CVPR
#1988

CVPR
#1988

CVPR 2013 Submission #1988. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

# Sparse Quantization for Patch Description
# Supplementary Material

Anonymous CVPR submission

Paper ID 1988

In the paper we left several derivations for the supplementary material. We detail here these derivations, following the same order of appearance as in the paper. We provide the derivations of the proof of Proposition 1, and Proposition 2. We also provide an extension of the relation of SQ with other encodings.

## 1. Sparse Quantization

**Proposition 1.** *Let $\hat{\mathbf{v}}^\star$ be the quantization into $\mathbb{T}_k^q$ of $\mathbf{v} \in \mathbb{R}^q$, which is $\hat{\mathbf{v}}^\star = \arg\min_{\hat{\mathbf{v}} \in \mathbb{T}_k^q} \|\hat{\mathbf{v}} - \mathbf{v}\|^2$. For $\|\mathbf{v}\|_2 \leq \|\mathbf{s}\|_2/\sqrt{k}$, where $\mathbf{s} \in \mathbb{T}_k^q$, $\hat{\mathbf{v}}^\star$ can be computed by*

$$\hat{v}_i^\star = \begin{cases} sign(v_i) & \text{if } i \in k\text{-Highest}(|\mathbf{v}|) \\ 0 & \text{otherwise} \end{cases}, \qquad (1)$$

*where $|\mathbf{v}|$ is the element-wise absolute value of $\mathbf{v}$, and $k$-Highest$(|\mathbf{v}|)$ is the set of dimension indices that indicate which are the $k$ elements of the vector $|\mathbf{v}|$ with the highest values.*

*Proof.* We first rewrite $\|\hat{\mathbf{v}} - \mathbf{v}\|^2$ as $\sum_i (\hat{v}_i - v_i)^2$. Since $\hat{\mathbf{v}} \in \mathbb{T}_k^q$ has $k$ elements set to 1 or $-1$ and $(q-k)$ set to 0, we can write the above summation as

$$\sum_i (\hat{v}_i - v_i)^2 =$$
$$\sum_{i:\hat{v}_i=(+1)} ((+1) - v_i)^2 + \sum_{i:\hat{v}_i=(-1)} ((-1) - v_i)^2 + \sum_{i:\hat{v}_i=0} (v_i)^2 \qquad (2)$$

We sort in descending order the absolute value of the set of values at each dimension of $\mathbf{v}$, *i.e.* we sort $\{|v_i|\}$, and we use a new indexing in this ordered set. We indicate so by using $\mathbf{v}'$, and we index it with $s$ instead of $i$, such that $|v'_{(s-1)}| > |v'_s|$. To see when (2) is minimum, note that

$$(v'_1)^2 > \ldots > (v'_{(s-1)})^2 > (v'_s)^2 > \ldots; \qquad (3)$$
$$(1-|v'_1|)^2 < \ldots < (1-|v'_{(s-1)}|)^2 < (1-|v'_s|)^2 < \ldots; \qquad (4)$$

where (4) is due to the assumption $\|\mathbf{v}\|_2 \leq \|\mathbf{s}\|_2/\sqrt{k}$, and it is equivalent to

$$\ldots < (sign(v'_{(s-1)}) - sign(v'_{(s-1)})|v'_{(s-1)}|)^2 <$$
$$< (sign(v'_s) - sign(v'_s)|v'_s|)^2 < \ldots. \qquad (5)$$

We rewrite Eq. (2):

$$\sum_i (\hat{v}_i - v_i)^2 =$$
$$\sum_{i:\hat{v}_i \neq 0} (sign(v_i) - sign(v_i)|v_i|)^2 + \sum_{i:\hat{v}_i=0} (v_i)^2 \qquad (6)$$

Therefore, to make the two terms in (6) minimum, we set the $k$ elements in $\hat{\mathbf{v}}$ to $sign(v_i)$ such that $(sign(v'_s) - sign(v'_s)|v'_s|)^2$ in (5) are minimum, and we set $(q-k)$ zeros in $\hat{\mathbf{v}}$ such that $(v'_s)^2$ in (3) are minimum. Thus, we set the $k$ highest values of $|\mathbf{v}|$ to $sign(v_i)$, and 0 to the other $(q-k)$ values. $\square$

## 2. SQ for Encoding in Patch Description

**Definition 1.** *Let $\boldsymbol{\alpha}^\star \in \mathbb{R}_k^Q$ be the encoding of $\mathbf{f} \in \mathbb{R}^q$ such that*

$$\boldsymbol{\alpha}^\star = \arg\min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} \|\boldsymbol{\alpha} - \Psi(\mathbf{f}, \bigcup_{0<p\leq q} \bar{\mathbb{T}}_p^q)\|^2. \qquad (7)$$

**Proposition 2.** *Let $q \leq 4$ and $k \leq 2$. Then, Algorithm 1 obtains the global minimum for $\boldsymbol{\alpha}^\star$ in Definition 1 with computational complexity $O(q^2)$.*

---

**Algorithm 1**: Sparse Quantization in Proposition 2

**Input**: $\mathbf{f} \in \mathbb{R}^q$
**Output**: $\boldsymbol{\alpha}^\star \in \mathbb{R}_k^Q$
**forall** $0 < p \leq q$ **do**
$\quad \boldsymbol{\beta}_p^\star = \arg\min_{\boldsymbol{\beta} \in \mathbb{T}_p^q} \|\boldsymbol{\beta} - \mathbf{f}\|_2$
**end**
$\boldsymbol{\alpha}^\star = \arg\min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} \|\boldsymbol{\alpha} - \tilde{\Psi}(\mathbf{f}, \{\boldsymbol{\beta}_p^\star\})\|_2$

---

CVPR
#1988

CVPR
#1988

CVPR 2013 Submission #1988. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

*Proof.* The Proposition is saying that if we constraint $k \leq 2$ and $q \leq 4$, we can assure that the minimum of

$$\boldsymbol{\alpha}^{\star} = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} \|\boldsymbol{\alpha} - \Psi(\mathbf{f}, \bigcup_{0 < p \leq q} \bar{\mathbb{T}}_p^q)\|^2 \qquad (8)$$

can be achieved by splitting the optimization into the optimization of each of the sets of the codebook $\bar{\mathbb{T}}_p^q$ independently, and then picking the elements that has higher similarity measure between $\mathbf{f}$ and the selected $\bar{\mathbb{T}}_p^q$. We use $\boldsymbol{\beta}_p^{\star}$ to denote the candidate selected for the set $\bar{\mathbb{T}}_p^q$, which is the solution of $\arg \min_{\boldsymbol{\beta} \in \mathbb{T}_p^q} \|\boldsymbol{\beta} - \mathbf{f}\|_2$. We define the *second* best solution of such SQ as $\boldsymbol{\beta}_p^{\star 2}$, which are the discarded candidates that are closer to $\boldsymbol{\beta}_p^{\star}$.

Let $\mathbf{f}'$ be the vector $\mathbf{f}$ such that the higher elements are at the beginning of the vector, *i.e.* $f_1' > \cdots > f_q'$. In Proposition 1, we showed that $\boldsymbol{\beta}_p^{\star}$ can be constructed selecting the first $p$ components of $\mathbf{f}'$. We can also see from the proof of Proposition 1, that $\boldsymbol{\beta}_p^{\star 2}$, consists on selecting the $p - 1$ first components of $\mathbf{f}'$ and also the $p + 1$ component. In this way, we maintain the $p - 1$ components with lower reconstruction error, and we only change the $p$-th term for the $(p+1)$-th, which keeps $p$ non-zero components and the reconstruction error is the closest to the one of the optimal SQ. When $p = q$, the second best solution is built by changing the sign of the $p$-th term, since it does not exist the $(p+1)$-th term, since $p \leq q$. Observe that the distance between $\boldsymbol{\beta}_p^{\star}$ and $\mathbf{f}'$ is

$$d(\boldsymbol{\beta}_p^{\star}, \mathbf{f}') = \sum_{j=1}^p (\frac{1}{\sqrt{p}} - |f_j'|)^2 + \sum_{j=p+1}^q (f_j')^2. \qquad (9)$$

The distance between the $\boldsymbol{\beta}_p^{\star 2}$ and $\mathbf{f}'$, in case $p < q$, is

$$d(\boldsymbol{\beta}_p^{\star 2}, \mathbf{f}') =$$
$$\sum_{j=1}^{p-1} (\frac{1}{\sqrt{p}} - |f_j'|)^2 + (\frac{1}{\sqrt{p}} - |f_{p+1}'|)^2 + (f_p')^2 + \sum_{j=p+2}^q (f_j')^2, \qquad (10)$$

and for $p = q$ is

$$d(\boldsymbol{\beta}_q^{\star 2}, \mathbf{f}') = \sum_{j=1}^{q-1} (\frac{1}{\sqrt{q}} - |f_j'|)^2 + (\frac{1}{\sqrt{q}} + |f_q'|)^2. \qquad (11)$$

The proof of the Proposition consists on verifying that the $k$ components of the set $\bigcup_{0 < p \leq q} \bar{\mathbb{T}}_p^q$, that have higher similarity measure with $\mathbf{f}$, are always in $\{\boldsymbol{\beta}_p^{\star}\}$ and never in $\{\boldsymbol{\beta}_p^{\star 2}\}$. Note that showing that the closest elements to $\mathbf{f}$ are never in $\{\boldsymbol{\beta}_p^{\star 2}\}$, means that they necessarily are in $\{\boldsymbol{\beta}_p^{\star}\}$. This is equivalent to show that

$$d(\boldsymbol{\beta}_a^{\star}, \mathbf{f}') \leq d(\boldsymbol{\beta}_b^{\star 2}, \mathbf{f}'), \qquad (12)$$

for all $a, b \leq q$. This condition is to verify that in general Algorithm 1 obtains the global maximum for $k \leq q$. In the following, we are only able to show that for $k \leq 2$ and $q \leq 4$.

When $k = 1$ it is trivial to proof the Proposition, since one of the elements in $\{\boldsymbol{\beta}_p^{\star}\}$ is necessarily the closest element to $\mathbf{f}$. For $k = 2$, this might not be the case, because for a certain $p$, $\boldsymbol{\beta}_p^{\star}$ and $\boldsymbol{\beta}_p^{\star 2}$ can be the closest elements to $\mathbf{f}$, rather than two different $\boldsymbol{\beta}_p^{\star}$ with different $p$'s. We use $\boldsymbol{\beta}_o^{\star}$ to denote the closest element to $\mathbf{f}$ in $\bigcup_{0 < p \leq q} \mathbb{T}_p^q$, which is the solution for $k = 1$, and we denote $\boldsymbol{\beta}_o^{\star 2}$ as the discarded candidate which is closest to $\boldsymbol{\beta}_o^{\star}$ in $\bar{\mathbb{T}}_o^q$. Thus, the proof for $k = 2$, consists on validating that the second closest element to $\mathbf{f}$ is not $\boldsymbol{\beta}_o^{\star 2}$, and that it is in $\{\boldsymbol{\beta}_o^{\star}\}$. We develop Eq. (12) using the distances previously calculated in Eq. (9) and (10):

$$d(\boldsymbol{\beta}_a^{\star}, \mathbf{f}') \leq d(\boldsymbol{\beta}_o^{\star 2}, \mathbf{f}') \iff \sqrt{\frac{o}{a}} \geq \frac{\sum_{j=1}^{o-1} f_j' + f_{o+1}}{\sum_{j=1}^a f_j'}. \qquad (13)$$

Thus, if it always exist a $\boldsymbol{\beta}_a^{\star}$ that verifies Eq. (13), we prove that for $k \leq 2$, Algorithm 1 finds the optimal SQ. We show that either $\boldsymbol{\beta}_{o+1}^{\star}$ and $\boldsymbol{\beta}_{o-1}^{\star}$ always fulfill such condition, for $q \leq 4$. For notation simplicity, we define $K = \sum_{j=1}^{o-1} f_j' + f_{o+1}$. Thus, for $\boldsymbol{\beta}_{o-1}^{\star}$, Eq. (13) becomes:

$$d(\boldsymbol{\beta}_{o-1}^{\star}, \mathbf{f}') \leq d(\boldsymbol{\beta}_o^{\star 2}, \mathbf{f}') \iff 1 - \frac{f_{o+1}'}{K} \geq \sqrt{1 - \frac{1}{o}}, \qquad (14)$$

and for $\boldsymbol{\beta}_{o+1}^{\star}$ (for $o < q$):

$$d(\boldsymbol{\beta}_{o+1}^{\star}, \mathbf{f}') \leq d(\boldsymbol{\beta}_o^{\star 2}, \mathbf{f}') \iff 1 + \frac{f_o'}{K} \geq \sqrt{1 + \frac{1}{o}}. \qquad (15)$$

From Eq. (14) and (15), and taking into account that $f_1' > \cdots > f_q'$, we can verify algebraically, that for $o \leq 3$, either Eq. (14) or (15) are fulfilled, or both. For $o = 4$, and if $q = 4$, it can be verified in the same way that $d(\boldsymbol{\beta}_3^{\star}, \mathbf{f}') \leq d(\boldsymbol{\beta}_4^{\star 2}, \mathbf{f}')$, where $d(\boldsymbol{\beta}_4^{\star 2}, \mathbf{f}')$ takes the form in Eq. (11), since for this case $p = q$. $\square$

## 3. Relation with Other Encodings

**Sparse Coding.** The formulation for the kernelized sparse coding [1] is

$$\boldsymbol{\alpha}^{\star} = \arg \min_{\alpha_i \in \mathbb{R}_k^Q} = \|\phi(\mathbf{x}) - \sum_i \alpha_i \phi(\mathbf{b}_i)\|^2, \qquad (16)$$

CVPR
#1988

CVPR
#1988

CVPR 2013 Submission #1988. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

where $\phi(\mathbf{x})$ is a non-linear mapping of $\mathbf{x}$. Eq. (16) can be decomposed in the following terms:

$$K(\mathbf{x}, \mathbf{x}) - 2\sum_i \alpha_i K(\mathbf{x}, \mathbf{b}_i) + \sum_i \sum_j \alpha_i \alpha_j K(\mathbf{b}_i, \mathbf{b}_j),$$
(17)

in which $K(\mathbf{x}, \mathbf{b}_i) = \phi(\mathbf{x})^T \phi(\mathbf{b}_i)$. $K(\mathbf{x}, \mathbf{x})$ can be treated as a constant because it does not influence on the optimization problem. Thus, the optimization becomes

$$\arg \min_{\alpha_i \in \mathbb{R}_k^Q} -2\sum_i \alpha_i K(\mathbf{x}, \mathbf{b}_i) + \sum_i \sum_j \alpha_i \alpha_j K(\mathbf{b}_i, \mathbf{b}_j).$$
(18)

Recall that the encoding with SQ can be formulated as (Eq. (3) in the paper)

$$\boldsymbol{\alpha}^\star = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} \|\boldsymbol{\alpha} - \Psi(\mathbf{f}, \{\mathbf{b}_i\})\|^2,$$
(19)

which decomposes into

$$\boldsymbol{\alpha}^T \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^T \Psi(\mathbf{f}, \{\mathbf{b}_i\}) + \Psi(\mathbf{f}, \{\mathbf{b}_i\})^T \Psi(\mathbf{f}, \{\mathbf{b}_i\}),$$
(20)

in which $\boldsymbol{\alpha}^T \boldsymbol{\alpha}$ is constant because of the constraint $\boldsymbol{\alpha} \in \mathbb{R}_k^Q$, and $\Psi(\mathbf{f}, \{\mathbf{b}_i\})^T \Psi(\mathbf{f}, \{\mathbf{b}_i\})$ can also be dropped because it does not depend on $\boldsymbol{\alpha}$, and hence, it does not influence in the minimization. Thus, the optimization becomes

$$\arg \min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} -\boldsymbol{\alpha}^T \Psi(\mathbf{f}, \{\mathbf{b}_i\}).$$
(21)

Noting that each entry of $\Psi$ corresponds to the similarity measure $K(\mathbf{f}, \mathbf{b}_i)$, we can rewrite Eq. (20) as:

$$\arg \min_{\boldsymbol{\alpha} \in \mathbb{R}_k^Q} -\sum_i \alpha_i K(\mathbf{f}, \mathbf{b}_i).$$
(22)

We can see that the main difference between SQ and the kernelized version of Sparse Coding lies in the term $\sum_i \sum_j \alpha_i \alpha_j K(\mathbf{b}_i, \mathbf{b}_j)$, which is a regularization term.

**Convolutional Networks.** Let $\mathbf{W} \in \mathbb{R}^{q \times m}$ be the matrix containing the filters we use in our formulation to extract the features, and let $\mathbf{x} \in \mathbb{R}_+^m$ be the raw image where $\mathbf{W}$ is applied. Thus, $\mathbf{f} = \mathbf{W}\mathbf{x} \in \mathbb{R}^q$. Here we show that

$$\Psi(\mathbf{f}, \{\mathbf{b}_i\}) \propto \Psi(\mathbf{x}, \{\frac{1}{w}\mathbf{W}^T \mathbf{b}_i\}) \in \mathbb{R}^Q,$$
(23)

where

$$\Psi(\mathbf{f}, \{\mathbf{b}_i\}) = \frac{1}{Z}(K(\mathbf{f}, \mathbf{b}_1) \ldots K(\mathbf{f}, \mathbf{b}_Q)) \in \mathbb{R}^Q.$$
(24)

$\mathbf{f} = \mathbf{W}\mathbf{x} \in \mathbb{R}^q$, $\mathbf{W} \in \mathbb{R}^{q \times m}$ and $\mathbf{x} \in \mathbb{R}_+^m$, $w$ is a normalization factor, and we assume that $\mathbf{f}$ and $\mathbf{b}_i$ are $\ell_2$-normalized.

First we decompose the left hand side of Eq. (23) for $\mathbf{b}_i$, assuming that we use the Gaussian kernel similarity, which is the one we use in the paper. Thus,

$$K(\mathbf{f}, \mathbf{b}_i) = \exp\left(-\frac{\|\mathbf{f} - \mathbf{b}_i\|^2}{\sigma^2}\right) =$$
(25)

$$\exp\left(-\frac{\|\mathbf{f}\|^2 + \|\mathbf{b}_i\|^2}{\sigma^2}\right)\exp\left(-\frac{-2\mathbf{f}^T \mathbf{b}_i}{\sigma^2}\right) =$$
(26)

$$K_1 \exp\left(\frac{2\mathbf{f}^T \mathbf{b}_i}{\sigma^2}\right),$$
(27)

wher $K_1$ is a constant since $\mathbf{f}$ and $\mathbf{b}_i$ are normalized.

Now we develop the right hand side of Eq. (23) with the same assumptions. This is

$$\Psi(\mathbf{x}, \frac{1}{w}\mathbf{W}^T \mathbf{b}_i) = K(\mathbf{x}, \frac{1}{w}\mathbf{W}^T \mathbf{b}_i) =$$
(28)

$$\exp\left(-\frac{\|\mathbf{x} - \frac{1}{w}\mathbf{W}^T \mathbf{b}_i\|^2}{\sigma^2}\right) =$$
(29)

$$\exp\left(\frac{\|\mathbf{x}\|^2 + \|\frac{1}{w}\mathbf{W}^T \mathbf{b}_i\|^2}{\sigma^2}\right)\exp\left(-\frac{-2\mathbf{x}^T \mathbf{W}^T \mathbf{b}_i}{w\sigma^2}\right) =$$
(30)

$$K_2 \exp\left(\frac{2\mathbf{f}^T \mathbf{b}_i}{w\sigma^2}\right),$$
(31)

in which we use the normalization factor $w$ and the equivalence $\mathbf{f} = \mathbf{W}\mathbf{x}$ to go from Eq. (30) to (31). Then, since $K_1$ and $K_2$ are two constants, we can recover the proportion of Eq. (23), *i.e.*

$$K_1 \exp\left(\frac{2\mathbf{f}^T \mathbf{b}_i}{\sigma^2}\right) \propto K_2 \exp\left(\frac{2\mathbf{f}^T \mathbf{b}_i}{w\sigma^2}\right),$$
(32)

which can be extended to all the set $\{\mathbf{b}_i\}$.

## References

[1] M. T. Harandi, C. Sanderson, R. Hartley, and B. C. Lovell. Sparse coding and dictionary learning for symmetric positive definite matrices: A kernel approach. In *ECCV*, 2012.