

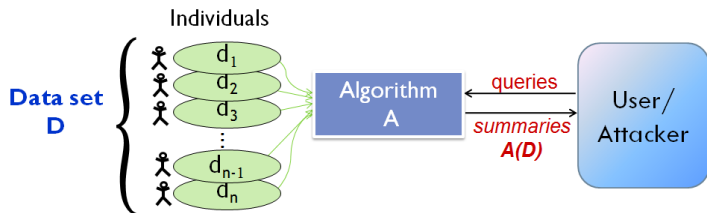
Private Model Selection via Stability Arguments and the Robustness of Lasso

Adam Smith¹ Abhradeep Thakurta²

¹Pennsylvania State University

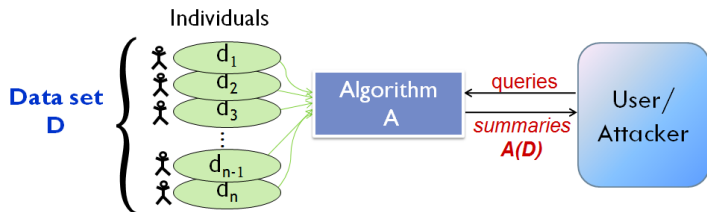
²Stanford University and Microsoft Research SVC

What Can We Learn Privately?



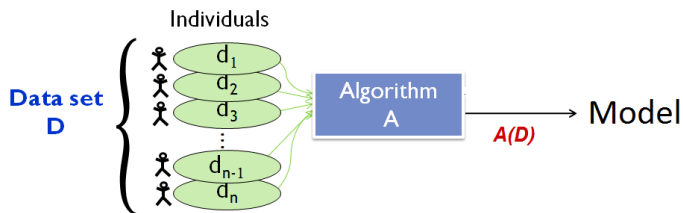
- Consider data set of sensitive individual data (e.g, medical records, purchase history)
- **Two conflicting goals:**
 - **Utility:** Release global information about \mathcal{D}
 - **Privacy:** Protect the privacy of individual entries

What Can We Learn Privately?



- Consider data set of sensitive individual data (e.g, medical records, purchase history)
- **Difficult problem:** Many attacks on “anonymized” data [NS08,Korolova11,CKN11,...]
- **Active research area:** Databases, Learning Theory, Programming Languages, Cryptography, Algorithms...

Model Selection



Models: Discrete collection of probability distribution families

$$M_1, \quad , M_2, \quad \dots$$

Goal: Select a model that best “fits” the data set \mathcal{D}

Example: Model could be a set of features

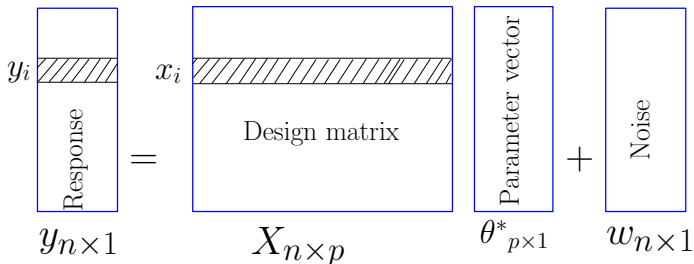
Running Example: Sparse Linear Regression

- Data point $\mathbf{d}_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$

Running Example: Sparse Linear Regression

- Data point $\mathbf{d}_i = (\mathbf{x}_i, \mathbf{y}_i)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $\mathbf{y}_i \in \mathbb{R}$
- **Model assumption:** Data generated by noisy linear system

i -th user data: (x_i, y_i)



Running Example: Sparse Linear Regression

- Data point $\mathbf{d}_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$
- **Model assumption:** Data generated by noisy linear system

i -th user data: (x_i, y_i)

The diagram shows the equation: $y_{n \times 1} = X_{n \times p} \theta_{p \times 1}^* + w_{n \times 1}$. The $y_{n \times 1}$ vector has a hatched row labeled y_i and the word "Response" written vertically. The $X_{n \times p}$ matrix has a hatched row labeled x_i and the words "Design matrix" written horizontally. The $\theta_{p \times 1}^*$ vector has three blue horizontal bars. The $w_{n \times 1}$ vector is labeled "Noise".

- **Sparsity:** θ^* has $s < n$ **non-zero** entries
- **High-dimensionality:** $p \gg n$

Running Example: Sparse Linear Regression

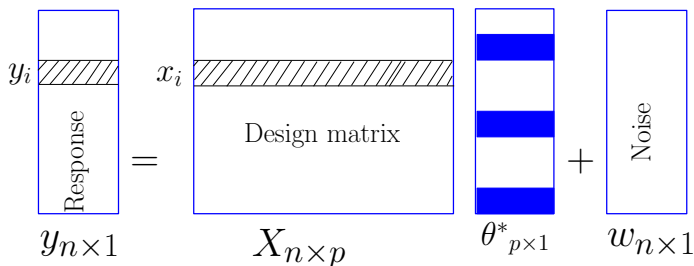
- Data point $\mathbf{d}_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$
- **Model assumption:** Data generated by noisy linear system
 i -th user data: (x_i, y_i)

The diagram illustrates the linear regression model equation: $y_{n \times 1} = X_{n \times p} \theta_{p \times 1}^* + w_{n \times 1}$. On the left, a vertical rectangle labeled "Response" has a hatched horizontal band at the top labeled y_i . Below it is the dimension $y_{n \times 1}$. In the middle, an equals sign is followed by a larger rectangle labeled "Design matrix" with a hatched horizontal band at the top labeled x_i . Below it is the dimension $X_{n \times p}$. To the right of the design matrix is a plus sign, followed by a vertical rectangle labeled $\theta_{p \times 1}^*$ with three blue horizontal bands. Below it is the dimension $\theta_{p \times 1}^*$. To the right of $\theta_{p \times 1}^*$ is another plus sign, followed by a vertical rectangle labeled "Noise" with a hatched horizontal band at the top. Below it is the dimension $w_{n \times 1}$.

- **Sparsity:** θ^* has $s < n$ **non-zero** entries
- **High-dimensionality:** $p \gg n$
- **Model selection problem:** Find the **support** of θ^*

Running Example: Sparse Linear Regression

- Data point $\mathbf{d}_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$
- **Model assumption:** Data generated by noisy linear system
 i -th user data: (x_i, y_i)



- **Sparse linear regression:** LASSO estimator

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{2n} \|y - X\theta\|_2^2 + \frac{\Lambda}{n} \|\theta\|_1$$

This Paper:

- Differentially private model selection
 - **Application:** Sparse linear regression
- Stability analysis of LASSO

Why is Privacy a Concern in Model Selection?

Many model selection procedures use convex optimization

Examples of convex optimization problems:

- LASSO, mean, median and support vector machines (SVM)

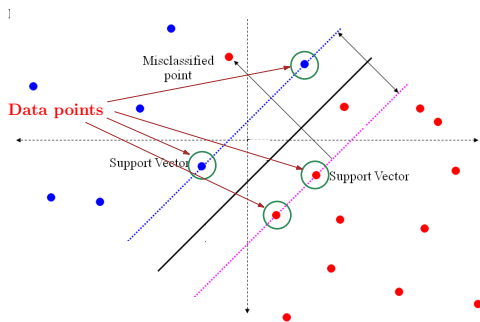
Why is Privacy a Concern in Model Selection?

Many model selection procedures use convex optimization

Examples of convex optimization problems:

- LASSO, mean, median and support vector machines (SVM)

Privacy breach: SVM example



- Reveals some data points exactly

Why is Privacy a Concern in Model Selection?

Many model selection procedures use convex optimization

Examples of convex optimization problems:

- LASSO, mean, median and support vector machines (SVM)

Question: How can we ensure privacy?

Why is Privacy a Concern in Model Selection?

Many model selection procedures use convex optimization

Examples of convex optimization problems:

- LASSO, mean, median and support vector machines (SVM)

Question: How can we ensure privacy?

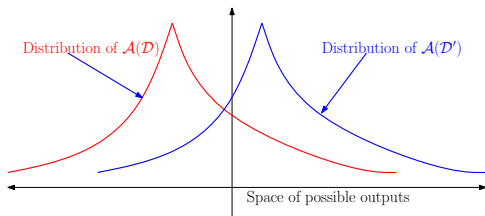
Answer: Design algorithms that satisfy differential privacy

- **Intuition:** Users learn roughly the same things about you whether or not your data is in the data set
- It is a **restriction** on the algorithm \mathcal{A}
 - Guarantee privacy even with **arbitrary** side information

- **Intuition:** Users learn roughly the same things about you whether or not your data is in the data set
- It is a **restriction** on the algorithm \mathcal{A}
 - Guarantee privacy even with **arbitrary** side information

Requirement

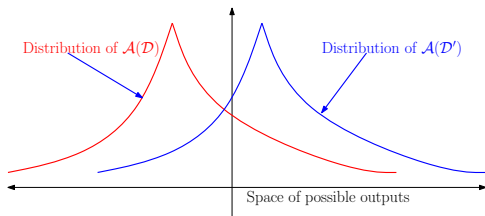
*Any two data sets \mathcal{D} and \mathcal{D}' that differ in one element induce **close** distributions on the space of outputs of \mathcal{A}*



- **Intuition:** Users learn roughly the same things about you whether or not your data is in the data set
- It is a **restriction** on the algorithm \mathcal{A}
 - Guarantee privacy even with **arbitrary** side information

Requirement

*Any two data sets \mathcal{D} and \mathcal{D}' that differ in one element induce **close** distributions on the space of outputs of \mathcal{A}*



- Closeness is measured by parameters ϵ and δ
 - Think $\epsilon = \text{small constant}$ and $\delta = 1/\text{poly}(n)$

Prior Works on Learning and Privacy

- Lots of work on differentially private learning
[KLNRS08,CM08,ZLW09,CMS11,KST12,JKT12,⋯]
- Almost all the results apply only to low-dimensional regime (i.e., $p < n$)
- In the context of sparse regression:
 - [ZLW09] works in the regime when $p < n$
 - [Kifer, Smith, T.'12] is the only result in $p \gg n$ regime

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies** **private** model selection
 - Two notions of stability
 - Subsampling stability [Shao96,Bach08,MB10,...]
 - Perturbation stability [BE02, SSSS10,DL09,...]
- ② Algorithm based on subsampling always efficient

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies** **private** model selection
 - Two notions of stability
 - Subsampling stability [Shao96,Bach08,MB10,...]
 - Perturbation stability [BE02, SSSS10,DL09,...]
- ② Algorithm based on subsampling always efficient

Stability results for the LASSO estimator

- ① Consistency assumptions imply stability
- ② **Efficient** test of perturbation stability
- ③ Private algorithm for sparse linear regression

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies** **private** model selection
 - Two notions of stability
 - Subsampling stability
 - Perturbation stability [This talk]
- ② Algorithm based on subsampling always efficient

Stability results for the LASSO estimator

- ① Consistency assumptions imply stability
- ② **Efficient** test of perturbation stability [This talk]
- ③ Private algorithm for sparse linear regression

Our Results for Sparse Linear Regression

Stochastic setting [in this talk]:

- $\|\theta^*\|_\infty \leq 1$, $\text{supp}(\theta^*) \leq s$, min. non zero entry of $\theta^* = \Omega(1)$
- $\mathbf{x}_i \sim \mathcal{N}(0, 1)^p$, $w \sim \mathcal{N}(0, \sigma^2)^p$

Our Results for Sparse Linear Regression

Stochastic setting [in this talk]:

- $\|\theta^*\|_\infty \leq 1$, $\text{supp}(\theta^*) \leq s$, min. non zero entry of $\theta^* = \Omega(1)$
- $\mathbf{x}_i \sim \mathcal{N}(0, 1)^p$, $w \sim \mathcal{N}(0, \sigma^2)^p$

Algorithm	Sample complexity
1. Subsampling based [Kifer, Smith, T.'12]	$O\left(\frac{1}{\epsilon} s^2 \log^2 p\right)$
2. Subsampling based	$O^*\left(\frac{\log(1/\delta)}{\epsilon} s \log p\right)$
3. Perturbation stability based	$O^*\left(s \log p + \frac{\log(1/\delta)}{\epsilon} \lambda\right)$ λ : low-order terms in s and $\log p$ for small s

O^* hides poly log factors

Our Results for Sparse Linear Regression

Stochastic setting [in this talk]:

- $\|\theta^*\|_\infty \leq 1$, $\text{supp}(\theta^*) \leq s$, min. non zero entry of $\theta^* = \Omega(1)$
- $\mathbf{x}_i \sim \mathcal{N}(0, 1)^p$, $\mathbf{w} \sim \mathcal{N}(0, \sigma^2)^p$

Algorithm	Sample complexity
1. Subsampling based [Kifer, Smith, T.'12]	$O\left(\frac{1}{\epsilon} s^2 \log^2 p\right)$
2. Subsampling based	$O^*\left(\frac{\log(1/\delta)}{\epsilon} s \log p\right)$
3. Perturbation stability based	$O^*\left(s \log p + \frac{\log(1/\delta)}{\epsilon} \lambda\right)$ λ : low-order terms in s and $\log p$ for small s

O^* hides poly log factors

- Non-private optimal sample complexity: $s \log p$

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies** **private** model selection
 - Two notions of stability
 - Subsampling stability
 - **Perturbation stability**
- ② Algorithm based on subsampling always efficient

Stability results for the LASSO estimator

- ① Consistency assumptions imply stability
- ② **Efficient** test of perturbation stability
- ③ Private algorithm for sparse linear regression

Perturbation Stability based Model Selection

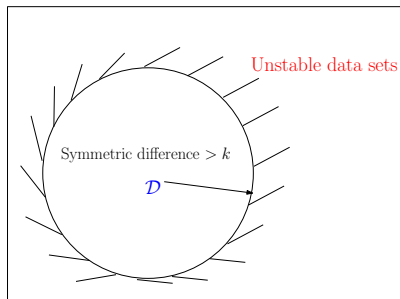
Distance to Instability Framework:

- A function $f : \mathcal{U}^* \rightarrow \mathcal{R}$ is **k-stable** at \mathcal{D} if
 - For any data set $\mathcal{D}' \in \mathcal{U}^*$, with $|\mathcal{D} \Delta \mathcal{D}'| \leq k$

$$f(\mathcal{D}) = f(\mathcal{D}')$$

- **Objective:** output $f(\mathcal{D})$ while preserving privacy

All data sets



Perturbation Stability based Model Selection

- 1 $\text{dist} \leftarrow \max_k (f(\mathcal{D}) \text{ is } k\text{-stable})$
- 2 $\widetilde{\text{dist}} \leftarrow \text{dist} + \text{Lap}(\frac{1}{\epsilon})$ [Think of the noise as constant]
- 3 If $\widetilde{\text{dist}} > \frac{\log(1/\delta)}{\epsilon}$, then **return** $f(\mathcal{D})$, o.w. **return** \perp

Perturbation Stability based Model Selection

- 1 $\text{dist} \leftarrow \max_k (f(\mathcal{D}) \text{ is } k\text{-stable})$
- 2 $\widetilde{\text{dist}} \leftarrow \text{dist} + \text{Lap}(\frac{1}{\epsilon})$ [Think of the noise as constant]
- 3 If $\widetilde{\text{dist}} > \frac{\log(1/\delta)}{\epsilon}$, then **return** $f(\mathcal{D})$, o.w. **return** \perp

Privacy guarantee

Theorem (variants in [DL09, KRSY11])

The algorithm is (ϵ, δ) -differentially private

Utility guarantee

Theorem

*If f is $\left(\frac{2 \log(1/\delta)}{\epsilon}\right)$ -stable w.r.t. \mathcal{D} , **then** the algorithm outputs $f(\mathcal{D})$*

Issue:

- For a given f , **distance to instability** might not efficiently be computable

Issue:

- For a given f , **distance to instability** might not efficiently computable

Fix: Obtain an efficiently computable proxy distance \hat{d}

- **Safety:** $\forall \mathcal{D}, \hat{d}(\mathcal{D}) \leq \text{distance to instability of } f$
- **\hat{d} is stable:** $\forall \mathcal{D}, \mathcal{D}' \text{ s.t. } |\mathcal{D} \Delta \mathcal{D}'| = 1,$

$$|\hat{d}(\mathcal{D}) - \hat{d}(\mathcal{D}')| = 1$$

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies private** model selection
 - Two notions of stability
 - Subsampling stability
 - Perturbation stability
- ② Algorithm based on subsampling always efficient

Stability results for the LASSO estimator

- ① Consistency assumptions imply stability
- ② **Efficient test of perturbation stability**
- ③ Private algorithm for sparse linear regression

Perturbation Stability of LASSO

Recall: LASSO estimator

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{2n} \|y - X\theta\|_2^2 + \frac{\Lambda}{n} \|\theta\|_1$$

Theorem (Consistency [...Wainwright06])

In the *stochastic setting*, under suitable choice of Λ and $n = \omega(s \log p)$, w.h.p. $\text{supp}(\hat{\theta}) = \text{supp}(\theta^*)$

Perturbation Stability of LASSO

Recall: LASSO estimator

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{2n} \|y - X\theta\|_2^2 + \frac{\Lambda}{n} \|\theta\|_1$$

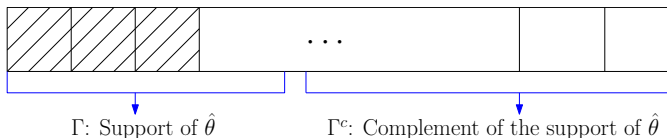
Theorem (Stability of LASSO)

Consistency assumptions on LASSO are sufficient to guarantee perturbation stability

Proof Idea:

- Analyze the KKT condition at $\hat{\theta}$
- Show that $\text{supp}(\hat{\theta})$ is stable by constructing “dual certificate” for nearby instances

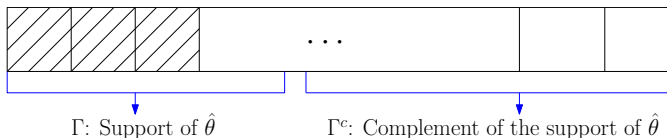
Stability Test for LASSO



Test for the following (real test is more complicated)

- **Restricted Strong Convexity (RSC)**: Restricted to subspace induced by Γ , the least-squared loss is strongly convex
 - **Note:** We only need to test strong convexity at Γ

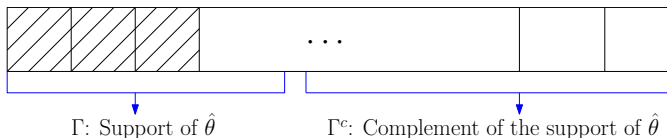
Stability Test for LASSO



Test for the following (real test is more complicated)

- **Restricted Strong Convexity (RSC)**: Restricted to subspace induced by Γ , the least-squared loss is strongly convex
 - **Note**: We only need to test strong convexity at Γ
- **“Strong stability”**: The coordinates of the gradient of the least-squared loss in Γ^c are $\ll \Lambda$

Stability Test for LASSO



Test for the following (real test is more complicated)

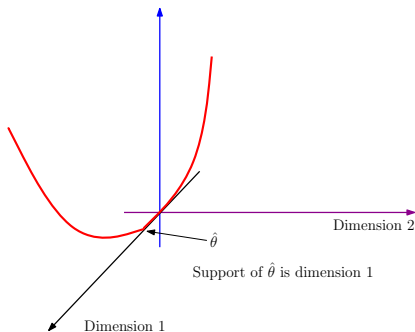
- **Restricted Strong Convexity (RSC)**: Restricted to subspace induced by Γ , the least-squared loss is strongly convex
 - **Note**: We only need to test strong convexity at Γ
- **“Strong stability”**: The coordinates of the gradient of the least-squared loss in Γ^c are $\ll \Lambda$

These conditions hold w.h.p. in the stochastic setting

Geometry of LASSO

Intuition:

- Strong sparsity and RSC implies stability of $\text{supp}(\hat{\theta})$
 - i.e., $\hat{\theta}$ does not shift *too much* by changing one entry in \mathcal{D}

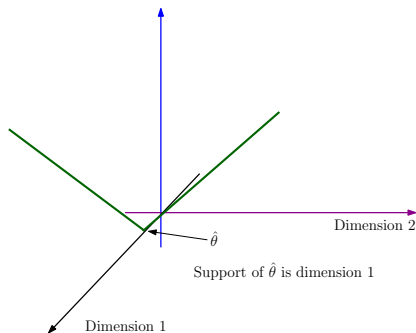


Strong convexity gives stability along *dimension 1*

Geometry of LASSO

Intuition:

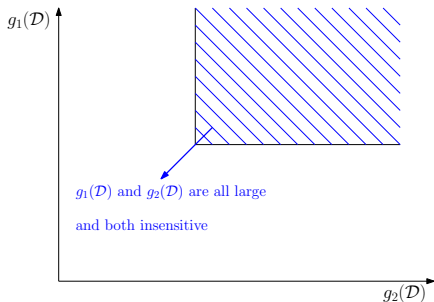
- Strong sparsity and RSC implies stability of $\text{supp}(\hat{\theta})$
 - i.e., $\hat{\theta}$ does not shift *too much* by changing one entry in \mathcal{D}



L_1 -regularization gives stability along *dimension 2*

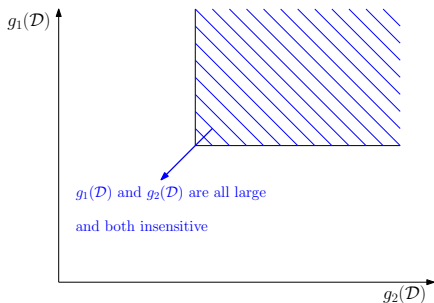
Making the Test Private (Simplified)

- Test for sparsity g_1 : Negative of the absolute value of $(s + 1)$ -th entry of $X^T(y - X\hat{\theta})$
- Test for strong convexity g_2 : Min. eigenvalue of $X_\Gamma^T X_\Gamma$



Making the Test Private (Simplified)

- Test for sparsity g_1 : Negative of the absolute value of $(s + 1)$ -th entry of $X^T(y - X\hat{\theta})$
- Test for strong convexity g_2 : Min. eigenvalue of $X_\Gamma^T X_\Gamma$



Our result:

- Can privately **test** if g_1 and g_2 are both large
- Serves as a proxy for the distance to instability \hat{d}

Putting the Pieces Together

- 1 Compute \hat{a} = function of $g_1(\mathcal{D})$ and $g_2(\mathcal{D})$
- 2 $\tilde{a} \leftarrow \hat{a} + \text{Lap}(\frac{1}{\epsilon})$ [Think of the noise as constant]
- 3 If $\hat{a} > \frac{\log(1/\delta)}{\epsilon}$, then **return** $f(\mathcal{D})$, o.w. **return** \perp

Privacy guarantee

Theorem (variants in [DL09, KRSY11])

The algorithm is (ϵ, δ) -differentially private

Utility guarantee

Theorem

*If \hat{a} is at least $\frac{2\log(1/\delta)}{\epsilon}$, **then** w.h.p. the algorithm outputs the support of the underlying parameter θ^**

New connections between stability and differential privacy

- ① **Stable** nonprivate model selection **implies** **private** model selection
 - Two notions of stability
 - Subsampling stability
 - Perturbation stability
- ② Algorithm based on subsampling always efficient

Stability results for the LASSO estimator

- ① Consistency assumptions imply stability
- ② **Efficient** test of perturbation stability
- ③ Private algorithm for sparse linear regression

Future work

- Analyze privacy/stability properties of iterative sparse estimation algorithms (e.g., [least angle regression](#))
- Privacy for high-dimensional learning where “good” model is not unique
 - Low rank matrix approximation
 - Sparse representation with many “good” sparse vectors