# Hypothesis Testing Details

We used the hypothesis testing to verify the significance of the improvements provided by EGA optimized quantization thresholds. Hereafter, we explain the details of this test.

**Part I. Hypothesis Testing**

Our major concern was to test, against zero, the difference $P(n,\Theta_{opt}) - P(n,\Theta_{man})$ between the average precisions of the wavelet correlogram indexing algorithm with the EGA-optimized ($\Theta_{opt}$) and original ($\Theta_{man}$) quantization thresholds. The improvements of other evaluation measures including average weighted precision, average recall, and average rank were verified in the same manner.

To decide about the closeness of the two average precisions, we tested for the following hypotheses:

$$
\begin{aligned}
H_1: & \quad \Delta P = P(n,\Theta_{opt}) - P(n,\Theta_{man}) \neq 0 \\
H_0: & \quad \Delta P = P(n,\Theta_{opt}) - P(n,\Theta_{man}) = 0
\end{aligned}
$$

(i)

To this end, we computed the following test statistic (Theodoridis *et al* 2003[1]):

$$
q_P = \left[ P_{ave}(n,\Theta_{opt}) - P_{ave}(n,\Theta_{man}) \right] \Bigg/ \sqrt{\frac{P_{std}(n,\Theta_{opt})^2 + P_{std}(n,\Theta_{man})^2}{|\mathbf{D}|}}
$$

(ii)

where $|\mathbf{D}|$ is the cardinality of the reference imagebase. $P_{ave}(n,\Theta_{opt})$ and $P_{ave}(n,\Theta_{man})$ are computed by (63), and $P_{std}(n,\Theta_{opt})$ and $P_{std}(n,\Theta_{man})$ are given by (64) in the revised paper. The random variable $q_P$ follows the *t*-distribution with $2|\mathbf{D}|-2$ degrees of freedom.

---

[1] S. Theodoridis and K. Koutroumbas, *Pattern Recognition*. San Diego: Academic Press, 2003.

For a specified significance level $\rho$, if the following criterion is satisfied, the hypothesis $H_0$ will be accepted which means $P_{\text{ave}}(n,\Theta_{\text{opt}})$ and $P_{\text{ave}}(n,\Theta_{\text{man}})$ are not statistically different with significance level $\rho$:

$$q_P \in C_\rho = [-c_\rho, c_\rho] \qquad \text{(iii)}$$

where $C_\rho = [-c_\rho, c_\rho]$ is the confidence interval at level $1-\rho$:

$$P(-c_\rho \leq q_P \leq c_\rho) = 1-\rho \qquad \text{(iv)}$$

where $P$ indicates the probability. Otherwise, the hypothesis $H_1$ will be accepted which means $P_{\text{ave}}(n,\Theta_{\text{opt}})$ and $P_{\text{ave}}(n,\Theta_{\text{man}})$ are statistically different with significance level $\rho$.

**Part II. Validating the Improvements**

As stated in Subsection II-A (in the revised paper), we have $|\mathbf{D}|=1000$ therefore, the confidence intervals for significance levels $\rho = 0.01$ and $\rho = 0.05$ are respectively as follows:

$$C_{0.01} = [-2.58, 2.58] \qquad \text{(v)}$$
$$C_{0.05} = [-1.96, 1.96] \qquad \text{(vi)}$$

Using (ii) and the results demonstrated in Tables I and II, we obtained:

$$q_P = 3.06 \notin C_{0.01} \qquad \text{(vii)}$$
$$q_{\bar{P}} = 2.42 \notin C_{0.05} \qquad \text{(viii)}$$
$$q_R = 3.17 \notin C_{0.01} \qquad \text{(ix)}$$
$$q_C = 2.10 \notin C_{0.05} \qquad \text{(x)}$$

According to (iii), the above test confirmed that the average precision and recall were improved with statistically significance level $\rho$=0.01. It also demonstrated that the average weighted precision and rank were improved with statistically significance level $\rho$=0.05.