

COMP 790-125: Goals for today

- ▶ Matching pursuit
- ▶ Learning dictionaries
 - ▶ batch learning
 - ▶ online learning
- ▶ Sparse representations, compressive sensing

Dictionary based representations

Last time we talked about sparsely representing a patch

$$\operatorname{argmin}_{\alpha} \frac{1}{2} \left\| \underbrace{\mathbf{R}_{i,j} \mathbf{y}}_{\text{extracted patch}} - \underbrace{\mathbf{D}\alpha}_{\substack{\text{combination of} \\ \text{dictionary elements}}} \right\|_2^2 + \lambda \underbrace{\|\alpha\|_1}_{\text{sparsity penalty}}$$

We talked about using dictionaries to obtain denoised images

$$\operatorname{argmin}_{\alpha, \mathbf{x}} \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \sum_{i,j} \left\{ \frac{1}{2} \|\mathbf{R}_{i,j} \mathbf{x} - \mathbf{D}\alpha_{i,j}\|_2^2 + \lambda \|\alpha_{i,j}\|_1 \right\}$$

All the while, we assumed that we know \mathbf{D} , and we looked at some standard (orthonormal) dictionaries.

Matching pursuit

One of the benefits of having an orthogonal dictionary \mathbf{D} is that

$$\operatorname{argmin}_{\alpha} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\alpha\|_2^2 + \|\alpha\|_1$$

has a closed form solution

$$\alpha_j = \mathbf{d}_j^T \mathbf{y}$$

But, if the dictionary is overcomplete, and consequently cannot possibly be orthonormal, we cannot guarantee that this update is optimal.

Matching pursuit – dictionary coherence

Recall that orthonormal dictionaries are those for which $\mathbf{d}_i^T \mathbf{d}_j = 0, i \neq j$.

Departure from orthonormality called coherence is defined as

$$\mathcal{M}(\mathbf{D}) = \max_{i,j} |\mathbf{d}_i^T \mathbf{d}_j|.$$

If coherence is small then update

$$\alpha_j = \mathbf{d}_j^T \mathbf{y}$$

might still be approximately correct.

Matching pursuit algorithm [2] is based on this idea

Initialize: $\alpha_i = 0, \forall i, \mathbf{r} = \mathbf{y}$

while $\|\mathbf{r}\|_2 >$ threshold **do**

$$k = \operatorname{argmax}_j |\mathbf{d}_j^T \mathbf{r}|$$

$$\alpha_k = \mathbf{d}_k^T \mathbf{r}$$

$$\mathbf{r} = \mathbf{r} - \alpha_k \mathbf{d}_k$$

end while

Learning dictionaries

In order to make story simpler we will assume that there is a patchset¹ Data = $\{\mathbf{y}_i : i = 1..n\}$

$$\underset{\alpha, \mathbf{D}}{\text{minimize}} \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$$

We want to simultaneously learn weights for each patch and dictionary.

¹We do not need to bother with details of patch extraction, since we are not focusing on reconstruction of clean signal. We did that last time.

Learning dictionaries

A problem

$$\mathbf{D}\alpha_i = \left(\frac{1}{3}\mathbf{D}\right)(3\alpha_i)$$

so we could learn dictionary \mathbf{D} and weights α_i or just scale the dictionary down and weights up.

More broadly, matrix factorization $\mathbf{A} = \mathbf{BC}$ is not unique²

For some \mathbf{T} such that $\mathbf{TT}^T = \mathbf{I}$ we can obtain a different factorization $\mathbf{A} = (\mathbf{BT})(\mathbf{T}^T \mathbf{C})$.

Q: What could we possibly do?

²Be on the lookout for such problems when your optimization problem has product terms.

Learning dictionaries

Just constrain the dictionary elements to have norm less or equal to 1

$$\begin{aligned} & \underset{\alpha, \mathbf{D}}{\text{minimize}} && \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha\|_1 \\ & \text{subject to} && \|\mathbf{d}_i\|_2 \leq 1, \forall i \end{aligned}$$

Q: Can you argue that with this constraint optimal dictionary consists of elements with norm exactly 1?

Learning dictionaries

$$\begin{aligned} \underset{\alpha, \mathbf{D}}{\text{minimize}} \quad & \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha\|_1 \\ \text{subject to} \quad & \|\mathbf{d}_j\|_2 \leq 1, \forall j \end{aligned}$$

Even though we eliminated spurious weight and dictionary combinations with unit norm constraint a more insidious problem persists.

There exist local minima, solutions that locally look optimal but are not global optima.

The problem is convex in α given \mathbf{D} and vice versa.

$$\begin{aligned} \underset{\alpha}{\operatorname{argmin}} \quad & \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha\|_1 \\ \underset{\mathbf{D}, \forall j \|\mathbf{d}_j\|_2 \leq 1}{\operatorname{argmin}} \quad & \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha\|_1 \end{aligned}$$

Learning dictionaries – an algorithm

```
for t = 1,...MAXITERATIONS do
    for i = 1,...T do
         $\alpha_i = \operatorname{argmin}_{\alpha_i} \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$ 
    end for
     $\mathbf{D} = \operatorname{argmin}_{\mathbf{D}, \forall j \|\mathbf{d}_j\|_2 \leq 1} \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2$ 
end for
```

Updating dictionaries

A number of different algorithms to solve

$$\alpha_i = \operatorname{argmin}_{\alpha_i} \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$$

- ▶ Coordinate descent
- ▶ Matching Pursuit
- ▶ LARS (Least Angle Regression)[3]
- ▶ FISTA (Fast Iterative Shrinkage and Thresholding Algorithm)[1].

What about dictionary update

$$\mathbf{D} = \operatorname{argmin}_{\mathbf{D}, \forall j \|\mathbf{d}_j\|_2 \leq 1} \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2?$$

We will use *projected* gradient descent.

Projected gradient descent

Like gradient descent, but with projection back onto the constrained set.

$$\underset{\mathbf{x} \in C}{\text{minimize}} f(\mathbf{x})$$

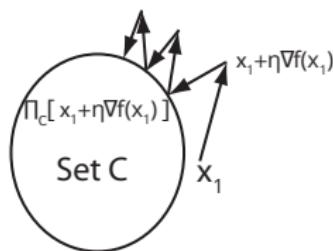
The algorithm alternates two steps:

for $i = 1, \dots, \text{MAXIT}$ **do**

$$\mathbf{x}^0 = \mathbf{x} - \eta \nabla f(\mathbf{x})$$

$$\mathbf{x} = \Pi_C[\mathbf{x}^0]$$

end for



Projection step for unit norm ball

Projection is defined as

$$\Pi_C[\mathbf{x}] = \underset{\mathbf{z} \in C}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{z}\|_2^2$$

a point in set C closest to point \mathbf{x} .

In our case for a single column of dictionary $S = \{\mathbf{z} : \|\mathbf{z}\|_2 \leq 1\}$ and

$$\Pi_S[\mathbf{x}] = \frac{\mathbf{x}}{\max(1, \|\mathbf{x}\|_2)}$$

and it amounts to normalizing a vector only if its norm is greater than 1.

For full dictionary $C = \{\mathbf{A} : \forall j \|\mathbf{a}_j\|_2 \leq 1\}$ and

$$\Pi_C[\mathbf{D}] = [\Pi_S[\mathbf{d}_1] \ \Pi_S[\mathbf{d}_2] \ \cdots \ \Pi_S[\mathbf{d}_n]]$$

So, just normalize columns whose norm is larger than 1.

Dictionary learning update

Hence a single dictionary update, for given α amounts to running this code

```
for i = 1,.., MAXIT do
    D0 = D + η/n ∑i(yi + ∑i(yi - Dαi)αiT
    D = ΠC[D0]
end for
```

where n is the number of samples (or patches) and η is a step size.

Full learning algorithm

```
for t = 1,...MAXITERATIONS do
    for i = 1,...T do
         $\alpha_i = \operatorname{argmin}_{\alpha_i} \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$ 
    end for
    for i = 1,..,MAXIT do
         $\mathbf{D}^0 = \mathbf{D} + \frac{\eta}{n} \sum_i (\mathbf{y}_i + \sum_i (\mathbf{y}_i - \mathbf{D}\alpha_i)\alpha_i^T)$ 
         $\mathbf{D} = \Pi_C[\mathbf{D}^0]$ 
    end for
end for
```

Illustrative applications of dictionary learning – inpainting

Remember that the Sabine Indians to the east of here were light gray mountain Indians and
had a kind of vegetation so that you wanted to climb into their nests because as you went to
climb into the top of a beehive, they were backbiting mountain with a howling, jaws like. The Santa
Clara would go against the sky in the west and take the valley from the open sea, and they were dark and
bewitchingly dark and dangerous. I always found in myself a desire over all the peaks of east where I never
went, but I could almost say, unless it could be that the morning over all the peaks of the Gila and the
Gila River back from the edges of the Santa Lucia. It may be that the birth and death of the day had some
power in my feelings about the two ranges of mountains.

The floor of the Salinas Valley, between the ranges and below the foothills, is level because this valley used to be the bottom of a hundred-mile inlet from the sea. The river mouth at Moss Landing was centuries ago the entrance to this long inland water. Once, fifty miles down the valley, my father bored a well. The drill came up bare with topsoil and then with gravel and then with white sea sand full of shells and even pebbles.



Illustrative applications of dictionary learning – completion



Online dictionary learning

Batch algorithms, pass through whole dataset before updating parameters

Online algorithms work on smaller subsets of data at a time.

for $t = 1, \dots, \text{MAXITERATIONS}$ **do**

randomly choose a sample \mathbf{z}_t

$$\boldsymbol{\alpha}_t = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{z}_t - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$$

$$\mathbf{D}^t = \underset{\mathbf{D}, \forall j \|\mathbf{d}_j\|_2}{\operatorname{argmin}} \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \sum_i \|\mathbf{z}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2$$

end for

A dictionary learning implementation

Use SPArse Modeling Software package.

It is available from:

<http://spams-devel.gforge.inria.fr/downloads.html>

It has interfaces for R, python and Matlab.

It is fairly efficient, and you can try out your ideas in a matter of hours.

Sparse representations

We implicitly assumed that given a patch \mathbf{y} and a dictionary \mathbf{D} we can find a *sparse* $\boldsymbol{\alpha}$ to encode that patch so that

$$\mathbf{y} \approx \mathbf{D}\boldsymbol{\alpha}.$$

We are going to look into when we can guarantee that such sparse representation exists, and in some cases when it is guaranteed to be exact, that is

$$\mathbf{y} = \mathbf{D}\boldsymbol{\alpha}.$$

Johnson-Lindenstrauss Lemma

Lemma

For a set V points in \mathbf{R}^n and some $\epsilon \in (0, 1)$ and integer $k = O(1/\epsilon^2 \log |V|)$ there exists a map $f : \mathbf{R}^n \rightarrow \mathbf{R}^k$ such that for all $u, v \in V$

$$(1 - \epsilon) \|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon) \|u - v\|^2$$

Johnson-Lindenstrauss Lemma

The number of dimensions needed to capture the distances between datapoints with a small distortion is logarithmic in the size of the set.

- ▶ No mention of distributions – **any set of points**
- ▶ No mention of the source space dimension n just the number of data points **$\log |V|$**

Note that the projection f need not preserve the feature identities.

A coordinate in the low dimensional space may combine apples and oranges, e.g. $1/2*\text{height} + 1/3*\text{weight} - 1/4*\text{age}$

Practical construction of mapping³

We will use a simple linear mapping $f(\mathbf{u}) = \mathbf{Ax}$

First, randomly generate matrix $\mathbf{U} : k \times n$

- ▶ rows uniformly sampled basis vectors of length n
- ▶ rows sampled from a spherical Gaussian $\mathcal{N}(0, I_n)$
- ▶ each entry is randomly sampled from $+1/-1$

$$p(u_{i,j} = 1) = 1/2 \quad p(u_{i,j} = -1) = 1/2$$

- ▶ each entry is randomly sampled from $\{-1, 0, +1\}$

$$p(u_{i,j} = \frac{1}{\sqrt{q}}) = \frac{1}{2q} \quad p(u_{i,j} = -\frac{1}{\sqrt{q}}) = \frac{1}{2q} \quad p(u_{i,j} = 0) = 1 - \frac{1}{q}$$

From matrix \mathbf{U} construct final projection matrix

$$\mathbf{A} = \frac{1}{\sqrt{k}} \mathbf{U}$$

³Johnson and Lindenstrauss 1984, Indyk and Motwani 1998, Achlioptas 2003, Matoušek 2008.

Sparse mapping

The last choice ensured that there are $1 - \frac{1}{q}$ zero entries in the matrix **A**.

This results in a cheap projection operation with all of the benefits of the JL lemma.

Locality sensitive hashing

In our exploration of JL lemma we came up with a scheme

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x}$$

that did dimensionality reduction but preserved the distances between points.

It turns out you can go even further and construct a mapping

$$h_{\mathbf{a}}(\mathbf{x}) = \text{sign}(\mathbf{a}^T \mathbf{x})$$

Each row of matrix \mathbf{A} defines a single hash function.

We will denote the group of hashing functions

$$h_A : \mathbf{R}^N \rightarrow \{-1, +1\}^k$$

which produces a k long sign vector. Note that storing this vector takes k bits!

Sparse representations

The framework is very similar to the regression setup.

We are given a set of vectors $(\mathbf{x}_1, \dots, \mathbf{x}_p)$ that form a basis of some space S or alternatively a frame.⁴

We wish to find a representation of a vector \mathbf{y}

$$\mathbf{y} = \sum_i w_i \mathbf{x}_i = \mathbf{Xw}$$

where the number of nonzero w_i s

$$\|\mathbf{w}\|_0 = \sum_i [w_i \neq 0]$$

is as small as possible.

⁴A frame is a set of possibly linearly dependent vectors that span a space with a constraint that there exist $0 < A \leq B < \infty$ such that

$$A \|\mathbf{v}\|_2 \leq \sum_i \langle \mathbf{x}_i, \mathbf{v} \rangle \leq B \|\mathbf{v}\|_2$$

ℓ_0 is hard but ℓ_1 is not

Solving optimization problem (P_0)

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\mathbf{w}\|_0 \\ & \text{subject to} && \mathbf{y} = \mathbf{X}\mathbf{w} \end{aligned}$$

is NP-hard, but its convex relaxation (P_1)

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\mathbf{w}\|_1 \\ & \text{subject to} && \mathbf{y} = \mathbf{X}\mathbf{w} \end{aligned}$$

is not. This problem is named Basis Pursuit.

Note that unlike the regression we require exact equality between \mathbf{y} and its representation $\mathbf{X}\mathbf{w}$

Basis Pursuit Denoising is a lasso regression analog.

When does solving (P_1) give solution to (P_0) ?

We can define coherence property for our basis/frame

$$\mathcal{M}(\mathbf{X}) = \max_{i \neq j} |\mathbf{x}_i^T \mathbf{x}_j|$$

Theorem

Suppose sparsest solution w^*

$$\|\mathbf{w}^*\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mathcal{M}(\mathbf{X})} \right)$$

then solving (P_1) gives the same solutions as the original (P_0) problem.

Restricted Isometry Property

Definition

A vector $\mathbf{w} \in \mathbf{R}^n$ is k -sparse if it has at most k non-zero entries.

Definition

If for any k -sparse vector \mathbf{w} it holds that

$$(1 - \delta_k) \|\mathbf{w}\|_2^2 \leq \|\mathbf{X}\mathbf{w}\|_2^2 \leq (1 + \delta_k) \|\mathbf{w}\|_2^2$$

then X has k -Restricted Isometry Property (RIP) with δ_k .

Put differently, using \mathbf{X} to map \mathbf{w} into a different space gives a vector of similar length as \mathbf{w} .

Contrast with Johnson Lindenstrauss for $f(\mathbf{u}) = \mathbf{A}\mathbf{u}$

$$(1 - \epsilon) \|\mathbf{u} - \mathbf{v}\|^2 \leq \|\mathbf{A}(\mathbf{u} - \mathbf{v})\|^2 \leq (1 + \epsilon) \|\mathbf{u} - \mathbf{v}\|^2$$

If matrix \mathbf{A} satisfies JL Lemma then RIP can be shown to hold for this matrix as well.

When does solving (P_1) give solution to (P_0) ?

$$(1 - \delta_k) \|\mathbf{w}\|_2^2 \leq \|\mathbf{Xw}\|_2^2 \leq (1 + \delta_k) \|\mathbf{w}\|_2^2$$

Theorem

If \mathbf{X} is k -RIP with $\delta_{2k} < \sqrt{2} - 1$, then for all k -sparse vectors \mathbf{w} such that $\mathbf{Xw} = \mathbf{y}$ the solution for (P_1) is equal to solution of (P_0) .

To summarize, if you construct a matrix that is k -RIP then you can represent a vector \mathbf{y} in this basis with some small distortion. Further, recovery of this representation can be done solving a convex optimization problem (P_1) .

Sensing

If a high-dimensional signal \mathbf{y} is sparsely representable in some basis \mathbf{X}

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

how much of \mathbf{y} do we need to know to reconstruct \mathbf{w} and hence the rest of \mathbf{y} .

We will split the issues of representation and sensing:

- ▶ representation: $\mathbf{y} = \mathbf{X}\mathbf{w}$ where \mathbf{w} is k-sparse vector
- ▶ sensing: $\mathbf{z} = \mathbf{R}\mathbf{y} = \mathbf{R}\mathbf{X}\mathbf{w}$ is measured

Sensing

Assuming we only measured some entries of vector $\mathbf{z} = \mathbf{R}\mathbf{y} = \mathbf{R}\mathbf{X}\mathbf{w}$ indexed by set M the recovery problem is

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\mathbf{w}\|_1 \\ & \text{subject to} && z_i = \langle \mathbf{r}_i, \mathbf{X}\mathbf{w} \rangle, \forall i \in M \end{aligned}$$

How large does set M have to be?

Recovery from compressive sensing measurements

We define coherence between basis

$$\mu(X, R) = \sqrt{n} \max_{i,j} |\langle \mathbf{x}_i, \mathbf{r}_j \rangle|$$

and for any two bases $\mu(\mathbf{A}, \mathbf{B}) \in [1, \sqrt{n}]$

Theorem

Given a representation basis \mathbf{X} and sensing basis \mathbf{R} and a vector $\mathbf{y} = \mathbf{X}\mathbf{w}$ with \mathbf{w} k -sparse then if

$$m \geq C\mu^2(\mathbf{X}, \mathbf{R})k \log \frac{n}{\delta}$$

and we randomly select m values to sense from $\mathbf{z} = \mathbf{R}\mathbf{y}$ we can recover \mathbf{w} with probability $1 - \delta$.

Practicalities

You don't need to know the k up front, you simply solve the problem

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && \|\mathbf{w}\|_1 \\ & \text{subject to} && z_i = \langle \mathbf{r}_i, \mathbf{Xw} \rangle, \forall i \in M \end{aligned}$$

The choice of M is random so any subset of vectors in \mathbf{R} is fine.

Note that the problem is not the same as Lasso but you can derive an ADMM algorithm for example.

Compressive sensing further reading

Compressive sensing is a very active field.

Your best bet in exploring applications is to start at
<http://dsp.rice.edu/cs>

Interesting ideas for projects relating to compressive sensing would involve learning dictionaries using prior knowledge and possibly giving strong theoretical guarantees about ability to reconstruct signals.

-  Amir Beck and Marc Teboulle.
A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems.
SIAM J. Img. Sci., 2(1):183–202, March 2009.
-  S.G. Mallat and Z. Zhang.
Matching pursuits with time-frequency dictionaries.
Signal Processing, IEEE Transactions on, 41(12):3397–3415, Dec 1993.
-  Bradley Efron Trevor, Trevor Hastie, Lain Johnstone, and Robert Tibshirani.
Least angle regression.
Annals of Statistics, 32:407–499, 2002.