

Probabilistic Sparse Recovery: Approximate Message Passing Algorithm

Geethu Joseph

SPC Lab, IISc

June 4, 2016

Sparse Signal Recover Problem

$$\mathbf{y} = \mathbf{A} \mathbf{x} + \mathbf{w}$$

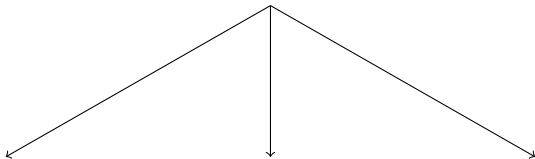
$m \times 1$ $m \times N$ $N \times 1$ $N \times 1$

$\|\mathbf{x}_k\|_0 < m < N$

$\sim \mathcal{N}(0, \chi\mathbf{I})$

Goal: recover the signal \mathbf{x} from observation \mathbf{y}

Sparse Recovery



Geometric¹

LASSO, l_p norm-based

Combinatorial²

IHT, CoSaMP, SP, OMP

Probabilistic³

AMP, SBL

¹[Donoho '04],[Candes-Tao '04, '06], [Rudelson-Vershynin'06]

²[Needell et al. '08] [Rubinstein et al. '09]

³[Zhou et al. '09] [Donoho et al. '09]

Probabilistic Approach

- 1 Impose a prior distribution on \mathbf{x}
- 2 Noise distribution gives $p(\mathbf{y}|\mathbf{x})$
- 3 Find the **MAP/MMSE estimate** using conditional pdf

$$\hat{\mathbf{x}}_{\text{MAP}} = \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y})$$

$$\hat{\mathbf{x}}_{\text{MMSE}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}\}$$

Motivation

- Impose Laplacian prior on \mathbf{x} , x_i are iid

$$p(\mathbf{x}) \cong \exp \left\{ -\frac{\lambda}{\chi} \|\mathbf{x}\|_1 \right\}$$

- Noise distribution

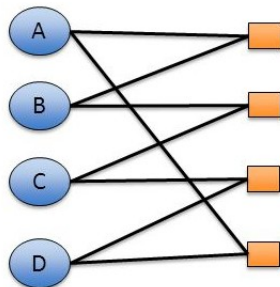
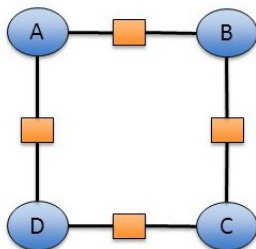
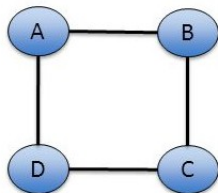
$$p(\mathbf{y}|\mathbf{x}) \cong \exp \left\{ -\frac{1}{2\chi} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 \right\}$$

- Laplacian MAP = LASSO

$$\begin{aligned} \hat{\mathbf{x}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \\ &= \arg \min_{\mathbf{x}} -\log p(\mathbf{x}, \mathbf{y}) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_1 \end{aligned}$$

Factor Graphs

- Joint probability of many variables factors into little pieces
- Nodes: random variables, factors
- No normalization: determine all probabilities



MAP/MMSE using factor graphs

- Construction of bipartite graph
 - ▶ Variable nodes: x_1, x_2, \dots, x_N
 - ▶ Factor nodes: $p(y_1|\mathbf{x}), p(y_2|\mathbf{x}), \dots, p(y_m|\mathbf{x})$
 - ▶ An edge between a variable node i and a factor node a if $A_{ai} \neq 0$
- **Goal:** Obtain marginal $p(x_i|\mathbf{y})$ or MAP estimate, $\arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y})$

Solution: Belief propagation

Belief Propagation

- **Intuition:** **Peer pressure** - A node determines a final belief distribution by listening to its neighbors.
- Keep passing messages (functions of the variable) until a stable belief state is reached
- Two versions
 - ▶ **Sum-product:** to compute marginals (MMSE), $p(\mathbf{x})$
 - ▶ **Max-product:** to compute maximizer (MAP), $\arg \max_{\mathbf{x}} p(\mathbf{x})$

Message Passing

- From a variable node to factor node

$$m_{i \rightarrow a} = p(x_i) \prod_{b \in \partial i \setminus a} m_{b \rightarrow i}(x_i)$$

- From factor node to variable node:

- 1 **Sum-product:** estimate marginals

$$m_{a \rightarrow i} = \int_{\mathbf{x}_{-i}} p(y_a | \mathbf{x}) \prod_{j \in \partial a \setminus i} m_{j \rightarrow a}(x_j) d\mathbf{x}_{-i}$$

★ **Marginal:** $p(x_i | \mathbf{y}) = p(x_i) \prod_{b \in \partial i} m_{b \rightarrow i}(x_i)$

- 2 **Max-product:** scoring functions whose maxima are most likely states

$$m_{a \rightarrow i} = \max_{x_i} p(y_a | \mathbf{x}) \prod_{j \in \partial a \setminus i} m_{j \rightarrow a}(x_j)$$

★ **Maximum:** $\hat{x}_i = \arg \max_{x_i} p(x_i) \prod_{b \in \partial i} m_{b \rightarrow i}(x_i)$

AMP Algorithms

- AMP: Belief propagation with carefully constructed approximations
 - ▶ The original AMP⁴: Laplacian MAP, with *iid* Gaussian noise, and *iid* matrix \mathbf{A}
 - ▶ The Bayesian AMP
 - ▶ The generalized AMP
 - ▶ Expectation-Maximization Bernoulli-Gaussian AMP
- Very fast: ≈ 15 iterations, simple operations

⁴[Donoho, Maleki, Montanari '09]

Design: Laplacian MAP estimate

- 1 Construct a joint distribution on data and write down the corresponding sum-product algorithm
- 2 Approximate the sum-product messages by the families with two scalar parameters
- 3 Use law of large numbers to approximate messages for the large system limit

Step 1: Basic Messages

- Recall: With $\mathbf{z} = \mathbf{Ax}$,

$$p(\mathbf{x}) \cong \prod_{i=1}^N \exp\{-\sigma^{-2} \lambda |x_i|\}$$
$$p(\mathbf{y}|\mathbf{x}) \cong \prod_{a=1}^m \exp\left\{-\frac{1}{2\chi} \|y_a - z_a\|\right\}$$

- Messages in negative log domain are

$$m_{i \rightarrow a}^{(t+1)}(x_i) = \lambda |x_i| + \sum_{b \in [m] \setminus a} m_{b \rightarrow i}^{(t)}(x_i)$$
$$m_{a \rightarrow i}^{(t)}(x_i) = \min_{x_j \in [M] \setminus i} \frac{1}{2} (y_a - z_a)^2 + \sum_{j \in [M] \setminus i} m_{j \rightarrow a}^{(t)}(x_j)$$

- Each iteration: $2Nm$ functions on real line

Assumptions

- $x_j \sim o(1)$
- $A_{ai}, iid \sim O(1/\sqrt{m})$
- \mathbf{A} is normalized
 - ▶ Column sum: $\sum_{a=1}^m A_{ai} = 0$
 - ▶ l_2 norm of column $\sum_{a=1}^m A_{ai}^2 = 1$
- m scales with N

Step 2: Approximation to Family of 2 Scalars

- Approximation to quadratic function
 - ▶ Expand the messages using Taylor series, around 0

$$\begin{aligned}m_{a \rightarrow i}^{(t)}(x_i) &= -\alpha_{a \rightarrow i}^{(t)} A_{ai} x_i + \beta_{a \rightarrow i}^{(t)} (A_{ai} x_i)^2 \\m_{i \rightarrow a}^{(t+1)}(x_i) &= \frac{1}{2\gamma_{i \rightarrow a}^{(t+1)}} \left(x_i - x_{i \rightarrow a}^{(t+1)} \right)^2\end{aligned}$$

- Each iteration: **$2Nm$ real numbers** functions of $\{\mathbf{y}, \mathbf{A}\}$

$$m_{i \rightarrow a}^{(t+1)}(x_i) \rightarrow \left[x_{i \rightarrow a}^{(t+1)} \quad \gamma_{i \rightarrow a} \right]$$

$$m_{a \rightarrow i}^{(t)}(x_i) \rightarrow \left[\alpha_{a \rightarrow i}^{(t)} \quad \beta_{a \rightarrow i}^{(t)} \right]$$

- Memory requirements $\sim O(mN)$

Further Simplification!

- Relation between messages:

$$r_{a \rightarrow i}^{(t)} = \frac{\alpha_{a \rightarrow i}^{(t)}}{\beta_{a \rightarrow i}^{(t)}} = y_a - \sum_{j \in [M] \setminus i} A_{aj} x_{j \rightarrow a}^{(t)}$$
$$x_{i \rightarrow a}^{(t+1)} = \eta \left(\sum_{b \in [m] \setminus a} A_{bi} r_{b \rightarrow i}^{(t)}; \theta_t \right)$$

- Law of large numbers: $\frac{\lambda}{\sum_{b \in [m] \setminus a} \beta_{b \rightarrow i}^{(t)} A_{bi}^2} \approx \theta_t$
- $\eta(a; \theta)$ is **soft thresholding function**

$$\eta(a; \theta) = \begin{cases} a - \theta & a > \theta \\ 0 & -\theta \leq a \leq \theta \\ a + \theta & a < -\theta \end{cases} \quad (1)$$

Step 3: Large System Limit

- Approximation:

$$r_{a \rightarrow i}^{(t)} \approx r_a^{(t)} + \delta r_{a \rightarrow i}^{(t)} + \mathcal{O}(1/N)$$

$$x_{i \rightarrow a}^{(t)} \approx x_i^{(t)} + \delta x_{i \rightarrow a}^{(t)} + \mathcal{O}(1/N)$$

- Approximation steps involve
 - 1 Taylor series expansion of $x_{i \rightarrow a}^{(t+1)}$
 - 2 Law of large numbers
- Memory requirements $\sim \mathcal{O}(m + N)$
- On convergence, declare estimate as $x_i^{(t)}$

AMP Algorithm

- After simplifications,

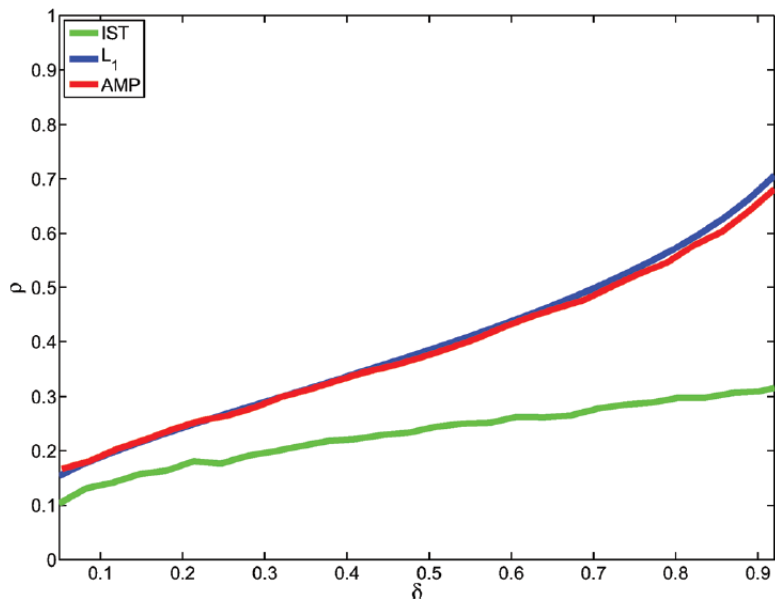
$$\begin{aligned}\mathbf{x}^{t+1} &= \eta\left(\mathbf{x}^t + \mathbf{A}^T \mathbf{r}^{(t)}; \theta_t\right) \\ \mathbf{r}^t &= \mathbf{y} - \mathbf{A}\mathbf{x} + b_t \mathbf{r}^{(t-1)}\end{aligned}$$

where $b_t \approx \frac{1}{m} \sum_{i=1}^N \eta'(\mathbf{x}_i^{(t-1)} + [\mathbf{A}^T \mathbf{r}^{(t-1)}]_i; \theta_{t-1})$

- Similar to iterative soft thresholding:

$$\begin{aligned}\mathbf{x}^{t+1} &= \eta\left(\mathbf{x}^t + \mathbf{A}^T \mathbf{r}^{(t)}; \theta_t\right) \\ \mathbf{r}^t &= \mathbf{y} - \mathbf{A}\mathbf{x}\end{aligned}$$

Phase Transition Curve



Covergence

Theorem

^a For Laplacian prior, and \mathbf{A} is Gaussian with iid entries, then (for n large enough) AMP converge within relative distance ϵ from a minimizer in $t = O(\log(1/\epsilon))$ iterations.

^a[Bayati, Montanari '11]