# Approximate Message-Passing for Convex Optimization with Non-Separable Penalties

Andre Manoel[†,‡], Florent Krzakala[♯], Bertrand Thirion[†], Gaël Varoquaux[†] and Lenka Zdeborová[‡]

† Parietal Team, Inria, Neurospin, CEA, ‡ Institut de Physique Théorique, CEA and ♯ Laboratoire de Physique Statistique, ENS

## Convex optimization with non-separable penalties

▷ **Problem statement**— we consider the minimization of an objective consisting of a quadratic loss and a non-separable penalty

$$\arg \min_{\mathbf{x}} \quad \frac{1}{2}\|\mathbf{y} - A\mathbf{x}\|_2^2 + \lambda \sum_{k=1}^{R} f\big((K\mathbf{x})_k\big)$$

for given $\mathbf{y} \in \mathbb{R}^N$, $A \in \mathbb{R}^{N \times P}$ and $K \in \mathbb{R}^{R \times P}$. Prominent examples are the total variation (TV) penalty and the cosparse analysis model

▷ **Better algorithms?**— proximal algorithms such as FISTA and ADMM are the state-of-the-art for performing this minimization. However, both have issues

- in FISTA, convergence is slow due to the inner loop requiring more and more iterations
- in ADMM, the behavior is highly dependent on the stepsize, which is hard to set

We thus look for alternative approaches that are hopefully faster and/or require less parameters

▷ **Our approach**— promising new class of algorithms: *approximate message-passing* (AMP)

Idea: adapt the vector approximate-message passing (VAMP) algorithm [Rangan et al. 2017] and the expectation-consistent (EC) approximation [Opper and Winther 2005] for non-separable penalties

We rederive the iteration from scratch and benchmark it on standard datasets: promising results!

## Approximate message-passing

▷ **Probabilistic framework**— given the following probability distribution

$$P(\mathbf{x}|A, \mathbf{y}) = \frac{1}{\mathcal{Z}} e^{-\frac{1}{2}\|\mathbf{y}-A\mathbf{x}\|_2^2} \prod_{j=1}^{P} e^{-\lambda f(x_j)}$$

the AMP algorithm [Donoho et al. 2009] is able to compute the MAP estimator by means of the following iteration

$$\mathbf{x}^{t+1} = \eta_{\lambda\sigma^t}(\mathbf{x}^t + A^T\mathbf{z}^t)$$
$$\mathbf{z}^t = \mathbf{y} - A\mathbf{x}^t + \frac{1}{\alpha}\mathbf{z}^{t-1}\big\langle \nabla\eta_{\lambda\sigma^t}(\mathbf{x}^t + A^T\mathbf{z}^t)\big\rangle$$

where $\eta_\lambda(v) = \text{prox}_{\lambda f}(v)$. For a $\ell_1$ penalty, $f(x) = |x|$: soft-thresholding

A lot like ISTA, but w/ an additional term and an adaptive stepsize based on the variance of $\mathbf{x}$. Usually faster, however: <u>convergence issues!</u>

▷ **The vector AMP (VAMP) algorithm** [Rangan et al. 2017]— more robust than AMP; MAP estimator comes from

$$\mathbf{x}^t = (A^T A + \rho^t I_N)^{-1}(A^T\mathbf{y} + \mathbf{u}^t), \qquad \mathbf{z}^t = \eta_{\lambda\sigma_x^t/(1-\sigma_x^t\rho^t)}\left(\frac{\mathbf{x}^t - \sigma_x^t\mathbf{u}^t}{1 - \sigma_x^t\rho^t}\right),$$
$$\mathbf{u}^{t+1} = \mathbf{u}^t + (\mathbf{z}^t/\sigma_z^t - \mathbf{x}^t/\sigma_x^t), \qquad \rho^{t+1} = \rho^t + (1/\sigma_z^t - 1/\sigma_x^t).$$

A lot like ADMM (more precisely, the Peaceman-Rachford splitting) with, once again, an <u>adaptive stepsize</u> based on variance of $\mathbf{x}$. Upon convergence, $\mathbf{x} = \mathbf{z}$

▷ **Some facts about VAMP**—

▷ as in ADMM, $A^T A$ can be replaced by its eigendecomposition, so that matrix inverse is not necessary (however: extra preprocessing costs)

▷ is not able in principle to deal with losses other than quadratic (but can be adapted [He et al. 2017]), nor <u>non-separable penalties</u>
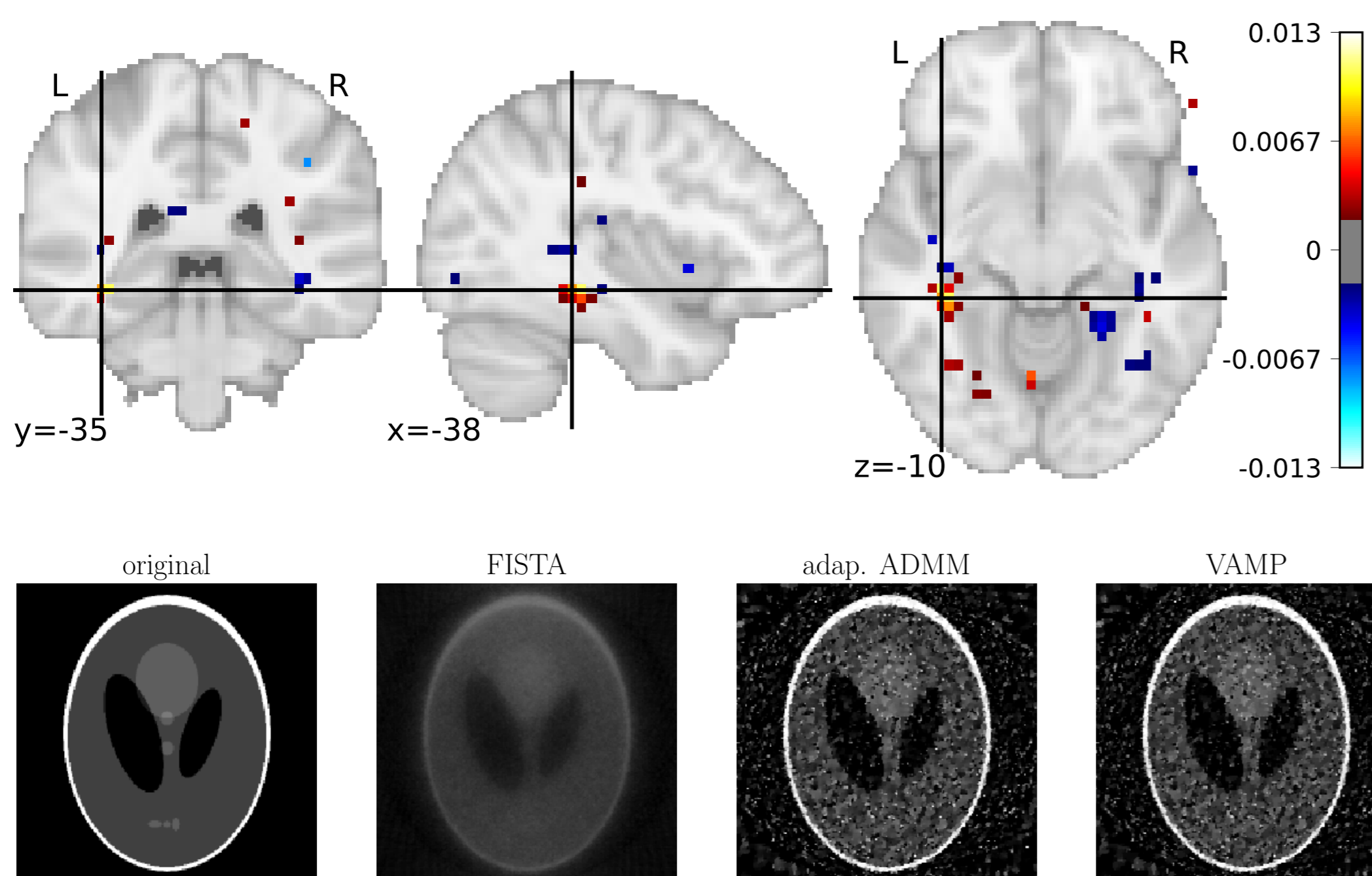




Figure: *Left:* **Sample of results obtained using proposed iteration** with TV penalties. Top: classification on Haxby, final result. Bottom: tomography on the Shepp-Logan phantom (bottom), 10s after preprocessing
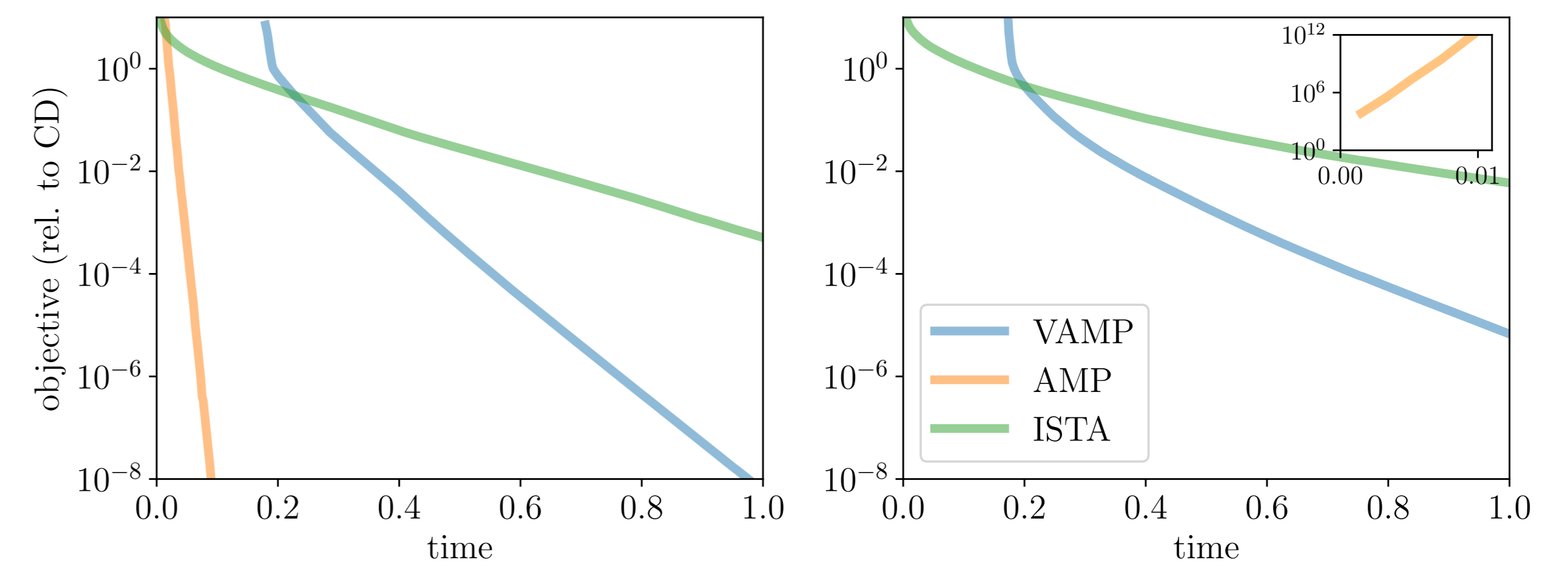


Figure: **Comparison between AMP and VAMP using $\ell_1$ regularization** on synthetic data: for i.i.d. Gaussian matrices (left) AMP is faster, but already for products of Gaussian i.i.d. matrices it diverges (right)

## Adapting VAMP to non-separable penalties

▷ **The expectation-consistent (EC) approximation** [Opper and Winther 2005]— given a probability distribution

$$P(\mathbf{x}) = \frac{1}{\mathcal{Z}} P_\ell(\mathbf{x}) \, P_r(\mathbf{x})$$

the (negative) log-partition function $-\log \mathcal{Z}$ is approximated by

$$\mathcal{F}[Q_\ell, Q_r] = -\log \int d\mathbf{x} \, P_\ell(\mathbf{x}) Q_r(\mathbf{x}) - \log \int d\mathbf{x} \, Q_\ell(\mathbf{x}) P_r(\mathbf{x}) + \log \int d\mathbf{x} \, Q_\ell(\mathbf{x}) Q_r(\mathbf{x})$$

for a tractable choice of $Q_{\ell,r}$, typically

$$Q_{\ell,r}(\mathbf{x}) = \exp\left(-\frac{1}{2}\rho_{\ell,r}\mathbf{x}^T\mathbf{x} + \mathbf{u}_{\ell,r}^T\mathbf{x}\right)$$

One must then optimize over $\mathbf{u}_{\ell,r}$ and $\rho_{\ell,r}$. <u>VAMP performs this optimization via a fixed-point iteration</u>

If one considers instead $[P_\ell(\mathbf{x}) P_r(\mathbf{x})]^\beta$ and takes the limit $\beta \to \infty$, the MAP estimate is recovered from $\mathbb{E}\mathbf{x} = -\nabla_{\mathbf{u}_{\ell,r}}\mathcal{F}$

▷ **Adapting to TV**— we introduce a new variable $\mathbf{z} = K\mathbf{x}$ and use the EC approximation not on $P(\mathbf{x}|A, \mathbf{y})$ but on

$$P(\mathbf{z}|A, \mathbf{y}) \propto \int d\mathbf{x} \, e^{-\frac{1}{2}\|\mathbf{y}-A\mathbf{x}\|_2^2} \delta(\mathbf{z} - K\mathbf{x}) \prod_{k=1}^{R} e^{-\lambda f(z_k)}.$$

The following VAMP-like iteration can be derived in this case

$$\mathbf{x}^t = (A^T A + \rho^t K^T K)^{-1}(A^T\mathbf{y} + K^T\mathbf{u}^t) \qquad \mathbf{z}^t = \eta_{\lambda\sigma_x^t/(1-\sigma_x^t\rho^t)}\left(\frac{K\mathbf{x}^t - \sigma_x^t\mathbf{u}^t}{1 - \sigma_x^t\rho^t}\right)$$
$$\mathbf{u}^{t+1} = \mathbf{u}^t + (\mathbf{z}^t/\sigma_z^t - K\mathbf{x}^t/\sigma_x^t) \qquad \rho^{t+1} = \rho^t + (1/\sigma_z^t - 1/\sigma_x^t)$$



$$P(\mathbf{y}|A\mathbf{x}) \qquad \delta(K\mathbf{x} - \mathbf{z}) \qquad P(\mathbf{z})$$

## Benchmarking the proposed iteration

▷ **Benchmarks**— we use a TV penalty ($K = \nabla$) and approach two problems:

▷ one vs. all <u>classification</u> on task fMRI on the Haxby dataset ($N = 1452$, $P = 136840$), 3 labels: "face", "house" and "chair"

▷ <u>tomography</u> on noisy projections of the Shepp-Logan phantom ($P = 40000$), 1% SNR noise

▷ **Tricks to speed up iteration**— in the $N \ll P$ setting: Woodbury formula, $K^T K = \Delta$ diagonal in Fourier basis (use FFT instead)

▷ **Some inconvenients**— as with PRS, convergence is not assured: relaxation parameter in the updates of $\mathbf{u}$ and $\rho$

▷ **Perspectives**— losses other than quadratic, imposing monotonicity, confidence interval from variance estimates, more experiments

## References

[1] Rangan, S., Schniter, P. & Fletcher, A. K. (2017). Vector approximate message passing. In IEEE International Symposium on Information Theory (pp. 1588-1592).

[2] Opper, M. & Winther, O. (2005). Expectation consistent approximate inference. Journal of Machine Learning Research, 6 (Dec), 2177-2204.

[3] Donoho, D. L., Maleki, A. & Montanari, A. (2009). Message-passing algorithms for compressed sensing. *PNAS*.

[4] Manoel, A., Krzakala, F., Thirion, B., Varoquaux, G. & Zdeborová, L. (2018). Approximate message-passing for convex optimization with non-separable penalties. In preparation.
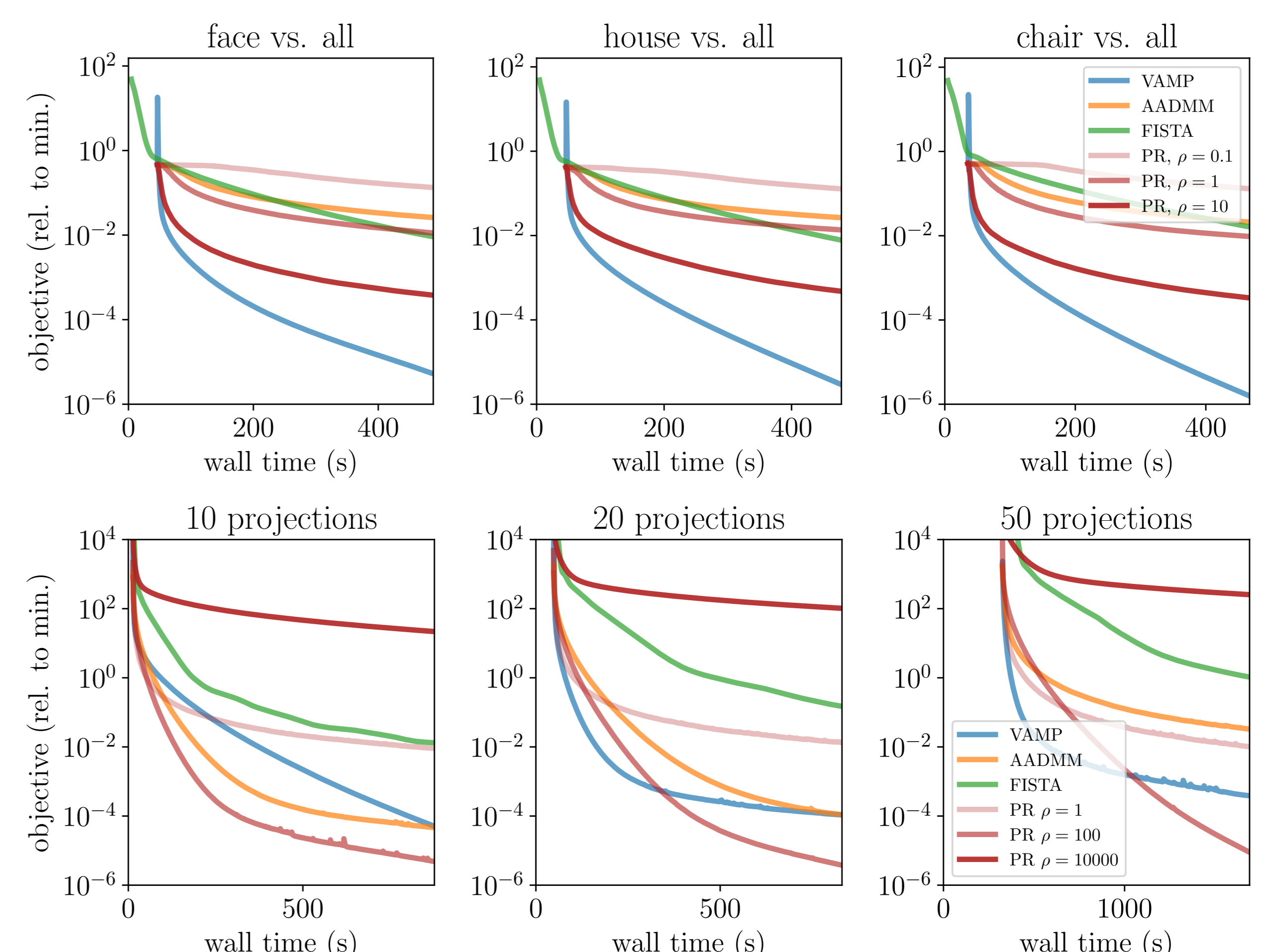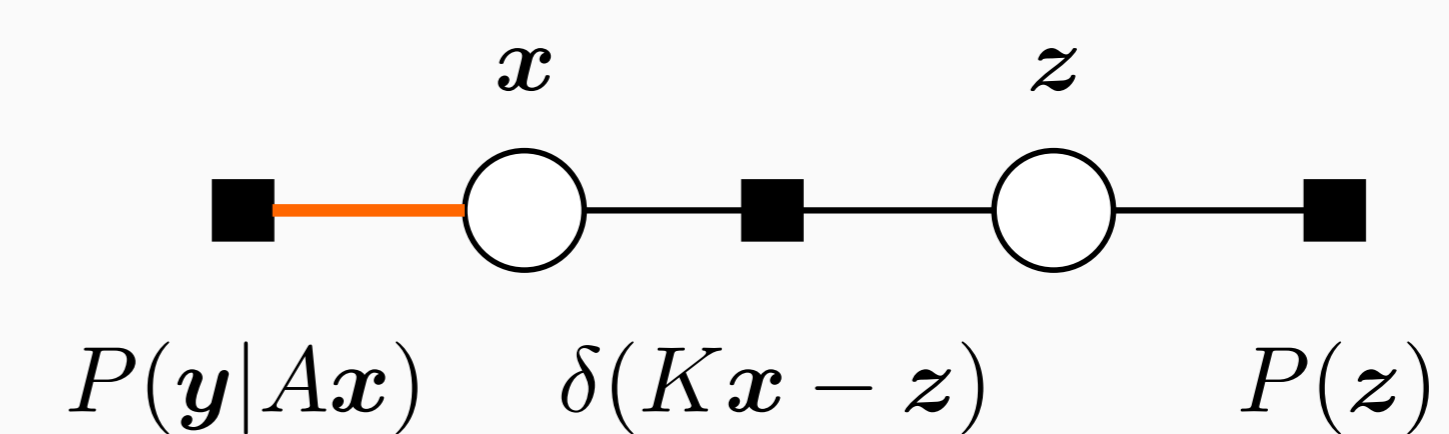
Figure: **Comparison between different approaches using TV penalties**, for classification (top) and tomography (bottom). VAMP is competitive, and often faster than other approaches