Medical Principles and Practice

# Evaluating the Ability of the Bedside Index for Severity of Acute Pancreatitis Score to Predict Severe Acute Pancreatitis: A Meta-Analysis

Yu-Xia Yang    Li Li

Department of Emergency Medicine, The First Affiliated Hospital of Zhengzhou University, Zhengzhou, China

## Abstract

**Objective:** To evaluate the diagnostic performance of the bedside index for severity in acute pancreatitis (BISAP) score in predicting severe acute pancreatitis (SAP). **Materials and Methods:** A systematic search was conducted using PubMed, Cochrane library and EMBASE databases up to May 2014, and 9 related studies, including 1,972 subjects, were reviewed. Pooled sensitivity, specificity, positive likelihood ratio (PLR), negative likelihood ratio (NLR), diagnosis of odds ratio (DOR) and hierarchic summary receiver-operating characteristic (HSROC) curves, as well as the area under the HSROC curve (AUC), were assessed using the HSROC and bivariate mixed effects models. Moreover, a subgroup analysis stratified by cutoff value was performed to measure the effect of the diagnostic threshold on the performance of the BISAP score. Finally, publication bias was assessed using Deeks' funnel plot asymmetry test. Statistical analyses were performed using the STATA 12.0 software. **Results:** The pooled sensitivity, specificity, PLR, NLR and DOR of the BISAP for predicting SAP were 64.82% (95% CI: 54.47–73.74%), 83.62% (95% CI: 70.03–91.77%), 3.96 (95% CI: 2.27–6.89), 0.42 (95% CI: 0.34–0.52) and 9.41 (95% CI: 5.38–16.45), respectively. The AUC was 0.77 (95% CI: 0.73–0.80). Moreover, the subgroup analysis results demonstrated that the BISAP cutoff point at 3 had a higher specificity and greater accuracy than at 2 to predict SAP. No significant publication bias was detected across the studies (p = 0.359). **Conclusion:** The BISAP score showed low sensitivity but high specificity for assessing the severity of acute pancreatitis.

© 2015 S. Karger AG, Basel

## Introduction

Acute pancreatitis (AP) is an inflammatory condition of the pancreas with a clinical course that varies from mild to severe and is characterized by activation of pancreatic enzymes to cause self-digestion of the pancreas [1]. Generally, AP is mild, self-limiting and requires no special treatment but 20–30% of patients would develop a severe disease that can progress to systemic inflammation and cause pancreatic necrosis, multiorgan failure and potentially death [1, 2]. Therefore, it is important to choose early, quick and accurate risk stratification for AP patients, which would permit evidence-based early initiation of intensive care therapy for patients with severe AP (SAP) to prevent adverse outcomes and possible complications.

Li Li, MD
Department of Emergency Medicine
The First Affiliated Hospital of Zhenzhou University
Jianshe Donglu, No. 1, Zhengzhou 450052, Henan (China)
E-Mail lili197212@163.com

Currently, there is a variety of scoring systems for the early detection of SAP, such as Ranson's score [3], acute physiology and chronic health examination (APACHE) II [4] and the computed tomography severity index (CTSI) [5]. Moreover, many inflammation markers such as C-reactive protein or interleukin-6 are also used in the clinic [6]. In 2008, Wu et al. [7] proposed a new prognostic scoring system, the bedside index of severity in acute pancreatitis (BISAP), which was a simple and accurate method that can predict the clinical severity of AP within 24 h after admission. However, the diagnostic value of the BISAP for the diagnosis of SAP was limited by the small sample size.

In 2014, a systematic literature review showed that the BISAP was one of the best predictors of persistent organ failure for AP [8]. However, a pooled clinical value of BISAP for the diagnosis of SAP was not obtained. Therefore, the purpose of this meta-analysis was to aggregate the reported data across the different studies and to estimate summary diagnostic test accuracy measures of the BISAP score using hierarchic summary receiver-operating characteristic (HSROC) and bivariate mixed effects models.

## Materials and Methods

### Literature Search

A systematic search was performed using PubMed, Cochrane library and EMBASE databases of publications up to May 2014. The meta-analysis was done using the search terms 'acute pancreatitis' AND ('BISAP' OR 'bedside index of severity in acute pancreatitis'). Moreover, the obtained bibliographies of the enrolled studies were further hand-searched for additional references.

### Inclusion and Exclusion Criteria

Inclusion criteria were: (a) studies evaluated the BISAP score for predicting SAP; (b) the subjects were diagnosed with AP; (c) the trial design was a prospective cohort study; (d) the absolute numbers of true-positive (TP), false-negative (FN), false-positive (FP) and true-negative (TN) test results were available or derivable from the article; (e) the clinical result of patients was indicated as SAP according to the Atlanta classification.

Exclusion criteria were: (a) the numbers of TP, FN, FP and TN test results were not derivable from the article; (b) the trial was designed as cross-sectional study; (c) studies were nonoriginal items, such as review, meeting abstract, case report and comment; (d) studies had been reported in previous publications. Two reviewers (Y.-X.Y. and L.L.) independently judged the study eligibility while screening the citations. Disagreements were resolved by discussion and reached a consensus.

In total, 32 studies were originally obtained. Among them, 14 were excluded after screening abstracts or titles, then the remaining 18 articles were full-text reviewed. Of these 18 articles, 9 were excluded: 4 non-English language articles, 3 reviews and 2 studies without sufficient data for calculations. Finally, 9 studies [9–17] that involved 1,972 subjects were included in the meta-analysis.

**Table 1.** Characteristics of included studies

| Ref. No. | Patients, n | Evaluation time, h | Study design | TP | FN | TN | FP | Cutoff point |
|---|---|---|---|---|---|---|---|---|
| 9 | 310 | <24 | cohort | 21 | 11 | 178 | 100 | 2 |
| 10 | 303 | <24 | cohort | 22 | 9 | 231 | 41 | 2 |
| 14 | 50 | <24 | cohort | 19 | 5 | 23 | 3 | 2 |
| 11 | 72 | <24 | cohort | 23 | 8 | 28 | 13 | 2 |
| 12 | 497 | <24 | cohort | 62 | 39 | 329 | 67 | 2 |
| 13 | 155 | <24 | cohort | 24 | 3 | 64 | 64 | 3 |
| 15 | 299 | <24 | cohort | 10 | 12 | 274 | 3 | 3 |
| 16 | 51 | <24 | cohort | 21 | 8 | 13 | 9 | 3 |
| 17 | 185 | <24 | cohort | 15 | 25 | 134 | 11 | 3 |

### Data Extraction and Quality Assessment

The two authors (Y.-X.Y. and L.L.) independently extracted data using the predefined information sheet, and disagreements were resolved by discussion and agreed by consensus. The following characteristics were collected from each study: the first author, year of publication, source, experiment design, sample size, the reference standard (gold standard), and the numbers of TP, FN, FP and TN results. The evidence-based quality assessment tool Quality Assessment of Diagnostic Accuracy Studies (QUADAS) criteria based on 14 items was used to assess the quality of diagnostic accuracy studies included in this meta-analysis [18, 19].

### Statistical Analyses

The meta-analysis was performed by STATA 12.0 (Stata Corp., College Station, Tex., USA) software using the program 'metandi'. Sensitivity, specificity, positive likelihood ratio (PLR), negative likelihood ratio (NLR) and diagnostic odds ratio (DOR) with their 95% confidence intervals (CI) were all calculated [20]. In addition, an HSROC curve was obtained, and the area under the curve (AUC) was calculated by bivariate mixed effects models to assess the predictive accuracy of the BISAP scoring system [21]. An AUC of a perfect test was 1.0, whereas an AUC of 0.5 represents a test that performs no better than chance [22]. The HSROC curve for individual studies was generated and analyzed to explore the influence of threshold effects.

Subgroup analysis was further performed to measure the effect of diagnostic threshold on the performance of the BISAP score. Publication bias was assessed using Deeks' funnel plot asymmetry test [23]. Funnel plots for publication bias were made by a regression of the diagnostic log odds ratio against 1/square root of effective sample size, weighting by effective sample size. If a funnel plot was symmetric, publication bias was neglected, and some mechanism that links to study results with sample size was present.

## Results

### Characteristics of the Studies and Quality Assessment

Of the 32 studies accessed, 9 (28.12%) were selected. The characteristics of the 9 studies are shown in table 1, and they all adequately describe the performance charac-

**Table 2.** Quality assessment tool for diagnostic accuracy systematic review of quality criteria of included studies

| Criterion No. | Reference No. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 13 | 9 | 10 | 14 | 11 | 15 | 12 | 16 | 17 |
| 1 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 2 | Y | N | Y | Y | Y | Y | Y | Y | Y |
| 3 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 4 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 5 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 6 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 7 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 8 | Y | N | Y | Y | Y | Y | Y | Y | Y |
| 9 | N | N | Y | Y | Y | Y | Y | Y | Y |
| 10 | U | U | U | U | U | U | U | U | U |
| 11 | U | U | U | U | U | U | U | U | U |
| 12 | Y | Y | Y | Y | Y | Y | Y | Y | Y |
| 13 | U | U | U | U | U | U | U | U | U |
| 14 | Y | Y | Y | U | U | U | U | U | U |

Y = Yes, represents certain answer for the corresponding question; N = no, represents negative answer for the corresponding question; U = unclear, i.e. the information provided in the individual studies was insufficient to answer the corresponding question. QUADAS criteria: 1 = the spectrum of patients representative of the patients who received the test in practice; 2 = selection criteria clearly described; 3 = the reference standard is likely to correctly classify the target condition; 4 = the time period between reference standard and index test is short enough to be reasonably sure that the target condition did not change between the two tests; 5 = the whole sample or a random selection of the sample received verification using a reference standard of diagnosis; 6 = patients received the same reference standard regardless of the index test result; 7 = the reference was standard independently of the index test (i.e. the index test did not form part of the reference standard); 8 = the execution of the index test was described in sufficient detail to permit replication of the test; 9 = the execution of the reference was standard described in sufficient detail to permit its replication; 10 = the index test results were interpreted without knowledge of the results of the reference standard; 11 = the reference standard results were interpreted without knowledge of the results of the index test; 12 = the same clinical data were available when test results were interpreted as would be available when the test is used in practice; 13 = uninterpretable/intermediate test results were reported; 14 = withdrawals from the study were explained.

**Table 3.** Meta-analysis results for the diagnostic performance of the BISAP in predicting SAP

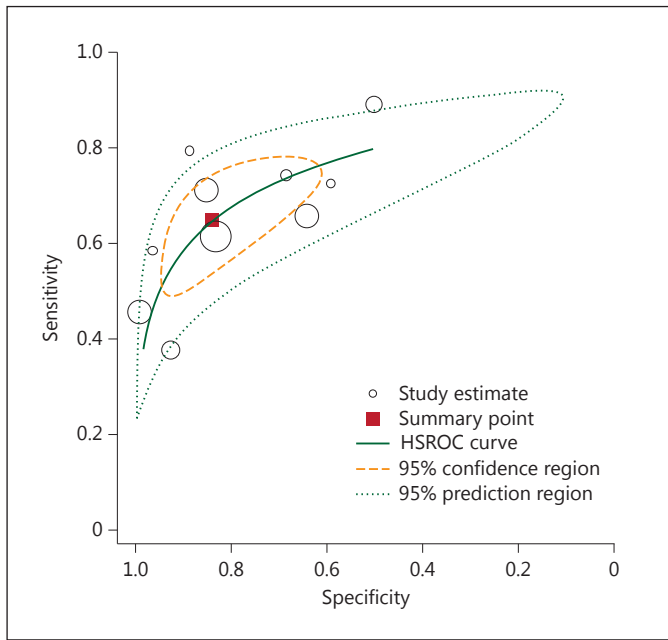| Analysis | Coefficient | 95% CI |
|---|---|---|
| Sensitivity | 0.65 | 0.54–0.74 |
| Specificity | 0.84 | 0.70–0.92 |
| DOR | 9.41 | 5.38–16.45 |
| PLR | 3.96 | 2.27–6.89 |
| NLR | 0.42 | 0.34–0.52 |
| AUC | 0.77 | 0.73–0.80 |

**Table 4.** Diagnostic performance of the BISAP in predicting SAP at different cutoff values

| Subgroup analysis | Cutoff = 2 | | Cutoff = 3 | |
|---|---|---|---|---|
| | coefficient | 95% CI | coefficient | 95% CI |
| Sensitivity | 67.30 | 60.53–73.42 | 61.18 | 41.20–78.00 |
| Specificity | 78.28 | 68.86–85.46 | 88.64 | 63.88–97.18 |
| DOR | 7.42 | 4.39–12.54 | 12.30 | 4.44–34.03 |
| PLR | 3.10 | 2.12–4.52 | 5.39 | 1.80–16.12 |
| NLR | 0.42 | 0.34–0.51 | 0.44 | 0.30–0.64 |
| AUC | 0.70 | 0.66–0.74 | 0.78 | 0.75–0.82 |

The assessment of study-specific quality according to QUADAS criteria is summarized in table 2. Overall, the enrolled studies were suitable for the meta-analysis with high quality except 3 unclear items, which were the tenth quality indicator (the index test results were interpreted without knowledge of the results of the reference), the eleventh quality indicator (the reference standard results were interpreted without knowledge of the results of the index test) and the thirteenth indicator (uninterpretable/ intermediate test results were reported) [18]. In addition, there were some studies that did not describe the details for elimination and exit objects.

*Diagnostic Value of the BISAP for SAP*
The results of the HSROC model are shown in table 3 and figure 1. The pooled sensitivity of BISAP testing for the diagnosis of SAP was 64.82% (95% CI: 54.47–73.74%), and the specificity was 83.62% (95% CI: 70.03–91.77%). The pooled DOR was 9.41 (95% CI: 5.38–16.45), the PLR was 3.96 (95% CI: 2.27–6.89), and the NLR was 0.42 (95% CI: 0.34–0.52). The AUC of the HSROC was 0.77 (95% CI: 0.73–0.80; fig. 1). The inversed and symmetry shape for the overall analysis showed that there was no significant publication bias (p = 0.359) as shown in figure 2.

teristics of the BISAP score for predicting SAP. Among these studies, Kim et al. [14] reported the results with the cutoff values at 2 and 3, respectively. BISAP scores for patients in all the studies were calculated using data within the first 24 h after admission. All included citations were designed as prospective cohort studies. The absolute numbers of TP, FN, FP and TN results were calculated based on sample size and the degree of sensitivity and specificity.

**Fig. 1.** HSROC curve of the sensitivity versus specificity of the BISAP score for the diagnosis of SAP. The curve is represented by the straight line; each of the analyzed studies is represented by a circle; the point estimate to which summary sensitivity and specificity correspond is represented by the square and the respective 95% CI by the dashed line, whereas the 95% confidence area in which a new study will be located is represented by the dotted line.
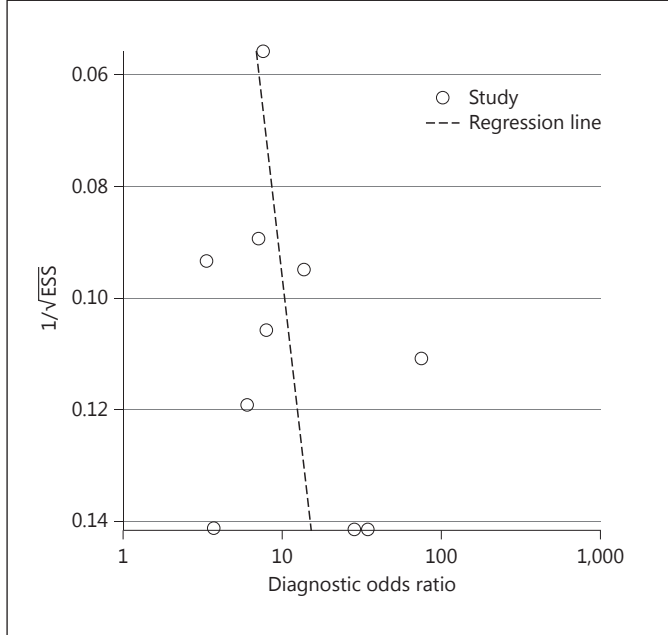


**Fig. 2.** Publication bias for identifying the diagnostic value of the BISAP score for the diagnosis of SAP: Deeks' funnel plot asymmetry test, p = 0.36. ESS = Effective sample size.

## Subgroup Analyses

The results for subgroup analyses stratified by cutoff value are shown in table 4, and a negative correlation between the logits of sensitivity and specificity was observed. When the cutoff value of the BISAP was set at 2, the pooled sensitivity, specificity, PLR, NLR and DOR were 67.30% (95% CI: 60.53–73.42%), 78.28% (95% CI: 68.86–85.46%), 3.10 (95% CI: 2.12–4.52), 0.42 (95% CI: 0.34–0.51) and 7.42 (95% CI: 4.39–12.54), respectively. The AUC of the HSROC was 0.70 (95% CI: 0.66–0.74). However, when the cutoff value was set at 3, the pooled sensitivity, specificity, PLR, NLR and DOR were 61.18% (95% CI: 41.20–78.00%), 88.64% (95% CI: 88–97.18%), 5.39 (95% CI: 1.80–16.12), 0.44 (95% CI: 0.30–0.64) and 12.30 (95% CI: 4.44–34.03), respectively. The AUC of the HSROC was 0.78 (95% CI: 0.75–0.82).

## Discussion

The pooled sensitivity and specificity values were 64.82 and 83.62%, respectively, from the meta-analysis of the diagnostic performance of the BISAP in 1,972 individuals from 9 research studies [9–17], indicating a good specificity but moderate sensitivity in diagnosing SAP. These sensitivity and specificity values confirmed the findings of several other studies that the BISAP score was a reliable and accurate means for predicting the severity of AP in the early phase [10, 12, 13]. However, when the BISAP is compared to Ranson's, APACHE II and CTSI scoring systems, it has a higher specificity but a lower sensitivity [12, 13, 17, 24]. A previous study [25] which compared the accuracy of the scoring system for SAP diagnosis based on clinical data from 2 prospective cohorts revealed that all the scoring systems in pancreatitis were cumbersome to use. Hence, the clinical application of the BISAP might be limited by its low sensitivity, and a new diagnosis scoring system for the early prediction of SAP is urgently needed.

Given the wide range of geographical distribution, we could not exclude the influence of patient selection bias and population differences on the sensitivity assessment. The DOR is a single indicator of test accuracy that combines the sensitivity and specificity data into a single number, with higher values indicating better discriminatory test performance (higher accuracy). A DOR of 1.0 indicates that a test does not discriminate between patients with the disorder and those without [25]. In our study, the AUC of the HSROC was 0.77; the result revealed that the BISAP had a relatively good discrimina-

tion to assess the severity of disease, which was similar to other reports [12, 13, 15]. For example, Khanna et al. [11] had reported that there was no significant difference between the BISAP and other scoring systems in predicting SAP in terms of AUC (BISAP 0.80, APACHE II 0.88 and Ranson's 0.85).

Since likelihood ratios are considered to be more clinically meaningful, we also presented both PLR and NLR as our measures of diagnostic accuracy. Likelihood ratios >10 or <0.1 are considered to provide strong evidence to rule in or out diagnoses, respectively, in most circumstances (indicating high accuracy) [26]. The PLR of 3.96 and NLR of 0.42 in the current study were similar to those of traditional scoring systems in predicting SAP. Zhang et al. [13] reported that the PLR and NLR values of the BISAP in predicting SAP were 1.778 and 0.222, those of APACHE II were 2.321 and 0.233, and those of Ranson's score were 4.625 and 0.264. The results of the BISAP for predicting SAP suggest that the accuracy still needs to be improved. However, the BISAP has several important advantages: it is simple to calculate and might be able to predict the severity of AP in the first 24 h after admission, and hence it is a promising method to predict SAP as previously reported [24, 27]. Furthermore, it could be used in medical decision-making at the extreme of the prediction range, such as enrollment criteria for clinical trials, and as triaging intensive care unit admission [28, 29].

Our meta-analysis had several limitations. Firstly, the definition of SAP in all enrolled studies was based on the Atlanta classification, which defined the persistence of organ failure for more than 48 h as SAP. Recently, SAP has more widely been recognized as persistent organ failure. Equally important, the Atlanta classification has some limitation that includes an uncomplicated clinical course in most patients with pseudocysts. Therefore, further meta-analyses could be needed using the newly developed SAP definition. Secondly, not all data were obtained from the transferred patients, such as their mental status, systemic inflammatory response syndrome, and the presence or absence of pleural effusion on imaging. Moreover, some complications are usually found to accompany AP, such as hepatic artery pseudoaneurysm and pneumonia, which caused a complex individual background in the studies. Further research is needed to enable a comprehensive reassessment of the pathological mechanisms of AP with attention to the effects of pre-existing risk factors (e.g. age, obesity, genetics) and well-defined end points, as well as an identification of accurate biomarkers to assess activity on these pathways that have strong predictive accuracy.

### Conclusion

The BISAP was not an ideal single method for assessing the severity of AP because of low sensitivity but high specificity. For the early prediction of AP severity, a new diagnosis strategy is needed to be developed in the future for the combination of different predictive rules.

### References

1 Forsmark CE, Baillie J: AGA Institute technical review on acute pancreatitis. Gastroenterology 2007;132:2022–2044.
2 Fagenholz PJ, Castillo CF-D, Harris NS, et al: Increasing United States hospital admissions for acute pancreatitis, 1988–2003. Ann Epidemiol 2007;17:491–497.
3 Ranson JH, Pasternack BS: Statistical methods for quantifying the severity of clinical acute pancreatitis. J Surg Res 1977;22:79–91.
4 Yeung YP, Lam B, Yip A: APACHE system is better than Ranson system in the prediction of severity of acute pancreatitis. Hepatobiliary Pancreat Dis Int 2006;5:294–299.
5 Balthazar EJ, Robinson DL, Megibow AJ, et al: Acute pancreatitis: value of CT in establishing prognosis. Radiology 1990;174:331–336.
6 Gross V, Leser H, Heinisch A, et al: Inflammatory mediators and cytokines – new aspects of the pathophysiology and assessment of severity of acute pancreatitis? Gastroenterology 1993;40:522–530.
7 Wu BU, Johannes RS, Sun X, et al: The early prediction of mortality in acute pancreatitis: a large population-based study. Gut 2008;57:1698–1703.
8 Yang CJ, Chen J, Phillips AR, et al: Predictors of severe and critical acute pancreatitis: a systematic review. Dig Liver Dis 2014;46:446–451.
9 Wang A, Xu S, Hong J, et al: The comparison of different clinical scoring systems for predicting prognosis in acute pancreatitis based on the revised Atlanta classification (in Chinese). Zhonghua Nei Ke Za Zhi/Chin J Intern Med 2013;52:668–671.
10 Park JY, Jeon TJ, Ha TH, et al: Bedside index for severity in acute pancreatitis: comparison with other scoring systems in predicting severity and organ failure. Hepatobiliary Pancreat Dis Int 2013;12:645–650.
11 Khanna AK, Meher S, Prakash S, et al: Comparison of Ranson, Glasgow, MOSS, SIRS, BISAP, APACHE-II, CTSI scores, IL-6, CRP, and procalcitonin in predicting severity, organ failure, pancreatic necrosis, and mortality in acute pancreatitis. HPB Surg 2013;2013:367581.
12 Chen L, Lu G, Zhou Q, et al: Evaluation of the BISAP score in predicting severity and prognoses of acute pancreatitis in Chinese patients. Int Surg 2013;98:6–12.
13 Zhang J, Shahbaz M, Fang R, et al: Comparison of the BISAP scores for predicting the severity of acute pancreatitis in Chinese patients according to the latest Atlanta classification. J Hepatobiliary Pancreat Sci 2014;21:689–694.

14 Kim BG, Noh MH, Ryu CH, et al: A comparison of the BISAP score and serum procalcitonin for predicting the severity of acute pancreatitis. Korean J Intern Med 2013;28:322–329.

15 Cho Y-S, Kim H-K, Jang E-C, et al: Usefulness of the bedside index for severity in acute pancreatitis in the early prediction of severity and mortality in acute pancreatitis. Pancreas 2013; 42:483–487.

16 Bezmarević M, Kostić Z, Jovanović M, et al: Procalcitonin and BISAP score versus C-reactive protein and APACHE II score in early assessment of severity and outcome of acute pancreatitis. Vojnosanit Pregl 2012;69:425–431.

17 Papachristou GI, Muddana V, Yadav D, et al: Comparison of BISAP, Ranson's, APACHE-II, and CTSI scores in predicting organ failure, complications, and mortality in acute pancreatitis. Am J Gastroenterol 2009;105:435–441.

18 Whiting P, Rutjes AW, Reitsma JB, et al: The development of QUADAS: a tool for the quality assessment of studies of diagnostic accuracy included in systematic reviews. BMC Med Res Methodol 2003;3:25.

19 Whiting PF, Weswood ME, Rutjes AW, et al: Evaluation of QUADAS, a tool for the quality assessment of diagnostic accuracy studies. BMC Med Res Methodol 2006;6:9.

20 Chu H, Cole SR: Bivariate meta-analysis of sensitivity and specificity with sparse data: a generalized linear mixed model approach. J Clin Epidemiol 2006;59:1331–1332.

21 Ferrin TE, Huang CC, Jarvis LE, et al: The MIDAS display system. J Mol Graph 1988;6: 13–27.

22 Rodday AM, Triedman JK, Alexander ME, et al: Electrocardiogram screening for disorders that cause sudden cardiac death in asymptomatic children: a meta-analysis. Pediatrics 2012;129:e999–e1010.

23 Deeks JJ, Macaskill P, Irwig L: The performance of tests of publication bias and other sample size effects in systematic reviews of diagnostic test accuracy was assessed. J Clin Epidemiol 2005;58:882–893.

24 Bollen TL, Singh VK, Maurer R, et al: A comparative evaluation of radiologic and clinical scoring systems in the early prediction of severity in acute pancreatitis. Am J Gastroenterol 2011;107:612–619.

25 Glas AS, Lijmer JG, Prins MH, et al: The diagnostic odds ratio: a single indicator of test performance. J Clin Epidemiol 2003;56:1129–1135.

26 Deeks JJ: Systematic reviews of evaluations of diagnostic and screening tests. BMJ 2001;323: 157–162.

27 Singh VK, Wu BU, Bollen TL, et al: A prospective evaluation of the bedside index for severity in acute pancreatitis score in assessing mortality and intermediate markers of severity in acute pancreatitis. Am J Gastroenterol 2009;104:966–971.

28 Zimmerman JE, Rousseau DM, Duffy J, et al: Intensive care at two teaching hospitals: an organizational case study. Am J Crit Care 1994;3:129–138.

29 Papachristou GI: Prediction of severe acute pancreatitis: current knowledge and novel insights. World J Gastroenterol 2008;14:6273–6275.