

# The Career Atlas: Mathematical Notation

Cao Bittencourt

<sup>a</sup>*B. Sc. in Economics from EPGE (FGV), RJ, Brazil.*

<sup>b</sup>*Statistician at Atlas Career Guide Inc., FL, USA.*

---

## Abstract

This is a brief document to define statistical methods for data-drive career choice and development. It deals with topics such as: career matching (i.e. vocational choice); estimation of competence, or overall skill level; estimation of skill set generality; versatility; skill set profitability; employability; labor market competitiveness; labor market taxonomy; optimal human resources acquisition and allocation; and so on and so forth. Each concept shall be explained at length in separate articles.

**Keywords:** Career choice; Career development; Matching algorithms; Competence; Similarity.

---

## 1. Basic Definitions

### 1.1. Skill Sets

The  $i$ -th professional attribute, or competency, of a person  $k$  is defined as:

$$a_i^k \in [0, 100], \quad (1)$$

where the interval  $[0, 100]$  determines the bounds for every competency.<sup>1</sup>

The skill set, or career profile, of a person  $k$  is the vector of their  $m$  attributes:

$$\mathbf{a}_k = (a_1^k, \dots, a_m^k). \quad (2)$$

The skill set matrix, or career profile matrix, is the collection of all  $n$  skill sets in the economy:

$$\mathbf{A} = \begin{bmatrix} a_1^1 & \dots & a_m^1 \\ \vdots & \ddots & \vdots \\ a_1^n & \dots & a_m^n \end{bmatrix}. \quad (3)$$

---

<sup>1</sup>More generally, these could be defined as  $a_{lb}$  (the lower bound) and  $a_{ub}$  (the upper bound). Here, the interval  $[0, 100]$  is used because of its ease of interpretation.

### 1.2. Skill Set Normalization

Normalization by the scale bounds is defined by the tilde operator:

$$\tilde{a}_i^k = \frac{a_i^k - 0}{100 - 0} = \frac{a_i^k}{100} \in [0, 1]; \quad (4)$$

$$\tilde{\mathbf{a}}_{\mathbf{k}} = (\tilde{a}_1^k, \dots, \tilde{a}_m^k); \quad (5)$$

$$\tilde{\mathbf{A}} = \begin{bmatrix} \tilde{a}_1^1 & \dots & \tilde{a}_m^1 \\ \vdots & \ddots & \vdots \\ \tilde{a}_1^n & \dots & \tilde{a}_m^n \end{bmatrix}. \quad (6)$$

Normalization by a skill set's highest attribute is defined by the hat operator:

$$\hat{a}_i^k = \frac{a_i^k}{\max_j a_j^k} \in [0, 1]; \quad (7)$$

$$\hat{\mathbf{a}}_{\mathbf{k}} = (\hat{a}_1^k, \dots, \hat{a}_m^k); \quad (8)$$

$$\hat{\mathbf{A}} = \begin{bmatrix} \hat{a}_1^1 & \dots & \hat{a}_m^1 \\ \vdots & \ddots & \vdots \\ \hat{a}_1^n & \dots & \hat{a}_m^n \end{bmatrix}. \quad (9)$$

## 2. Basic Skill Set Models

### 2.1. Skill Set Generality

The generality of a skill set is the mean of its maxima-normalized attributes:

$$\gamma_k = \left( \frac{1}{m} \right) \sum_{i=1}^m \hat{a}_i^k \in [0, 1]. \quad (10)$$

People with high  $\gamma_k$  scores are called *generalists*. Conversely, those with low  $\gamma_k$  scores are called *specialists*. Career profiles that are neither broad nor specialized are said to be *balanced*.

The generality vector of all  $n$  skill sets in the economy is:

$$\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_n). \quad (11)$$

### 2.2. Attribute Equivalence

The attribute equivalence of a particular attribute in a skill set measures the importance of that attribute relative to the skill set's highest attribute. It is calculated using skill set generality as both a midpoint and scaling parameter in the following linear-logistic classification function:

$$\text{aeq}(\hat{a}_i^k, \gamma_k) = \hat{a}_i^k \left[ 1 + \gamma_k (1 - \hat{a}_i^k) \exp \left( \frac{\hat{a}_i^k - \gamma_k}{\gamma_k - 1} \right) \right]^{-\frac{\gamma_k}{\hat{a}_i^k}} \in [0, 1]. \quad (12)$$

For a short-hand notation, the attribute equivalence can be denoted by the umlaut operator:

$$\ddot{a}_i^k = \text{aeq}(\hat{a}_i^k, \gamma_k). \quad (13)$$

This is not to be confused with the Newtonian dot notation for partial derivatives, which we do employ, instead preferring the more explicit  $\frac{\partial x}{\partial y}$  derivative operator of Leibniz.

At any rate, attributes with high levels of  $\ddot{a}_i^k$  are said to be equivalent to the skill set's most important attribute. These are a career profile's *core* competencies. The remaining competencies are classified as either *important*, *auxiliary*, *minor*, or *unimportant*.

The attribute equivalence vector of a skill set is given by the collection of their  $m$  umlauted attributes:

$$\ddot{\mathbf{a}}_k = (\ddot{a}_1^k, \dots, \ddot{a}_m^k). \quad (14)$$

Finally, the attribute equivalence matrix is the collection of all attribute equivalence vectors in the economy:

$$\ddot{\mathbf{A}} = \begin{bmatrix} \ddot{a}_1^1 & \dots & \ddot{a}_m^1 \\ \vdots & \ddots & \vdots \\ \ddot{a}_1^n & \dots & \ddot{a}_m^n \end{bmatrix}. \quad (15)$$

### 2.3. Skill Set Competence

The overall competence of a skill set is the mean of its scale-normalized attributes, weighted by each attribute's importance (i.e. its attribute equivalence):

$$c_k = \frac{\sum_{i=1}^m \ddot{a}_i^k \tilde{a}_i^k}{\sum_{i=1}^m \ddot{a}_i^k} \in [0, 1]. \quad (16)$$

Career profiles with high  $c_k$  are said to be competent. However, this adjective can be seen as offensive to some people; and, most importantly, it could also be misleading, because competence is often defined relative to the specific requirements of a particular job. Therefore, we opt for the more generic competence classification of *high level*, *mid level*, and *low level*, which is somewhat less ambiguous.

The competence vector of all  $n$  skill sets in the economy is:

$$\mathbf{c} = (c_1, \dots, c_n). \quad (17)$$

## 3. Comparative Models

### 3.1. Matching Models

The most basic comparative model is that of Euclidean matching with linear weights:

$$s_{k,q} = s(\mathbf{a}_k, \mathbf{a}_q) = 1 - \tilde{d}(\mathbf{a}_k, \mathbf{a}_q) \in [0, 1], \quad (18)$$

where

$$\tilde{d}_{k,q} = \tilde{d}(\mathbf{a}_k, \mathbf{a}_q) = \sqrt{\frac{\sum_{i=1}^m a_i^q (a_i^k - a_i^q)^2}{\sum_{i=1}^m a_i^q \max(100 - a_i^q, a_i^q)^2}} \in [0, 1]. \quad (19)$$

In this model, we compare a skill set  $\mathbf{a}_k$  to a skill set  $\mathbf{a}_q$  by calculating the weighted Euclidean distance from  $\mathbf{a}_k$  to  $\mathbf{a}_q$  normalized by the maximum theoretical distance to  $\mathbf{a}_q$ .

Other weighting systems can be employed in this type of matching model. We could, for instance, substitute the linear weights with either quadratic weights,

$$a_i^{q^2} \in [0, 1], \quad (20)$$

or speciality-root weights,

$$a_i^{q^{\frac{1}{1-\gamma_k}}} \in [0, 1]. \quad (21)$$

But the best and most interpretable results are obtained using attribute equivalence as the weighting function:

$$\tilde{d}_{k,q} = \tilde{d}(\mathbf{a}_k, \mathbf{a}_q) = \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q (a_i^k - a_i^q)^2}{\sum_{i=1}^m \ddot{a}_i^q \max(100 - a_i^q, a_i^q)^2}} \in [0, 1]. \quad (22)$$

We could also employ other matching methods instead of the “baseline” weighted Euclidean model. [detail each method later]:

1. logit regression matching
2. probit regression matching
3. bvls regression matching
4. tobit regression matching
5. pearson correlation matching
6. kendal nonparametric correlation matching
7. spearman nonparametric correlation matching

At last, similarity and normalized distance metrics determine the respective vectors and matrices, as follows:

$$\mathbf{s}_k = (s_{k,1}, \dots, s_{k,n}); \quad (23)$$

$$\tilde{\mathbf{d}}_k = (\tilde{d}_{k,1}, \dots, \tilde{d}_{k,n}); \quad (24)$$

$$\mathbf{S} = \begin{bmatrix} s_{1,1} & \dots & s_{n,1} \\ \vdots & \ddots & \vdots \\ s_{1,n} & \dots & s_{n,n} \end{bmatrix} = \begin{bmatrix} 1 & \dots & s_{k,1} & \dots & s_{n,1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ s_{1,k} & \dots & 1 & \dots & s_{n,k} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ s_{1,n} & \dots & s_{k,n} & \dots & 1 \end{bmatrix}; \quad (25)$$

$$\mathbf{D} = \begin{bmatrix} \tilde{d}_{1,1} & \dots & \tilde{d}_{n,1} \\ \vdots & \ddots & \vdots \\ \tilde{d}_{1,n} & \dots & \tilde{d}_{n,n} \end{bmatrix} = \begin{bmatrix} 0 & \dots & \tilde{d}_{k,1} & \dots & \tilde{d}_{n,1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{d}_{1,k} & \dots & 0 & \dots & \tilde{d}_{n,k} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{d}_{1,n} & \dots & \tilde{d}_{k,n} & \dots & 0 \end{bmatrix}. \quad (26)$$

### 3.2. Qualification Models

A closely related concept to matching is the qualification comparative model. In this family of functions, however, Euclidean matching is mandatory, as other matching methods do not make sense for this specific type of calculation. The reason for this is at that, here, we are not particularly interested in matching (i.e. a typical classification problem), but rather in the actual distances between comparison skill sets.

To define these models, we first have to define the gap function, which measures only positive competency gaps:

$$\delta_{k,q}^i = \delta(a_i^k, a_i^q) = \max(a_i^k - a_i^q, 0) \in [0, 100]. \quad (27)$$

Having defined the gap function, we can write the underqualification model:

$$\tilde{\delta}_{k,q}^< = \text{uqa}(\mathbf{a}_k, \mathbf{a}_q) = \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, a_i^k)^2}{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, 0)^2}} = \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, a_i^k)^2}{\sum_{i=1}^m \ddot{a}_i^q a_i^{q^2}}}. \quad (28)$$

And, analogously, the overqualification model is given by:

$$\tilde{\delta}_{k,q}^{\geq} = \text{oqa}(\mathbf{a}_k, \mathbf{a}_q) = \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^k, a_i^q)^2}{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, 100)^2}} = \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^k, a_i^q)^2}{\sum_{i=1}^m \ddot{a}_i^q (100 - a_i^q)^2}}. \quad (29)$$

The final set of “sufficient qualification” is, evidently, the complement of the underqualification model:

$$s_{k,q}^{\geq} = \text{sqa}(\mathbf{a}_k, \mathbf{a}_q) = 1 - \sqrt{\frac{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, a_i^k)^2}{\sum_{i=1}^m \ddot{a}_i^q \delta(a_i^q, 0)^2}} = 1 - \text{uqa}(\mathbf{a}_k, \mathbf{a}_q). \quad (30)$$

As with the similarity and normalized distance statistics described above, these three qualification models are bounded to the  $[0, 1]$  interval. Likewise, they also determine qualification vectors:

$$\tilde{\delta}_k^< = (\tilde{\delta}_{k,1}^<, \dots, \tilde{\delta}_{k,n}^<); \quad (31)$$

$$\tilde{\delta}_k^{\geq} = (\tilde{\delta}_{k,1}^{\geq}, \dots, \tilde{\delta}_{k,n}^{\geq}); \quad (32)$$

$$s_k^{\geq} = (s_{k,1}^{\geq}, \dots, s_{k,n}^{\geq}); \quad (33)$$

and matrices

$$\tilde{\Delta}_{<} = \begin{bmatrix} \tilde{\delta}_{1,1}^{<} & \cdots & \tilde{\delta}_{n,1}^{<} \\ \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,n}^{<} & \cdots & \tilde{\delta}_{n,n}^{<} \end{bmatrix} = \begin{bmatrix} 0 & \cdots & \tilde{\delta}_{k,1}^{<} & \cdots & \tilde{\delta}_{n,1}^{<} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,k}^{<} & \cdots & 0 & \cdots & \tilde{\delta}_{n,k}^{<} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,n}^{<} & \cdots & \tilde{\delta}_{k,n}^{<} & \cdots & 0 \end{bmatrix}; \quad (34)$$

$$\tilde{\Delta}_{\geq} = \begin{bmatrix} \tilde{\delta}_{1,1}^{\geq} & \cdots & \tilde{\delta}_{n,1}^{\geq} \\ \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,n}^{\geq} & \cdots & \tilde{\delta}_{n,n}^{\geq} \end{bmatrix} = \begin{bmatrix} 0 & \cdots & \tilde{\delta}_{k,1}^{\geq} & \cdots & \tilde{\delta}_{n,1}^{\geq} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,k}^{\geq} & \cdots & 0 & \cdots & \tilde{\delta}_{n,k}^{\geq} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{\delta}_{1,n}^{\geq} & \cdots & \tilde{\delta}_{k,n}^{\geq} & \cdots & 0 \end{bmatrix}; \quad (35)$$

$$\mathbf{S}_{\geq} = \begin{bmatrix} s_{1,1}^{>} & \cdots & s_{n,1}^{>} \\ \vdots & \ddots & \vdots \\ s_{1,n}^{>} & \cdots & s_{n,n}^{>} \end{bmatrix} = \begin{bmatrix} 1 & \cdots & s_{k,1}^{>} & \cdots & s_{n,1}^{>} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ s_{1,k}^{>} & \cdots & 1 & \cdots & s_{n,k}^{>} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ s_{1,n}^{>} & \cdots & s_{k,n}^{>} & \cdots & 1 \end{bmatrix}. \quad (36)$$