

Fuzzy Techniques Provide a Theoretical Explanation for the Heuristic ℓ_p -Regularization of Signals and Images

Fernando Cervantes, Bryan Usevitch
Department of Electrical and Computer Engineering
University of Texas at El Paso
500 W. University
El Paso, TX 79968, USA
fcervantes@miners.utep.edu
usevitch@utep.edu

Leobardo Valera, Vladik Kreinovich,
and Olga Kosheleva
Computational Science Program
University of Texas at El Paso
El Paso, TX 79968, USA
leobardovalera@gmail.com, vladik@utep.edu
olgak@utep.edu

Abstract—One of the main techniques used to de-noise and de-blur signals and images is *regularization*, which is based on the fact that signals and images are usually smoother than noise. Traditional Tikhonov regularization assumes that signals and images are differentiable, but, as Mandelbrot has shown in his fractal theory, many signals and images are not differentiable. To de-noise and de-blur such images, researchers have designed a heuristic method of ℓ^p -regularization.

ℓ^p -regularization leads to good results, but it is not used as widely as should be, because it lacks a convincing theoretical explanation – and thus, practitioners are often reluctant to use it, especially in critical situations. In this paper, we show that fuzzy techniques provide a theoretical explanation for the ℓ^p -regularization.

Fuzzy techniques also enables us to come up with natural next approximations to be used when the accuracy of the ℓ^p -based de-noising and de-blurring is not sufficient.

I. INTRODUCTION

Measurement results are, in general, somewhat different from the actual values of the measured quantities. Our information about the physical world comes from measurements. Measurement results are, in general, different from the actual values of the corresponding quantities.

For example, when we measure the values of the quantity of interest at different moments of time, then the measurement results y_1, y_2 , etc., are, in general, different from the actual values x_1, x_2, \dots , of the physical quantity at the corresponding moments of time.

Two main reasons for measurement errors. There are two main reasons for this difference between the actual values x_k and the observed values y_k :

- first, there is usually an additive noise n_k ;
- second, there is inertia (“blur”): even if the actual value changes abruptly, it takes some time for the sensor to capture this change.

As a result of the noise, the value y_k is different from the value x_k . As a result of inertia, we have what is called a blur: the measurement result y_k depends not only on the current

value x_k of the physical quantity, but also on the previous values x_{k-1}, x_{k-2}, \dots . Usually, the dependence of y_k on x_i is linear, so we conclude that

$$y_k = h_{k,0} \cdot x_k + h_{k,1} \cdot x_{k-1} + \dots + n_k,$$

for some coefficients $h_{k,i}$. The coefficients $h_{k,0}, h_{k,1}, \dots$ describe how the measuring instrument distorts the signal in the k -th moment of time.

Need for signal reconstruction. Once we get the measurement results y_k , we would like to reconstruct the actual values x_k from y_k as accurately as possible.

Continuous approximation. In practice, we only measure finitely many values x_k . However, in many practical situations, the time difference between the two consequent measurement is so small that, in effect, we have a continuous dependence $y(t)$ of the measured values on time t .

In this continuous approximation, the sum turns into an integral, so we have the dependence

$$y(t) = \int h(t, s) \cdot x(s) ds + n(t).$$

Need for image reconstruction. A similar problem can be formulated for image processing. When we observe an image, we usually observe the image intensity values $s(i, j)$ at different locations (i, j) on a rectangular grid, i.e., at a spatial location $(u_0 + i \cdot \Delta u, v_0 + j \cdot \Delta v)$, where (u_0, v_0) is the starting point and Δu and Δv are distances between the neighboring pixels in the u - and v -directions.

Similarly to signals, each observed value $s(i, j)$ is, in general, different from the actual (desired) value $I(i, j)$ of the corresponding intensity. First, there is noise (measurement error), and second, the image is blurred, in the sense that the observed signal $s(i, j)$ reflects not only the actual intensity $I(i, j)$ at the same spatial location (i, j) , but also the intensities

$I(i', j')$ at nearby locations. Under the assumption that the dependence of $s(i, j)$ on I is linear, we conclude that

$$s(i, j) = \sum_{i', j'} h(i, j, i', j') \cdot I(i', j') + n(i, j)$$

for some coefficients $h(i, j, i', j')$, where $n(i, j)$ denotes the (additive) noise.

Based on the observed image $s(i, j)$, we need to reconstruct the original image $I(i, j)$.

In many practical situations, there is a big need for signal and image reconstruction. When we take a photo of a friend with a modern sophisticated cell phone camera, this blur is barely visible – and does not constitute a serious problem. However, when a spaceship takes a photo of a distant planet, the blur is very visible – and needs to be eliminated. In such situations, we need to reconstruct the original image $I(x, y)$ from the blurred image $s(x, y)$.

Continuous approximation. In the continuous approximation, we need to reconstruct the intensity $I(x, y)$ at different spatial locations (x, y) from the observed signal $s(x, y)$:

$$s(x, y) = \int h(x, y, x', y') \cdot I(x', y') dx' dy',$$

for appropriate weights $h(x, y, x', y')$.

Regularization as a de-noising technique. One of the main ideas behind de-noising is that:

- a signal or an image is usually rather smooth, in the sense that the intensities $x(t)$ and $x(t')$ (or $I(x)$ and $I(x')$) at neighboring points t and t' are usually close to each other, while
- the noise is usually not smooth: the effects of noise on two neighboring points may be drastically different.

It is therefore reasonable, when we reconstruct an image from observation, to impose an additional constraint that the resulting image should be, in some reasonable sense, smooth. This introduction of the additional constraint is known as *regularization*.

Traditional regularization. The usual regularization – first introduced by Tikhonov – imposes the smoothness constraint, which in case of signals has the form

$$\int \left(\frac{dx}{dt} \right)^2 dt \leq C$$

and in case of an image has the form

$$\int |\nabla I(x)|^2 dx = \int \left(\left(\frac{\partial I}{\partial x} \right)^2 + \left(\frac{\partial I}{\partial y} \right)^2 \right) dx dy \leq C$$

for some constant C , where ∇I denotes the gradient of the image $I(x)$; see, e.g., [10].

In practice, the signal is given as a 1-D array of values x_i at different moments of time

$$t_i = t_0 + i \cdot \Delta t,$$

and the image is given as a 2-D array of values $I_{i,j}$ at different points

$$(x_i, y_j) = (x_0 + h_x \cdot i, y_0 + h_y \cdot j)$$

on a rectangular grid. In this case, we should use a discrete approximation to the derivatives, i.e., impose a discretized constraint, which for signals takes the form

$$\sum_i (d_i)^2 \leq C,$$

where

$$d_i \stackrel{\text{def}}{=} \frac{x_i - x_{i-1}}{\Delta t},$$

and for images takes the form

$$\sum_{i,j} |\nabla I_{i,j}|^2 \leq C,$$

where

$$|\nabla I_{i,j}|^2 \stackrel{\text{def}}{=} ((\Delta_x I_{i,j})^2 + (\Delta_y I_{i,j})^2) \leq C,$$

$$\Delta_x I_{i,j} \stackrel{\text{def}}{=} \frac{I_{i,j} - I_{i-1,j}}{h_x} \quad \text{and} \quad \Delta_y I_{i,j} \stackrel{\text{def}}{=} \frac{I_{i,j} - I_{i,j-1}}{h_y}.$$

Tikhonov regularization is not always adequate. Tikhonov regularization assumes that the actual signals and images are differentiable. Real-life signals images are often rather smooth, but not differentiable: this was one of the main discoveries of Benoit Mandelbrot, the father of fractals; see, e.g., [6].

For such signals and images, Tikhonov regularization distorts them, by making them too smooth.

How to make regularization more adequate: a heuristic ℓ_p -idea. To make regularization more adequate, researchers proposed to replace Tikhonov's term with a slightly different term

$$\sum_i |d_i|^p$$

for signals or

$$\sum_{i,j} |\nabla I_{i,j}|^p$$

for images, for an appropriate value $p < 2$; see, e.g., [2], [4], [5].

Advantages and limitations of the ℓ_p idea. the main advantage of the above ℓ_p -idea is that it works: in many real-life cases, we get a much better de-noising and de-blurring than with the Tikhonov regularization.

However, this method also has two major limitations. The first limitation is that this method is a heuristic, it does not have a convincing theoretical justification – and, as a result, practitioners are not very willing to use it in critical situations.

The second related limitation is that we do not know what to do when the ℓ_p -method does not work well. For theoretically justified methods, the next is often clear; for example:

- linear models are justified by the possibility of Taylor expansion,

- so if a linear model is not adequate enough, we can try quadratic models, cubic, etc.

For the ℓ_p -regularization, however, we do not have a theoretical explanation and, as a result, there is no clear next approximation.

What we do in this paper. In this paper, we show that fuzzy techniques – a known methodology for translating imprecise (“fuzzy”) expert knowledge into precise terms (see, e.g., [3], [8], [11]) – leads to a theoretical explanation for the ℓ_p -heuristic.

We also show that this theoretical explanation leads to a natural next approximation.

II. FORMALIZING THE PROBLEM

Towards formalizing the problem. We want to describe the requirement that the neighboring values are close, e.g., that the values x_i and x_{i-1} (or $I_{i,j}$ and $I_{i-1,j}$) are close to each other. In other words, we want to describe the requirement that the difference $d \stackrel{\text{def}}{=} x_i - x_{i-1}$ between the neighboring values is small.

“Small” is a relative notion: a small building is much taller than a small dog. So, to properly formalize this notion, we need to explicitly take into account the corresponding scale σ . In other words, since it is not possible to provide a single description for smallness, we would like to have descriptions of the notion “small, of size σ ” corresponding to different scales σ .

In fuzzy logic, each property is characterized by its *membership function* that assigns, to each possible value of the corresponding quantity x , the degree $\mu(x) \in [0, 1]$ to which, in the expert’s opinion, the value x satisfies the given property. Thus, we need, for each scale σ , to come up with a function $\mu_\sigma(d)$ to which the value d is small of size σ .

Let us describe reasonable restrictions of these functions.

Monotonicity. The larger the difference d , the smaller our degree of confidence that this difference d is small. Thus, for each σ , the function $\mu_\sigma(d)$ should be a decreasing function of d .

Continuity. Very small changes in d and σ should not affect our degree of belief $\mu_\sigma(d)$ that d is small of size σ . Thus, the function $\mu_\sigma(d)$ should be a continuous function of both its variables σ and d .

Scale-invariance. The numerical values of all physical quantities depend on the choice of the measuring unit. For example, if, instead of meters, we start using centimeters to describe distances, the distances will not change but their numerical values will all multiply by 100.

In general, if we replace the original measuring unit with a new unit which is λ times smaller, all the numerical values are multiplied (“re-scaled”) by this factor λ .

Since changing the units does not change the physics, it makes sense to require that all our conclusions should also not change if we simply change the measuring unit. In other words, all our conclusions should be scale-invariant.

In our case, this means that the value $\mu_\sigma(d)$ should not change if we use a different unit for measuring intensity. Under a different unit, instead of the difference d , we have $d' = \lambda \cdot d$, and instead of the scale σ , we have the new value $\sigma' = \lambda \cdot \sigma$. Thus, we must have

$$\mu_\sigma(d) = \mu_{\sigma'}(d') = \mu_{\lambda \cdot \sigma}(\lambda \cdot d).$$

In particular, for $\lambda = \sigma^{-1}$, when we use the original scale σ as the measuring unit, we get

$$\mu_\sigma(d) = \mu_1\left(\frac{d}{\sigma}\right).$$

Let us start with the simplest 1-D case. Let us start with the simplest case, when all the membership functions describing closeness form a 1-parametric family.

Since the functions $\mu_1\left(\frac{d}{\sigma}\right)$ corresponding to different scales σ already form a 1-D family, this means that we will consider only the functions from this family.

We usually have several experts. Our goal is to describe the expert’s opinion re what is small. Usually, we have several experts, so we need to combine their knowledge.

Different experts may provide different scales σ_i of smallness. So, for the same difference d , we may get different degrees $\mu_1\left(\frac{d}{\sigma_i}\right)$ to which d is small. We want to take into account the opinion of all these experts. In other words, we want to say that d is small in the opinion of the first expert (i.e., of size σ_1) and in the opinion of the second expert (i.e., of size σ_2), etc.

In fuzzy logic, once we know the degrees of confidence s_1, s_2, \dots , in different statements S_1, S_2, \dots , to estimate our degree of confidence s in the corresponding “and”-statement $S_1 \& S_2 \& \dots$, we need to use an appropriate “and”-operation (a.k.a. t-norm) $f_{\&}(s_1, s_2, \dots)$.

Thus, in general, to describe the expert’s opinion about smallness, we need to use membership functions of the type

$$f_{\&}\left(\mu_1\left(\frac{d}{\sigma_1}\right), \mu_1\left(\frac{d}{\sigma_2}\right), \dots\right).$$

We have made a simplifying assumption that all membership functions should belong to a 1-parametric family $\mu_1\left(\frac{d}{\sigma}\right)$. Thus, for every set of values $\sigma_1, \sigma_2, \dots$, there should exist a single value σ for which

$$f_{\&}\left(\mu_1\left(\frac{d}{\sigma_1}\right), \mu_1\left(\frac{d}{\sigma_2}\right), \dots\right) = \mu_1\left(\frac{d}{\sigma}\right).$$

What we do next. We have described reasonable constraints on the membership function. Now, we need to find membership functions that satisfy these constraints.

III. ANALYSIS OF THE PROBLEM AND THE MAIN RESULT

Reduction to the product t-norm. In general, there are many different t-norms. It is known, however (see, e.g., [7]) that each continuous t-norm can be approximated, with any given accuracy, by a so-called strict Archimedean t-norm, i.e., a t-norm that has the form $f_{\&}(a, b) = g^{-1}(g(a) \cdot g(b))$ for some strictly increasing continuous function $g(x)$.

Thus, for all practical purposes, we can safely assume that the actual t-norm is strictly Archimedean. For this t-norm, if we use “re-scaled” degrees of confidence $m(x) \stackrel{\text{def}}{=} g(\mu_1(x))$, the “and”-operation turns into a product and thus, the above requirement takes the following form.

Definition. We say that a continuous strictly decreasing function $m(d)$ describes closeness if for every tuple $\sigma_1 > 0, \sigma_2 > 0, \dots$, there exists a value σ for which, for all $d > 0$, we have

$$m\left(\frac{d}{\sigma_1}\right) \cdot m\left(\frac{d}{\sigma_2}\right) \cdot \dots = m\left(\frac{d}{\sigma}\right).$$

Main Result. A membership functions describes closeness if and only if it has the form $m(d) = \exp(-A \cdot d^p)$ for some $A > 0$ and $p > 0$.

Discussion. The requirement that all the differences d_i are small means that the difference d_1 is small, and the difference d_2 is small, etc. So the degree of confidence that all the differences are small is equal to the result of applying “and”-operation to the corresponding degrees.

In our scale, “and”-operation is a product, so this degree is equal to the product

$$\prod_{i=1}^n \exp(-A \cdot |d_i|^p) = \exp\left(-A \cdot \sum_{i=1}^n |d_i|^p\right).$$

The constraint is that this degree of confidence should be larger than or equal to some threshold t :

$$\exp\left(-A \cdot \sum_{i=1}^n |d_i|^p\right) \geq t.$$

After taking negative logarithm of both sides, and dividing both sides by A , we can get an equivalent inequality

$$\sum_{i=1}^n |d_i|^p \leq C \stackrel{\text{def}}{=} -\frac{\ln(t)}{A}.$$

This is exactly the ℓ^p -approach, so fuzzy logic indeed leads to a theoretical explanation for this approach.

Proof of Proposition 1. Since the function $m(d) \in [0, 1]$ is strictly decreasing, it is positive for all d , and thus, we can take logarithms of both sides of the desired equality. The logarithm of the product is equal to the sum of the logarithms. Thus, if we denote $M(d) \stackrel{\text{def}}{=} -\ln(m(d))$, then the above equality takes the following form:

$$M\left(\frac{d}{\sigma_1}\right) + M\left(\frac{d}{\sigma_2}\right) + \dots = M\left(\frac{d}{\sigma}\right).$$

In particular, if we take n terms

$$\sigma_1 = \sigma_2 = \dots = \sigma_n = 1,$$

we conclude that

$$n \cdot M(d) = M(k(n) \cdot d)$$

for some value

$$k(n) \stackrel{\text{def}}{=} \frac{1}{\sigma}.$$

This equality can be represented in the following equivalent form

$$\frac{1}{n} \cdot M(k(n) \cdot d) = M(d).$$

This equality should be true for all d , in particular, for $d' = k(n) \cdot d$. For this d' , we get

$$\frac{1}{n} \cdot M(d') = M\left(k\left(\frac{1}{n}\right) \cdot d'\right),$$

where we denoted

$$k\left(\frac{1}{n}\right) \stackrel{\text{def}}{=} \frac{1}{k(n)}.$$

So, for every two integers m and n , we have

$$\begin{aligned} \frac{m}{n} \cdot M(d) &= m \cdot \left(\frac{1}{n} \cdot M(d)\right) = m \cdot M\left(k\left(\frac{1}{n}\right) \cdot d\right) = \\ &= M\left(k(m) \cdot k\left(\frac{1}{n}\right) \cdot d\right). \end{aligned}$$

Thus,

$$\frac{m}{n} \cdot M(d) = M\left(k\left(\frac{m}{n}\right) \cdot d\right),$$

where we denoted

$$k\left(\frac{m}{n}\right) \stackrel{\text{def}}{=} k(m) \cdot k\left(\frac{1}{n}\right).$$

In other words, for every rational number r , we have

$$r \cdot M(d) = M(k(r) \cdot d).$$

For $d = 1$, we get $M(k(r)) = r \cdot M(1)$. Thus, if we denote $s \stackrel{\text{def}}{=} k(r)$, we get

$$r = \frac{M(s)}{M(1)}$$

and so, the above equality takes the form

$$\frac{M(s)}{M(1)} \cdot M(d) = M(s \cdot d).$$

This equality has been proven only for values $k(r)$ for rational r , but since the function $M(r)$ is continuous, it can be extended to all s .

In particular, for

$$D(x) \stackrel{\text{def}}{=} \frac{M(x)}{M(1)},$$

we get

$$D(s) \cdot D(d) = D(s \cdot d).$$

For continuous functions $D(d)$, all solutions to this functional equation are known (see, e.g., [1]), they all have the form $D(d) = d^p$ for some p . Thus,

$$M(d) = M(1) \cdot D(d) = A \cdot d^p,$$

where $A \stackrel{\text{def}}{=} M(a)$, and so, for $m(d) = \exp(-M(d))$, we have the desired expression.

The main result is proven.

IV. WHAT NEXT?

Discussion. In the previous section, we considered the simplest case, when all the membership functions form a 1-D family. A natural next step is to consider situations when they form a 2-D family, then a 3-D family, etc.

Analysis of the problem. In the above proof, we showed that the fact the set of the corresponding membership functions is closed under multiplication, we can conclude that the set of its logarithms forms a linear space.

In general, each n -dimensional space is formed by linear combinations of n basis functions $f_1(x), \dots, f_n(x)$. Scale-invariance means for each of these functions, the re-scaled function $f_i(\lambda \cdot x)$ belongs to the same linear spaces, i.e., that

$$f_i(\lambda \cdot x) = \sum_{j=1}^n c_{ij}(\lambda) \cdot f_j(x)$$

for functions $c_{ij}(\lambda)$.

Differentiating both sides of this equality relative to λ and taking $\lambda = 1$, we conclude that

$$x \cdot \frac{df_i(x)}{dx} = \sum_{j=1}^n c'_{ij}(1) \cdot f_j(x).$$

Here, $\frac{dx}{x} = dz$ for $z = \ln(x)$. Thus, if we express all the functions $f_i(x)$ in terms of z , i.e., consider $f_i(x) = F_i(\ln(x))$, with $F_i(z) \stackrel{\text{def}}{=} f_i(\exp(z))$, then for the new functions $F_i(z)$, we get a system of linear differential equations with constant coefficients:

$$\frac{dF_i(z)}{dz} = \sum_{j=1}^n c'_{ij}(1) \cdot F_j(z).$$

Solutions to such systems are known (see, e.g., [9]): they are linear combinations of functions of the type

$$z^k \cdot \exp(a \cdot z) \cdot \sin(\omega \cdot z + \varphi),$$

where $k \geq 0$ is a natural number and $a + \omega \cdot i$ is an eigenvalue of the corresponding matrix.

Thus, the functions $F_i(x)$ are linear combinations of the functions of the type

$$z^k \cdot \exp(a \cdot z) \cdot \sin(\omega \cdot z + \varphi).$$

Substituting $z = \ln(x)$ into this formula, we arrive at the following conclusion.

Result. The function $f_i(x) = -\ln(\mu(x))$ is a linear combination of functions the type

$$(\ln(x))^k \cdot x^a \cdot \sin(\omega \cdot \ln(x) + \varphi).$$

Thus, each membership function takes the form $\exp(-f_i(x))$ for such functions $f_i(x)$.

1-D and 2-D cases. For a 1-D real-valued matrix, the eigenvalue is a real number, so $\omega = 0$, $k = 0$, and we have $f(x) = x^a$, which is exactly what we showed in our main result.

In the 2-D case, we can have two different real eigenvalues, or we can have double real value, or we can have two mutually conjugate complex eigenvalues. For the complex eigenvalues, we do not have monotonicity, so this case has to be dismissed. Thus, for the 2-D case, only two options are left:

- the case of two different eigenvalues, when the membership function is equal to $\exp(-a \cdot |d|^p - a' \cdot |d|^{p'})$ and thus, regularization is equivalent to the constraint

$$\sum_i |d_i|^p + a \cdot \sum_i |d_i|^{p'} \leq C$$

for some a and p' , and

- the case of a double eigenvalue, when the membership function is equal to $\exp(-a \cdot |d|^p - a' \cdot |d|^p \cdot \ln(|d|))$ and thus, regularization is equivalent to the constraint

$$\sum_i |d_i|^p + a \cdot \sum_i |d_i|^p \cdot \ln(|d_i|) \leq C.$$

ACKNOWLEDGMENT

This work was supported in part by the US National Science Foundation grants HRD-0734825, HRD-1242122, and DUE-0926721.

The authors are greatly thankful to the anonymous referees for valuable suggestions.

REFERENCES

- [1] J. Aczél and J. Dhombres, *Functional Equations in Several Variables*, Cambridge University Press, 2008.
- [2] B. Amizic, L. Spinoulas, R. Molina, and A. K. Katsaggelos, "Compressive blind image decomposition", *IEEE Transactions on Image Processing*, 2013, Vol. 22, No. 10, pp. 3994–4006.
- [3] G. Klir and B. Yuan, "Fuzzy Sets and Fuzzy Logic", Prentice Hall, Upper Saddle River, New Jersey, 1995.
- [4] D. Krishnan and R. Fergus, "Fast image deconvolution using hyper-Laplacian priors", *Proc. Adv. Neural Inf. Processing Systems*, 2009, pp. 1033–1041.
- [5] A. Levin, R. Rergus, F. Durand, and W. T. Freeman, "Image and depth from a conventional camera with a coded aperture", *ACM Transactions on Graphics*, 2007, Vol. 26, No. 3, pp. 1–8.
- [6] B. Mandelbrot, *The Fractal Geometry of Nature*, Freeman, San Francisco, California, 1983.
- [7] H. T. Nguyen, V. Kreinovich, and P. Wojciechowski, "Strict Archimedean t-Norms and t-Conorms as Universal Approximators", *International Journal of Approximate Reasoning*, 1998, Vol. 18, Nos. 3–4, pp. 239–249.
- [8] H. T. Nguyen and E. A. Walker, *A First Course in Fuzzy Logic*, Chapman and Hall/CRC, Boca Raton, Florida, 2006.
- [9] J. C. Robinson, *An Introduction to Ordinary Differential Equations*, Cambridge University Press, Cambridge, UK, 2004.
- [10] A. N. Tikhonov and V. Y. Arsenin, *Solutions of Ill-Posed Problems*, V. H. Winston & Sons, Washington, DC, 1977.
- [11] L. A. Zadeh, "Fuzzy sets", *Information and Control*, 1965, Vol. 8, pp. 338–353.