

# Power Law Regularization in Probabilistic Inverse Problems: a Theoretical Justification

Kanita Petcharat, Piyaphon Paichit, Sa-aat Niwitpong\*

*Department of Applied Statistics, Faculty of Applied Science  
King Mongkut's University of Technology North Bangkok,  
1518 Pibulsongkram Rd., Bangsue, Bangkok 10800 Thailand*

Received 9 January, 2009; Revised 30 January, 2010

## Abstract

In many practical situations, it is desirable to reconstruct the probability distribution  $\rho(x)$  of a quantity  $x$  from the noisy measurements of this quantity. This reconstruction often results in heavier tails, so to get a more accurate reconstruction, it is reasonable to decrease these tails. An empirical comparison of several possible tail-decreasing strategies has shown that the best possible strategy is to raise all the values of the reconstructed density  $\tilde{\rho}(x)$  to a power  $a > 1$ , and then normalize the resulting distribution.

In this paper, we provide a theoretical explanation for this empirical fact.

©2011 World Academic Press, UK. All rights reserved.

**Keywords:** inverse problem, power law regularization, theoretical justification

## 1 Formulation of the Problem

**Reconstructing probability distributions is a practically important problem.** In many practical situations, we need to find the probability density  $\rho(x)$  of a physical quantity  $x$ .

**Reconstructing probability distributions: traditional statistical approach.** Traditional statistical approach assumes that we have a sample of values  $x_1, \dots, x_n$  distributed according to this distribution. Based on this sample, we can then reconstruct the desired probability density function  $\rho(x)$ ; see, e.g., [4]

**In practice, we must take measurement errors into account.** In practice, the values  $x_i$  come from measurement, and measurement is never absolutely accurate: what we actually observe are values  $m_i = x_i + e_i \neq x_i$ , where  $e_i = m_i - x_i$  are measurement errors.

Thus, what we reconstruct by using the traditional statistical techniques is *not* the probability density for the desired quantity  $x$ , but rather the probability density  $\rho_m(x)$  of the sum  $m = x + e$ .

**How measurement errors can be taken into account.** Measurement errors are usually independent from the desired quantity, so it is reasonable to assume that the corresponding random variables  $x$  and  $e$  are independent. Under this assumption, the observed probability density  $\rho_m(x)$  is related to the probability densities  $\rho(x)$  and  $\rho_e(x)$  by the known convolution relationship:

$$\rho_m(x) = \int \rho(y) \cdot \rho_e(x - y) dy. \quad (1)$$

By performing measurements in the absence of the signal (or for a well-calibrated signal), we can determine the probability distribution  $\rho_e(x)$  of the measurement errors. Thus, we know  $\rho_m(x)$  and  $\rho_e(x)$  and so, in principle, we can solve the equation (1) and reconstruct the desired probability distribution  $\rho(x)$ ; see, e.g., [2] and references therein.

\*Corresponding author. Email: snw@kmutnb.ac.th (S. Niwitpong).

**Need for further improvement.** Because of the noise (= measurement errors), the reconstructed distribution  $\tilde{\rho}(x)$  is somewhat different from the actual distribution  $\rho(x)$ . It is therefore desirable to decrease this difference.

**How to improve the reconstructed distribution: analysis of the problem.** The difference between the actual and the reconstructed distributions is especially large for the values  $x$  which occur rarely – and for which, therefore, there are few data points.

In other words, the reconstructed values  $\tilde{\rho}(x)$  may be reasonably accurate when  $\rho(x)$  is large, but the values  $\tilde{\rho}(x)$  may be corrupted when  $\rho(x)$  is small – i.e., at the tails of the original distribution.

When the actual values  $\rho(x)$  are close to 0, the only drastic difference is when the reconstructed values  $\tilde{\rho}(x)$  are much larger. So, due to the measurement errors, the tails of the reconstructed distribution are heavier than of the actual one.

Thus, to improve the reconstructed distribution  $\tilde{\rho}(x)$ , it is desirable to decrease its tails.

**How to decrease the tails: a general idea.** A reasonable idea is to select an appropriate tail-decreasing function  $f(x)$  and then replace each value  $\tilde{\rho}(x)$  with a new value  $f(\tilde{\rho}(x))$  in such a way that when  $\tilde{\rho}(x)$  is small, the new value  $f(\tilde{\rho}(x))$  is even smaller.

This replacement may change the overall probability value from

$$\int \tilde{\rho}(x) dx = 1 \quad (2)$$

to

$$\int f(\tilde{\rho}(x)) dx \neq 1. \quad (3)$$

So, to get a probability density function, we must normalize the values  $f(\tilde{\rho}(x))$ , i.e., take the new probability density function

$$\rho_c(x) = \frac{f(\tilde{\rho}(x))}{\int f(\tilde{\rho}(y)) dy}. \quad (4)$$

**Empirical fact: power-law functions are the best for tail decrease.** An empirical analysis described in [2] showed that the most adequate reconstruction occurs when we use power-law functions  $f(x) = x^a$  for some  $a > 1$ .

**What we do in this paper.** In this paper, we provide a theoretical explanation for the above empirical fact.

## 2 A New Theoretical Justification

**The main idea behind our explanation.** The numerical value of a physical quantity depends on the choice of a measuring unit. For example, 1 inch is exactly the same distance as 2.54 cm. It is natural to require that the results of the tail decrease procedure not change if we simply change the measuring unit. We will show that this reasonable requirement uniquely determines the power law functions.

**Towards formalizing the main idea.** If we change the original measuring unit to a new one which is  $\lambda$  times larger, then the quantity that is described by the value  $x$  in the old units will have a numerical value  $x' = x/\lambda$  in the new units.

The probability density  $\tilde{\rho}(x)$  means that the probability to find the value between  $x$  and  $x + dx$  is equal to  $\tilde{\rho}(x) \cdot dx$ . In terms of the new units  $x' = x/\lambda$ , this same value  $\tilde{\rho}(x) \cdot dx$  is the probability to find the value  $x'$  between  $x' = x/\lambda$  and  $x' + dx'$ , where  $dx' = dx/\lambda$ . From  $x' = x/\lambda$ , we conclude that  $x = x' \cdot \lambda$  and  $dx = dx' \cdot \lambda$ . Thus, the probability to find a value between  $x'$  and  $x' + dx'$  is equal to

$$\tilde{\rho}(x) \cdot dx = \tilde{\rho}(\lambda \cdot x') \cdot \lambda \cdot dx'. \quad (5)$$

On the other hand, by definition of the probability density  $\tilde{\rho}'(x')$ , this same probability is equal to

$$\tilde{\rho}'(x) \cdot dx'. \quad (6)$$

By equating the expressions (5) and (6), we conclude that

$$\tilde{\rho}'(x') = \tilde{\rho}(\lambda \cdot x') \cdot \lambda. \quad (7)$$

**From the main idea to the justification of the power law.** If we apply the tail decrease procedure with a function  $f(x)$  to the probability density function  $\tilde{\rho}(x)$  described in the original units, we get a new function

$$\rho_c(x) = c \cdot f(\tilde{\rho}(x)) \quad (8)$$

for an appropriate normalizing constant  $c$ . In terms of the new units  $x'$ , this new probability density function has the form

$$\rho'_c(x') = \rho_c(\lambda \cdot x') \cdot \lambda = c \cdot f(\tilde{\rho}(\lambda \cdot x')). \quad (9)$$

On the other hand, if we directly apply the same tail decrease procedure to the expression (7) describing the probability density in the new units, we get an expression

$$\rho''_c(x') = c' \cdot f(\tilde{\rho}'(x')) \quad (10)$$

for an appropriate normalizing constant  $c'$ . Substituting the expression (7) for  $\rho'(x')$  into this formula, we conclude that

$$\rho''_c(x') = c' \cdot f(\tilde{\rho}(\lambda \cdot x') \cdot \lambda). \quad (11)$$

Our requirement is that the result of tail decrease should not depend on whether we apply this procedure in the original units or in the new units, i.e., that the resulting probability density functions  $\rho'_c(x')$  and  $\rho''_c(x')$  must coincide. By equating the right-hand sides of the expressions (9) and (11), we get

$$c' \cdot f(\tilde{\rho}(\lambda \cdot x') \cdot \lambda) = c \cdot f(\tilde{\rho}(\lambda \cdot x')). \quad (12)$$

This must be true for all values  $\tilde{\rho}$ , so we conclude that

$$c' \cdot f(\tilde{\rho} \cdot \lambda) = c \cdot f(\tilde{\rho}). \quad (13)$$

Dividing both sides of this equality by  $c'$  and taking into account that the ratio  $r \stackrel{\text{def}}{=} c/c'$  may depend on  $\lambda$ , we deduce that for every two numbers  $\tilde{\rho} > 0$  and  $\lambda$ , the following equality holds:

$$f(\tilde{\rho} \cdot \lambda) = r(\lambda) \cdot f(\tilde{\rho}). \quad (14)$$

It is known that every continuous function  $f(x)$  satisfying this property has the form  $f(x) = C \cdot x^a$  for some  $a$ ; see, e.g., [1], Section 3.1.1. (This result was first proven in [3].)

The coefficient  $C$  does not matter since we normalize the probability density function anyway, so we can simply conclude that  $f(x) = x^a$ .

**Conclusion.** The only tail decrease procedures for which the results do not depend on the choice of the measuring units are the power-law procedures, with  $f(x) = x^a$ .

**Comment.** For differentiable functions  $f(x)$ , the result about power functions is easy to prove. Indeed, if we differentiate both sides of (14) by  $\lambda$  and take  $\lambda = 1$ , we get

$$\tilde{\rho} \cdot \frac{df}{d\tilde{\rho}} = a \cdot f, \quad (15)$$

where  $a \stackrel{\text{def}}{=} r(1)$ . By moving all the terms containing  $f$  into one side and all the terms containing  $\tilde{\rho}$  to the other side, we conclude that

$$\frac{df}{f} = a \cdot \frac{d\tilde{\rho}}{\tilde{\rho}}. \quad (16)$$

Integrating both sides, we get

$$\ln(f) = a \cdot \ln(\tilde{\rho}) + c, \quad (17)$$

hence

$$f = e^{a \cdot \ln(\tilde{\rho}) + c} = e^c \cdot \left(e^{\ln(\tilde{\rho})}\right)^a = C \cdot \tilde{\rho}^a \quad (18)$$

for  $C = e^c$ .

### 3 Observation: Tail Decrease Procedure is Related to the Notion of Entropy

When the measurements are precise, we get the exact values of the probabilities, i.e., we should take  $f(x) = x$ , with  $a = 1$ . When the measurements are very accurate, we get the probabilities close to the actual ones, so the tail decrease procedure  $f(x) = x^a$  should minimally change the probability values – and thus, we should have  $a = 1 + \varepsilon$  for some small  $\varepsilon$ .

For a small  $\varepsilon$ , we can expand the expression  $\rho^{1+\varepsilon}$  in Taylor series in  $\varepsilon$  and keep only linear terms in this expansion (thus ignoring quadratic and higher order terms). As a result, we get

$$\tilde{\rho}^{1+\varepsilon} = \tilde{\rho} + \varepsilon \cdot \tilde{\rho} \cdot \ln(\tilde{\rho}). \quad (19)$$

Indeed,

$$\tilde{\rho}^{1+\varepsilon} = \tilde{\rho} \cdot \tilde{\rho}^\varepsilon = \tilde{\rho} \cdot (e^{\ln(\tilde{\rho})})^\varepsilon = \tilde{\rho} \cdot e^{\varepsilon \cdot \ln(\tilde{\rho})} \approx \tilde{\rho} \cdot (1 + \varepsilon \cdot \ln(\tilde{\rho})). \quad (20)$$

Due to (19), the normalizing constant is equal to

$$\int f(\tilde{\rho}(y)) dy = \int \tilde{\rho}(y) dy + \varepsilon \cdot \int \tilde{\rho}(y) \cdot \ln(\tilde{\rho}(y)) dy. \quad (21)$$

The first integral is the full probability, i.e., 1, the second integral differs only by the sign from the entropy

$$S(\tilde{\rho}) = - \int \tilde{\rho}(y) \cdot \ln(\tilde{\rho}(y)) dy \quad (22)$$

of the probability distribution, thus,

$$\int f(\tilde{\rho}(y)) dy = 1 - \varepsilon \cdot S(\tilde{\rho}), \quad (23)$$

and the resulting normalized distribution (4) takes the form

$$\rho_c(x) = \frac{(\tilde{\rho}(x))^{1+\varepsilon}}{1 - \varepsilon \cdot S(\tilde{\rho})}. \quad (24)$$

## Acknowledgments

The authors are thankful to Vladik Kreinovich for his help and encouragement.

## References

- [1] Aczel, J., *Lectures on Functional Equations and Their Applications*, Dover Publ., New York, 2006.
- [2] Ferson, S., Using approximate deconvolution to estimate cleanup targets in probabilistic risk analyses, *Hydrocarbon Contaminated Soils*, edited by P.T. Kostecki, E.J. Calabrese and M. Bonazountas, Amherst Scientific Publishers, Amherst, Massachusetts, vol.5, pp.245–254, 1995.
- [3] Pexider, J., Notiz uber Funktionaltheoreme, *Monatsch. Math. Phys.*, vol.14, pp.293–301, 1903.
- [4] Sheskin, D.J., *Handbook of Parametric and Nonparametric Statistical Procedures*, Chapman & Hall/CRC Press, Boca Raton, Florida, 2007.