

K-Means Based Prediction of Transcoded JPEG File Size and Structural Similarity

Steven Pigeon, Université de Montréal, Canada

Stéphane Coulombe, École de Technologie Supérieure, Canada

ABSTRACT

The problem of efficiently adapting JPEG images to satisfy given constraints such as maximum file size and resolution arises in a number of applications, from universal media access for mobile browsing to multimedia messaging services. However, optimizing for perceived quality (user experience) commands a non-negligible computational cost which in the authors work, they aim to minimize by the use of low-cost predictors. In previous work, the authors presented predictors and predictor-based systems to achieve low-cost and near-optimal adaption of JPEG images under given constraints of file size and resolution. In this work, they extend and improve these solutions by including more information about images to obtain more accurate predictions of file size and quality resulting from transcoding. The authors show that the proposed method, based on the clustering of transcoding operations represented as high-dimensional vectors, significantly outperforms previous methods in accuracy.

Keywords: *FileSize Reduction, Image Adaptation, Image Quality, Image Transcoding, Joint-Photographic-Experts-Group (JPEG), Resolution Reduction, Structural Similarity Index (SSIM)*

1. INTRODUCTION

The need for efficient image adaptation arises in a number of contexts, ranging from universal media access with varying browsing conditions (Han et al., 1998; Mohan, Smith, & Li, 1999), to multimedia messaging services (MMS) (Coulombe & Grassel, 2004). In the case of universal access, one uses a mobile device, either a smart-phone, PDA, or a tablet, to access resources or services on the Web. The traditional response has been to use rather crude

adaptation strategies (Han et al., 1998) such as simply preparing a single “mobile” version of the resource (Fling, 2009), but this one-size-fits-all solution will leave users at both ends of the device capability spectrum dissatisfied: some will find the mobile version exceeding (or cumbersome for) their devices’ capabilities, while others will find it inadequate and lacking.

In the context of MMS, for another example, a receiving terminal is characterized by its capabilities—or more exactly its *limitations*—such as the maximum resolution of images it can display, the formats it can decode, and the maximum message size it can receive

DOI: 10.4018/jmdem.2012040103

and interpret correctly (Open Mobile Alliance, 2010). Interoperability between MMS users will require server-side adaptation, as the sender's device may be more capable than the receiver's, and the receiving device will be unable to display correctly, if at all, a message that exceeds its capabilities. In this context, adaptation will require that the sender's images are converted to comply with the receiver's device capabilities, that is, changing the file size (by altering the compression parameters) and resolution of images (by scaling them). Adaptation can also include the case where the compression format itself needs to be changed. But this is seldom a problem since MMS image traffic is mostly composed of JPEG images taken from the devices' cameras. Accordingly, we will neglect the case where the format also needs to be adapted (for example, from PNG to GIF) and concentrate on the prevalent problem of JPEG to JPEG image adaptation subject to changes in compression parameters (e.g., the quality factor) and scaling (resolution).

Therefore, whether in the context of universal access or multimedia messaging, the challenge is to adapt images to fit given constraints, dictated by the network conditions and the receiving device capabilities, while simultaneously maximizing the user experience and minimizing the computational cost of adaptation. In the context of high-volume service providing, whether for MMS or universal media access, only the fastest adaptation algorithms yielding the best perceived quality can be considered.

Of course, previous studies have addressed the problem of efficient image adaptation, but the solutions they propose are either still computationally expensive (and extensive modifications to existing JPEG manipulation libraries) (Ridge, 2003; Shu & Chau, 2005) or overly rigid, focusing on unrealistically constrained transformations such as scaling by powers of two (Lei & Georganas, 2002; Ratnakar & Ivashin, 2001; Ridge, 2003), or using a small, fixed, number of possible adaptations, without real consideration for the perceived quality resulting from adaptation. For example,

Ridge's method is accurate, but requires the JPEG image to be partly decompressed so that the DCT coefficients are available, on which successive re-quantization passes are performed until the quality factor yielding the largest file not exceeding the constraint is found (Ridge, 2003). Other methods exploit the structure of the DCT to yield fast scaling algorithms in the (partially) compressed domain by manipulating the DCT coefficients directly, but such methods also require the image to be partly decoded so that the DCT coefficients are available, and they are constrained to scaling by powers of two. Furthermore, it is unclear what the expected speed-ups are, as the DCT-coefficient based scaling algorithms are still relatively complex and may compare to an efficient implementation scaling using space-domain filters in terms of computational complexity. But, in our opinion, the principal shortcoming of previous methods is that they do not consider joint changes in compression parameters and scaling as a means of adaptation maximizing perceived quality.

In previous work, we have proposed low-cost predictor-based adaptation systems for the prediction of the file size and perceived quality resulting from an image subjected to simultaneous changes in compression parameters and scaling (Coulombe & Pigeon, 2009, 2010; Pigeon & Coulombe, 2008). These lookup-table based constant time predictors, which we will refer to as JQSP1 (for JPEG Quality and Size Predictor) and JQSP2 in this work, are described in Section 3, *Prediction Algorithms*. These predictors, unlike the solutions discussed in the previous paragraph, use only information that is readily available *without* decompressing the images, such as the original file size and the original quality factor (the parameter that controls the aggressiveness of compression) which can be accessed by reading only the file's header. These predictors are used in combination with an adaptation system (Coulombe & Pigeon, 2010; Pigeon & Coulombe, 2011, 2012) to predict the best transcoding parameters for adapting a given JPEG image subject to receiving terminal constraints, where "best" is defined as most likely to minimize perceived distortion under

the considered viewing conditions as defined by the characteristics of the receiving device.

These predictors, although shown to perform well, do not make use of all the readily available information about the images such as resolution and bits per pixel, both of which are very likely to help formulate even more accurate predictions about file size and quality resulting from a given transcoding operation. However, the table-based schemes in our earlier work do not lend themselves easily to a larger number of parameters, and we believe that the uniform quantization of parameters (which was a key component of the methods' computational efficiency) is also an unwanted restriction for the problem at hand.

Therefore, in this work, we propose to extend and improve the solutions previously presented by the authors by using a method capable of including more information about images in order to help formulate more accurate predictions of file size and quality resulting from transcoding, and by lifting the restrictions of the previous methods, in particular the uniform quantization of parameters. To do so, we propose a method based on K-Means, which can formulate predictions using the clustering of transcoding operations represented as high-dimensional vectors.

The work is structured as follows: the next section, Section 2, presents the proposed solution. Section 3 details the validation methodology as well as the algorithms of previous work. Section 4 presents the results from the proposed method as well as the results from the previously presented methods. In Section 5, we discuss the results, accuracy, algorithmic complexity, and the memory usage of the algorithms considered. The Appendix discusses the efficient implementation of K-Means and prediction. We conclude in Section 6.

2. PROPOSED CLUSTERING-BASED SOLUTION

In this section, we detail the Enhanced JPEG Quality and Size Predictor (EJQSP), the solu-

tion we are proposing for the prediction of relative file size and quality resulting from a JPEG image adaptation based on clustering (Hastie, Tibshirani, & Friedman, 2009). We first describe the general problem of clustering, and then we describe its application to our particular objective.

The general clustering problem is as follows: we have n points in \mathbb{R}^d (or other metric space), the $\{x_j\}_{j=1}^n$, which we want to partition into $C = \{C_1, C_2, C_3, \dots, C_m\}$, m disjoint subsets. The partition C of X , per definition, is such that $\bigcup_{i=1}^m C_i = X$ and $C_i \cap C_j = \emptyset$, for $1 \leq i \neq j \leq m$. That is, the union of the disjoint subsets forms C . For each subset C_i , we elect a value, the prototype, denoted \bar{x}_i , deemed representative (under a given metric) of the elements of C_i . There are many sensible metrics one can use, but the metric considered here is the usual Euclidean distance, and the prototype \bar{x}_i for subset C_i is given by

$$\bar{x}_i = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j,$$

the L_2 centroid, and where $|C_i|$ denotes the cardinality of C_i . Let $\bar{X} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m\}$ be the set of the prototypes. The goodness of a partition C (and of its corresponding set of prototypes \bar{X}), is assessed using the error function

$$E(C) = \sum_{i=1}^m \sum_{x_j \in C_i} \|x_j - \bar{x}_i\|^2, \quad (1)$$

and the goal is to find the optimal partition C^* that minimizes $E(C)$, that is,

$$C^* = \arg \min_C E(C). \quad (2)$$

Applying clustering to our problem of predicting the resulting (relative) file size and quality of an image subjected to changes in quality factor and scaling commands that we represent our exemplars as d -dimensional vectors encoding information about the original image, the transformation applied, and the quantities to be predicted, for example, relative file size and resulting quality. The transformations applied to the image, the change in compression parameters and in resolution, will be referred to as the *transcoding operation*. In this work, the transcoding operation describes only the change in quality factor and scaling; but it could also include more parameters such as, say, the chroma sub-sampling strategy (Pennebaker & Mitchell, 1993).

Let us define more notations. The original compressed image I_j is represented by a tuple (QF_j, w_j, h_j, f_j) , where QF_j is the quality factor with which the image was originally compressed (and in this work we suppose that the quality factor complies with the Independent JPEG group definition, that is, it is an integer between 0 and 100) (IJG, 2012), w_j and h_j are its width and height in pixels, respectively, and f_j , the original compressed file size of image I_j (possibly including extraneous data such as EXIFs and comments) (ISO/IEC, 2011). A transcoding operation (QF_{out}, z) describes the desired quality factor with which to recompress the image as well as $0 < z \leq 1$, a scaling factor. Applying a transcoding operation (QF_{out}, z) to an image I_j yields an image with resolution $zw_j \times zh_j$, quality factor QF_{out} , with perceived quality $q(I_j, QF_{out}, z)$ and resulting *relative* file size $f(I_j, QF_{out}, z)$, expressed as a ratio of the original file size f_j (let us note that in the original work of (Pigeon & Coulombe, 2008), relative file size rather than absolute file size was used in an effort to abstract the effect of absolute file size from the predictor; and in this work we will adhere to this convention). Both

$f(I_j, QF_{out}, z)$ and $q(I_j, QF_{out}, z)$ are measured after the actual transcoding of image I_j . The transcoding operations are constrained to be such that $QF_{out} \in \{10, 20, \dots, 100\}$ and $z \in \{0.1, 0.2, \dots, 1.0\}$. While this is not required for the solution proposed in this work, we used these restrictions for computational reasons in previous work (Pigeon & Coulombe, 2008) and they are included here for a fairer comparison of the methods.

The resulting quality measured between the original and the transcoded image is assessed by a quality metric. Ideally, it would be estimated using a MOS-like measure, but an accurate subjective evaluation of quality is difficult; we will ordinarily rely on simpler measures. In this work, we chose the structural similarity index, or SSIM (Wang, Bovick, Sheikh, & Simoncelli, 2004), to assess the perceived quality of resulting images. While PSNR has been a *de facto* standard for measuring image quality for a long time, SSIM is more robust to transformations that do not affect perceived quality, and is therefore deemed a better estimation of the user experience. Should the transcoded image resolution differ from the original (whenever $z \neq 1$), the transcoded image is scaled back to the original resolution for comparison. In all cases, scaling is performed using the Blackman filter, chosen for its spectral characteristics (Blackman & Tukey, 1959). In previous work, we assessed quality using a more sophisticated approach that depended on the receiving device's screen resolution and canvas (the maximum image size it can manipulate; not necessarily related to screen resolution); but in this work, for the sake of simplicity, we will omit such complications and consider only the case where the images are compared at the original image resolution.

To the original data from the image characteristics and the transcoding operation, we will further add *features*. The added features will consist of transformations of the original data used to put forward characteristics that would be otherwise impossible to discover

using K-Means alone. Such features can be created randomly (for example a random linear combination of the original data), but they can also be constructed from *a priori* knowledge. The first feature we will use is the bits per pixel of the original image I_j , denoted b_j , which gives a measure of the image complexity. The second feature is the difference of quality factors, $QF_{out} - QF_j$, which gives information on the expected drop in quality and file size resulting from the change in quality factor. We will discuss features and their selection further in Section 5.

The exemplars considered are therefore 9-dimensional vectors. The vector associated to an image $I_j = (QF_j, w_j, h_j, f_j)$ being applied a transcoding operation (QF_{out}, z) is therefore

$$x_j = (QF_j, w_j, h_j, b_j, QF_{out}, z, QF_{out} - QF_j, f, q), \quad (3)$$

where,

$$b_j = \frac{f_j}{w_j h_j},$$

and f and q stand for $f(I_j, QF_{out}, z)$ and $q(I_j, QF_{out}, z)$, respectively.

Solving eq. (2) exactly is an NP-Hard problem (Aloise, Deshpande, Hansen, & Popat, 2010; Mahajan, Nimhorkar, & Varadarajan, 2009; Vattani, 2009) and therefore we will seek approximate algorithms such as K-Means (Lloyd, 1982), an algorithm very similar to LBG (Linde, Buzo, & Gray, 1980), which, while stochastic and in general sub-optimal, was shown to converge rapidly to good solutions under most circumstances (Bottou & Bengio, 1995).

The K-Means is an iterative algorithm and proceeds as follows: the initialization consists in randomly picking, uniformly and without

replacement, m vectors from X to serve as initial prototypes, the \bar{X}_0 . At iteration $t = 1, 2, \dots$, for each vector $x_j \in X$, we find the closest prototype $\bar{x}_{t-1,i} \in \bar{X}_{t-1}$ amongst the prototypes of the previous iteration. That is, we find

$$i = \arg \min_i \|x_j - \bar{x}_{t-1,i}\|^2.$$

We then assign the vector x_j to the subset

$C_{t,i}$. All vectors assigned to a same subset are then used to compute the prototype. Since in our case the metric is the Euclidean distance, the prototype is the average vector within set $C_{t-1,i}$, that is, $\bar{x}_{t,i} = |C_{t-1,i}|^{-1} \sum_{x_j \in C_{t-1,i}} x_j$ —had

we used an L_1 metric, the prototype would have been the vector median, considerably more troublesome to compute (Barni, 1997). The algorithm iterates until satisfactory convergence is obtained, that is, the error $E(C_t)$ is not significantly smaller than the error $E(C_{t-1})$. In our experiments, we set the threshold to a relative difference of $\alpha = 10^{-6}$ or less. Algorithm 1 details the procedure in pseudo-code, and the Appendix discusses the computational complexity and parallelization of Algorithm 1.

The L_2 norm considered for the minimization of eq. (2) and the lack of a distance matrix in the problem formulation suggests that for best results, the exemplars must lie in an isotropic space; in other words, all dimensions should be spread along similar scales. If exemplars do not lie in such a space, the usual approach is to use principal component analysis or similar techniques to project the exemplars onto a vector space that provides isotropism (Hastie et al., 2009). Results, however, suggest that dimension-wise standardization suffices to provide satisfactory isotropy.

Algorithm 1. K-Means

```

Inputs:    $X$ , the exemplars
            $m$ , the number of prototypes
            $\alpha$ , the convergence threshold

 $t \leftarrow 0$ 
 $E_0 \leftarrow BIG\_NUM \{ "infinity" \}$ 
 $\bar{X}_0 \leftarrow m$  vectors from  $X$  (uniform random, without replacement)
 $A \leftarrow \{0,0, \dots, 0\}$  { the initial prototype assignment }
for all  $x_j \in X$  do
     $a_j \leftarrow \arg \min_i \|x_j - \bar{x}_{0,i}\|^2$  { compute membership }
end for
repeat
     $t \leftarrow t + 1$ 
     $n_t \leftarrow \{0,0, \dots, 0\}$  {  $m$  elements, the  $n_{t,i}$  }
     $\bar{X}_t \leftarrow \{0,0, \dots, 0\}$  {  $m$  elements, the  $\bar{x}_{t,i}$  }
    for all  $x_j \in X$  do
         $i \leftarrow a_j$  { reuse membership }
         $\bar{x}_{t,i} \leftarrow \bar{x}_{t,i} + x_j$  { update prototype }
         $n_{t,i} \leftarrow n_{t,i} + 1$  { update the number of exemplars in this subset }
    end for
    for  $i = 1$  to  $m$  do
         $\bar{x}_{t,i} \leftarrow \bar{x}_{t,i} / n_{t,i}$  { normalize prototype }
    end for
     $E_t \leftarrow 0$ 
    for all  $x_j \in X$  do
         $i \leftarrow a_j \leftarrow \arg \min_k \|x_j - \bar{x}_{t,k}\|^2$  { update membership }
         $E_t \leftarrow E_t + \|x_j - \bar{x}_{t,i}\|^2$  { update total error }
    end for
until  $\text{converges}(E_t, E_{t-1}, \alpha)$ 
outputs:  $\bar{X}_t$ , the prototypes after  $t$  iterations.

```

It only remains to decide on the number of prototypes, m . While allowing m to grow arbitrarily large reduces the training error (with zero error when $m = n$, for example), it is not desirable to have a very large m as the table holding the prototypes becomes very large. The number of prototypes has to be only as large as necessary so that at the local scale, the manifold onto which the exemplars lie appears approximately isotropic. The value of m is therefore subject to a number of trade-offs between prediction accuracy (that we should call *generalization error*, which differ from minimizing eq.

(2), since eq. (2) is only concerned with the training exemplars, and not with all possible exemplars (especially the infinite number of exemplars we have yet to observe) and it is possible that “all” the exemplars are distributed somewhat differently than the training set exemplars), the cost of storing the table, and the time for searching it. We discuss these issues in the Appendix. Fortunately, we will show that m need not be very large to outperform previously proposed predictors (Pigeon & Coulombe, 2008).

3. PREDICTION ALGORITHMS AND SIMULATIONS

The data set used in all experiments is composed of approximately 73000 JPEG images collected from the Web using a crawler, with high-profile Web sites as origination points (Pigeon & Coulombe, 2008). The data set was split into two disjoint parts, 90% and 10%, for the training and test sets respectively. As with previous experiments (Pigeon & Coulombe, 2008), each image was subjected to 100 different transcodings (corresponding to all combinations $(QF_{out}, z) \in \{10, 20, \dots, 100\} \times \{0.1, 0.2, \dots, 1.0\}$), for which we observed resulting file size and perceived quality, yielding approximately 6570000 training exemplars and 730000 test exemplars, each modeled after eq. (3).

The first predictor we presented in Pigeon and Coulombe (2008), denoted here JQSP1, uses a table look-up scheme to formulate its predictions. First, the original quality factor QF_{in} (corresponding to QF_j in this work), the desired output quality factor QF_{out} , and the scaling z are uniformly quantized to the desired precision, giving \widetilde{QF}_{in} , \widetilde{QF}_{out} , and \widetilde{z} (the tilde notation denotes quantized values throughout this paper). In Pigeon and Coulombe (2008), we constrained both \widetilde{QF}_{in} and \widetilde{QF}_{out} to $\{10, 20, \dots, 100\}$, and \widetilde{z} to $\{0.1, 0.2, \dots, 1.0\}$, effectively indexing a $10 \times 10 \times 10$ array containing, in each cell, the relative file size and quality predictions corresponding to the tuple $(\widetilde{QF}_{in}, \widetilde{QF}_{out}, \widetilde{z})$. The predictions are simply formulated as the centroid (the average) of all training exemplars whose quantized parameters fall into a same cell.

The predictor introduced in Coulombe and Pigeon (2010), denoted here JQSP2, unlike the JQSP1, does not predict resulting file size nor quality, but the transcoding parameters maximizing quality under the constraint of file size. This allows the predictor to formulate finer transcoding operations than the coarser JQSP1,

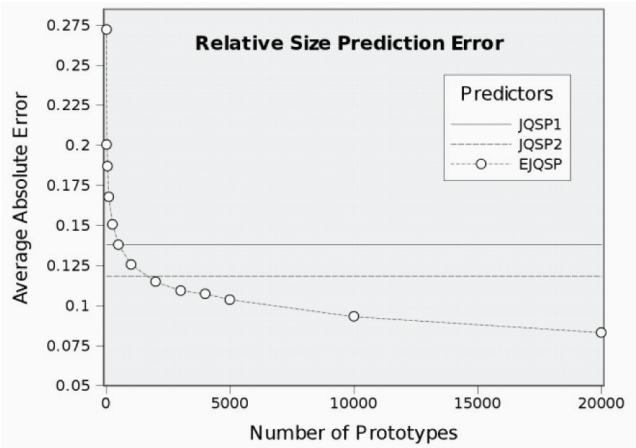
despite being optimized on the same training set. JQSP2 formulates predictions as follows:

given an original quality \widetilde{QF}_{in} , a scaling factor \widetilde{z} and a target file size \widetilde{f}_{max} , the algorithm formulates the prediction of the transcoding parameter \widetilde{QF}_{out} and resulting quality \hat{q} (the hat notation will denote predictions) as the centroid of, for each unique image in the training set with original quality factor of \widetilde{QF}_{in} , the transcoding parameters maximizing file size without exceeding the constraint. The size of the table can be adjusted by varying the quantization on \widetilde{QF}_{in} , \widetilde{z} , and \widetilde{f} . In Coulombe and Pigeon (2010), we have that $\widetilde{QF}_{in} \in \{10, 20, \dots, 100\}$, $\widetilde{z} \in \{0.1, 0.2, \dots, 1.0\}$, and \widetilde{f} varies from 0.001 to 1.0 by increments of 0.001, thus minimizing errors introduced by the quantization on f —since f represents *relative* file size, even a small fraction of the *original* file size may actually correspond to a large portion of the *target* file size.

The variant of K-Means described by Algorithm 1 is not very sensitive to the initial conditions (the randomly chosen \bar{X}_0) but it is not impervious to them either; its stochastic nature will require, to obtain a truly satisfying minimum, several independent optimizations. In our experiments, we opted for 30 independent optimizations (using different initial conditions but the same training set) and chose the optimization with the smallest error, as defined by eq. (1). The efficient implementation of K-Means is discussed in the Appendix.

All algorithms share the same training set (90% of the exemplars) and the same test set (the remaining 10%). Even if the exemplars are modeled after eq. (3), algorithms JQSP1 and JQSP2 use only QF_j (or QF_{in} in the original papers), QF_{out} , z , f , and q , ignoring the resolution information, w_j and h_j , and the features, b_j , $QF_{out} - QF_j$. Algorithm EJQSP, based on K-Means, will compute m clusters from the training set using the vectors as given by eq. (3). Tests were conducted on the remain-

Figure 1. Average absolute error for relative size prediction



ing 10% of the exemplars using the three prediction algorithms, and we compared the actual resulting file size f and observed quality q against the predictions \hat{f} and \hat{q} . The differences are reported as average absolute error as the quantities lie in $[0, 1]$, using a mean square error would yield exceedingly small quantities (as, for example, $0.1^2 = 0.01$), thus exaggerating the method's performance.

4. RESULTS

The JQSP1, JQSP2, and the proposed EQSP predictors are compared, as described in the previous section, using the same test exemplars. For each test exemplar (formed according to eq. (3)), each predictor was asked to compute the predicted file size \hat{f} and quality \hat{q} , and the resulting errors $\hat{f} - f$ and $\hat{q} - q$ were mea-

Figure 2. Average absolute error for quality prediction

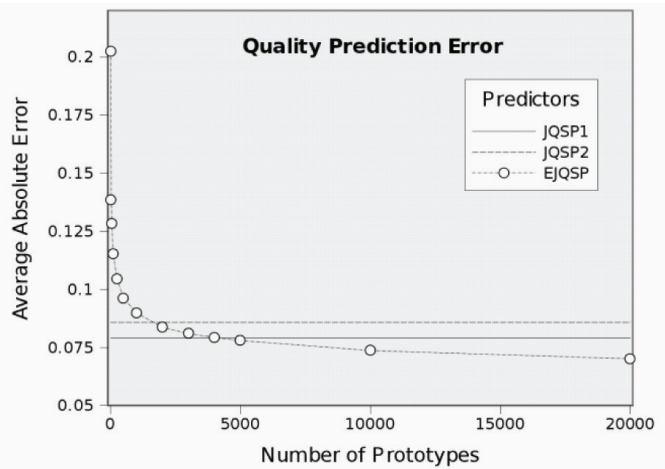


Table 1. Average Absolute Errors for the different algorithms

Prototypes	Average Absolute Size Error			Average Absolute Quality Error		
	JQSP1	JQSP2	EJQSP	JQSP1	JQSP2	EJQSP
—	0.1379	0.1183	—	0.7090	0.0858	—
10			0.2721			0.2023
25			0.2004			0.1385
50			0.1869			0.1282
100			0.1678			0.1153
250			0.1506			0.1044
500			0.1379			0.0961
1000			0.1255			0.0899
2000			0.1148			0.0836
3000			0.1093			0.0809
4000			0.1073			0.0783
5000			0.1037			0.0780
10000			0.0931			0.0737
20000			0.0831			0.0701

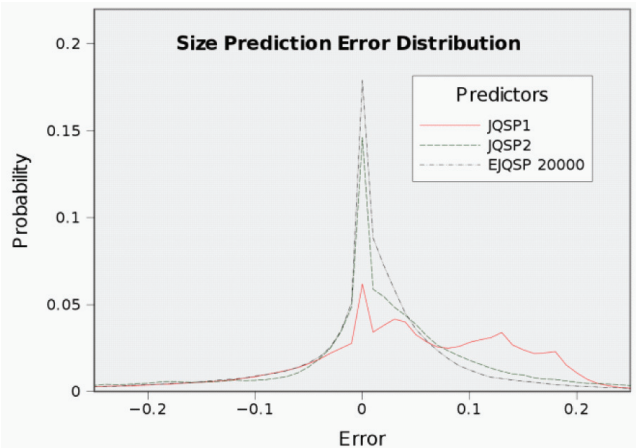
sured. The EJQSP predictor was trained using arbitrary, but not unlikely, values of m , namely, 10, 25, 50, 100, 250, 500, 1000, 2000, 3000, 4000, 5000, 10000, and 20000. The values are chosen not only to show that the prediction accuracy increases as the number of prototypes increases, but also to show the graceful behavior of algorithm EJQSP as the number of allowed prototypes is *reduced*, and that EJQSP is amenable to trade-offs.

Figure 1 and Figure 2 presents the different average errors over all test exemplars for the different predictors and for varying number of prototypes. The errors of algorithms JQSP1 and JQSP2 appear as straight lines as they are unaffected by the number of prototypes of algorithm EJQSP. Examining Figure 1 and Table 1, we see that the EJQSP predictor breaks even, on the accuracy of the prediction of the resulting file size, with predictor JQSP1 using only 500 prototypes and with JQSP2 using 2000. The performance difference continues to increase as the number of prototypes grows, to a point where, at 20000 prototypes, EJQSP has an error $\approx 40\%$ smaller than JQSP1, and $\approx 27\%$

smaller than JQSP2. While the prediction error is $\approx 10\%$ smaller with 20000 prototypes than with 10000, the gain is obtained at the cost of doubling the run-time, as we will discuss further in Section 5. Figure 2 shows similar behavior for quality prediction. EJQSP breaks even with JQSP2 using approximately 1500 prototypes and with JQSP1 at 4000; however, EJQSP ultimately yields an error that is $\approx 20\%$ smaller than JQSP2, and $\approx 12\%$ less than JQSP1. Both Figures 1 and 2 show EJQSP accuracy increases smoothly with the number of prototypes (and the large number of exemplars allows the use of a model with a rather high capacity without risks of over-fitting) (Hastie et al., 2009).

Figure 3 shows the distribution of errors $\hat{f} - f$ (the error of the predicted relative file size \hat{f} against the observed file size, after transcoding, f) for the different predictors in this study. Examining Figure 3, we can see that the distribution of errors from the JQSP1 predictor, despite peaking near zero, exhibits a strong skewness resulting in a definite propen-

Figure 3. Distribution of errors on file size prediction

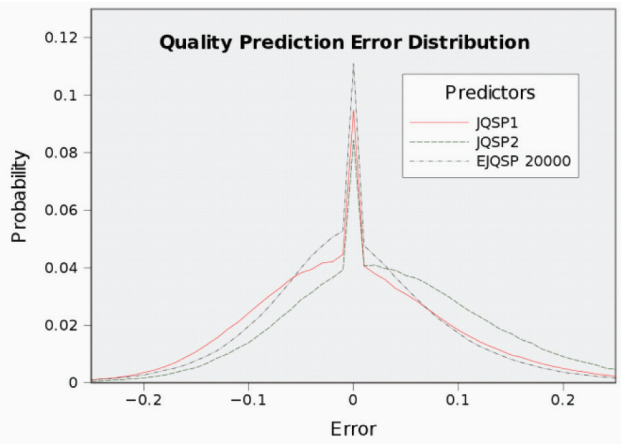


sity to overestimate the resulting file size. An algorithm using JQSP1 as its predictor will therefore tend to take (*a posteriori*) conservative decisions, and will likely fail to use the full file size budget, leading to images with a lower quality than might have been actually possible. Predictor JQSP2 also exhibits such a skewness, but to a much lesser degree, and shows a stronger peak around zero, explaining its greater accuracy (as shown in Figure 1). Lastly, EJQSP, with 20000 prototypes, shows a general behavior comparable to JQSP2, but again with a much

stronger peak at zero than either JQSP1 or JQSP2.

If Figure 3 shows that the individual relative file size predictors behave differently, Figure 4 shows that for the resulting distributions of quality prediction errors $\hat{q} - q$ (the error of the predicted quality \hat{q} against the observed quality q) are quite similar. For all predictors, the distribution of the errors shows two distinct components: a component similar to a skewed Gaussian (Azzalini, 1985) and a strong peak around zero. While JQSP2 formu-

Figure 4. Distribution of error on quality prediction



lates better relative file size prediction than JQSP1, the situation is reversed for quality prediction; a fact that is reflected in Figure 2, where the overall error on quality prediction of JQSP1 is $\approx 8\%$ smaller than that of JQSP2. This can be explained by the fact that the predictor JQSP2 was designed to avoid overshooting the predicted file size but has no special provision regarding quality prediction. EJQSP, like JQSP1, tends to underestimate resulting quality, however, the central peak and a more compact distribution of errors around that peak show that it yields, overall, better predictions than previous predictors.

5. DISCUSSION

For all algorithms, training and test exemplars must be obtained. In our experiments, it meant subjecting every image from the database obtained by crawling to 100 different transcodings—varying over combinations of \widetilde{QF}_{out} and \widetilde{z} —to yield a sufficiently large (and dense) pool of transcoding examples. This process is of course very expensive as one has to perform the actual transcoding and assess resulting quality, but the operation can be performed off-line (and incrementally) as trends in image characteristics will vary over time, at application-specific pace.

Training the JQSP1 predictor is $O(n)$ in the number of exemplars, n . Training consists, for each exemplar $x_j \in X$, in quantizing the QF_j , QF_{out} , and z , into \widetilde{QF}_j , \widetilde{QF}_{out} , and \widetilde{z} to index an entry in the array and accumulate the partial sums to compute the centroids. The finalization of the centroid computation is proportional to the number of entries in the array, that is, $O\left(\left|\widetilde{QF}_j\right|\left|\widetilde{QF}_{out}\right|\left|\widetilde{z}\right|\right)$, where, by abuse of notation, $\left|\widetilde{QF}_j\right|$, $\left|\widetilde{QF}_{out}\right|$, and $\left|\widetilde{z}\right|$ denote the number of distinct values each can take. The cost of the normalization is therefore negligible compared to the cost of computing the partial sums, since the number of entries in

the array will be very small compared to the number of exemplars used for training; furthermore, the array cannot become very large, not only because of memory consumption, but also to avoid the problem of context dilution (Hastie et al., 2009), where it becomes increasingly likely that only a few exemplars (or even none at all) are mapped to a given entry.

The operating principle of JQSP2 is quite different as it formulates, given an original quality factor QF_j , a scaling z , a prediction on the QF_{out} needed to meet, without exceeding, a target file size (it also predicts resulting file size and prediction as a by-product), and has linear-time complexity for training. JQSP2's design makes it less likely to overshoot significantly on file size (Coulombe & Pigeon, 2010). The training phase constructs the table by going through all the exemplars in the training set and for each exemplar, it updates partial sums indexed by \widetilde{QF}_j , \widetilde{z} , and \widetilde{f} ; which is performed in $O(n)$. The partial sums are then normalized at the cost of $O\left(\left|\widetilde{QF}_j\right|\left|\widetilde{z}\right|\left|\widetilde{f}\right|\right)$, which is again negligible compared to the scanning of the training set and the updates of the partial sums.

EJQSP formulates its prediction by clustering, as described in Section 2. However, solving eq. (2) exactly is NP-Hard (Aloise et al., 2010; Mahajan et al., 2009; Vattani, 2009), and one has to revert to an approximate algorithm such as the Linde-Buzo-Gray (Linde et al., 1980) or K-means (Lloyd, 1982). An iteration for K-means is $O(dmn)$ for m prototypes and n exemplars in \mathbb{R}^d . Since $O(\lg n)$ iterations seem to be sufficient to bring K-Means to a converging solution, even for moderate m and n , we get an over-all complexity of $O(dmn \lg n)$.

Predicting resulting file size and quality or transcoding operations from algorithm JQSP1 and JQSP2 is a constant-time process as it suffices to quantize the appropriate parameters and use the quantized versions to index a table

containing the desired prediction. In both cases, we assume that quantizing the parameters is also a constant time operation; or at least, constant in the sense that its complexity does not depend on the number of exemplars used for training and only loosely on the table density. One can think of a quantization that is essentially a “rounding” of values, which certainly can be performed in constant-time. Other, possibly PDF-optimized (Graf & Luschgy, 2000), quantization schemes could be used, and in this case the cost would be at most proportional to the entropy of the quantized values. The clustering method proposed in this work does not require quantization (at least, not explicitly as a preprocessing stage) but was, for the sake of comparability with previous work, trained using the same training exemplars which have quality factors and scalings constrained to a small set of possible values (Section 3). The clustering-based predictor, EJQSP, therefore simply takes the original data, computes the features (in constant time) and searches for the prototype closest to the test exemplar (ignoring quantities f and q , as those we want to predict and are unknown in a prediction-time exemplar); which is essentially nearest neighbor search between one point, the exemplar, and the m prototypes. Although with some preprocessing (Mahajan et al., 2009) or approximate search (Indyk & Motwani, 1998) nearest neighbor search can be made sub-linear, we consider it requires linear time; therefore, if the prototypes lie in \mathbb{R}^d (including features), the search is $O(md)$.

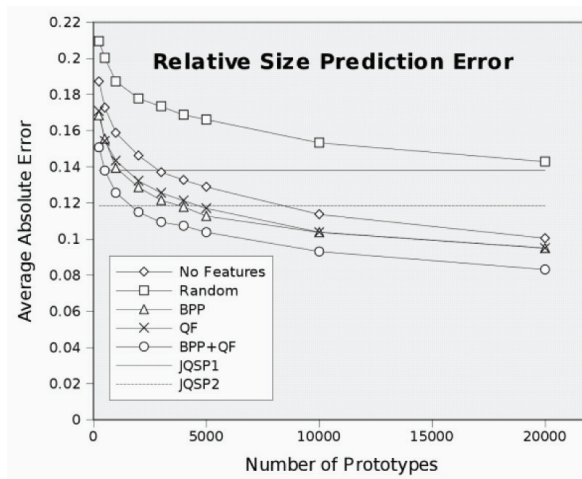
Let us note that while the table size in either algorithm JQSP1 or JQSP2 is limited upwards by the problem of context dilution—where, at some point, there are not enough exemplars to fill all the entries in the table or with sufficiently low variance—algorithm EJQSP degrades gracefully because even if it is presented an exemplar that resembles no other exemplar it has seen during the training phase, it can still formulate a prediction using the closest prototypes. In the worst case, JQSP1 and JQSP2 could fail because the corresponding cell in the

array is empty. To avoid this problem, arrays in algorithms JQSP1 and JQSP2 must be kept small enough so that each cell is susceptible to receive a sufficiently large number of exemplars.

The input can be transformed or augmented using features to make information available to the prediction algorithm, information that would be otherwise hard to discover (in a machine-learning sense). Finding good transformations or good features is not, in general, a trivial task. *A priori* knowledge can help us add a very small number of very effective features. In our experiments, we devised two such features. The first feature is the number of bits per pixel of the image, denoted b_j for image I_j , will help distinguish images of the same resolution but of different file sizes. The rationale is that the ratio of the file size to the resolution is indicative of the intrinsic complexity of the image. At equal resolution, an image representing only a featureless blue sky will have a rather small file size, while an image representing a complex natural scene—say a picture taken in a forest—will likely have a much larger file size. They will also behave differently under transformations, and this knowledge will help the predictor formulate more accurate predictions. In the same way, the quality factor difference feature, $QF_{out} - QF_j$, encodes the drop in quality incurred by the transcoding operation, and will help extract information about transcodings that have a similar drop in quality onto similar hyperplanes, thus helping prediction, which is the primary criterion for retaining a feature over another.

Indeed, features can be heuristic, formulated from *a priori* knowledge, or even random; but in all cases, features are to be evaluated, and added only if they ameliorate prediction. Figure 5 and Figure 6 present our experiments for feature selection. In these experiments, we compare clustering using different vector lengths. First, we compare with vectors representing only the images’ basic features (that is, using only QF_j , w_j , h_j , z , and QF_{out}), then we compare using vectors formed from the

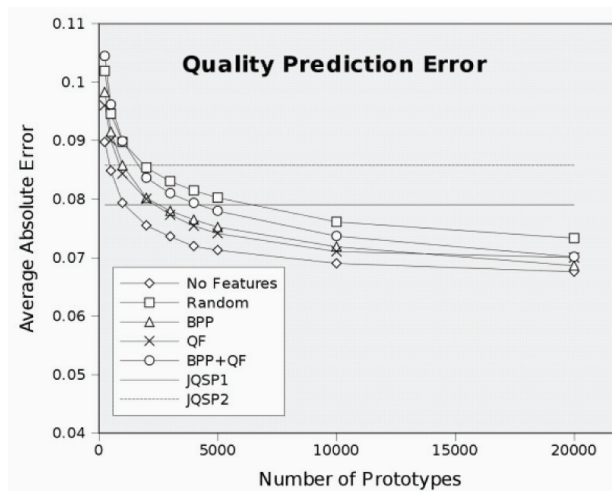
Figure 5. Features' influence on relative file size prediction (relative size prediction error)



basic features plus a uniform random feature (a random variable uniformly drawn on $[0, 1]$, different for each vector), then with the basic vectors augmented with b_j , then with the basic vectors augmented by $QF_{out} - QF_j$, and finally with vectors augmented by both b_j and $QF_{out} - QF_j$. In Figure 5 and Figure 6, the experiments are respectively denoted as no features, random, BPP, QF, and BPP + QF. In

Figure 5, we compare the effects of the different features (or combinations of features) on the relative file size prediction error. We see that no feature is preferable to adding a random feature (thus showing it is not very useful to select random features) but that successively adding the proposed features helps the clustering achieving better relative file size prediction. Indeed, it is with both added features that prediction error is minimized, at 0.083, using 20000

Figure 6. Features' influence on quality prediction (quality prediction error)



prototypes. The situation is different when we consider quality prediction, as shown in Figure 6. Again, the random feature is worse; but using no features is slightly better than using both features. However, using a large number of prototypes, the difference in quality prediction error between using both features and no feature is likely not significant (in the order of 0.003) whereas the difference in relative file size prediction is more important (in the order of 0.01). The experiments show that feature selection, especially when predicting more than one quantity, may lead to application-specific trade-offs. In our feature selection, we explicitly favored precision on relative file size prediction in view of applications such as Pigeon and Coulombe (2011).

The memory usage is also worth discussing. As JQSP1 and JQSP2 training procedures consist in streaming in exemplars one after the other and discarding them after usage, the memory usage is limited to the table needed to store intermediate results. The size of the table is determined by the quantization imposed on the parameters. For JQSP1, for example, the quantization on QF_j , QF_{out} , and z , as used in previous works, yields a $10 \times 10 \times 10$ table with a small number of values stored in each entry (the cumulated file sizes, qualities, and number of exemplars mapping to this entry), which is unlikely to pose problems in a server-type environment. The memory needed by JQSP2 is also determined by the quantization of its input, QF_j , z , and desired file size f_{max} , and so its memory usage can also be quite moderate, even if the target file sizes are rather finely quantized. For EQJSP, the memory usage during training can be made $O(md)$ for m prototypes if one streams the n exemplars from external storage, but one would likely keep all exemplars in memory for faster iterations, yielding $O((m+n)d)$ storage. During run-time, where only prediction is needed, the storage is brought back to $O(md)$, which is also likely to be essentially negligible in a server-type environment.

Of course, there are trade-offs between storage, algorithmic complexity, and prediction accuracy to consider. Since one would think that storage, especially in a server-type environment, will be negligible even for large m , the main trade-off for EJQSP will be between algorithmic complexity and accuracy. Accuracy plays an important role in applications such as MMS adaptation where the task is not to adapt a single image, but a series of images which are part of the same message (Pigeon & Coulombe, 2011, 2012). In this type of optimization problem, errors propagate and do not necessarily cancel out, and it is therefore preferable to favor accuracy over the price of slightly increased complexity (which we will mitigate in the Appendix) so that the final adaptation quality is not jeopardized.

6. CONCLUSION

Despite formulating predictions in (at most) linear time in the number of prototypes, algorithm EJQSP accuracy outperforms the constant-time JQSP1 and JQSP2 algorithms we presented in previous work. Algorithm EJQSP, using 20000 prototypes, yields significantly better prediction of resulting file size and quality of JPEG images subject to transcoding operations. It yields $\approx 40\%$ smaller prediction error on file size and $\approx 12\%$ on quality than algorithm JQSP1, while yielding $\approx 27\%$ smaller prediction error on file size and $\approx 20\%$ on quality than algorithm JQSP2. The new predictor, with its reduced error, can be combined with systems such as those presented in Coulombe and Pigeon (2009, 2010) and Pigeon and Coulombe (2011, 2012) to yield more efficient and more precise transcoding systems, whether for universal media access, mobile browsing, or multimedia messaging services.

ACKNOWLEDGMENT

This work was funded by Vantrix Corporation and by the Natural Sciences and Engineering Research Council of Canada under the Col-

laborative Research and Development Program (NSERC-CRD 326637-05).

REFERENCES

- Aloise, D., Deshpande, A., Hansen, P., & Papat, P. (2010). NP-hardness of Euclidean sum-of-squares clustering. *Machine Learning*, 75(2), 245–248. doi:10.1007/s10994-009-5103-0
- Azzalini, A. (1985). A class of distributions that includes the normal ones. *Scandinavian Journal of Statistics*, 12, 171–178.
- Barni, M. (1997). A fast algorithm for l-norm vector median filtering. *IEEE Transactions on Image Processing*, 6(10), 1452–1455. doi:10.1109/83.624972
- Blackman, R. B., & Tukey, J. W. (1959). *The measurement of power spectra, from the point of view of communications engineering*. Mineola, NY: Dover.
- Bottou, L., & Bengio, Y. (1995). Convergence properties of the k-means algorithms. In Tesauro, G., & Touretzky, D. S. (Eds.), *Advances in neural information processing systems* (Vol. 7, pp. 585–592). Cambridge, MA: MIT Press.
- Carver, R. H., & Tai, K.-C. (2005). *Modern multithreading implementing, testing, and debugging multithreaded Java and C++/Pthreads/Win32 programs*. Hoboken, NJ: Wiley-Interscience. doi:10.1002/0471744174
- Coulombe, S., & Grassel, G. (2004). Multimedia adaptation for the multimedia messaging service. *IEEE Communications Magazine*, 42(7), 120–126. doi:10.1109/MCOM.2004.1316543
- Coulombe, S., & Pigeon, S. (2009). Quality-aware selection of quality factor and scaling parameters in JPEG image transcoding. In *Proceedings of the IEEE International Conference on Computational Intelligence for Multimedia, Signal, and Video Processing* (pp. 68–74).
- Coulombe, S., & Pigeon, S. (2010). Low-complexity transcoding of JPEG images with near-optimal quality using a predictive quality factor and scaling parameters. *IEEE Transactions on Image Processing*, 19(3), 712–721. doi:10.1109/TIP.2009.2036716
- Fling, B. (2009). *Mobile design and development: Practical concepts and techniques for creating mobile sites and Web apps - Animal guide*. Sebastopol, CA: O'Reilly Media.
- Graf, S., & Luschgy, H. (2000). *Foundations of quantization for probability distributions (Lecture Notes in Mathematics)*. New York, NY: Springer.
- Han, R., Bhagwat, P., LaMaire, R., Mummert, T., Perret, V., & Rubas, J. (1998). Dynamic adaptation in an image transcoding proxy for mobile Web browsing. *IEEE Personal Communications Magazine*, 5(6), 8–17. doi:10.1109/98.736473
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning*. New York, NY: Springer.
- IJG. (2012). *The independent JPEG group*. Retrieved March 10, 2012, from <http://www.ijg.org/>
- Indyk, P., & Motwani, R. (1998). Approximate nearest neighbor: Towards removing the curse of dimensionality. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing* (pp. 604–613).
- ISO/IEC. (2011). *10918-5: Information technology, digital compression and coding of continuous-tone still images: JPEG File Interchange Format (JFIF)*. Geneva, Switzerland: ISO/IEC.
- Lei, Z., & Georganas, N. D. (2002). Accurate bit allocation and rate control for DCT domain video transcoding. In *Proceedings of the IEEE Canadian Conference on Electrical and Computer Engineering* (pp. 968–973).
- Linde, Y., Buzo, A., & Gray, R. M. (1980). An algorithm for vector quantizer design. *IEEE Transactions on Communications*, 28(1), 84–95. doi:10.1109/TCOM.1980.1094577
- Lloyd, S. P. (1982). Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2), 129–137. doi:10.1109/TIT.1982.1056489
- Mahajan, M., Nimbhorkar, P., & Varadarajan, K. (2009). The planar k-means problem is NP-hard. In *Proceedings of the 3rd International Workshop on Algorithms and Computation* (pp. 274–285).
- Mattson, T. G., Sanders, B. A., & Massingill, B. (2005). *Patterns for parallel programming*. Reading, MA: Addison-Wesley.
- Mohan, R., Smith, J. R., & Li, C.-S. (1999). Adapting multimedia internet content for universal access. *IEEE Transactions on Multimedia*, 1(1), 104–114. doi:10.1109/6046.748175
- Open Mobile Alliance. (2010). *Enabler test specification (for conformance) for MMS Candidate Version 1.3*. San Diego, CA: Author.

- Pennebaker, W. B., & Mitchell, J. L. (1993). *JPEG still image data compression standard*. New York, NY: Van Nostrand Reinhold.
- Pigeon, S., & Coulombe, S. (2008). Computationally efficient algorithms for predicting the file size of JPEG images subject to changes of quality factor and scaling. In *Proceedings of the 24th Queen's University Biennial Symposium on Communications* (pp. 378-382).
- Pigeon, S., & Coulombe, S. (2011). Optimal quality-aware predictor-based adaptation of multimedia messages. In *Proceedings of the IEEE 6th International Conference on Intelligent Data Acquisition and Advanced Computing Systems* (pp. 496-499).
- Pigeon, S., & Coulombe, S. (2012). Optimal quality-aware predictor-based adaptation of multimedia messages. In Duro, R. (Ed.), *Digital image, signal and data processing*. Keswick, UK: Rivers. doi:10.1109/IDAACS.2011.6072803
- Ratnakar, V., & Ivashin, V. (2001). *U.S. Patent No. 6,233,359: File size bounded JPEG transcoder*. Washington, DC: United States Patent and Trade-mark Office.
- Ridge, J. (2003). Efficient transform-domain size and resolution reduction of images. *Signal Processing Image Communication*, 18(8), 621–639. doi:10.1016/S0923-5965(03)00056-0
- Shu, H., & Chau, L.-P. (2005). Frame size selection in video downsizing transcoding application. In *Proceedings of the International Symposium on Circuits and Systems* (pp. 896-899).
- Vattani, A. (2009). K-means require exponentially many iterations even in the plane. In *Proceedings of the 25th Symposium on Computational Geometry* (pp. 324-332).
- Wang, Z., Bovick, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4), 600–612. doi:10.1109/TIP.2003.819861
- Williams, A. (2007). *Boost threads*. Retrieved March 10, 2012, from http://www.boost.org/doc/libs/1_48_0/doc/html/thread.html

APPENDIX

Efficient Implementation of K-Means

Translating Algorithm 1 directly to a programming language such as C would yield a rather straightforward implementation of K-means, where the program scans the exemplars one by one, finds the nearest prototypes, assigns them to the corresponding partition, then uses the assignment to update the prototypes, repeating the procedure until error ceases to decrease significantly. A careful, but sequential, implementation of this procedure will fail to exploit the inherent parallelism of the algorithm. For example, searching for the nearest prototype can be performed in parallel for every exemplar very efficiently because there are no dependencies between the exemplars: either the values are read-only (the x_j , the $\bar{x}_{t,i}$) or write-only (the a_j), thus dispensing us entirely of the need for mutual exclusion mechanisms. The complexity of one

iteration can therefore be reduced from $O(dmn)$ to $O\left(dm\frac{n}{c}\right)$, where c is the number of available cores (or hardware threads). Other parts can also be parallelized, such as the update of the prototypes, but not as easily; one would partition the problem by groups of n/c exemplars and hold c series of m partial sums to be combined once all the exemplars are processed.

Parallelism can be obtained by various means, mostly depending on the programming language chosen, in our case C++, such as POSIX Threads (pthreads) (Carver & Tai, 2005), Boost threads (Williams, 2007), but the simplest mean by far is to use OpenMP (Mattson, Sanders, & Massingill, 2005). OpenMP is a compiler extension that reduces the parallelization of a classical C (or C++) program to little more than the addition of a few well-placed compiler-specific #pragmas that specify parallel for-loops, memory fences, and reductions.

A parallel implementation can mitigate the complexity of prediction as well. Reducing the (sequential) complexity from $O(md)$ to $O\left(\frac{m}{c}d\right)$ may mean a significant speed-up especially that c can be quite large in modern server-type shared-memory CPUs. Furthermore, as the operation is read-only, there is no need for mutual exclusion, and synchronization is limited to waiting for all sub-problems to terminate and combine (sequentially) the c sub-answers into the final answer.