

Information k -means, fragmentation and syntax analysis: A new approach to unsupervised Machine Learning

Gautier Appert

PhD supervisor : Olivier Catoni

October 29, 2020

Given a sample X_1, \dots, X_n

of n independent copies of a signal $X \in \mathbb{R}^d$ (we will focus on digital images), we want to produce a better representation of its content.

We want to

- 1 identify significant patterns
- 2 represent their interactions

We propose to

- 1 perform some vector quantization of image fragments
- 2 compute syntax trees for the quantized fragments

both steps being based on a data compression criterion (performing lossy coding in the first step and lossless coding in the second one).

In the large sample limit we can recover the signal distribution

- Let $\bar{X} = (X_1, \dots, X_n) \in \mathbb{R}^{nd}$ and let $\theta(\bar{X})$ be a (lossy) binary prefix code for \bar{X} (we code the whole sample \bar{X} , not a single random image X).
- Let $\bar{Y} = f[\theta(\bar{X})]$ be the lossy decoding of \bar{X} from its binary representation $\theta(\bar{X})$.
- Assume that $\|\bar{Y} - \bar{X}\|^2 \leq nd\alpha$, for some distortion level $\alpha > 0$.
- Introduce a coding distribution Q_θ . Assume that $Q_{\bar{Y}} = Q_\theta \circ f^{-1} = Q_{f(\theta)}$ is exchangeable.
- Consider a blurred version of the sample \bar{X} defined as $\bar{V} = (V_1, \dots, V_n) = \bar{X} + \bar{W}$, where \bar{W} is independent of the sample \bar{X} and distributed as $\mathcal{N}(0, \sigma^2 I_{nd})$.
- Define $Q_{\bar{V}}$ from $Q_{\bar{Y}}$, setting $Q_{\bar{V}|\bar{Y}} = \mathcal{N}(\bar{Y}, \sigma^2 I_{nd})$.

Lemma : The progressive estimator

$$\bar{Q}_{V_n|V_1,\dots,V_{n-1}} = \frac{1}{n} \sum_{i=1}^n Q_{V_n|V_1,\dots,V_{i-1}}$$

is such that

$$\frac{1}{d} \mathbb{P}_{V_1,\dots,V_{n-1}} [\mathcal{K}(\mathbb{P}_{V_n}, \bar{Q}_{V_n|V_1,\dots,V_{n-1}})] \leq \frac{1}{nd} \mathcal{K}(\mathbb{P}_\theta, Q_\theta) + \frac{\alpha}{2\sigma^2}.$$

The estimator \bar{Q} is successful at approaching \mathbb{P}_V when Q_θ is a successful data compression scheme for the source $\mathbb{P}_\theta = \mathbb{P}_{\theta(\bar{X})}$. Moreover \bar{Q} is a function of Q_θ , so that Q_θ contains in this case enough information to characterize \mathbb{P}_V that is itself an approximation of \mathbb{P}_X .

- $\{ \text{statistical models} \} \subsetneq \{ \text{compression models} \}$... and both can produce estimators, according to the previous lemma.
- Namely, if we use a **parametric model** $\{Q_{\theta|\beta}\}$, hoping that $\inf_{\beta} \mathcal{K}(\mathbb{P}_{\theta}, Q_{\theta|\beta})$ is small and **if we set** $Q_{\theta} = \int Q_{\theta|\beta} d\mu(\beta)$ for some prior probability measure μ on the parameter β , the lemma produces a **statistical estimator** \bar{Q} from a **statistical model** $\{Q_{\theta|\beta}\}$.
- But we can choose Q_{θ} otherwise, and in particular we can use a **grammar based code**. The general principle is to **define small binary indexes for repeated large patterns** and to gather those possibly nested indexation rules into a **grammar**. A widely used grammar based code is the **Lempel Ziv algorithm**. We will propose a more evolved choice of grammar producing **syntax trees** related to the conditional probability of observing some selected patterns in a given context.

- 1 Grammar based compression can **process** in a meaningful way **infrequent patterns** : even seen in the sample only twice, a large pattern is worth being indexed by a small index and described only once. On the other hand a statistical estimator, based on the estimation of expectations, cannot take into account in a reliable way an event that occurs only twice in the sample.
- 2 A parametric model $Q_{\theta|\beta}$ that accounts for elaborate conditional independence assumptions will presumably depend on a **high-dimensional parameter** β , and estimation in high dimension requires huge samples unless the complexity is restricted in some other way.

For those two reasons, we hope to **make sense of small samples** compared to a more traditional statistical approach.

First lossy compression step : fragmentation

- Let $\mathcal{C} \subset \mathbb{R}$ be finite (or more generally countable). For instance $\mathcal{C} = \llbracket 0, 255 \rrbracket$.
- We code the sample $\bar{X} = (X_1, \dots, X_n) \in \mathbb{R}^{nd}$ by

$$\theta = [(A_i \subset \llbracket 1, k \rrbracket, 1 \leq i \leq n), (B_j \subset \llbracket 1, d \rrbracket, C_j \in \mathcal{C}^{B_j}, 1 \leq j \leq k)],$$

where $A_i \in \mathcal{T}_B = \left\{ A \subset \llbracket 1, k \rrbracket : \bigsqcup_{j \in A} B_j = \llbracket 1, d \rrbracket \right\}$,

- the lossy decoding function being

$$f(\theta) = \bar{Y} = (C_{A_i}, 1 \leq i \leq n), \text{ where } C_A = \sum_{j \in A} C_j,$$

and where C_j is set to 0 outside its support B_j .

- The distortion $\mathcal{D}(\bar{X}, \theta) = (nd)^{-1} \|\bar{X} - \bar{Y}\|^2$ is minimized for a given set of fragments (B, C) when $A_i \in \arg \min_{A \in \mathcal{T}_B} \|X_i - C_A\|$ and is then equal to

$$\mathcal{D}(\bar{X}, B, C) = \inf_{A \in \mathcal{T}_B^n} \mathcal{D}(\bar{X}, \underbrace{(A, B, C)}_{\theta}) = d^{-1} \bar{\mathbb{P}}_X \left(\min_{A \in \mathcal{T}_B} \|X - C_A\|^2 \right),$$

where $\bar{\mathbb{P}}_X = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ is the empirical measure.

- We see that the distortion is a generalization of the k -means empirical criterion, that we get when $B_j = \llbracket 1, d \rrbracket$, $j \in \llbracket 1, k \rrbracket$.

- Create a sequence of codes $\theta_k = (A_k, B_k, C_k), k \geq n$.
- Start with one fragment per image: $k = n, A_{n,i} = \{i\}, i \in \llbracket 1, n \rrbracket, B_{n,j} = \llbracket 1, d \rrbracket, C_{n,j,s} = \arg \min_{c \in \mathcal{C}} |X_{j,s} - c|, j \in \llbracket 1, k \rrbracket$.
- Maintain through iterations on k the two properties
 - $\mathcal{D}(\bar{X}, \theta_k) \leq \alpha$
 - and $\llbracket 1, d \rrbracket = \bigsqcup_{j \in A_{k,i}} B_{k,j}, i \in \llbracket 1, n \rrbracket$,while minimizing the code length of θ_k at each step.
- Stop when it is no more possible to decrease the code length of θ_k or earlier for a partial fragmentation.

The fragmentation loop

Iterate on $k \geq n$ the following.

- Choose a pair $J_k \subset \llbracket 1, k \rrbracket$,
- set $A_{k+1,i} = A_{k,i} \cup \{k+1\}$ if $A_{k,i} \cap J_k \neq \emptyset$, $A_{k+1,i} = A_{k,i}$ otherwise,
- choose a support $B_{k+1,k+1} \subset \bigcap_{j \in J_k} B_{k,j}$,
- set $B_{k+1,j} = B_{k,j} \setminus B_{k+1,k+1}$ if $j \in J_k$, and $B_{k+1,j} = B_{k,j}$, if $j \in \llbracket 1, k \rrbracket \setminus J_k$.
- set $C_{k+1,j,s} \in \arg \min_{c \in \mathcal{C}} \{ |c - \bar{\mathbb{P}}(X_{l,s} | j \in A_{k+1,l})| \}$,

Choice of $B_{k+1,k+1}$, ensuring that $\mathcal{D}(\bar{X}, \theta_{k+1}) \leq \alpha$:

$$B_{k+1,k+1} = \left\{ s \in \bigcap_{j \in J_k} B_{k,j} : \mathbf{Var}(\bar{\mathbb{P}}_{X_{l,s} | k+1 \in A_{k+1,l}}) + \min_{c \in \mathcal{C}} [c - \bar{\mathbb{P}}(X_{l,s} | k+1 \in A_{k+1,l})]^2 \leq \alpha \right\},$$

where $\bar{\mathbb{P}}_l = \frac{1}{n} \sum_{i=1}^n \delta_i$. $s \in B_j$.

- Let $|A_k| = \sum_{i=1}^n |A_{k,i}|$ and $|B_k| = \sum_{j=1}^k |B_{k,j}|$.
- Assuming that everything is coded on L bits,
 $Q(\theta_k) = 2^{-L(|A_k| + 2|B_k| + (n+2k))}$ is a sub-probability measure.
- Thus

$$\begin{aligned} \xi(k) &= L^{-1} \log_2(Q(\theta_{k+1})/Q(\theta_k)) \\ &= 2|B_{k+1,k+1}| - \sum_{i=1}^n \mathbb{1}(k+1 \in A_{k+1,i}) - 2 \end{aligned}$$

can be maximized in J_k to find the pair giving the best data compression.

- A faster less optimal choice consists in choosing $J_k = \{j_{k,1}, j_{k,2}\}$, where
 $j_{k,1} \in \arg \max_j 2|B_{k,j}| - \sum_{i=1}^n \mathbb{1}(j \in A_{k,i})$ and $j_{k,2} \in \arg \max_{j_{k,2}} \xi(k)$.

Purpose

We get from the first fragmentation step a lossy code $\theta = (A, B, C)$ made of a fragment codebook $(B, C) = (B_j, C_j)_{j=1}^k$ described by their supports B_j and their contents C_j , and n sets of labels $A = (A_i)_{i=1}^n$, where $A_i \subset \llbracket 1, k \rrbracket$, describing the sample images. **The aim of the syntax analysis step is to perform lossless data compression on the sample description A .**

Starting point

The enumeration of A as a list of words with separators between the sets

$$w_{1,1} \dots w_{1,r_1} \wedge w_{2,1} \dots, w_{2,r_2} \wedge \dots \wedge w_{n,1} \dots w_{n,r_n} \wedge .$$

This representation follows the syntax $\{\{w\}\wedge\}$ where $w \in \llbracket 1, k \rrbracket$ is of type fragment label and where $\{\}$ denotes repetition, according to Extended Backus Naur specifications.

First step : grammar based compression

- Starting from $A_{0,i} = A_i$, $i \in \llbracket 1, n \rrbracket$, iterate the following for $m \geq 0$.
 - Choose a pair $J_m \subset \llbracket 1, k+m \rrbracket$ and define $A_{m+1,i} = (A_{m,i} \setminus J_m) \cup \{k+m+1\}$ when $J_m \subset A_{m,i}$ and $A_{m+1,i} = A_{m,i}$ otherwise.
 - Choose J_m to maximize $\sum_{i=1}^n \mathbb{1}(J_m \subset A_{m,i})$ as long as it is greater than 2.
- Appending the description of the pairs, we get a representation of type $\{\{w|p\} \wedge\} \{ab\}$ where $w \in \llbracket 1, k \rrbracket$, $p \in \llbracket k+1, k+m \rrbracket$ and $a, b \in \llbracket 1, k+m \rrbracket$.
- Its length decreases at each step.

We appended to the sample representation $R ::= \{\{w|p\} \wedge\}$ containing non terminal symbols of type p , the description of a grammar $G ::= \{ab\}$ represented by the right-hand side of context free rewriting rules $p \rightarrow ab$.

Second step : grammar based grammar compression

- We can reindex the grammar G to get $a_1 b_{1,1} a_1 b_{1,2} \dots a_1 b_{1,q_1} \dots a_m b_{m,1} \dots a_m b_{m,q_m}$ where the a_i are distinct and q_1, \dots, q_m are maximal.
- We can then factorize G into $a_1 \dots a_m \wedge b_{1,1} \dots b_{1,q_1} \wedge \dots \wedge b_{m,1} \dots b_{m,q_m} \wedge$ that is of the type $G ::= \{a\} \wedge \{\{b\} \wedge\}$.
- The second part of the description of G , describing the contexts $\{\{b\} \wedge\}$ of the pairs first elements $\{a\}$, is of the same type as the original representation and can be compressed again.
- Thus $\{\{b\} \wedge\}$ becomes $\{\{b|s\} \wedge\} \{cd\}$, where s are new non terminal syntax symbols and where $\{cd\}$ is a second grammar.
- The syntax labels s represent sets of configurations that appear in the same contexts $\{a\}$ and therefore perform a certain kind of context analysis.
- The whole representation follows the syntax $\{\{w|p\} \wedge\} \{a\} \wedge \{\{b|s\} \wedge\} \{cd\}$.

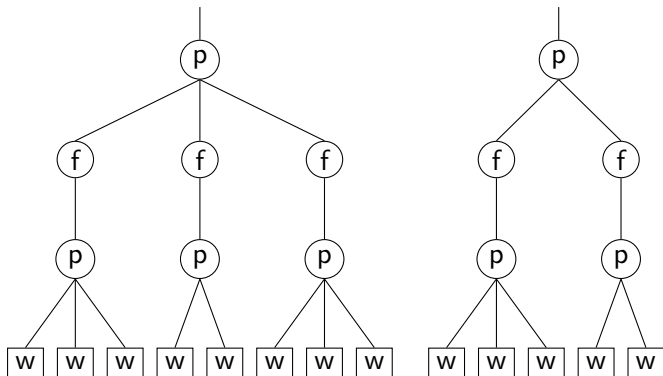
The syntax labels of step 2 induce a classification of the symbols of step 1. We define

- $f(w) = w$ when $w \in \llbracket 1, k \rrbracket$,
- $f(p) = s$ when $p \rightarrow ab$ and b is recovered from the compressed representation of ab by rewriting s
- and $f(p) = b$ when $p \rightarrow ab$ has not been compressed in the second step.

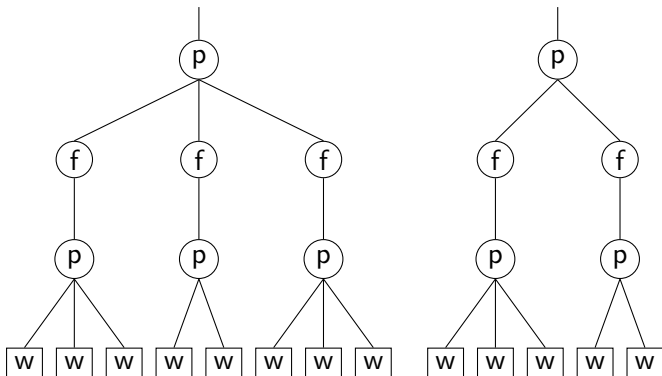
We obtain $f : \llbracket 1, k + m \rrbracket \rightarrow \llbracket 1, k + m + m' \rrbracket$ and we can recode $j \in \llbracket 1, k + m \rrbracket$ as $f(j)h(j)$ where $h(j)$ is the rank of $j \in f^{-1}[f(j)]$.

The sample representation R becomes of type $\{\{fh\}^\wedge\}$. It can be split into $\{\{f\}^\wedge\}\{\{h\}^\wedge\}$ where the values of h corresponding to each value of f have been gathered on the right. Each $\{f\}$ describes the syntax of a sample image at the first level.

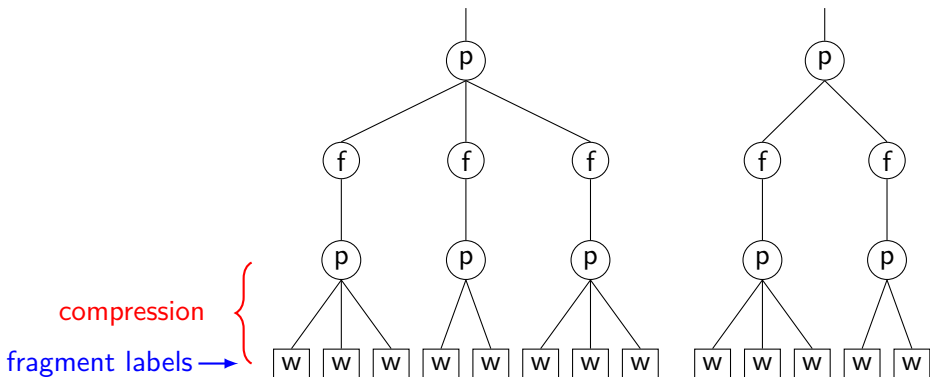
The whole compression process can be repeated on the first level sample syntax $\{\{f\}^\wedge\}$ to get several levels of syntax organized into a syntax tree.



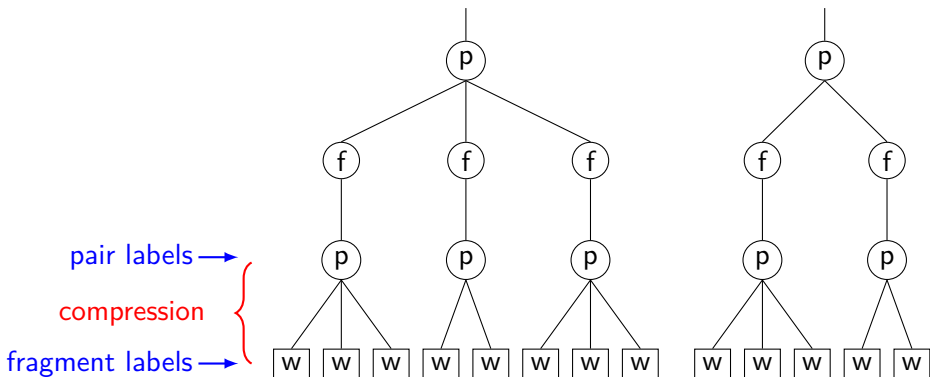
fragment labels →



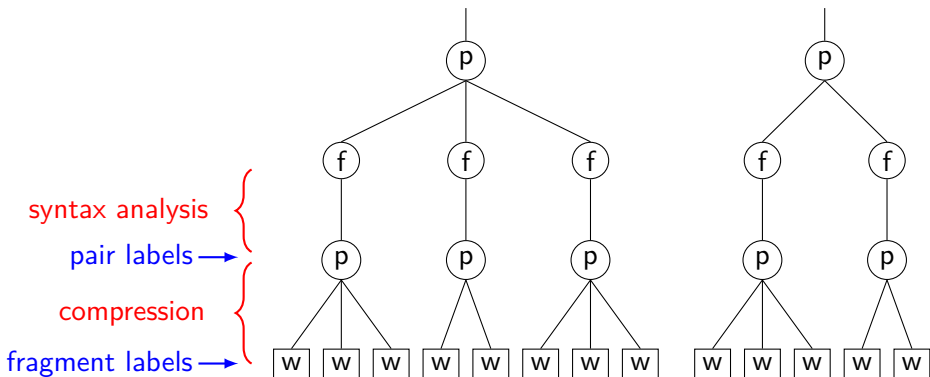
Syntax tree



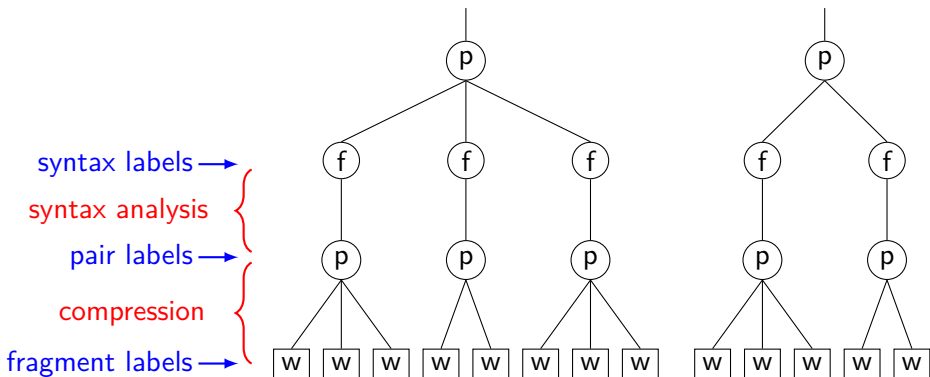
Syntax tree



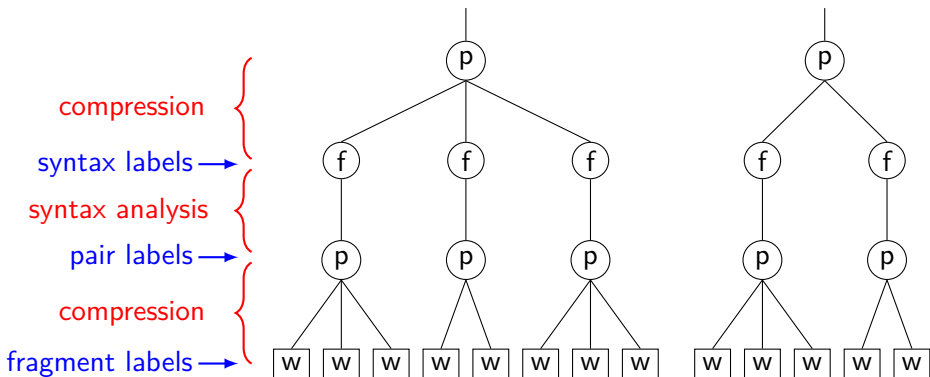
Syntax tree



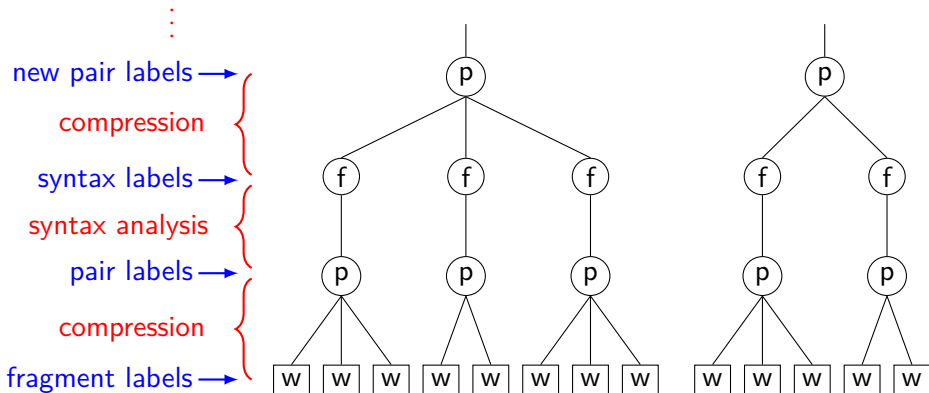
Syntax tree



Syntax tree



Syntax tree



Lemma for linear k -means in a separable Hilbert space

- Let (W_1, \dots, W_n) be n independent copies of $W \in \ell^2$.
- Assume that $\|W\|_\infty = \text{ess sup } \|W\| < +\infty$.
- Let $\Theta \subset (\ell^2)^k$ be bounded.
- Assume that $\mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \in [a, b] \text{ for all } \theta \in \Theta \right) = 1$.
- For any $k \geq 2$, any $n \geq 2k$ and any $\delta \in]0, 1[$, with probability at least $1 - \delta$, for any $\theta \in \Theta$,

$$\begin{aligned} \mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \right) &\leq \bar{\mathbb{P}}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \right) \\ &\quad + \left(\frac{\log(n/k)}{\log(2)} \sqrt{\frac{8 \log(k)}{n}} + 2 \sqrt{\frac{\log(k)}{n}} \right) \|\Theta\| \|W\|_\infty \\ &\quad + \sqrt{\frac{(\sqrt{2} + 1)(k(b-a)^2 + 2 \log(ek) \|W\|_\infty^2 \|\Theta\|^2)}{n}} + \sqrt{\frac{\log(\delta^{-1})}{2n}} (b-a). \end{aligned}$$

Lemma on the excess risk

- If $\theta^* \in \Theta$ is non random, with probability at least $1 - \delta$,

$$\begin{aligned} & \sup_{\theta \in \Theta} (\mathbb{P}_W - \bar{\mathbb{P}}_W) \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle - \min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j^*, W \rangle \right) \\ & \leq \left(\frac{\log(n/k)}{\log(2)} \sqrt{\frac{8 \log(k)}{n}} + 2 \sqrt{\frac{\log(k)}{n}} \right) \|\Theta\| \|W\|_\infty + \\ & \sqrt{\frac{(\sqrt{2} + 1)(k(b-a)^2 + 2 \log(ek) \|W\|_\infty^2 \|\Theta\|^2)}{n}} + \sqrt{\frac{2 \log(\delta^{-1})}{n}} (b-a). \end{aligned}$$

- If $\hat{\theta} \in \arg \min_{\theta \in \Theta} \bar{\mathbb{P}} \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \right)$, we have the same bound for

$$\mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \hat{\theta}_j, W \rangle \right) - \inf_{\theta \in \Theta} \mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \right).$$

Lemma

In expectation with respect to the sample distribution

$$\begin{aligned} & \mathbb{P}_{W_1, \dots, W_n} \left[\mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \hat{\theta}_j, W \rangle \right) - \inf_{\theta \in \Theta} \mathbb{P}_W \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle \right) \right] \\ & \leq \left(\frac{\log(n/k)}{\log(2)} \sqrt{\frac{8 \log(k)}{n}} + 2 \sqrt{\frac{\log(k)}{n}} \right) \|\Theta\| \|W\|_\infty \\ & \quad + \sqrt{\frac{(\sqrt{2} + 1)(k(b-a)^2 + 2 \log(ek)) \|W\|_\infty^2 \|\Theta\|^2}{n}}. \end{aligned}$$

Proposition

- Let (X_1, \dots, X_n) be n independent copies of $X \in H$, some separable Hilbert space.
- Assume that $\mathbb{P}(\|X\| \leq B) = 1$, $n \geq 2k$ and $k \geq 2$.
- Consider an estimator $\hat{C} \in \arg \min_{C \in H^k} \bar{\mathbb{P}}_X \left(\min_{j \in \llbracket 1, k \rrbracket} \|X - C_j\|^2 \right)$.

$$\begin{aligned} & \mathbb{P}_{X_1, \dots, X_n} \left[\mathbb{P}_X \left(\min_{j \in \llbracket 1, k \rrbracket} \|X - \hat{C}_j\|^2 \right) \right] \\ & \leq \underbrace{\inf_{C \in H^k} \mathbb{P}_X \left(\min_{j \in \llbracket 1, k \rrbracket} \|X - C_j\|^2 \right)}_{\leq 1} + 16 B^2 \log \left(\frac{n}{k} \right) \sqrt{\frac{k \log(k)}{n}}. \end{aligned}$$

Improves on [Biau Devroye and Lugosi, 2008] $\mathcal{O}(k/\sqrt{n})$.

Proof: Remark that $\|C_j\|^2 - 2\langle C_j, X \rangle = \langle \theta_j, W \rangle$, where $\theta_j = (C_j, \gamma^{-1} \|C_j\|^2 B^{-1}) \in H \times \mathbb{R}$ and $W = (-2X, \gamma B)$ and apply the previous lemma. Take $\gamma_{\text{opt}} = \sqrt{2}$.

Lemma. PAC-Bayesian inequality

- Consider $h : \mathcal{T} \times \mathcal{W} \rightarrow \mathbb{R}$.
- Consider a prior $\pi \in \mathcal{M}_+^1(\mathcal{T})$ on the parameter space \mathcal{T} .
- Let $\lambda > 0$ be a positive exponent.
- Let (W_1, \dots, W_n) be n independent copies of $W \in \mathcal{W}$.

$$\mathbb{P}_{W_1, \dots, W_n} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\mathcal{T})} \sup_{\eta \in \mathbb{N}} \left\{ \int \min \left\{ \eta, -\lambda \sum_{i=1}^n h(\theta, W_i) \right. \right. \right. \right. \right. \\ \left. \left. \left. \left. - n \log \left[\mathbb{P}_W \exp \left[-\lambda h(\theta, W) \right] \right] \right\} d\rho(\theta) - \mathcal{K}(\rho, \pi) \right\} \right] \right\} \leq 1.$$

- Let $\rho_{\theta'|\theta} = \bigotimes_{j=1}^k \left(\bigotimes_{i \in \mathbb{N}} \mathcal{N}(\theta_{j,i}, \beta^{-1}) \right) : (\mathbb{R}^{\mathbb{N}})^k \rightarrow \mathcal{M}_+^1((\mathbb{R}^{\mathbb{N}})^k)$.

- Let

$$\langle \theta, w \rangle = \begin{cases} \lim_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i, & \text{when } \overline{\lim}_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i = \underline{\lim}_{s \rightarrow +\infty} \sum_{i=0}^s \theta_i w_i \in \mathbb{R}, \\ 0, & \text{otherwise} \end{cases}$$

be a non bilinear but measurable extension of the scalar product from ℓ^2 to $\mathbb{R}^{\mathbb{N}}$.

- Introduce $f(\theta, w) = \min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, w \rangle$, $\theta \in (\mathbb{R}^{\mathbb{N}})^k, w \in \mathbb{R}^{\mathbb{N}}$
- and the centered loss function $\bar{f}(\theta, w) = f(\theta, w) - \mathbb{P}_W(f(\theta, W))$.

- Write

$$\begin{aligned}
 (\mathbb{P}_W - \bar{\mathbb{P}}_W)f(\theta, W) &= (\mathbb{P}_W - \bar{\mathbb{P}}_W) \underbrace{(\delta_{\theta'|\theta} - \rho_{\theta'|\theta})}_{\text{small perturbation}} f(\theta', W) \\
 &+ \sum_{p=1}^H (\mathbb{P}_W - \bar{\mathbb{P}}_W) \underbrace{(\rho_{\theta'|\theta}^{2^{p-1}} - \rho_{\theta'|\theta}^{2^p})}_{\text{chain of intermediate scales}} f(\theta', W) \\
 &+ (\mathbb{P}_W - \bar{\mathbb{P}}_W) \underbrace{\rho_{\theta'|\theta}^{2^H}}_{\text{big perturbation}} f(\theta', W).
 \end{aligned}$$

- Remark that

$$\begin{aligned}
 (\delta_{\theta'|\theta} - \rho_{\theta'|\theta})f(\theta', W) &= \rho_{\theta'|\theta} \left(\min_{j \in \llbracket 1, k \rrbracket} \langle \theta_j, W \rangle - \min_{j \in \llbracket 1, k \rrbracket} \langle \theta'_j, W \rangle \right) \\
 &\leq \rho_{\theta'|\theta} \underbrace{\left(\max_{j \in \llbracket 1, k \rrbracket} \langle \theta_j - \theta'_j, W \rangle \right)}_{\text{Gaussian}/\rho} \leq \sqrt{2 \log(k)/\beta} \|W\|_\infty.
 \end{aligned}$$

- From the PAC-Bayesian inequality applied to $h(\theta, w) = (\delta_{\theta'|\theta} - \rho_{\theta'|\theta}) \bar{f}(\theta', w)$,

$$\mathbb{P}_{W_1, \dots, W_n} \left\{ \exp \sup_{\theta \in (\ell^2)^k} \left[n\lambda (\mathbb{P}_W - \bar{\mathbb{P}}_W) (\rho_{\theta'|\theta} - \rho_{\theta''|\theta}^2) f(\theta', W) - n\rho_{\theta'|\theta} \left[\log \left(\mathbb{P}_W \left[\exp \left(-\lambda (\delta_{\theta''|\theta'} - \rho_{\theta''|\theta'}) \bar{f}(\theta'', W) \right) \right] \right) \right] - \frac{\beta \|\theta\|^2}{2} \right] \right\} \leq 1.$$

- This gives

$$\mathbb{P}_{W_1, \dots, W_n} \left[\sup_{\theta \in \Theta} (\mathbb{P}_W - \bar{\mathbb{P}}_W) (\rho_{\theta'|\theta} - \rho_{\theta''|\theta}^2) f(\theta', W) \right] \leq \frac{4\lambda}{\beta} \log(k) \|W\|_\infty^2 + \frac{\beta \|\Theta\|^2}{2n\lambda} \underset{\lambda_{\text{opt}}}{=} \|W\|_\infty \|\Theta\| \sqrt{\frac{8 \log(k)}{n}}.$$

- Consider $\psi(x) = \begin{cases} \log(1+x+x^2/2), & x \geq 0, \\ -\log(1-x+x^2/2), & x \leq 0, \end{cases}$
- and $\tilde{f}(\theta, W) = f(\theta, W) - \frac{a+b}{2}$.
- Remark that

$$\begin{aligned} (\mathbb{P}_W - \bar{\mathbb{P}}_W) \rho_{\theta'|\theta} f(\theta', W) &= \\ \rho_{\theta'|\theta} \left[\mathbb{P}_W \tilde{f}(\theta', W) - \bar{\mathbb{P}}_W \left(\lambda^{-1} \psi \left[\lambda \tilde{f}(\theta', W) \right] \right) \right] \\ + \underbrace{\rho_{\theta'|\theta} \bar{\mathbb{P}}_W \left[\lambda^{-1} \psi \left[\lambda \tilde{f}(\theta', W) \right] - \tilde{f}(\theta', W) \right]}_{\leq \frac{\lambda}{2(1+\sqrt{2})} \left[(b-a)^2/4 + 2\log(ek) \|W\|_\infty^2 / \beta \right]} \\ &\leq \frac{\lambda}{2(1+\sqrt{2})} \left[(b-a)^2/4 + 2\log(ek) \|W\|_\infty^2 / \beta \right] \\ &\quad \text{since } |x - \psi(x)| \leq \frac{x^2}{4(1+\sqrt{2})}, \quad x \in \mathbb{R}. \end{aligned}$$

- Take $h(\theta, w) = \lambda^{-1} \psi[\lambda \tilde{f}(\theta, w)]$ to obtain

$$\mathbb{P}_{W_1, \dots, W_n} \left\{ \sup_{\theta \in \Theta} \exp \left[-n \lambda \rho_{\theta' | \theta} \bar{\mathbb{P}}_W \left(\lambda^{-1} \psi[\lambda \tilde{f}(\theta', W)] \right) - n \rho_{\theta' | \theta} \left[\log \left(\mathbb{P}_W \left[\exp \left(-\psi[\lambda \tilde{f}(\theta', W)] \right) \right] \right) \right] - \frac{\beta \|\theta\|^2}{2} \right] \right\} \leq 1.$$

- Use $\psi(x) \leq \log(1 + x + x^2/2)$, $x \in \mathbb{R}$ to deduce

$$\begin{aligned} \mathbb{P}_{W_1, \dots, W_n} \left\{ \sup_{\theta \in \Theta} \rho_{\theta' | \theta} \left[\mathbb{P}_W \left(\tilde{f}(\theta', W) \right) - \bar{\mathbb{P}}_W \left(\lambda^{-1} \psi[\lambda \tilde{f}(\theta', W)] \right) \right] \right\} \\ \leq \lambda \left[(b - a)^2 / 4 + 2 \log(ek) \|W\|_\infty^2 / \beta \right] + \frac{\beta \|\Theta\|^2}{2n\lambda}. \end{aligned}$$

- For the biggest perturbation we get

$$\mathbb{P}_{W_1, \dots, W_n} \left\{ \sup_{\theta \in \Theta} \left(\mathbb{P}_W - \bar{\mathbb{P}}_W \right) \rho_{\theta' | \theta}^{2H} f(\theta', W) \right\} \\ \leq_{\lambda_{\text{opt}}} \sqrt{\frac{(\sqrt{2} + 1) \left(2^{-H} \beta (b - a)^2 + 8 \log(ek) \|W\|_{\infty}^2 \right) \|\Theta\|^2}{4n}}$$

- Putting all together

$$\mathbb{P}_{W_1, \dots, W_n} \left\{ \sup_{\theta \in \Theta} \left(\mathbb{P}_W - \bar{\mathbb{P}}_W \right) f(\theta, W) \right\} \leq 2 \sqrt{2 \log(k) / \beta} \|W\|_{\infty} \\ + \sqrt{\frac{(\sqrt{2} + 1) \left(2^{-H} \beta (b - a)^2 + 8 \log(ek) \|W\|_{\infty}^2 \right) \|\Theta\|^2}{4n}} \\ + H \|W\|_{\infty} \|\Theta\| \sqrt{\frac{8 \log(k)}{n}}$$

- Choose $\beta = 2n \|\Theta\|^{-2}$ and $H = \lfloor \log(n/k) / \log(2) \rfloor$.

According to the bounded difference inequality

With probability at least $1 - \delta$

$$\begin{aligned} & \sup_{\theta \in \Theta} (\mathbb{P}_W - \bar{\mathbb{P}}_W) f(\theta, W) \\ & \leq \mathbb{P}_{W_1, \dots, W_n} \left\{ \sup_{\theta \in \Theta} (\mathbb{P}_W - \bar{\mathbb{P}}_W) f(\theta, W) \right\} + \sqrt{\frac{2 \log(\delta^{-1})}{n}} (b - a). \end{aligned}$$

Generalization bounds for fragmentation

- Consider the fragmentation model

$$\mathcal{M}(\mathcal{S}) = \left\{ (B, C)_{j=1}^k : B_j \subset \llbracket 1, k \rrbracket, C_j \in [-a, a]^{B_j}, \sum_{j=1}^k \mathbb{P}_{\mathcal{S}}(B_j) \leq \mathcal{S}, |\mathcal{T}_{B,K}| \geq 2 \right\},$$

where $\mathcal{T}_{B,K} = \{A \subset \mathcal{T}_B : |A| \leq K\}$.

- Remark that $\log(|\mathcal{T}_{B,K}|) \leq K \log\left(\frac{ek}{K}\right)$.
- Consider the risk function

$$\begin{aligned} \mathcal{R}(B, C) &= \underbrace{\mathcal{D}(B, C)}_{\stackrel{\text{def}}{=} d^{-1} \mathbb{P}_X(\min_{A \in \mathcal{T}_{B,K}} \|X - C_A\|^2)} - d^{-1} \mathbb{P}_X(\|X\|^2) \\ &= d^{-1} \mathbb{P}_X \left[\min_{A \in \mathcal{T}_{B,K}} (\|C_A\|^2 - 2\langle X, C_A \rangle) \right]. \end{aligned}$$

- Assume that $\mathbb{P}\left(\max_{s \in \llbracket 1, d \rrbracket} |X_s| \leq a\right) = 1$.

Proposition

Assume that $k \geq 2$, $K \geq 1$, $S \in [1, k]$, $n \geq 2SK$ and $\delta \in]0, 1[$. With probability at least $1 - \delta$, for any $(B, C) \in \mathcal{M}(S)$,

$$\begin{aligned} \mathcal{R}(B, C) - \bar{\mathcal{R}}(B, C) &\leq a^2 \left(\frac{\sqrt{10} \log(nS^{-1}K^{-1})}{\log(2)} \sqrt{\frac{8SK \log(|\mathcal{T}_{B,K}|)}{n}} \right. \\ &\quad + 2\sqrt{10} \sqrt{\frac{SK \log(|\mathcal{T}_{B,K}|)}{n}} + \sqrt{\frac{4(\sqrt{2} + 1)(9 + 5 \log(|\mathcal{T}_{B,K}|))SK}{n}} \\ &\quad \left. + 2\sqrt{\frac{k(k-1) \log(2) + 2 \log(\delta^{-1})}{n}} \right) \\ &\leq a^2 \mathcal{O} \left(\log \left(\frac{n}{SK} \right) \sqrt{\frac{SK^2 \log(k/K)}{n}} + \sqrt{\frac{k^2 + \log(\delta^{-1})}{n}} \right). \end{aligned}$$

This is based on a lemma for linear fragmentation similar to the linear k -means lemma.

- For a given empirical distortion level $\mathcal{D}(\bar{X}, B, C)$, we should minimize $\mathcal{S} = d^{-1}|B|$, to get the best possible generalization bound, similarly to what was suggested by the compression approach.
- We can take advantage of this generalization bound to compute the fragmentation code book (B, C) on a subsample of the data base to be analyzed.

Proposition

$$\text{Let } (\widehat{B}, \widehat{C}) \in \arg \min_{(B,C) \in \mathcal{M}(S)} \overline{\mathcal{D}}(B, C) + \left(\frac{\log(nS^{-1}K^{-1})}{\log(2)} \sqrt{\frac{8SK \log(|\mathcal{T}_{B,K}|)}{n}} \right. \\ \left. + 2\sqrt{\frac{SK \log(|\mathcal{T}_{B,K}|)}{n}} \right) \sqrt{10}a^2 + \sqrt{\frac{4(\sqrt{2}+1)(4+5\log(e|\mathcal{T}_{B,K}|))SK}{n}} a^2.$$

With probability at least $1 - \delta$,

$$\mathcal{D}(\widehat{B}, \widehat{C}) \leq \inf_{(B,C) \in \mathcal{M}(S)} \mathcal{D}(B, C) + \left(\frac{\log(nS^{-1}K^{-1})}{\log(2)} \sqrt{\frac{8SK \log(|\mathcal{T}_{B,K}|)}{n}} \right. \\ \left. + 2\sqrt{\frac{SK \log(|\mathcal{T}_{B,K}|)}{n}} \right) \sqrt{10}a^2 + \sqrt{\frac{4(\sqrt{2}+1)(4+5\log(e|\mathcal{T}_{B,K}|))SK}{n}} \\ + 4\sqrt{\frac{k(k-1)\log(2) + 2\log(\delta^{-1})}{n}} a^2.$$

Application to images: a small example

Some images to analyze



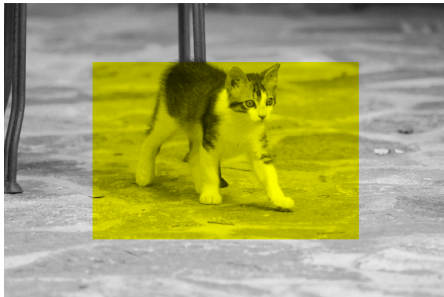
Application to images: a small example

Some images to analyze



Application to images: a small example

Some images to analyze



Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Application to images: extraction of the training sample



300 extracted patches of size 300×300 from the three previous yellow frames.

Fragmentation parameters

- Threshold for the distortion: $\alpha = 0.1 \times \frac{1}{d} \sum_{s=1}^d \mathbf{Var}(\bar{\mathbb{P}}_{X_s})$
- Number of fragments: 2000
- Fragmentation criterion: looking for indexes that maximize

$$2 \underbrace{|B_j|}_{\text{size of fragment } j} - \underbrace{\sum_{i=1}^n \mathbb{1}(j \in A_i)}_{\text{number of images sharing fragment } j} .$$

Fragmentation: image approximation



Figure: Image approximation resulting from the fragmentation vs original image.

Fragmentation: image approximation



Figure: Image approximation resulting from the fragmentation vs original image.

Fragmentation: image approximation



Figure: Image approximation resulting from the fragmentation vs original image.

Fragmentation: visualization of the fragments

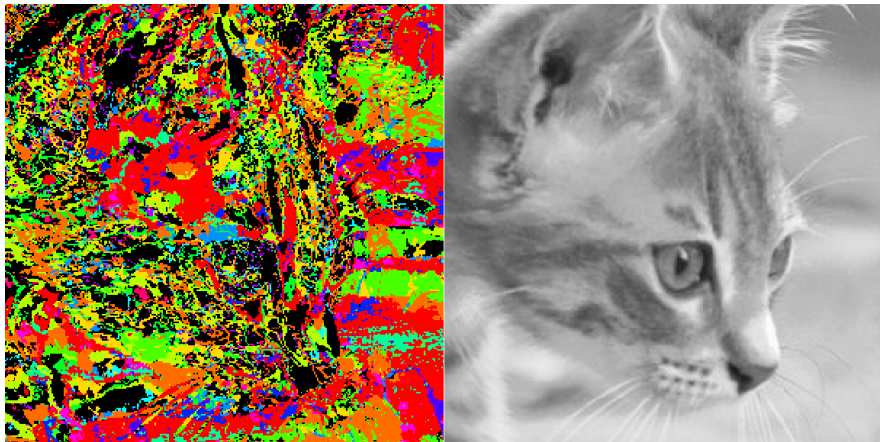


Figure: Visualization of the fragments.

Fragmentation: visualization of the fragments

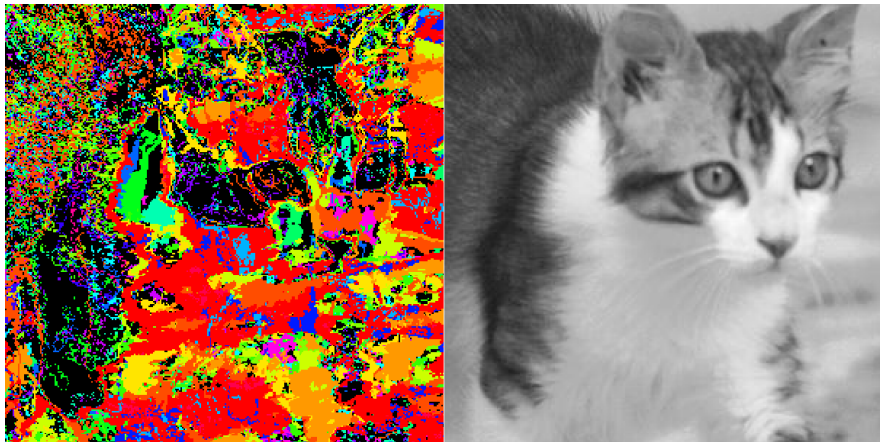


Figure: Visualization of the fragments.

Fragmentation: visualization of the fragments

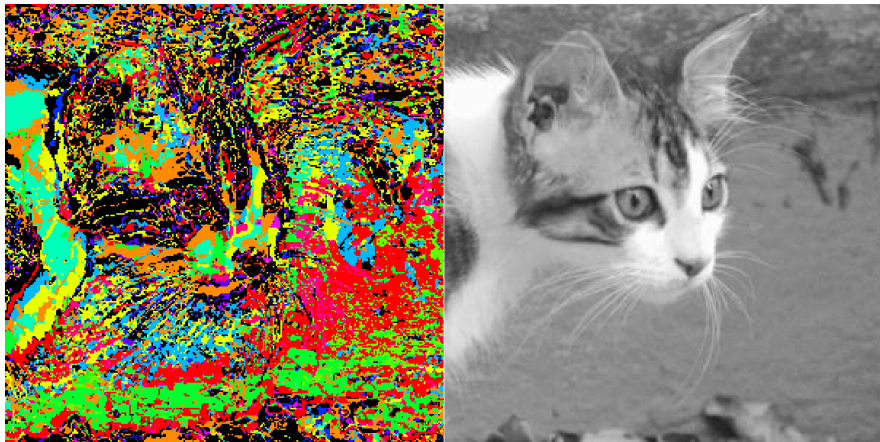
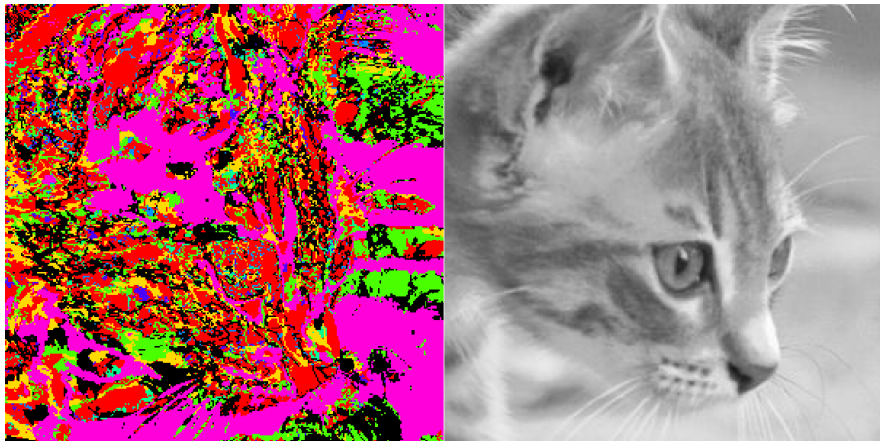


Figure: Visualization of the fragments.

Merging and Syntax parameters

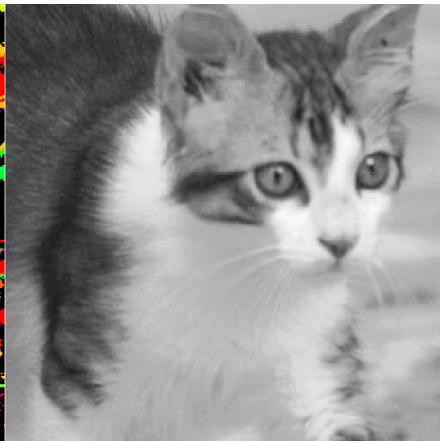
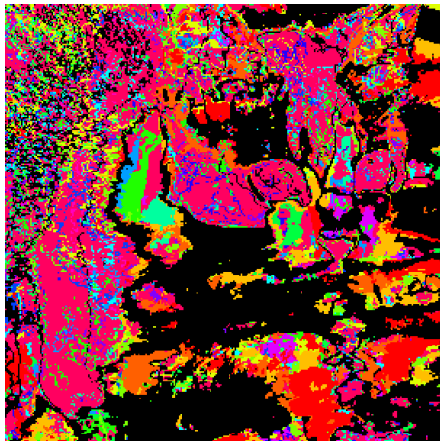
- Maximum number of merged/compressed fragments: 300
- Maximum number of syntax labels: 300
- Iterate the merging step and syntax step until no merging is performed: 14 syntax steps.

Syntax labels at highest level.



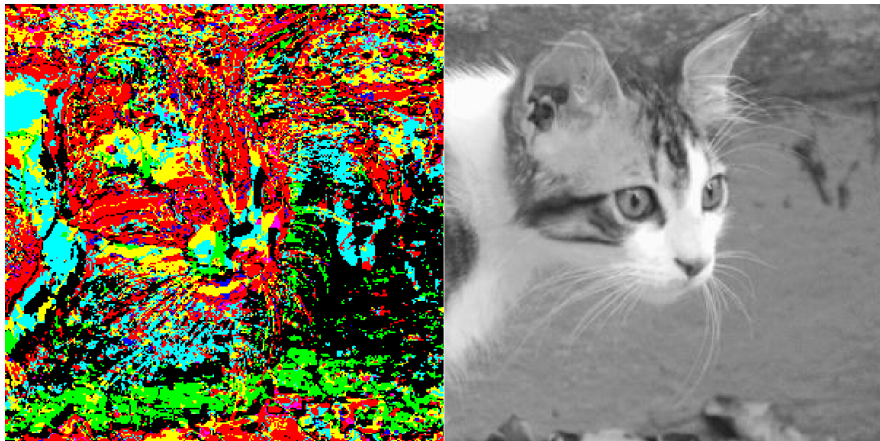
Compared to the result of fragmentation, some simplifications can be observed.

Syntax labels at highest level.



Compared to the result of fragmentation, some simplifications can be observed.

Syntax labels at highest level.



Compared to the result of fragmentation, some simplifications can be observed.

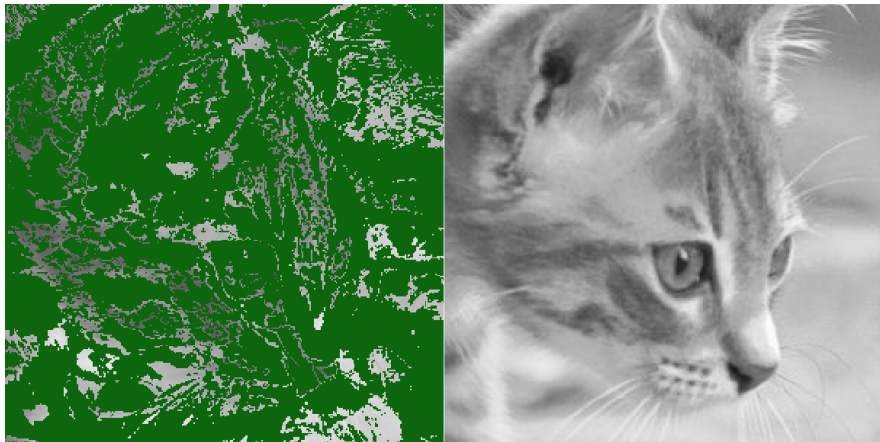


Figure: Visualization of a syntax label detecting some kind of translation.

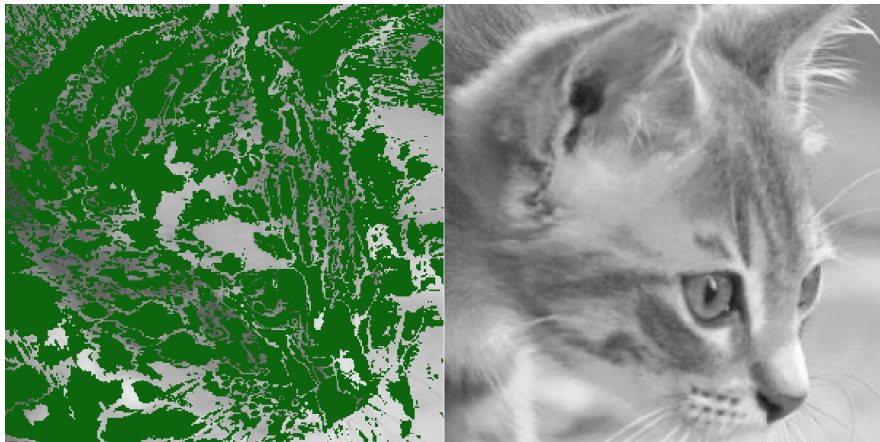


Figure: Visualization of a syntax label detecting some kind of translation.

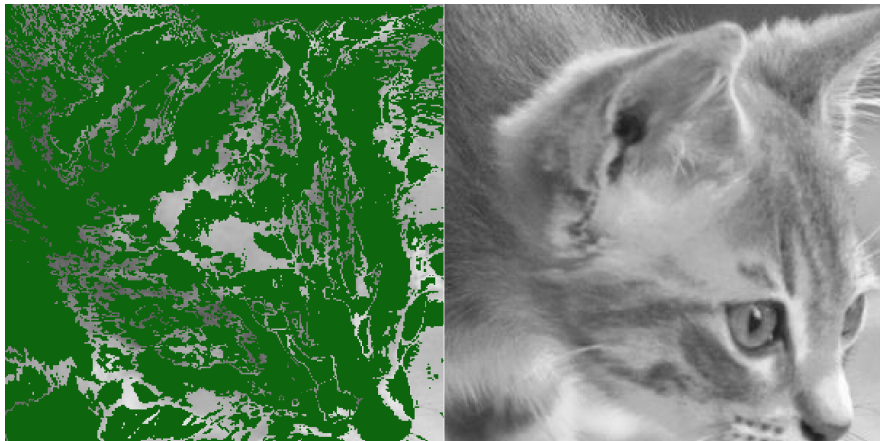


Figure: Visualization of a syntax label detecting some kind of translation.

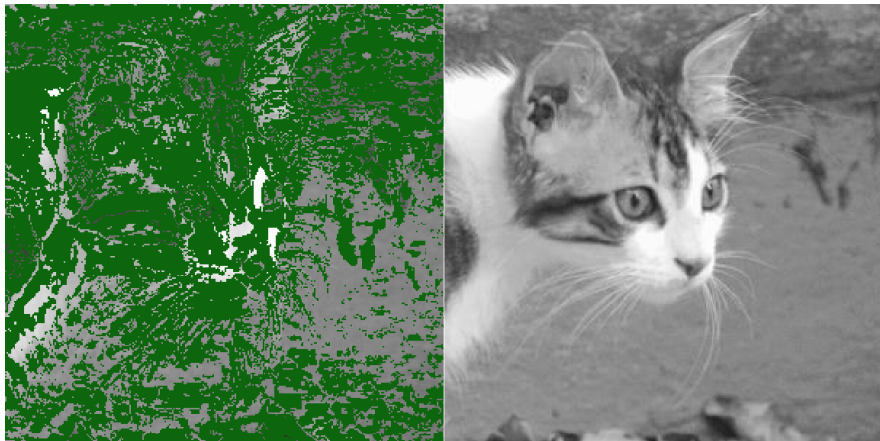


Figure: Visualization of a syntax label detecting some kind of translation.

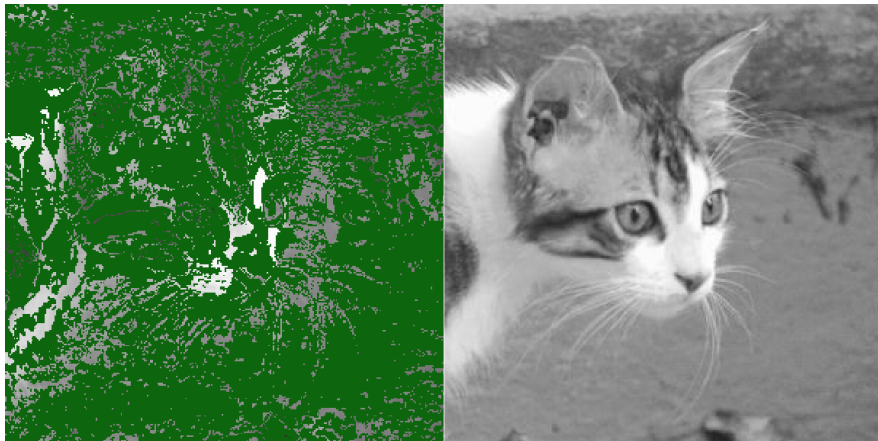


Figure: Visualization of a syntax label detecting some kind of translation.

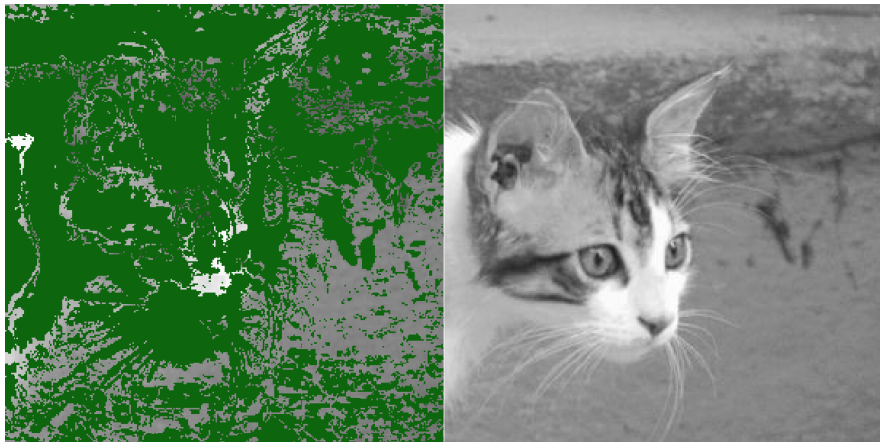


Figure: Visualization of a syntax label detecting some kind of translation.

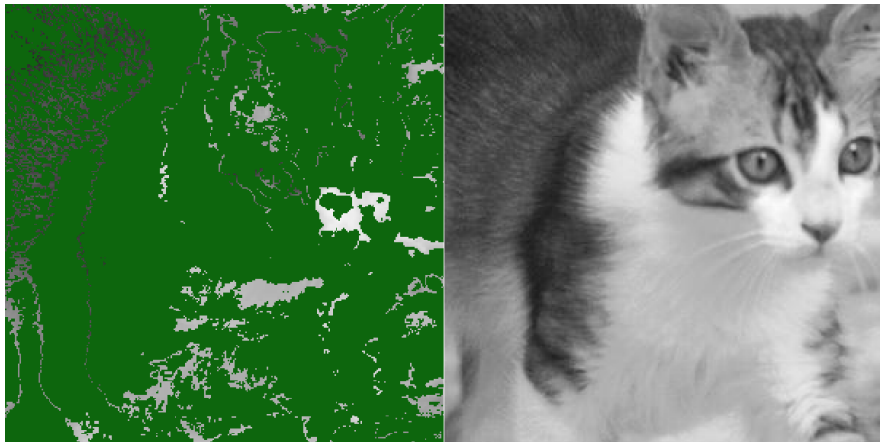


Figure: Visualization of a syntax label detecting some kind of translation.

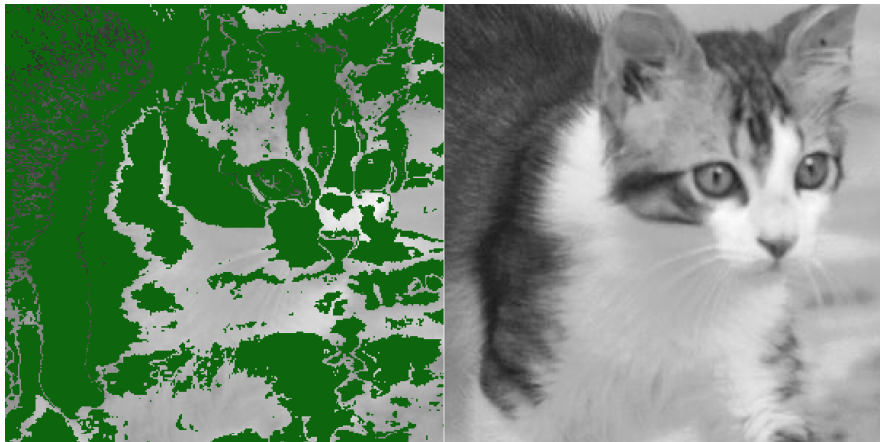


Figure: Visualization of a syntax label detecting some kind of translation.

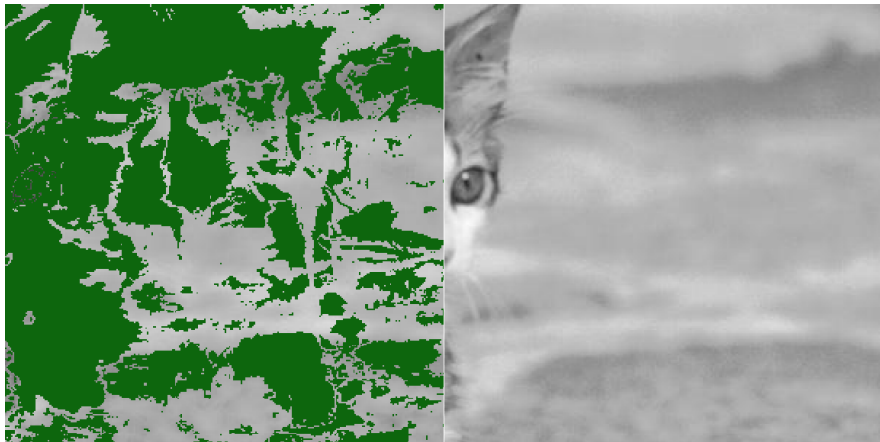


Figure: Visualization of a syntax label detecting some kind of translation.

First identification of the Kitten's head

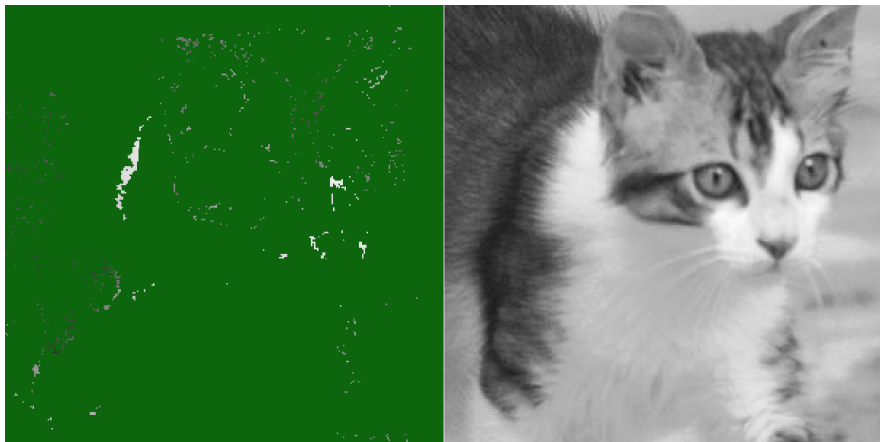


Figure: Visualization of a fixed syntax label recognizing related pattern.

First identification of the Kitten's head

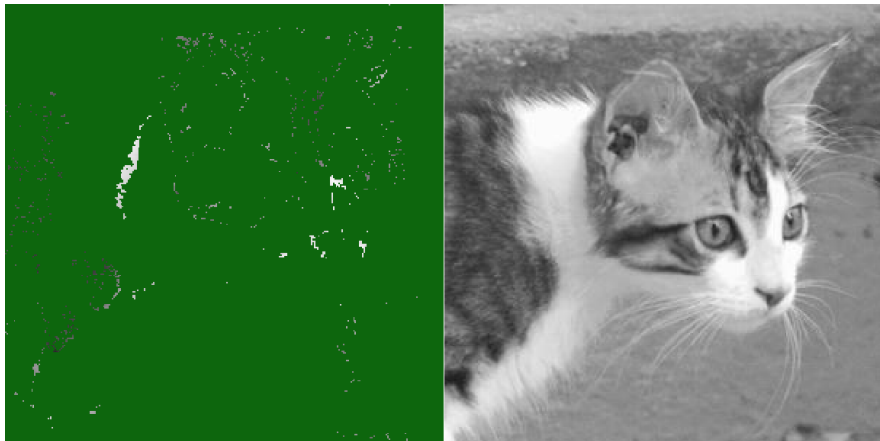


Figure: Visualization of a fixed syntax label recognizing related pattern.

Second identification of the Kitten's head

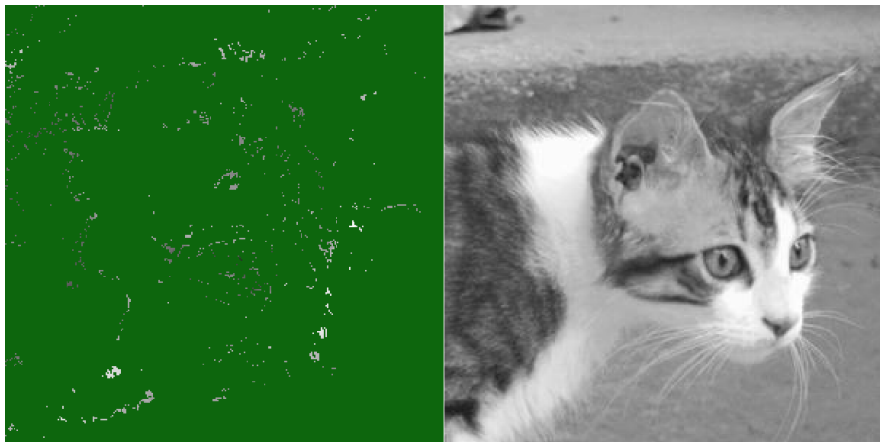


Figure: Visualization of a fixed syntax label recognizing related patterns.

Second identification of the Kitten's head

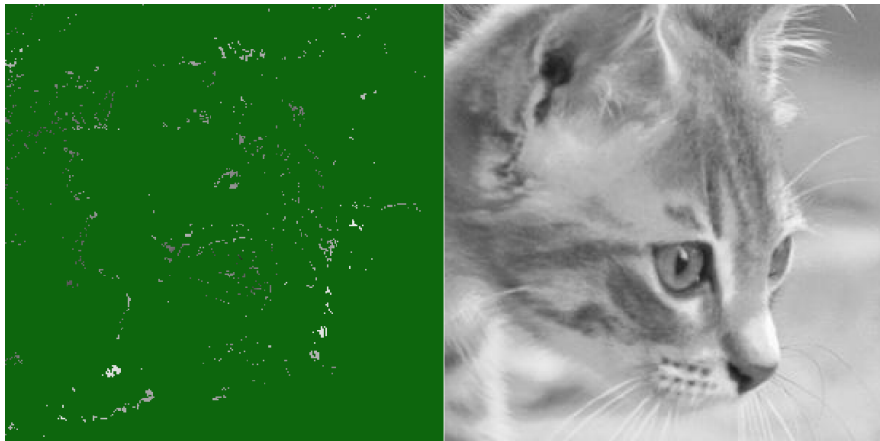


Figure: Visualization of a fixed syntax label recognizing related patterns.

Weak distortion and information fragmentation

- To get a weaker notion of distortion, introduce a probabilistic description of X as $\mathbb{P}_{S, V|X}$.
- Define the triplet (X, S, V) , where $\mathbb{P}_{S|X} = \mathbb{P}_S = \frac{1}{d} \sum_{s=1}^d \delta_s$ is uniform on pixel locations and where $\mathbb{P}_{V|X, S} = \mathcal{N}(X_S, \sigma^2)$ is the value X_S at pixel S blurred by a Gaussian noise.
- Describe the sample by $\bar{\mathbb{P}}_{X, S, V} = \bar{\mathbb{P}}_X \mathbb{P}_{S, V|X}$.
- Introduce the empirical weak distortion

$$\begin{aligned}\tilde{\mathcal{D}}(\bar{X}, B, C) &= 2\sigma^2 \left[1 - \exp\left(-\inf_{Q \in \mathcal{Q}} \mathcal{K}(Q_{X, S, V}, \bar{\mathbb{P}}_{X, S, V})\right) \right] \\ &= 2\sigma^2 \bar{\mathbb{P}}_X \left[\min_{A \in \mathcal{J}_{B, K}} \mathbb{P}_S \sum_{j \in A} \mathbb{1}(S \in B_j) [1 - \exp(-(X_S - C_{j, S})^2 / (2\sigma^2))] \right] \\ &\leq \mathcal{D}(\bar{X}, B, C), \text{ where}\end{aligned}$$

$$\mathcal{Q} = \left\{ Q_{X, S, V} : \text{for some } A : \mathbb{R}^d \rightarrow \mathcal{J}_{B, K}, Q_{V|X, S} = \mathcal{N}(C_{A(X), S}, \sigma^2) \right\},$$

and its expected value $\tilde{\mathcal{D}}(B, C)$.

Lemma

Minimizing the weak distortion in C , the fragment contents, gives the optimal approximation of $\bar{\mathbb{P}}_{X,S,V}$ under a conditional independence assumption only. Indeed

$$\inf_{C_j \in \mathbb{R}^{B_j}} \tilde{\mathcal{D}}(\bar{X}, B, C) = 2\sigma^2 \inf_{Q \in \mathcal{Q}} \left[1 - \exp\left(-\mathcal{K}(Q_{X,S,V}, \bar{\mathbb{P}}_{X,S,V})\right) \right]$$

where

$$\mathcal{Q} = \left\{ Q_{X,S,V} : \text{for some } A : \mathbb{R}^d \rightarrow \mathcal{T}_{B,K}, Q_{V|X,S} = Q_{V|S, \ell_{A,B}(X,S)} \right\}$$

and

$$\ell_{A,B}(X, S) = j \iff S \in B_j \text{ and } j \in A(X)$$

is the classification function defined by A and B .

Proposition

$$\text{Let } (\widehat{B}, \widehat{C}) \in \arg \min_{(B,C) \in \mathcal{M}(S)} \widetilde{\mathcal{D}}(\overline{X}, B, C) + 2\sigma^2 \left(\frac{\log(nS^{-1}K^{-1})}{\log(2)} \sqrt{\frac{8SK \log(|\mathcal{T}_{B,K}|)}{n}} \right. \\ \left. + 2\sqrt{\frac{SK \log(|\mathcal{T}_{B,K}|)}{n}} \right) + 2\sigma^2 \sqrt{\frac{(\sqrt{2}+1)(1+2\log(e|\mathcal{T}_{B,K}|))SK}{n}},$$

with probability at least $1 - \delta$,

$$\widetilde{\mathcal{D}}(\widehat{B}, \widehat{C}) \leq \inf_{(B,C) \in \mathcal{M}(S)} \widetilde{\mathcal{D}}(B, C) + 2\sigma^2 \left(\frac{\log(nS^{-1}K^{-1})}{\log(2)} \sqrt{\frac{8SK \log(|\mathcal{T}_{B,K}|)}{n}} \right. \\ \left. + 2\sqrt{\frac{SK \log(|\mathcal{T}_{B,K}|)}{n}} \right) + 2\sigma^2 \sqrt{\frac{(\sqrt{2}+1)(1+2\log(e|\mathcal{T}_{B,K}|))SK}{n}} \\ + 2\sigma^2 \sqrt{\frac{k(k-1)\log(2) + 2\log(\delta^{-1})}{n}}.$$