

EVALUATING MODEL FIDELITY IN AN AERIAL IMAGE ANALYSIS SYSTEM

F. Quint M. Sties
Institute for Photogrammetry and Remote Sensing
University of Karlsruhe
76128 Karlsruhe, Germany
quint@ipf.bau-verm.uni-karlsruhe.de
Commission III, Working Group 3

KEY WORDS: Aerial Image Understanding, Model, Knowledge Base, Semantic Networks

ABSTRACT

The purpose of the system MOSES is the automatic recognition of objects in aerial images. To direct the model based structural image analysis, one has to evaluate each state of the analysis process. One We present in this article the procedures used in MOSES to calculate a part of these valuations, the model fidelity, which is a measure for the goodness of match between the chosen image primitives and the specific model. Metrics defined on a parametric representation of the primitives are used to evaluate the model fidelity. The results of the image analysis process directed by these valuations are presented.

KURZFASSUNG

Das System MOSES dient der automatischen Erkennung von Objekten in Luftbildern. Zur Steuerung der modellbasierten, strukturellen Bildanalyse sind Bewertungen des aktuellen Analysezustandes anzugeben. In diesem Artikel werden die in MOSES verwendete Verfahren zur Berechnung eines Teils dieser Bewertungen, der Modelltreue, vorgestellt. Die Modelltreue ist ein Maß für die Übereinstimmung zwischen den gewählten Bildprimitiven und dem spezifischen Modell. Zu ihrer Berechnung werden Metriken auf einer parametrischen Darstellung der Primitiven verwendet. Ergebnisse des Bildanalyse unter Verwendung der vorgestellten Modelltreue werden erläutert.

1 INTRODUCTION

Understanding of aerial images is one of the most challenging tasks in computer vision. Due to its complexity, a model based analysis has been found to be mandatory since several years, see e.g. (McKeown et al., 1985), (Nicolin and Gabler, 1987), (Matsuyama and Hwang, 1990), (Sandakly and Giraudon, 1994), (Stilla, 1995). In our system MOSES (*Map Oriented SEMantic image underSTanding*) (Quint and Sties, 1995) we too perform a structural, model based analysis. We are interested in the recognition of objects in urban environment using large scale aerial images.

2 MOSES

One of the main characteristics of the system MOSES is that large scale topographical maps are used to automatically refine the models used for image analysis. The architecture of our system is shown in Fig. 1. The generative model contains domain independent, common sense knowledge the system designer has about the environment. The generic models in the map domain and in the image domain are specializations of the generative model and they reflect the particularities of the representations of our environment in the map and image respectively. The models contain both declarative knowledge, which describes the structure of the objects, and procedural knowledge, which contains the methods used during the map and image analysis process. As a repository for the models semantic networks (Findler, 1979) are used, as implemented by the system ERNEST (Kummert et al., 1993).

The generative model and the generic models are that part of the system which is build by the system developer. The models and scene descriptions described in the sequel are automatically build in analysis processes. Analysis takes place in three phases.

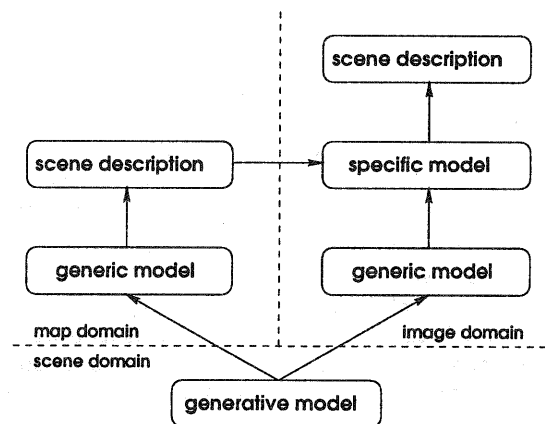


Figure 1: Architecture of the system MOSES

2.1 Map analysis

In the first phase, the generic model in the map domain is used to analyse the map, which is available as a list of digitized contours. The procedure by which map analysis is performed is similar to the one used in the image analysis process and will be described in a following section. The result of the map analysis is a description of the scene, as far as it can be constructed out of the map data. This scene description is also stored in a semantic network.

The nodes of the semantic network represent objects, parts and subparts of the scene. They are described with attributes, which in this case mainly contain the geometric properties of the scene objects. Links between the nodes represent relations between the corresponding objects or parts. Two typical relations are the *part-of* relation, which describes the structure of the scene objects and the *specialization* relation, along which properties of objects are inherited.

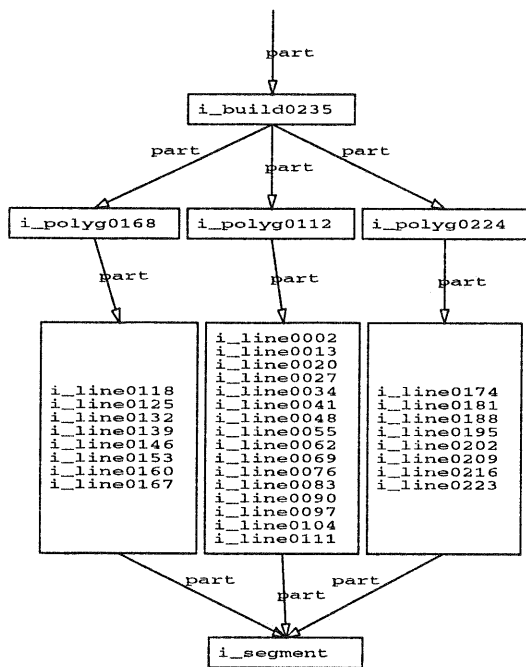


Figure 2: Detail of the part-of hierarchy of the specific model

2.2 Model building

In the second phase, the scene description obtained after the map analysis is combined with the generic model in the image domain and the specific model in the image domain is built. A detail of this specific model representing building 0235 and its parts as far as they are given in the map, is shown in Fig. 2.

For each node (instance) in the scene description we create a new node (concept) in the specific model. This new concept is a specialization of the corresponding concept in the generic model in the image domain and thus inherits its declarative and procedural knowledge. The values of the attributes in the scene description after map analysis are stored after a transformation as restrictions for the corresponding attributes of the newly created concepts. They serve as initial estimates for the calculation of the attribute values out of the image data.

The relations between the instances in the scene description are transferred accordingly into relations between the new concepts. Whilst the generic model in the image domain describes in a general form the representation of an arbitrary scene in an aerial image, the specific model in the image domain describes in a detailed manner that part of the world, which is subject to the current analysis. The grade of detail depends of course from the contents of the map.

2.3 Image primitives

Prior to the model based image analysis primitives are extracted from the image data. We work with large scale color aerial images, which after digitization have a pixel size of 30 cm x 30 cm on the ground. Line segments and regions serve as primitives. The line segments are extracted with a gradient based procedure (Quint and Bähr, 1994). The regions are gained by segmenting the aerial image using a Bayesian homogeneity predicate (Quint and Landes, 1996).

The regions and the line segments are combined in an at-

tributed undirected graph. The nodes of the graph are attributed with the regions. Nodes corresponding to neighbouring regions are connected with links. A link between two nodes is attributed with the line segment(s) which compose the border between the corresponding regions. This feature graph is the database on which the model based image analysis operates.

2.4 Image analysis

In the third phase, the specific model in the image domain is used to perform the actual image analysis. The aim of this phase is to verify in the image the objects found after the map analysis and to detect and describe other objects of the scene which are not represented in the map. For the later, the context gained through the verification of the map objects will be helpful.

The strategy followed in the analysis process is a general, problem independent strategy provided by the shell ERNEST. The analysis starts by creating a modified concept for the goal concept (expansion step). A modified concept is a preliminary result and it reflects constraints for the concept that have been determined out of the context of the current analysis state.

Following top-down the hierarchy in the semantic network, stepwise the concepts on lower hierarchical levels are expanded until a concept on the lowest level is reached. Since this concept does not depend from other concepts, its correspondence with a primitive in the database can be established and its attributes can be calculated. This is called instantiation.

Analysis now moves bottom-up to the concept at the next higher hierarchical level. If instances have been found for all parts of this concept, the concept itself can be instantiated. Otherwise the analysis continues with the next not yet instantiated concept on a lower level. After an instantiation, the acquired knowledge is propagated bottom-up and top-down to impose constraints and restrict the search space. Thus, in the analysis process top-down and bottom-up processing alternate. As well, expansion and instantiation alternate during the analysis.

Generally, while performing an instantiation it is possible to establish several correspondences between a concept and primitives in the data base. However, only one of these correspondences leads to the correct interpretation. Since it usually is not possible to ultimately decide at the lower levels which correspondence is correct, all possible correspondences have to be accounted for.

Thus, the image analysis is a search process, which can be graphically represented by a tree. Each node of the tree represents a state of the analysis process. If in a given state several correspondences are possible, the search tree is splitted: for each hypothesis a new node as successor of the current node is created.

The analysis process continues with that leaf node of the search tree which is considered to be the best according to a problem dependent evaluation. It is known that the problem of finding an optimal path in a search tree can be solved by the A^* -algorithm (Nilsson, 1982). Its application is possible if one can evaluate the path from the root node to the current node and if one can give an estimate for the valuation of the path from the current node to the (not yet known) terminal node containing the solution.

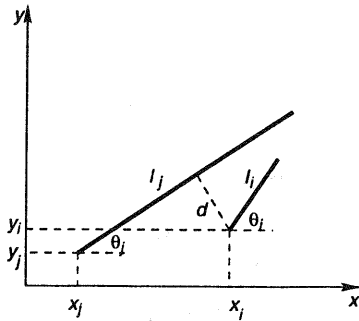


Figure 3: Parameters used to describe a line segment

3 VALUATIONS

The functions which evaluate the states of the analysis are very important since they are not only responsible for the efficiency of the search, but they are also decisive for the success or failure of the analysis. We relate the valuation of the search path to the valuation of the analysis goal in the given state of the analysis. The valuation of the goal is calculated considering the valuations of the instances and modified concepts already created and the estimates for the valuations of the instances and modified concepts which will be created in the path from the current node to the solution node.

When an instantiation is performed, implicitly a hypothesis of match is established between the concept under instantiation and the chosen primitives from the database. Since we can not ultimately decide at the moment the instantiation is performed, if it is the correct one, we are working under uncertainty and we have to quantify our uncertainty. At the level of each concept in the semantic network we have a dichotomous frame of discernment with the events: the chosen primitives

- match
- do not match

to the concept (i.e. model).

The valuations computed for the instances and modified concepts in each state of the analysis are measures of our subjective belief in these hypotheses. They take values between 0 and 1 and we interpret them as basic belief masses in the framework of the Dempster-Shafer theory of evidence (Shafer, 1976). The higher a valuation is, the stronger is our subjective belief in the corresponding hypothesis. Using the methods described in (Quint, 1995), the different valuations are combined and propagated in the hierarchy of the semantic network to result in the valuation of the analysis goal.

We evaluate two aspects for our hypotheses of match: the compatibility and the model fidelity. The compatibility evaluates an analysis state considering the principles of perceptual grouping. It is calculated based on geometric, topologic and radiometric properties of the image primitives only. In this category belong for example the goodness of fit of several line segments extracted from the image data to form an edge of an object, the goodness of fit of several edges to form a polygon, the compatibility of the polarity of edges to form a polygon etc.

The model fidelity measures the goodness of fit between the image primitives and the specific model gained through the

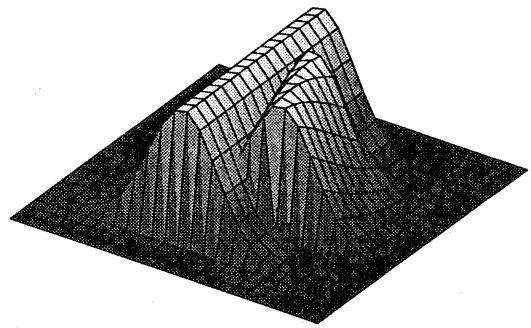


Figure 4: Neighbourhood function for the position of line segments

analysis of the map. Portraying it in simplified terms, one can say that the compatibility is a measure for the ability of the chosen image primitives to form an object of the generic model, whereas the model fidelity is a measure for the ability to form exactly that object, which is predicted by the map. We present in this article measures used for the evaluation of the model fidelity.

4 MODEL FIDELITY

4.1 Model fidelity for line segments

At the level of line segments we define the model fidelity with help of a distance function between the image primitives and the contours stored in the specific model. The distance function is part of a metric defined with help of a set of square integrable functions on a parametric space for line segments.

We describe a line segment with help of the coordinates of its starting point, its length and the angle between the line and the positive x -axis (see Fig. 3). Thus, a line segment s_i is represented in the space $S = (x, y, l, \theta)$ by the point $s_i = (x_i, y_i, l_i, \theta_i)$. The coordinates of a line segment are in the domain $(x, y) \in \mathbb{R}^2$, the length of a line is in $l \in \mathbb{R}_+$ and its angle is in $\theta \in (-\frac{\pi}{2}, \frac{\pi}{2}]$. The space (x, y, l, θ) is the Cartesian product of the enumerated domains and is different from \mathbb{R}^n . For this reason we do not use the Euclidean distance between two points in this space to calculate the distance between two line segments, but use instead a metric defined on an isomorphic space of functions.

We define an isomorphism by attaching each point s_i in the space S a function $n_i(x, y, l, \theta)$ from the space of square integrable functions $\mathcal{L}^2(S)$. We call this function *neighbourhood function*. As a distance between two line segments s_i and s_j we now use the distance defined on the family of functions n_i . It is well known that a distance function defined with the expression:

$$d_{ij} = \left[\int_S (n_i(x, y, l, \theta) - n_j(x, y, l, \theta))^2 dx dy dl d\theta \right]^{\frac{1}{2}} \quad (1)$$

satisfies the necessary properties for a metric on $\mathcal{L}^2(S)$. If we choose the functions $n_i(x, y, l, \theta)$ such that their norm in the induced metric is equal to 1, i.e.

$$\int_S (n_i(x, y, l, \theta))^2 dx dy dl d\theta \stackrel{!}{=} 1, \quad (2)$$

the expression (1) simplifies to:

$$d_{ij} = \left[2 - 2 \int_S n_i(x, y, l, \theta) n_j(x, y, l, \theta) dx dy dl d\theta \right]^{\frac{1}{2}} \quad (3)$$

The distance d_{ij} decreases when the integral in expression (3) increases. If the neighbourhood functions are positive functions, the integral in expression (3) takes values between 0 and 1.

We have formulated our search problem using as valuations of the nodes in the search tree merit functions and not cost functions. The reason for this is pragmatic: it is more natural to evaluate the goodness than the badness of a match. Thus, we will not use the distance as given by expression (3) but only the integral in expression (3) to define the model fidelity m_{ij} at the level of line segments:

$$m_{ij} = \int_S n_i(x, y, l, \theta) n_j(x, y, l, \theta) dx dy dl d\theta. \quad (4)$$

This integral equals to the cosinus of the angle between the two versors n_i and n_j in the vector space $\mathcal{L}^2(S)$ and can be thought of as a correlation measure between these two versors.

The neighbourhood functions are chosen regarding the physics of the image formation process and some heuristics motivated by experience. We construct the function $n_i(x, y, l, \theta)$ as a product of three functions defined on \mathbb{R}^2 , \mathbb{R}_+ and $(-\frac{\pi}{2}, \frac{\pi}{2}]$ respectively:

$$n_i(x, y, l, \theta) = f_i(x, y) g_i(l) h_i(\theta).$$

To define the function $f_i(x, y)$ we take advantage of the fact that the parameters of the camera and the position of the airplane at the moment the aerial image was taken are known. We can determinate the transformation between the image coordinates and the coordinates in the specific model (map coordinates) and transform the image primitives into the map coordinate system. Assuming that the corresponding contours are depicted in the map, there are several error sources which are responsible for the fact that the line segments extracted from the image will not overlap with the map contours. These are for example inaccuracies in:

- the extraction of line segments from the image,
- the determination of the transformation parameters,
- the acquisition and digitization of the map data.

Subsuming all these effects, we can safely assume that the position of the image primitives is normally distributed around their "true" position as given by the specific model.

For this reason we use as a neighbourhood function $f_i(x, y)$ for the position of the line segments a Gaussian shaped function. However, since we do not want to evaluate differently the situations when a short line segment lies in the middle of its model line or closer to the endpoints, our function is constant along the length of the line. We choose for the neighbourhood function $f_i(x, y)$:

$$f_i(x, y) = K_{xy} \exp \left(- \frac{((x - x_i) \sin \theta_i - (y - y_i) \cos \theta_i)^2}{2\sigma^2} \right)$$

for positions (x, y) between the endpoints of a line, i.e. $\{(x, y) \mid (x - x_i) \cos \theta_i + (y - y_i) \sin \theta_i \geq 0 \wedge (x - x_i) \cos \theta_i +$

$(y - y_i) \sin \theta_i \leq l_i\}$, and $f_i(x, y) = 0$ otherwise. The neighbourhood functions $f_i(x, y)$ and $f_j(x, y)$ for the constellation of line segments shown in Fig. 3 are displayed in Fig. 4. The variance of the Gaussian is chosen equal to the residual mean square error of the transformation.

For the part of the neighbourhood function, which depends from the length of the line, we choose a function which "inside" the line is proportional to the square root of the length and which is 0 "outside":

$$g_i(l) = \begin{cases} K_l \sqrt{l} & \text{if } l \in [0, l_i] \\ 0 & \text{otherwise.} \end{cases}$$

As we will see later, this choice penalizes image primitives in an amount proportional to the ratio of their length and the length of the model contour.

The considerations regarding the uncertainty in the position of line segments applies also for small deviations of the angle. Thus, the neighbourhood function for the angle is chosen following similar reflections. But because the domain of definition of the angle is an interval and because we want a stronger penalization of large deviations of the angle, we use a trigonometric function instead of the Gaussian shaped function:

$$h_i(\theta) = K_\theta \cos(\theta - \theta_i).$$

The constants K_{xy} , K_l and K_θ are calculated imposing normalization for each of the partial neighbourhood functions and we can thus assure the fulfillment of condition (2).

With this choice of neighbourhood functions, the integral for the model fidelity is separable into three terms: the position fidelity, the length fidelity and the angle fidelity. The integral over the product of the neighbourhood functions for the position, i.e. the position fidelity can generally not be expressed in a closed form. However, if the angle between the two lines is small or the parameter σ is in the same order of magnitude as the mean geometric distances between the two line segments then a good approximation is given by:

$$\int_{\mathbb{R}^2} f_i(x, y) f_j(x, y) dx dy = \frac{\sqrt{\pi} \sigma}{l_i \sin \Delta\theta} \times \left(\operatorname{erf} \left(\frac{u_1 \sin \Delta\theta - A}{\sigma \sqrt{2 + 2 \cos \Delta\theta^2}} \right) - \operatorname{erf} \left(\frac{u_2 \sin \Delta\theta - A}{\sigma \sqrt{2 + 2 \cos \Delta\theta^2}} \right) \right) \quad (5)$$

with $\Delta\theta = \theta_j - \theta_i$ and $A = -(x_i - x_j) \sin \theta_j + (y_i - y_j) \cos \theta_j$. The coordinates u_1 and u_2 are the coordinates of the start- and of the endpoint of line l_i in a coordinate system uOv with its origin in the starting point of line l_j and with the u -axis parallel to the line l_j . For a situation as shown in Fig. 3, when after a parallel displacement the perpendicular distance d between the two lines varies, the position fidelity varies in function of d as shown in Fig. 5.

The integrals over the neighbourhood functions for the length and the angle of the line segments can be expressed in closed form and result to:

$$\int_{\mathbb{R}_+} g_i(l) g_j(l) dl = \frac{\min(l_i, l_j)^2}{l_i l_j}$$

and

$$\int_{-\pi/2}^{\pi/2} h_i(\theta) h_j(\theta) d\theta = \cos(\theta_i - \theta_j).$$

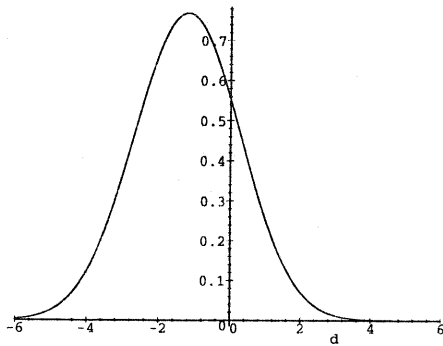


Figure 5: Position fidelity as a function of d (see also Fig. 3)

The length fidelity amounts thus to the ratio of the length of the shorter line to the length of the longer line. The angle fidelity is the cosine of the angle difference of the two lines. The total model fidelity for line segments is given by the product of the three components.

Usually, due to noise influence the visible contour of an object in the image is broken and thus several line segments will form the contour. In this case, the contour is constructed step by step by adding line segments until the contour is completed. The A^* -algorithm requires an optimistic estimate of the merit for future instantiations. To give an optimistic estimate for the future instantiations in the case of a partially estimated contour, we elongate the already instantiated line segments in order to simulate a virtual best fit with the model. The model fidelity for this "ideal" best fit is evaluated and serves as an optimistic estimate for the model fidelity of future instantiations.

4.2 Model fidelity for polygons

A different approach for the model fidelity is used at the hierarchical level of polygons. Whilst at the level of line segments the similarity in position and orientation between the selected image primitives and the model contour has been evaluated, we evaluate at the level of polygons the similarity between the shape of the polygon created by the image primitives and the shape of the model polygon.

The corner points of the polygon in the image domain are obtained as intersections of the chosen image primitives. In the case where several image primitives form an edge of an object, these primitives are replaced for the purpose of the corner point calculation with a regression line. The error produced by the approximation with the regression line is taken into account in the valuations of the compatibility. In the case where no correspondence could be established between an edge of an object and an image primitive we make a wildcard assignment to the current edge. In this case the corresponding corner points are chosen to be the end point of the image primitive assigned to the edge previous to and the starting point of the image primitive assigned to the edge after the wildcard-assigned edge. The wildcard assignments however lead to a penalization in the model fidelity of the line segments.

To not include position and orientation errors in our measure we first transform the polygon in the image domain on the model polygon. We take a similarity transformation between the corresponding corner points of the two polygons and calculate the transformation parameters such that the residual

mean square error is minimal. Since the scale of the image and the map are known, we fix the scale parameter in the similarity transformation to the known value.

The resulting minimal mean square error is a measure for the similarity of the shapes of the two polygons. We gain our subjective belief in the hypotheses of match between the image polygon and the model polygon with help of a fuzzy function:

$$p_{ij}(r) = \exp\left(-\frac{r^2}{\sigma_r^2}\right) \quad (6)$$

where r is the residual mean square error after the transformation.

4.3 Model fidelity for objects

The resulting model fidelity for an object of the scene is calculated by combining the model fidelities at the level of line segments and polygons. The model fidelities are interpreted as subjective beliefs in the corresponding hypotheses of match and treated in the framework of the Dempster-Shafer theory of evidence. With an extension (Quint, 1995) to approaches found in the literature we propagate the model fidelities calculated at a lower hierarchical level of the semantic network upwards. Model fidelities at the same hierarchical level are combined with Dempster's rule (Shafer, 1976).

The such computed model fidelity at the level of an object of the scene is used to decide whether an object represented in the map could have been verified in the image analysis process. Besides this, the model fidelity for an object is further propagated up to the goal concept of the analysis, which in our case represents the scene. At this level it is combined with the compatibility measures computed for the instances and contributes to the valuation of an analysis state. However, since we are not only interested in the verification of objects represented in the map, the model fidelity contributes in a smaller fraction to the valuation of the analysis state than the compatibility.

5 RESULTS

We present the results for the verification of buildings in the scene of Fig. 6. The line segments used in the image analysis process are overlaid in black color in the Figure. There were roughly 5000 line segments presented to the system. The line segments which are found after the analysis to compose the buildings of the scene are drawn in white color. For each building its identifier is also displayed in the Figure. The model fidelity for the recognized buildings is given in Table 1.

Excepting the house in the lower left corner of the image (i.house0106) all the other buildings in the image have been verified successfully. The main reason for the failure of the verification was that the position error for this building with respect to the specific model was twice as big as the position errors of the other buildings in the scene. In this experiment the parameters σ in expression (5) and σ_r in expression (6) were chosen such, that an absolute position error of 2m in the scene leads to a model fidelity of 0.5 (i.e. half of maximal value).

Those objects rejected by the verification process are marked and passed to the following phase of the analysis, the classification phase, where these image structures are interpreted regardless of a specific model gained from map analysis, but

