)

・ロト ・ 理 ト ・ ヨ ト ・

3

Sparse Coding

Sparse Coding

- Let $\mathcal{B} = B_1(\lambda) = \{ y \in \mathbb{R}^k : \|y\|_1 \le \lambda \}$
- $\mathcal{T} = \{ \text{ linear } \mathcal{T} : \|\mathcal{T}e_i\| \le \gamma, 1 \le i \le k \}$
- Sets are of the form $C = T(B_1(\lambda))$ (Typically sparse) linear combinations of *k* points

$$I_{\mathcal{S},k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(B_1(\lambda)))$$

• Free parameter: λ , controls sparsity.

・ロト ・聞 ト ・ヨト ・ヨト

Sparse Coding

- Let $\mathcal{B} = B_1(\lambda) = \{ y \in \mathbb{R}^k : \|y\|_1 \le \lambda \}$
- $\mathcal{T} = \{ \text{ linear } \mathcal{T} : \|\mathcal{T}e_i\| \le \gamma, 1 \le i \le k \}$
- Sets are of the form $C = T(B_1(\lambda))$ (Typically sparse) linear combinations of *k* points

$$I_{\mathcal{S},k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(B_1(\lambda)))$$

• Free parameter: λ , controls sparsity.

・ロト ・聞 ト ・ヨト ・ヨト

Sparse Coding

- Let $\mathcal{B} = B_1(\lambda) = \{ y \in \mathbb{R}^k : \|y\|_1 \le \lambda \}$
- $\mathcal{T} = \{ \text{ linear } \mathcal{T} : \|\mathcal{T}e_i\| \le \gamma, 1 \le i \le k \}$
- Sets are of the form $C = T(B_1(\lambda))$ (Typically sparse) linear combinations of *k* points

$$I_{\mathcal{S},k} = \min_{T \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, T(B_1(\lambda)))$$

• Free parameter: λ , controls sparsity.

く 同 ト く ヨ ト く ヨ ト

- From supervised learning to learning data representation
- Unsupervised Learning Algorithms
- Computations
- Final Comments

・ 同 ト ・ ヨ ト ・ ヨ ト

3

K-means problem: find set $C = \{m_1, \ldots, m_k\}$ of size *k* minimizing

$$I_{S}[C] = \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

- Non-linear, non-convex.
- NP-hard even for k = 2!

・ 同 ト ・ ヨ ト ・ ヨ ト

ъ

K-Means Problem

$$I_{S}[C] = \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

 Given means, optimal partition given by assignment function:

$$a(x_i) = \operatorname{argmin}_{j=1}^k d^2(x_i, m_j)$$

• Given a partition, optimal means given by centers of mass:

$$m_j = \sum_{i:a(x_i)=j} x_i / \sum_{i:a(x_i)=j} 1$$

個 とくき とくきと

K-Means with Lloyd's Algorithm

$$I_{S}[C] = \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

Lloyd's algorithm:

- Initialize $\{m_i : i = 1, ..., k\}$ randomly in $S = \{x_1, ..., x_n\}$ (without replacement)
- Repeat until convergence:
 - find partition given means,

$$a(x_i) = \operatorname{argmin}_{j=1}^k d^2(x_i, m_j)$$

• update means given the above partition,

$$m_j = \sum_{i:a(x_i)=j} x_i / \sum_{i:a(x_i)=j} 1$$

ヘロト ヘアト ヘビト ヘビト

K-Means with Lloyd's Algorithm

$$I_{S}[C] = \frac{1}{n} \sum_{i=1}^{n} \min_{j=1}^{k} d^{2}(x_{i}, m_{j})$$

Lloyd's algorithm:

- Initialize {m_i : i = 1,..., k} randomly in S = {x₁,..., x_n} (without replacement)
- Repeat until convergence:
 - find partition given means,

$$a(x_i) = \operatorname{argmin}_{j=1}^k d^2(x_i, m_j)$$

update means given the above partition,

$$m_j = \sum_{i:a(x_i)=j} x_i / \sum_{i:a(x_i)=j} 1$$

・ 同 ト ・ ヨ ト ・ ヨ ト …

K-Means with Lloyd's Algorithm (cont.)

Lloyd's algorithm:

- a) Given means \Rightarrow optimize partition.
- b) Given partition \Rightarrow optimize means.

Greedy block-coordinate descent. Does it converge?

- Steps a) and b) **strictly** decrease *I_S*[*C*], until convergence.
- \Rightarrow no partition is repeated (never twice in the same state).
- Number of different partitions of S is $\leq k^n$

 \Rightarrow must converge to **local minimum** in finite number of steps.

How close to **global optimum**?: we don't know.

K-Means with Lloyd's Algorithm (cont.)

Lloyd's algorithm:

- a) Given means \Rightarrow optimize partition.
- b) Given partition \Rightarrow optimize means.

Greedy block-coordinate descent. Does it converge?

- Steps a) and b) **strictly** decrease *I*_S[*C*], until convergence.
- \Rightarrow no partition is repeated (never twice in the same state).
- Number of different partitions of S is $\leq k^n$
- \Rightarrow must converge to **local minimum** in finite number of steps.

How close to global optimum?: we don't know.

< 回 > < 回 > < 回 >

K-Means with Lloyd's Algorithm (cont.)

Lloyd's algorithm:

- a) Given means \Rightarrow optimize partition.
- b) Given partition \Rightarrow optimize means.

Greedy block-coordinate descent. Does it converge?

- Steps a) and b) **strictly** decrease *I*_S[*C*], until convergence.
- \Rightarrow no partition is repeated (never twice in the same state).
- Number of different partitions of S is $\leq k^n$
- \Rightarrow must converge to **local minimum** in finite number of steps.

How close to global optimum?: we don't know.

Choose seeding adaptively:

- m_1 uniform random in $S = \{x_1, \ldots, x_n\}$.
- For *j* = 2, . . . , *k*, let

$$\mathbb{P}[m_j = i] \propto \min_{l=1,\ldots,j-1} d^2(x_i, m_l)$$

(pick means far from previously-inserted means)

イロト イポト イヨト イヨト 三日

Properties of K-Means++

K-Means++

- m_1 uniform random in $S = \{x_1, \ldots, x_n\}$.
- $\mathbb{P}[m_j = i] \propto \min_{l=1,...,j-1} d^2(x_i, m_l), \quad j = 2,...,k$

Guarantees:

 O(ln k)-Approximation randomized algorithm: after seeding it is

$$\mathbb{E}\left\{I_{S}[\{m_{1},\ldots,m_{k}\}]\right\} \leq 8(\ln k+2) \cdot \min_{|C|=k} I_{S}[C]$$

ヘロア ヘビア ヘビア・

Properties of K-Means++

K-Means++

- m_1 uniform random in $S = \{x_1, \ldots, x_n\}$.
- $\mathbb{P}[m_j = i] \propto \min_{l=1,...,j-1} d^2(x_i, m_l), \quad j = 2,...,k$

Guarantees:

 O(ln k)-Approximation randomized algorithm: after seeding it is

$$\mathbb{E}\left\{I_{\mathcal{S}}[\{m_1,\ldots,m_k\}]\right\} \leq 8(\ln k + 2) \cdot \min_{|\mathcal{C}|=k} I_{\mathcal{S}}[\mathcal{C}]$$

・ 同 ト ・ ヨ ト ・ ヨ ト

Complexity of K-Means++

K-Means++

- m_1 uniform random in $S = \{x_1, \ldots, x_n\}$.
- $\mathbb{P}[m_j = i] \propto \min_{l=1,\ldots,j-1} d^2(x_i, m_l), \quad j = 2, \ldots, k$

Complexity: O(kn).

- Choose m_1 randomly. Let $p_i \leftarrow d^2(x_i, m_1), i = 1, \dots, n$ $\Theta(n)$
- For j = 2, ..., k: $k \times$
 - Draw $m_j \propto [p_1, \ldots, p_{j-1}]$ $\Theta(n)$
 - $p_i \leftarrow \min\{p_i, d^2(x_i, m_j)\}$

ヘロア ヘビア ヘビア・

Complexity of K-Means++

K-Means++

- m_1 uniform random in $S = \{x_1, \ldots, x_n\}$.
- $\mathbb{P}[m_j = i] \propto \min_{l=1,...,j-1} d^2(x_i, m_l), \quad j = 2,...,k$

Complexity: O(kn).

• Choose m_1 randomly. Let $p_i \leftarrow d^2(x_i, m_1), i = 1, ..., n$ $\Theta(n)$ • For j = 2, ..., k: $k \times$ • Draw $m_j \propto [p_1, ..., p_{j-1}]$ $\Theta(n)$ • $p_i \leftarrow \min\{p_i, d^2(x_i, m_j)\}$ $\Theta(n)$

ヘロト ヘアト ヘビト ヘビト

Additionally:

- May still run Lloyd algorithm after.
- Incremental: computes all intermediate solutions (k' = 1, ..., k).
- Proof does not (trivially) extend to K-Flats.

・ 同 ト ・ ヨ ト ・ ヨ ト …

Can we quantify the generalization properties of the previous unsupervised learning algorithms?

Recall that we defined,

- Empirical reconstruction error: $I_{\mathcal{S}}[C] = \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, C)$
- Expected reconstruction error: $I[C] = \mathbb{E}[d^2(x, C)]$

Generalization Bounds (cont.)

A good algorithm is such that

generalization

$$\mathbb{P}(|I_{\mathcal{S}}[C_{\mathcal{S}}] - I[C_{\mathcal{S}}]| \geq \epsilon)$$

is small,or

consistency (excess risk bounds)

$$\mathbb{P}(I[C_{\mathcal{S}}] - \inf_{C \in \mathcal{H}_k} I[C] \geq \epsilon)$$

is small.

伺き くほき くほう

Many algorithms can be extended using kernels. Given a feature map $\Phi: X \to \mathcal{H}$, the idea is to consider

$$I_{\mathcal{S}}[C] = \frac{1}{n} \sum_{i=1}^{n} d^2(x_i, C)$$

where:

- C is a subset of \mathcal{H}
- d²(x_i, C) = min_{v∈C} ||Φ(x) − v||²_H is the square distance of the feature map Φ(x) of a point x to the set C

For example in the case of kernel k-means, the algorithms will compute (implicitly) means in the feature space

$$m_j = \frac{1}{\ell_j} \sum_{i=1}^{\ell_j} \Phi(x_i), \quad \ell_j \leq n.$$

Indeed the distance of a point to a such mean can be computed explicitly

$$\|\Phi(x) - m_j\|_{\mathcal{H}}^2 = K(x, x) - \frac{1}{\ell_j} \sum_{i=1}^{\ell_j} K(x, x_i) + \frac{1}{\ell_j^2} \sum_{i=1}^{\ell_j} K(x_i, x_i)$$

・ 同 ト ・ ヨ ト ・ ヨ ト

- Learning data representation can be described in a regularization framework.
- Different Hypotheses spaces induce different algorithms.
- Computations are considerably more complicated (lack of convexity).
- Sample/approximation tradef-offs?
- Partial supervision (semi-supervised, time, constraintssymmetry etc.)

・ 同 ト ・ ヨ ト ・ ヨ ト

L. Rosasco Regularization in Unsupervised Learning

くりょう 小田 マイボット 山下 シックション

K-Flats Problem

Modified Lloyd algorithm:

- Initialize flats "randomly" $\{\pi_i : i = 1, \ldots, k\}$.
- Repeat until convergence:
 - Given flats, optimal partition given by assignment function:

$$a(x_i) = \operatorname{argmin}_{j=1}^k d^2(x_i, \pi_j x_i)$$

• Given partition, optimal flat π_j given by top eigenvectors (PCA) of

$$S_j S_j^t = \sum_{i:a(x_i)=j} x_i x_i^t$$

• Similar guarantees: convergence to local minimum.

・ロ・ ・ 同・ ・ ヨ・ ・ ヨ・

K-Flats Problem

Modified Lloyd algorithm:

- Initialize flats "randomly" $\{\pi_i : i = 1, \ldots, k\}$.
- Repeat until convergence:
 - Given flats, optimal partition given by assignment function:

$$a(x_i) = \operatorname{argmin}_{j=1}^k d^2(x_i, \pi_j x_i)$$

• Given partition, optimal flat π_j given by top eigenvectors (PCA) of

$$S_j S_j^t = \sum_{i:a(x_i)=j} x_i x_i^t$$

• Similar guarantees: convergence to local minimum.

ヘロン 人間 とくほ とくほ とう