

Learning theory for matrix inverse problems with gaussian priors

Jonas Adler
KTH, Elekta

Olivier Verdier
KTH

1 Maximum a posteriori estimate

Suppose $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear forward operator (a matrix) and we have data of the form

$$\begin{aligned}x &\in \mathcal{N}(0, \Sigma) \\ y &\in \mathcal{N}(Ax, \Gamma)\end{aligned}$$

Theorem 1.1 ([Mu2012]). *The maximum a posteriori (MAP) estimate of x given y*

$$A_{MAP}^{-1}(y) = \arg \max_{x'} P(x'|y)$$

is

$$A_{MAP}^{-1}(y) = (A^T \Gamma^{-1} A + \Sigma^{-1})^{-1} A^T \Gamma^{-1} y$$

Proof. We have

$$P(x'|y) \propto \underbrace{P(y|x')}_{\mathcal{N}(Ax', I)} \underbrace{P(x')}_{\mathcal{N}(0, I)}$$

Instead of maximizing the posterior, we minimize the log of the posterior

$$\arg \max_{x'} P(x'|y) = \arg \min \log P(x'|y) = \arg \min (\log P(y|x') + \log P(x'))$$

where we let

$$\begin{aligned}f(x') &= \log P(y|x') + \log P(x') \\ &= \frac{1}{2} \|Ax - y\|_{\Gamma^{-1}}^2 + \frac{1}{2} \|x\|_{\Sigma^{-1}}^2\end{aligned}$$

we compute

$$\nabla f = A^T \Gamma^{-1} (Ax - y) + \Sigma^{-1} x$$

Setting this to zero gives the conditions for a maximum

$$(A^T \Gamma^{-1} A + \Sigma^{-1}) x = A^T \Gamma^{-1} y$$

which gives a maximum likelihood solution as

$$x = (A^T \Gamma^{-1} A + \Sigma^{-1})^{-1} A^T \Gamma^{-1} y$$

□

Connection to Tikhonov regularization: By identification, this is equivalent to solving the Tikhonov regularized problem:

$$\min_x \|\Gamma^{-1/2}(Ax - y)\| + \|\Sigma^{-1/2}x\|$$

2 Best function estimate

Theorem 2.1 ([Bi2006]). *We have*

$$\arg \min_B \mathbb{E}_{(x,y)} \|B(y) - x\|_2^2 = \mathbb{E}_x(x | y)$$

where the minimization is taken over all functions $B: \mathbb{R}^m \rightarrow \mathbb{R}^n$.

Proof. By the law of total probability

$$\mathbb{E}_{(x,y)} \|B(y) - x\|_2^2 = \mathbb{E}_y(\mathbb{E}_x(\|B(y) - x\|_2^2 | y))$$

By the monotonicity of the expectation we have

$$B(y) = \arg \min_z \mathbb{E}_x(\|z - x\|_2^2 | y)$$

we find this by differentiating and setting to zero

$$0 = \mathbb{E}_x(2(z - x) | y) = 2\mathbb{E}_x(z | y) - 2\mathbb{E}_x(x | y)$$

which gives

$$B(y) = \mathbb{E}_x(x | y)$$

□

Theorem 2.2. *A neural network, which we denote by A_{NN}^{-1} , will in the limit of infinite capacity and trained to solve the inverse problem equation 1 according to*

$$\arg \min_{A_{NN}^{-1}} \mathbb{E}_{(x,y)} \|A_{NN}^{-1}(y) - x\|_2^2$$

is given by the MAP estimate:

$$A_{NN}^{-1}(y) = (A^T \Gamma^{-1} A + \Sigma^{-1})^{-1} A^T \Gamma^{-1} y$$

Proof. By Theorem 2.1, we have that the optimum solution of the training problem is given by the conditional expectation. Since neural networks are universal function approximators, this will be attainable in the limit of infinite capacity.

Further, for Gaussian distributions the mean and the mode coincide. Thus the conditional expectation, being a gaussian, is the MAP estimate. □

Denoising Consider the case of denoising where instead of observing y as in 1, we first apply the pseudo-inverse

$$z = A^\dagger y$$

By standard arguments this is also a gaussian and we find

$$z \in \mathcal{N}(A^\dagger Ax, A^\dagger \Gamma (A^\dagger)^T)$$

Hence the MAP estimate is (with $A^\dagger A = P_{A^*}$, $AA^\dagger = P_A$ being projection onto the range of A^* and A respectively) via straightforward application of theorems 1.1 and 2.1:

$$\begin{aligned} A_{\text{denoising}}^{-1}(A^\dagger y) &= ((A^\dagger A)^T (A^\dagger \Gamma (A^\dagger)^T)^\dagger (A^\dagger A) + \Sigma^{-1})^\dagger (A^\dagger A)^T (A^\dagger \Gamma (A^\dagger)^T)^\dagger A^\dagger y \\ &= (P_{A^*}^T (A^T \Gamma^\dagger A) P_A + \Sigma^{-1})^\dagger P_{A^*}^T (A^T \Gamma^\dagger A) A^\dagger y \\ &= ((AP_{A^*})^T \Gamma^\dagger (AP_{A^*}) + \Sigma^{-1})^\dagger (AP_{A^*})^T \Gamma^\dagger P_A y \end{aligned}$$

Now via identification, $Az = AA^\dagger y$, so This shows that:

$$A_{\text{denoising}}^{-1}(A^\dagger y) = (AP_{A^*})_{NN}^{-1}(P_A y)$$