

# Urinary bladder cancer staging in CT urography using machine learning

Sankeerth S. Garapati, Lubomir Hadjiiski,<sup>a)</sup> Kenny H. Cha, Heang-Ping Chan, Elaine M. Caoili, and Richard H. Cohan

*Department of Radiology, The University of Michigan, Ann Arbor, MI 48109, USA*

Alon Weizer

*Department of Urology, Comprehensive Cancer Center, The University of Michigan, Ann Arbor, MI 48109, USA*

Ajjai Alva

*Department of Internal Medicine, Hematology-Oncology, The University of Michigan, Ann Arbor, MI 48109, USA*

Chintana Paramagul, Jun Wei, and Chuan Zhou

*Department of Radiology, The University of Michigan, Ann Arbor, MI 48109, USA*

(Received 10 February 2017; revised 4 July 2017; accepted for publication 30 July 2017; published 5 September 2017)

**Purpose:** To evaluate the feasibility of using an objective computer-aided system to assess bladder cancer stage in CT Urography (CTU).

**Materials and methods:** A dataset consisting of 84 bladder cancer lesions from 76 CTU cases was used to develop the computerized system for bladder cancer staging based on machine learning approaches. The cases were grouped into two classes based on pathological stage  $\geq T2$  or below  $T2$ , which is the decision threshold for neoadjuvant chemotherapy treatment clinically. There were 43 cancers below stage  $T2$  and 41 cancers at stage  $T2$  or above. All 84 lesions were automatically segmented using our previously developed auto-initialized cascaded level sets (AI-CALS) method. Morphological and texture features were extracted. The features were divided into subspaces of morphological features only, texture features only, and a combined set of both morphological and texture features. The dataset was split into Set 1 and Set 2 for two-fold cross-validation. Stepwise feature selection was used to select the most effective features. A linear discriminant analysis (LDA), a neural network (NN), a support vector machine (SVM), and a random forest (RAF) classifier were used to combine the features into a single score. The classification accuracy of the four classifiers was compared using the area under the receiver operating characteristic (ROC) curve ( $A_z$ ).

**Results:** Based on the texture features only, the LDA classifier achieved a test  $A_z$  of 0.91 on Set 1 and a test  $A_z$  of 0.88 on Set 2. The test  $A_z$  of the NN classifier for Set 1 and Set 2 were 0.89 and 0.92, respectively. The SVM classifier achieved test  $A_z$  of 0.91 on Set 1 and test  $A_z$  of 0.89 on Set 2. The test  $A_z$  of the RAF classifier for Set 1 and Set 2 was 0.89 and 0.97, respectively. The morphological features alone, the texture features alone, and the combined feature set achieved comparable classification performance.

**Conclusion:** The predictive model developed in this study shows promise as a classification tool for stratifying bladder cancer into two staging categories: greater than or equal to stage  $T2$  and below stage  $T2$ . © 2017 American Association of Physicists in Medicine [<https://doi.org/10.1002/mp.12510>]

**Key words:** bladder cancer staging, classification, computer-aided diagnosis, CT urography, feature extraction, machine learning, radiomics, segmentation

## 1. INTRODUCTION

Bladder cancer is one of the most common cancers affecting both men and women.<sup>1</sup> It can cause substantial morbidity and mortality among the patients with the disease. In 2017, it is estimated that there will be 79,030 new cases and 16,870 deaths from bladder cancer.<sup>1</sup> One in 42 Americans will be diagnosed with bladder cancer in their lifetime and 9 of 10 patients with this cancer are over the age of 55.<sup>1,2</sup> The average age of diagnosis is 73.<sup>1</sup> Approximately half of all bladder cancer cases are first found while the cancer is still confined to the inner wall of the bladder and has not invaded into deeper layers or distant parts of the body.<sup>1</sup> Bladder cancer has a recurrence rate of 50–80% and requires constant surveillance.

This makes it the most expensive cancer to treat, requiring a total of \$4.1 billion yearly, on a per patient basis in the United States.<sup>2</sup> Bladder cancer can be divided into three categories that include noninvasive, superficial, and invasive. The initial treatment for bladder cancer is transurethral resection of the bladder tumor (TURBT), which removes the tumor from the bladder and also helps provide information regarding the stage of the cancer.<sup>3–5</sup> Bladder cancer is staged in order to determine treatment options and estimate a prognosis for the patient. Accurate staging provides the physician with information about the extent of the cancer. The tumor stages  $T$  refer to the depth of the penetration of the tumor into the layers of the bladder.  $T0$  indicates no primary tumor,  $T1$  indicates that the tumor has invaded the connective tissue under

the epithelium, T2 indicates that the tumor has invaded the bladder muscle, T3 indicates that the tumor has invaded the fatty tissue around the bladder, and T4 indicates that the tumor has spread beyond the fatty tissue into other areas such as the pelvic wall, uterus, prostate or abdominal wall<sup>6</sup> (Fig. 1). An example of bladder cancer stage T2 is presented in Fig. 2.

The accurate staging of bladder cancer is crucial for providing proper treatment to the patient. Superficial diseases (under stage T2) can be managed with less aggressive treatment than invasive diseases (stage T2 and above).<sup>3–5</sup> There

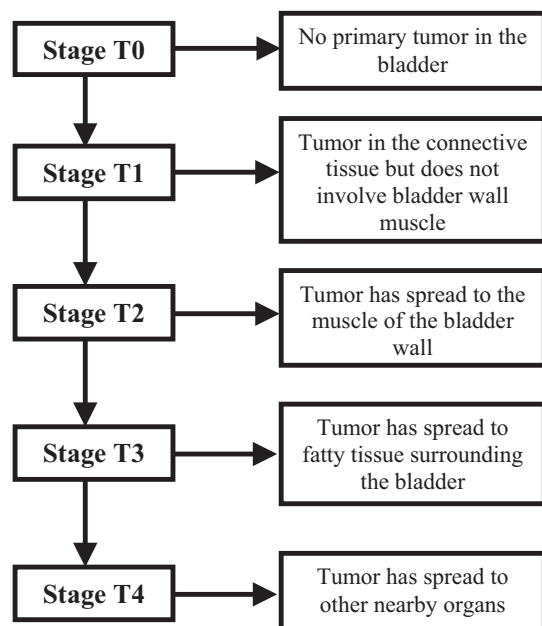


FIG. 1. Bladder cancer stage grading scale definition.

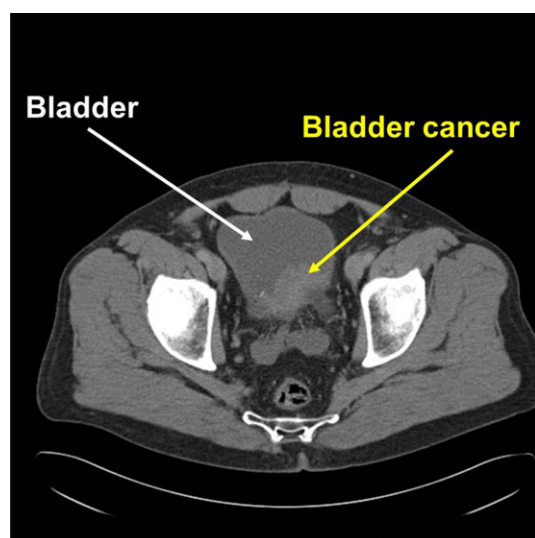


FIG. 2. Urinary bladder CT. The bladder cancer is marked and clearly visible. The cancer stage is T2. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

are two types of staging for bladder cancer — clinical and pathological. The clinical stage is the physicians' best estimate for the extent of the cancer based on physical exams and imaging. The pathological stage is determined by analysis of the tissue collected from the cancer after biopsy, tumor resection, or bladder cystectomy. The accuracy of the staging depends on the complete resection of the tumor. Incomplete resection of the tumor may reduce the reliability of the staging at the beginning of the tumor management process.<sup>7</sup> Bladder cystectomy ensures that the entire bladder tumor is present for pathological review; therefore, the pathological staging is based on the histological review of the cystectomy specimen.<sup>6</sup> Adjuvant chemotherapy is used in patients with locally advanced bladder cancer in order to reduce the chances of cancer recurrence following radical cystectomy.<sup>8</sup> Neoadjuvant chemotherapy is used prior to radical cystectomy in order to reduce the tumor size before surgical removal; for example, a cisplatin-based regimen has been shown to decrease the probability of finding extravesical disease and improve survival when compared to radical cystectomy alone.<sup>8–10</sup>

Correct staging of bladder cancer is crucial for the decision of neoadjuvant chemotherapy treatment and minimizing the risk of undertreatment or overtreatment. Patients with stage T2 to T4 carcinomas of the bladder are recommended for treatment with neoadjuvant chemotherapy. Studies found that up to 50% of the patients who are estimated to have a T1 disease at clinical staging are understaged and later upstaged after radical cystectomy.<sup>11–14</sup> This inaccuracy in staging can partly be attributed to the subjectivity and variability of clinicians in utilizing various diagnostic information. The purpose of this study is to develop an objective decision support system that can potentially reduce the risk of undertreatment or overtreatment by merging radiomic information in a predictive model using statistical outcomes and machine learning.

## 2. MATERIALS AND METHODS

### 2.A. Dataset

The data collection protocol was approved by our institutional review board and is HIPAA compliant. Patient informed consent was waived for this retrospective study. Our dataset consisted of 84 bladder cancer lesions from 76 bladder cancer CTU cases collected from patient files without additional imaging for research purpose. The CTU scans in this dataset were acquired at an image slice interval of 0.625–1.25 mm using 120 kVp and 120–280 mA. The dataset consisted of 22 noncontrast cases (22 lesions), 22 early phase contrast-enhanced cases (22 lesions), and 32 delayed-phase contrast-enhanced cases (40 lesions). Per imaging protocol, the early phase contrast-enhanced images are obtained 60 s following the initiation of a contrast injection. The delayed-phase contrast-enhanced images are obtained 12 min after the initiation of contrast injection. The type of scan a patient receives is determined by the protocol of the hospital

performing the scan. Our dataset includes patients referred to our hospital for treatment so that some scans were performed at outside hospitals and followed different scanning protocols, resulting in scans with inconsistent contrast-enhancement phase. A patient may also get a noncontrast scan due to risk factors, such as allergy to the contrast media, asthma, renal insufficiency, significant cardiac disease, or anxiety.<sup>15</sup>

For all cases, clinical and pathological staging were performed during the patient's clinical care. Cystectomy was performed after completing the course of neoadjuvant chemotherapy. The primary chemotherapy regimen used for the patients in our dataset were MVAC, which is a combination of four medications: Methotrexate, Vinblastine, Doxorubicin, and Cisplatin. Stage T2 is identified to be clinically important as a decision threshold for neoadjuvant chemotherapy treatment. The stage at the beginning of the tumor management process, based on the clinical staging and pathological staging was used as a reference standard of the tumor stage for our study.

In addition, for all bladder cancer lesions a radiologist measured the longest diameter on the pretreatment scans by using an electronic caliper provided by an in-house developed graphical user interface.

The 84 bladder cancer lesions were separated into two classes. The first class consisted of 41 cancers that were stage T2 or above and the patients were treated with neoadjuvant chemotherapy. The second class consisted of 43 cancers that were below stage T2 and patients were not referred to neoadjuvant chemotherapy treatment. The dataset was then split randomly by case into two sets with 42 cancers each while keeping the proportion of cancers between the two classes similar. The first set (Set 1) consisted of 22 cancers below stage T2 and 20 cancers stage T2 or above. The second set (Set 2) consisted of 21 cancers below stage T2 and 21 cancers stage T2 or above.

In Set 1, two patients had two lesions and one patient had three lesions. In Set 2, three patients had two lesions. In Set 1, the average tumor sizes (the longest diameters) of stage < T2

and  $\geq$  T2 were  $26.4 \pm 17.3$  and  $45.6 \pm 19.1$  mm, respectively [Fig. 3(a)]. In Set 2, the average tumor sizes (the longest diameters) of stage < T2 and  $\geq$  T2 were  $27.3 \pm 10.8$  mm and  $40.6 \pm 17.3$  mm, respectively [Fig. 3(b)].

## 2.B. Segmentation of bladder lesions on CT urography

Our previously developed method for bladder lesion segmentation using an auto-initialized cascaded level set (AI-CALS) was used.<sup>16</sup> Briefly, the system consists of three stages that include preprocessing, initial segmentation, and 3D level set segmentation (Fig. 4). The segmentation of bladder lesions is often difficult as some lesions are located in the noncontrast-enhanced region of the bladder such that contrast between the lesion and the surrounding background was low. Additionally, lesions often have irregular boundaries and can be very small and subtle. Each lesion in the dataset was marked by a bounding box as an input volume of interest (VOI). The lateral dimensions of the box were determined by an adjustable rectangle within the image slice that contains the best view of the lesion. The top and bottom slices are marked to completely enclose the lesion. The AI-CALS segmentation is then automatically performed in the VOI. In the preprocessing stage, image processing techniques including smoothing, anisotropic diffusion, gradient filters, and a rank transform of the gradient magnitude are used to generate sets of smoothed images, gradient magnitude images, and gradient vector images. The initial segmentation surface is obtained by combining information from these images. Three dimensional (3D) flood fill algorithm, morphological dilation filter, and morphologic erosion filter are applied to the initial segmentation surface to connect nearby components, which is then used to initialize the level set segmentation. The initial contour is propagated toward the lesion boundary using a bank of cascaded level sets. The level sets help refine the initial contour. The details of the AI-CALS method can be found in our previous paper.<sup>16</sup>

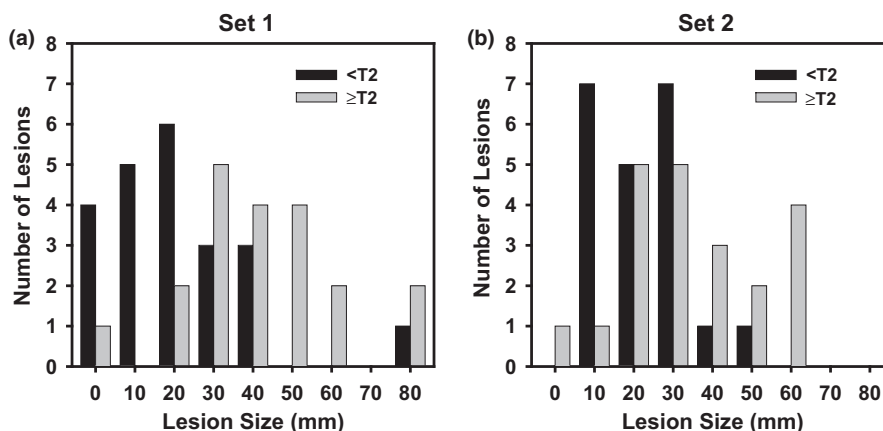


FIG. 3. Distribution of tumor sizes (the longest diameters) for Set 1 and Set 2. (a) Set 1: The average tumor sizes of stage < T2 and  $\geq$  T2 were  $26.4 \pm 17.3$  mm and  $45.6 \pm 19.1$  mm, respectively. (b) Set 2: The average tumor sizes of stage < T2 and  $\geq$  T2 were  $27.3 \pm 10.8$  mm and  $40.6 \pm 17.3$  mm, respectively.

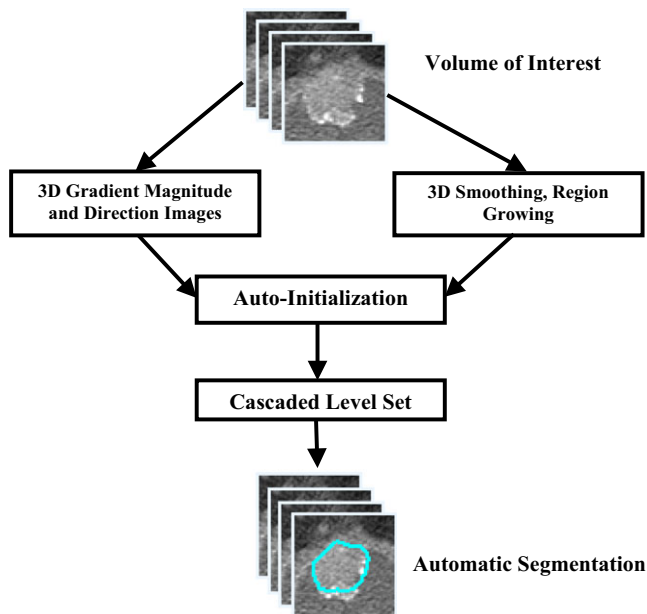


FIG. 4. Block diagram of the auto-initialized cascaded level sets (AI-CALS) method. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

### 3. CLASSIFICATION

#### 3.A. Feature extraction

Following automated computer segmentation, texture features and morphological features were extracted to characterize the lesion. The mass size was measured as its 3D volume. Five morphological features were extracted based on the normalized radial length (NRL). NRL is defined as the radial length normalized relative to the maximum radial length for the segmented object.<sup>17</sup> The NRL features extracted include zero crossing count, area ratio, standard deviation, mean, and entropy. In addition, 10 contrast features and a number of features including circularity, rectangularity, perimeter-to-area ratio, Fourier descriptor, gray level average, standard deviation of gray level, mean density, eccentricity, moment ratio, and axis ratio were extracted as shape descriptors.

The texture of the tumor margin can provide important information about its characteristics. We calculated texture features from the rubber band straightening transform (RBST) images<sup>18</sup> of the tumor margin including those from the run-length statistics matrices, filtered Dasarathy east-west direction, and filtered Dasarathy horizontal direction.<sup>19,20</sup> The texture feature set also included the gray level radial gradient direction features.

In total, 91 features were extracted to form the feature space, including 26 morphological features and 65 texture features.

#### 3.B. Feature selection/classification

A block diagram of the machine learning-based bladder cancer staging system is shown in Fig. 5. Stepwise feature selection was used to select the best subset of features to create an effective classifier.<sup>21</sup> A number of different

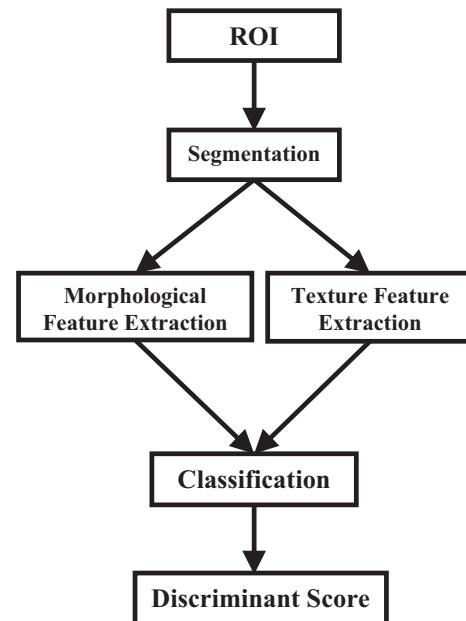


FIG. 5. Block diagram of our machine learning based staging system. We compared the linear discriminant analysis (LDA), back-propagation neural network (NN), Support vector machine (SVM), and Random forest classifiers (RAF) in the classification stage for this study.

classification experiments were performed to determine the best collection of input features. The classification performance was compared in three feature spaces: (a) morphological features only, (b) texture features only, and (c) morphological and texture features combined. A two-fold cross-validation was conducted by partitioning the dataset into Set 1 and Set 2. In the first fold, Set 1 was used for feature selection and classifier training. The trained classifier was then tested on Set 2. In the second fold, feature selection and classifier training were performed on Set 2 and then tested on Set 1.

When training on a given fold (for example, Set 1) a leave-one-case-out resampling scheme with stepwise feature selection was used to reduce the dimensionality of the feature space. In stepwise feature selection, one feature is entered or removed in alternate steps while their effect is analyzed using the Wilks' lambda criterion.<sup>21</sup> The significance of the change in the Wilks' lambda when a feature is included or removed was estimated by F statistics.  $F_{in}$ ,  $F_{out}$ , and tolerance are the parameters of the stepwise feature selection, which define the thresholds for inclusion or exclusion of a given feature. A range of  $F_{in}$ ,  $F_{out}$ , and tolerance values is evaluated by using an automated simplex optimization method. The set of  $F_{in}$ ,  $F_{out}$ , and tolerance values that lead to the highest classification result with the lowest number of features based on the training set are selected. A smaller number of features are preferred in order to reduce the chance of overfitting. Once the set of  $F_{in}$ ,  $F_{out}$ , and tolerance is selected, the stepwise feature selection with the selected parameter set is applied to the entire training fold to select a single set of features and train a single classifier. After the classifier is fixed it is applied to the test fold (e.g., Set 2) for performance evaluation.



Four different classifiers were evaluated in this study. The same partitioning of Set 1 and Set 2 was used for all classifiers. We compared the four classifiers for this classification task. The first classifier was linear discriminant analysis (LDA).<sup>22,23</sup> The LDA with the stepwise feature selection was used to determine the most effective features using the training set in each fold, as described above. The second classifier was a back-propagation neural network (NN)<sup>24</sup> with a single hidden layer and a single output node. The selected features from LDA were used for this classifier and they determined the number of input nodes to the NN. The parameters for the NN were adjusted using the training set, and the best performing network was applied to the test set. The third classifier was a support vector machine (SVM)<sup>25,26</sup> with a radial basis kernel. Using training data, a SVM determines a decision hyperplane to separate the two classes by maximizing the distance, or the margin, between the training samples of both classes and the hyperplane. The width of the SVM radial basis kernels  $\gamma$  was varied between 0.02 and 0.14 for the experiments. The best parameters for the SVM kernels for a specific experiment were selected using the training set, which were then applied to the test set. The LDA selected features were also used as the input to the SVM. The fourth one is the Random Forest (RAF) classifier.<sup>27</sup> We used the WEKA<sup>28</sup> implementation and selected 50–100 trees and five to seven features per tree for our classification task using the training set in each fold. The parameters for the random forest classifier were determined experimentally using the training sets. All 91 features were used as an input to the RAF.

### 3.C. Evaluation methods

Lesion segmentation performance was evaluated using radiologists' 3D hand-segmented contours as reference standards. The hand outlines of all 84 lesions were obtained from an experienced abdominal radiologist (RAD1). Hand outlines for a subset of 12 lesions were obtained from a second experienced abdominal radiologist (RAD2). The average distance and the Jaccard index<sup>29</sup> were calculated between the computer outlines and the hand outlines. The average distance, *AVDIST*, is defined as the average of the distances between the closest points of the two contours:

$$AVDIST(G, U) = \frac{1}{2} \left( \frac{\sum_{x \in G} \min\{d(x, y) : y \in U\}}{N_G} + \frac{\sum_{y \in U} \min\{d(x, y) : x \in G\}}{N_U} \right), \quad (1)$$

where  $G$  and  $U$  are two contours being compared.  $N_G$  and  $N_U$  denote the number of voxels on  $G$  and  $U$ , respectively. The function  $d$  is the Euclidean distance. For a given voxel along the contour  $G$ , the minimum distance to a point along the contour  $U$  is determined. The minimum distances obtained for all points along  $G$  are averaged. This process is repeated by switching the roles of  $G$  and  $U$ . *AVDIST* is then calculated as the average of the two average minimum distances.

The Jaccard index is defined as the ratio of the intersection between the reference volume and the segmented volume to the union of the reference volume and the segmented volume:

$$JACCARD^{3D} = \frac{V_G \cap V_U}{V_G \cup V_U}, \quad (2)$$

A value of 1 indicates that  $V_U$  completely overlaps with  $V_G$ , whereas a value of 0 implies  $V_U$  and  $V_G$  are disjoint.

To evaluate the classifier performance, the training and test scores output from the classifier were analyzed using the receiver-operating characteristic (ROC) methodology.<sup>30</sup> The classification accuracy was evaluated using the area under the ROC curve,  $A_z$ . The statistical significance of the differences between the different classifiers and feature spaces were estimated by the CLABROC program using ROC software by Metz et al.<sup>31,32</sup>

## 4. RESULTS

The lesion segmentation performance of the AI-CALS compared to the radiologist hand outlines for the 84 lesions are shown in Table I. Table II shows the computer segmentation performance compared to two different radiologists' hand outlines for a subset of 12 lesions.

The performance of the classifiers based on different machine learning techniques, the LDA, NN, SVM, and RAF, is summarized in Table III. Different feature spaces containing the morphological features, the texture features, and the combined set of both morphological and texture features were used for classification. The features selected with LDA were used in the SVM and NN classifiers. The LDA classifier with morphological features achieved a training  $A_z$  of 0.91 on Set 1 and a test  $A_z$  of 0.81 on Set 2. For training on Set 2 it achieved a  $A_z$  of 0.97 and a test  $A_z$  of 0.90 on Set 1. The selected features on the training sets included volume, a contrast feature, and gray level feature. The test  $A_z$  of the NN for Set 1 and Set 2 was 0.88 and 0.91, respectively. The SVM

TABLE I. Segmentation performance of the 84 lesions compared to hand outlines performed by radiologist 1 (RAD1).

| AI-CALS vs. RAD1                           |              |
|--|--------------|
| Average distance <i>AVDIST</i>             | 4.9 ± 2.7 mm |
| Jaccard index <i>JACCARD</i> <sup>3D</sup> | 43.5 ± 14.0% |

TABLE II. Segmentation performance for a subset of 12 lesions compared to hand outlines performed by two different radiologists (RAD1, RAD2).

|  | AI-CALS<br>vs. RAD1 | AI-CALS<br>vs. RAD2 | RAD1<br>vs. RAD2 |
|--|---------------------|---------------------|------------------|
| Average distance <i>AVDIST</i>             | 5.2 ± 2.5 mm        | 4.1 ± 1.5 mm        | 2.9 ± 1.1 mm     |
| Jaccard index <i>JACCARD</i> <sup>3D</sup> | 43.2 ± 13.2%        | 50.1 ± 14.7%        | 58.7 ± 11.1%     |

TABLE III. Summary results for LDA, NN, SVM, and RAF classifiers in morphological, texture, and combined feature spaces. The column “Number of Features” did not apply to the RAF classifier. All features were used for the RAF classifier. The differences in the  $A_z$  values between pair-wise comparison of the different classifiers did not achieve statistical significance after performing Bonferroni correction for the 18 comparisons ( $P > 0.0028$ ).

| Feature type                     | Number of features | LDA      |         | NN       |         | SVM      |         | RAF      |         |
|----------------------------------|--------------------|----------|---------|----------|---------|----------|---------|----------|---------|
|                                  |                    | Training | Testing | Training | Testing | Training | Testing | Training | Testing |
| Morphological features           |                    |          |         |          |         |          |         |          |         |
| Training (Set 1) testing (Set 2) | 4                  | 0.91     | 0.81    | 0.96     | 0.91    | 0.95     | 0.90    | 1        | 0.88    |
| Training (Set 2) testing (Set 1) | 4                  | 0.97     | 0.90    | 0.98     | 0.88    | 0.97     | 0.88    | 1        | 0.83    |
| Texture features                 |                    |          |         |          |         |          |         |          |         |
| Training (Set 1) testing (Set 2) | 2                  | 0.91     | 0.88    | 0.95     | 0.92    | 0.92     | 0.89    | 1        | 0.97    |
| Training (Set 2) testing (Set 1) | 7                  | 1        | 0.91    | 1        | 0.89    | 1        | 0.91    | 1        | 0.89    |
| Combined features                |                    |          |         |          |         |          |         |          |         |
| Training (Set 1) testing (Set 2) | 3                  | 0.92     | 0.90    | 0.97     | 0.95    | 0.92     | 0.89    | 1        | 0.96    |
| Training (Set 2) testing (Set 1) | 7                  | 1        | 0.89    | 1        | 0.91    | 1        | 0.92    | 1        | 0.86    |

achieved test  $A_z$  of 0.88 on Set 1 and test  $A_z$  of 0.90 on Set 2. The test  $A_z$  of the RAF for Set 1 and Set 2 was 0.83 and 0.88, respectively. The distribution of the discriminant scores from the four classifiers for testing on Set 1 and Set 2 in two-fold cross-validation in the morphological feature space are presented in Fig 6. It can be observed that most of the classifiers were able to provide a relatively good separation between the two classes.

By using the texture features the LDA classifier achieved a test  $A_z$  of 0.91 on Set 1 and a test  $A_z$  of 0.88 on Set 2. When trained on Set 1 or Set 2 the stepwise feature selection procedure selected subsets of the filtered Dasarathy east-west direction features, the filtered Dasarathy horizontal direction features and the gray level radial gradient direction features. The test  $A_z$  of the NN classifier for Set 1 and Set 2 was 0.89 and 0.92, respectively. The SVM classifier achieved test  $A_z$  of 0.91 on Set 1 and test  $A_z$  of 0.89 on Set 2. The test  $A_z$  of the RAF classifier for Set 1 and Set 2 was 0.89 and 0.97, respectively.

When the morphological and the texture features were combined, the LDA classifier achieved a test  $A_z$  of 0.89 on Set 1 and a test  $A_z$  of 0.90 on Set 2. When trained on Set 1 or Set 2 the stepwise feature selection procedure selected a contrast feature, subsets of the filtered Dasarathy horizontal direction features, and subsets of the gray level radial gradient direction features. The test  $A_z$  of the NN classifier for Set 1 and Set 2 was 0.91 and 0.95, respectively. The SVM classifier achieved test  $A_z$  of 0.92 on Set 1 and test  $A_z$  of 0.89 on Set 2. The test  $A_z$  of the RAF classifier for Set 1 and Set 2 was 0.86 and 0.96, respectively. The test ROC curves for all of the classifiers when tested on Set 1 and Set 2 in the two-fold cross-validation in the different feature spaces are shown in Fig. 7.

The differences in the  $A_z$  values between pairs of classifiers did not achieve statistical significance. The classifiers achieved slightly higher  $A_z$  values in the texture and combined feature spaces than in the morphological feature space; however, the differences did not achieve statistical significance after Bonferroni correction for the multiple comparisons ( $P$ -value  $< 0.05/18 = 0.0028$  to be considered significant).

## 5. DISCUSSION

The agreement between the AI-CALS lesion segmentation and the radiologists' manual segmentation was slightly lower than the agreement between two radiologists' hand outlines, indicating that the computer segmentation will need to be further improved. Both the morphological and the texture features were important for classifying the bladder cancer stage. When only morphological features were used in the classifier, volume, and contrast features were always selected. Volume was the primary feature used to describe lesion size. When the classifier used only the texture features, the features from the three main groups, the filtered Dasarathy east-west direction features, the filtered Dasarathy horizontal direction features, and the gray level radial gradient direction features were consistently selected. There was essentially no change in classification accuracy when the morphological features were added to the texture features in the combined set.

The LDA, SVM, and NN classifiers all led to relatively consistent results. There was no statistically significant difference in the performances between pairs of the classifiers. The best overall results for the two-fold cross-validation were obtained when a combined feature set was used with an NN classifier. Using Set 1 for training, the training  $A_z$  was 0.97 and the test  $A_z$  was 0.95. Using Set 2 for training, the training  $A_z$  was 1.00 and the test  $A_z$  was 0.91.

The RAF classifier showed greater imbalance between Set 1 and Set 2 than the other classifiers. When training was done on Set 2 and testing on Set 1, the  $A_z$  were substantially lower than the  $A_z$  values when training was done on Set 1 and testing on Set 2. For example, the test  $A_z$  decreased from 0.88 to 0.83 for morphological features, from 0.97 to 0.89 for texture features only, and from 0.96 to 0.86 for the combined features. This imbalance between the two sets could be due to the fact that RAF utilized all the features in the subspace, whereas the other three classifiers involved feature selection.

Examples of bladder cancers with stages  $\geq T2$  or  $< T2$  and the corresponding classifier scores are shown in Fig. 8. The reported scores are test scores for the LDA, SVM, NN, and RAF classifiers based on the morphological features. In

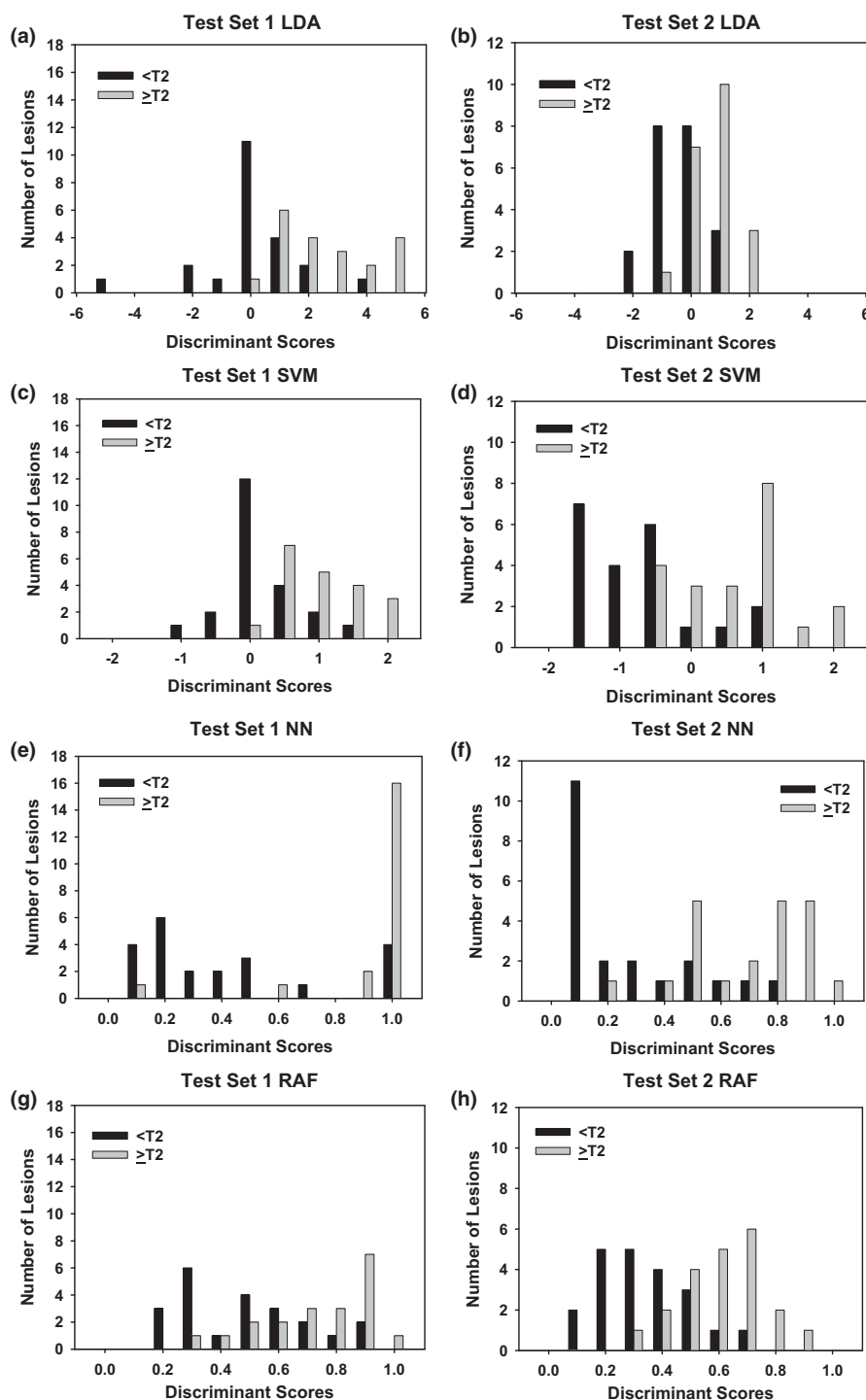


FIG. 6. Distribution of the classifiers discriminant scores for testing on Set 1 and Set 2 in two-fold cross-validation using the morphological features. (a) LDA (Set 1)  $A_z = 0.90$ , (b) LDA (Set 2)  $A_z = 0.81$ , (c) SVM (Set 1)  $A_z = 0.88$ , (d) SVM (Set 2)  $A_z = 0.90$ , (e) NN (Set 1)  $A_z = 0.88$ , (f) NN (Set 2)  $A_z = 0.91$ , (g) RAF (Set 1)  $A_z = 0.83$ , (h) RAF (Set 2)  $A_z = 0.88$ .

Figs. 8(a)–8(d) the T1 stage cancers of different sizes that were correctly classified with low scores by all classifiers are shown. Note that the output score ranges are different for different classifiers so that the score values should not be compared across classifiers. T3 stage and T2 stage cancers that were correctly classified with high scores from all classifiers are presented in Figs. 8(e)–8(h), respectively. A case that was

clinically identified as T1 stage pre-surgery but later was identified as a T2 stage cancer postsurgery is shown in Figs. 8(k) and 8(l). The classifiers classified the cancer as  $\geq T2$  with high scores. Figs. 8(m) and 8(n) show a T2 stage cancer that was incorrectly identified by the LDA, SVM, and NN classifiers with low scores, but correctly identified by the RAF with a high score.

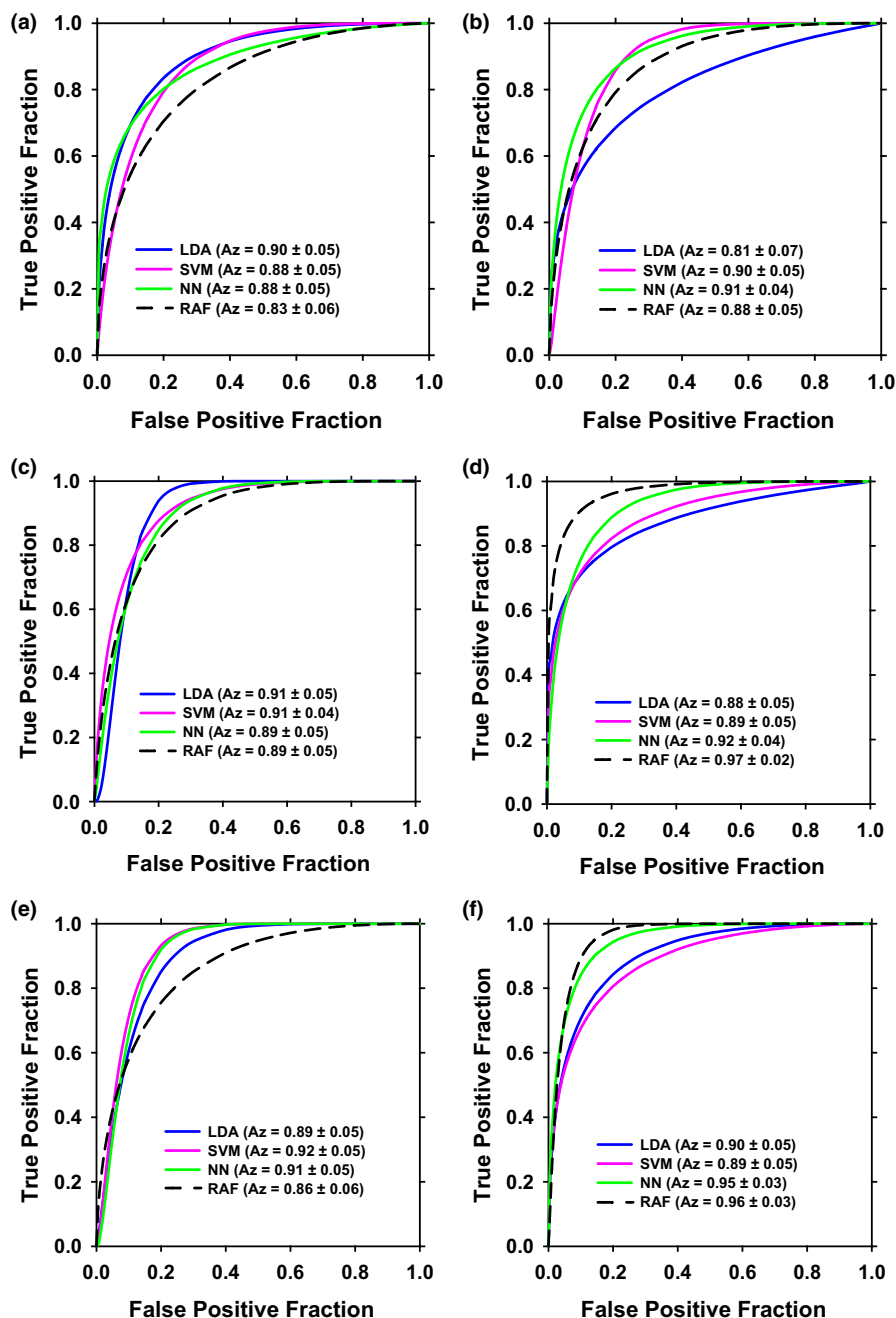


FIG. 7. ROC curves for testing on Set 1 and Set 2 in two-fold cross-validation for LDA, SVM, NN, and RAF classifiers: Left column: testing on Set 1, right column: testing on Set 2. (a) and (b) morphological features; (c) and (d) texture features; (e) and (f) combined features. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

We also have extracted features from the manually segmented bladder lesions and applied the four different types of classifiers with the different feature sets to the cancer stage prediction. The classifiers using features extracted from the manually segmented lesions performed similarly to the classifiers using features extracted from the AI-CALS segmented lesions. The test  $A_z$  values ranged from 0.77 to 0.95. For 6 of the 24 experiments the classifiers using features extracted from the manually segmented lesions performed better than classifiers using features extracted from the AI-CALS segmentations. However, the differences did not reach statistical

significance. Therefore, although the performance of the AI-CALS lesion segmentation was slightly lower than the radiologists' hand outlines the final classification results were similar.

The main limitation of the study is the small dataset. Another limitation is that we have not applied the deep learning convolution neural network (DLCNN) to this bladder cancer staging task. DLCNN has been shown to be superior to conventional classifiers in many classification tasks, especially the classification of natural scene images with millions of training samples. It also shows promise in number of



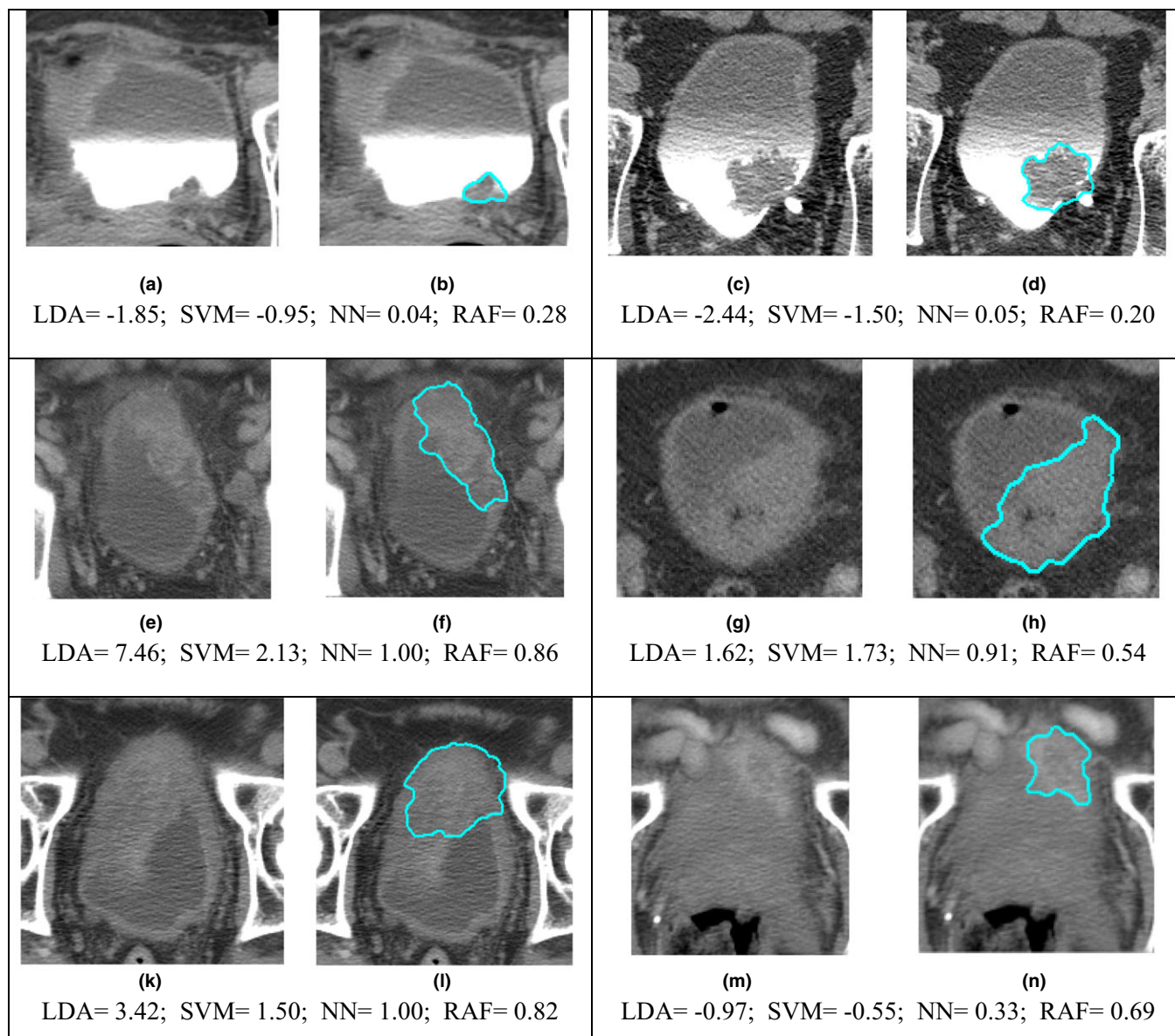


FIG. 8. Examples of bladder cancers with stages  $\geq T2$  or  $< T2$ . The outlines represent the AI-CALS segmentation. The reported scores are test scores for the LDA, SVM, NN, and RAF classifiers based on the morphological features. Note that the output score ranges are different for different classifiers so that the score values should not be compared across classifiers. The two cases in (a, b) and (c, d) both contained a T1 stage cancer that was properly classified with low scores from all classifiers. (e, f) was a T3 stage case that was properly classified with high scores from all classifiers. (g, h) was a T2 stage case that was properly classified with high scores from all classifiers. (k, l) was a case that was clinically identified as T1 presurgery but was identified as a T2 stage cancer postsurgery. The classifiers classified the cancer as  $\geq T2$  with high scores. (m, n) was T2 stage cancer that was incorrectly identified by the LDA, SVM, and NN classifiers with low scores and correctly identified by the RAF with a high score. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

medical imaging applications<sup>33,34</sup> including bladder segmentation<sup>35</sup> and bladder cancer treatment response monitoring.<sup>36</sup> However, our experience with DLCNN also indicates that it is not always the best, perhaps limited by the relatively small annotated training set in medical imaging, even with transfer learning. As the performances of the four conventional classifiers used in this study were quite high, it would not be a fair comparison for DLCNN if we do not have adequate training for the latter. We will continue to collect additional cases and compare the conventional classifiers with DLCNN for bladder cancer staging in a future study.

## 6. CONCLUSION

In this preliminary study we proposed machine learning methods for prediction of bladder cancer stage. It was found that the morphological features and texture features were useful for assessing the stage of bladder lesions. The LDA, SVM, and NN classifiers all led to relatively consistent results. There was a trend that the SVM and NN classifier slightly outperformed the LDA classifier. The best overall results for the two-fold cross-validation were obtained when a combined feature subspace was used with the NN

classifier. Further studies are under way to improve the staging of bladder cancer and test the classifier on a larger dataset, and to investigate the potential of improving the predictive model by combining imaging biomarkers with non-imaging biomarkers.

## ACKNOWLEDGMENTS

This work is supported by National Institutes of Health grant number U01CA179106.

<sup>a)</sup>Author to whom correspondence should be addressed. Electronic mail: lhadjisk@umich.edu; Telephone: (734) 647-7428.

## REFERENCES

- American Cancer Society. *Cancer Facts & Figures 2017*. Atlanta: American Cancer Society Inc.; 2017.
- Bladder Cancer Advocacy Network, [www.bcan.org/facts](http://www.bcan.org/facts); 2017, Bladder Cancer Facts (2017).
- Chang SS, Boorjian SA, Chou R, et al. Diagnosis and treatment of non-muscle invasive bladder cancer: AUA/SUO guideline. *J Urol*. 2016;196:1021–1029.
- Witjes JA, Comperat E, Cowan NC, et al. Guidelines on muscle-invasive and metastatic bladder cancer, European association of urology; 2016.
- Babjuk M, Bohle A, Burger M, et al. Guidelines on Non-muscle-invasive Bladder Cancer (Ta, T1 and CIS), European Association of Urology; 2016.
- AJCC Cancer Staging Handbook*, 8th ed. Chicago, IL: American Joint Committee on Cancer; 2016.
- Herr HW, Donat SM. Quality control in transurethral resection of bladder tumours. *BJU Int*. 2008;102:1242–1246.
- Meeks JJ, Bellmunt J, Bochner BH, et al. A systematic review of neoadjuvant and adjuvant chemotherapy for muscle-invasive bladder cancer. *Eur Urol*. 2012;62:523–533.
- Fagg SL, Dawsonedwards P, Hughes MA, Latief TN, Rolfe EB, Fielding JWL. CIS-Diamminedichloroplatinum (DDP) as initial treatment of invasive bladder cancer. *Br J Urol*. 1984;56:296–300.
- Raghavan D, Pearson B, Coorey G, et al. Intravenous CIS-platinum for invasive bladder cancer – safety and feasibility of a new approach. *Med J Aust*. 1984;140:276–278.
- Huguet J, Crego M, Sabate S, Salvador J, Palou J, Villavicencio H. Cystectomy in patients with high risk superficial bladder tumors who fail intravesical BCG therapy: pre-cystectomy prostate involvement as a prognostic factor. *Eur Urol*. 2005;48:53–59.
- Fritzsche HM, Burger M, Svatek RS, et al. Characteristics and outcomes of patients with clinical T1 grade 3 urothelial carcinoma treated with radical cystectomy: results from an international cohort. *Eur Urol*. 2010;57:300–309.
- Turker P, Bostrom PJ, Wroclawski ML, et al. Upstaging of urothelial cancer at the time of radical cystectomy: factors associated with upstaging and its effect on outcome. *BJU Int*. 2012;110:804–811.
- Shariat SF, Palapattu GS, Karakiewicz PI, et al. Discrepancy between clinical and pathologic stage: impact on prognosis after radical cystectomy. *Eur Urol*. 2007;51:137–151.
- ACR Manual on Contrast Media. ACR Committee on Drugs and Contrast Media; 2016.
- Hadjiiski LM, Chan H-P, Caoili EM, Cohan RH, Wei J, Zhou C. Auto-initialized cascaded level set (AI-CALS) segmentation of bladder lesions on multi-detector row CT urography. *Acad Radiol*. 2013;20:148–155.
- Hadjiiski LM, Sahiner B, Chan H-P, Petrick N, Helvie MA, Gurcan MN. Analysis of temporal change of mammographic features: computer-aided classification of malignant and benign breast masses. *Med Phys*. 2001;28:2309–2317.
- Sahiner B, Chan H-P, Petrick N, Helvie MA, Goodsitt MM. Computerized characterization of masses on mammograms: the rubber band straightening transform and texture analysis. *Med Phys*. 1998;25:516–526.
- Dasarathy BR, Holder EB. Image characterizations based on joint gray-level run-length distributions. *Pattern Recog Lett*. 1991;12:497–502.
- Way TW, Hadjiiski LM, Sahiner B, et al. Computer-aided diagnosis of pulmonary nodules on CT scans: segmentation and classification using 3D active contours. *Med Phys*. 2006;33:2323–2337.
- Chan H-P, Wei D, Helvie MA, et al. Computer-aided classification of mammographic masses and normal tissue: linear discriminant analysis in texture feature space. *Phys Med Biol*. 1995;40:857–876.
- Lachenbruch PA. *Discriminant Analysis*. New York: Hafner Press; 1975.
- Tatsuoka MM. *Multivariate Analysis, Techniques for Educational and Psychological Research*, 2nd edn. New York: Macmillan; 1988.
- Rumelhart DE, Hinton GE, Williams RJ. *Learning Internal Representation by Error Propagation, Parallel Distributed Processing*. Cambridge, MA: MIT Press; 1986.
- Vapnik VN. *Statistical Learning Theory*. New York: Wiley; 1998.
- Burges CJC. A tutorial on support vector machines for pattern recognition. *Data Min Knowl Disc*. 1998;2:121–167.
- Ho TK. The random subspace method for constructing decision forests. *IEEE Trans Pattern Anal Mach Intell*. 1998;20:832–844.
- Witten IH, Frank E, Hall MA, Pal CJ, The WEKA Workbench. *Online Appendix for Data Mining: Practical Machine Learning Tools and Techniques*. Burlington, MA: Morgan Kaufmann; 2016.
- Jaccard P. The distribution of the flora in the alpine zone. *New Phytol*. 1912;11:37–50.
- Metz CE. ROC methodology in radiologic imaging. *Invest Radiol*. 1986;21:720–733.
- Metz CE, Herman BA, Shen JH. Maximum-likelihood estimation of receiver operating characteristic (ROC) curves from continuously-distributed data. *Stat Med*. 1998;17:1033–1053.
- Metz ROC Software. University of Chicago Medical Center Department of Radiology, see <http://metz-roc.uchicago.edu/MetzROC/software>.
- Litjens G, Kooi T, Bejnordi BE, et al. A Survey on Deep Learning in Medical Image Analysis. arXiv:1702.05747; 2017.
- Greenspan H, van Ginneken B, Summers RM. Deep learning in medical imaging: overview and future promise of an exciting new technique. *IEEE Trans Med Imaging*. 2016;35:1153–1159.
- Cha KH, Hadjiiski L, Samala RK, Chan HP, Caoili EM, Cohan RH. Urinary bladder segmentation in CT urography using deep-learning convolutional neural network and level sets. *Med Phys*. 2016;43:1882–1896.
- Cha KH, Hadjiiski LM, Chan H-P, et al. Bladder cancer treatment response assessment using deep learning in CT with transfer learning. *Proc SPIE*. 2017;10134:101341–101346.