

Online Supplementary Material

Knottnerus BJ, Geerlings SE, Moll van Charante EP, ter Riet G. Towards a simple diagnostic index for acute uncomplicated urinary tract infections. *Ann Fam Med*. 2013;11(5):442-451.

<http://www.annfammed.org/content/11/5/442>

Supplemental Appendix. Detailed Description of Variable Selection Method

All analyses were performed in Stata/SE 10.1 (StataCorp LP).

Multiple Imputation

Missing values were imputed using multiple imputation by chained equations^{1,2} (using Stata's *ice* command), creating 45 complete data sets. For multiple imputation, 58 variables were used, including possible interactions (see below).

After multiple imputation, patients with missing values for the dependent variable were dropped.³

Variable Selection

To avoid overfitting by using too many candidate predictors, the 22 most relevant variables (based on literature and clinical usefulness) were used for the analysis, including 2 interaction variables (see below).

Logistic regression with bootstrapped backward elimination was used to derive a parsimonious diagnostic index. The result of the urine culture was the binary dependent variable, $\geq 10^3$ colony-forming units (CFUs) of a single uropathogen per milliliter (mL) being defined as a positive culture according to international guidelines.⁴ For nominal predictors, categorical coding was used to create binary dummy variables. For ordinal predictors, ordinal coding was used. For the only continuous variable "age," fractional polynomials were used to select the best fitting function, but only the linear function was selected.⁵

We selected 5 different sets of variables, based on clinical practice:

1. History = only patient history questions
2. History + dipstick = variables selected in the model above plus + urine dipstick variables
3. History + dipstick + sediment = variables selected in 2 models above + urinary sediment variables
4. History + dipstick + dipslide = variables selected in first 2 models + dipslide variables
5. History + dipstick + sediment + dipslide = all selected variables

For each of the 5 sets of variables, the following procedure was performed:

From each of the 45 imputed data sets, 600 bootstrap samples were generated. In each bootstrap sample, variables were eliminated backwards at a significance level of .05 using Stata's *mfpboot* command.

For each variable, a bootstrap inclusion fraction (BIF) was obtained in each imputed data set. Next, the mean BIF across all 45 imputed data sets was calculated for each variable. Variables with a mean BIF $\geq 66.67\%$ were retained in the final models.

After having obtained the final set of variables, parameter-wise shrinkage of their regression coefficients was performed in each imputed data set to correct for possible overoptimism,⁶ using Stata's *cval* command with parameter-wise option and 10-fold cross-validation. The resulting 45 shrunk regression coefficients and their standard errors were averaged. The averaged regression coefficients were used to compose risk scores as described by Sullivan et al.⁷

Variables

The 58 variables below were used for multiple imputation. The 22 variables in italics were used for the logistic regression analysis.

| History | Urine collection |
|--|--|
| Name of health centre / GP | Midstream urine |
| <i>Age in years</i> | Minutes from urine collection until placing into fridge |
| Ethnicity | |
| General health according to patient | |
| Marital status | Dipstick |
| <i>Duration of symptoms in days</i> | <i>Blood</i> |
| Lower abdominal pain | Glucose |
| Lower abdominal pain when? | <i>Leucocyte esterase</i> |
| Back pain | <i>Nitrite positive</i> |
| Back pain when? | Protein |
| Blood in urine according to patient | Sediment |
| Burning sensation during micturition | <i>Bacteria / HPF</i> |
| <i>Not able to empty bladder completely</i> | <i>Leucocytes / HPF</i> |
| <i>False urge to urinate</i> | Squamous epithelial cells / HPF |
| <i>More frequent micturition than usually</i> | |
| Urine incontinence: frequency | Dipslide |
| Urine incontinence: type | <i>Cystine Lactose Electrolyte Deficient (CLED) medium</i> |
| <i>Micturition of smaller amounts than usually</i> | MacConkey medium |
| <i>Pain during micturition</i> | |
| Bad smell of urine | Culture |
| Urge to urinate hard to control | $\geq 10^3$ CFU of a single uropathogen per milliliter (mL) ^a |
| <i>Vaginal discharge</i> | |
| <i>Vaginal irritation or itching</i> | Interactions |
| Currently menstruating? | <i>Patient thinks she has a UTI? \times UTIs in past year according to patient (n)</i> |
| <i>Sexual activity past week (n)</i> | <i>Patient thinks she has a UTI? \times ≥ 1 UTI ever diagnosed according to patient</i> |
| Voiding directly after sex | Patient thinks she has a UTI? \times Periods of UTI symptoms without consulting physician |
| <i>Patient thinks she has a UTI?</i> | Patient thinks she has a UTI? \times $\geq 10^3$ CFU/mL of uropathogen in culture |
| Bother at social activities | UTIs in past year according to patient (n) \times $\geq 10^3$ CFU/mL of uropathogen in culture |
| Bother at work/school | ≥ 1 UTI ever diagnosed according to patient \times $\geq 10^3$ CFU/mL of uropathogen in culture |
| Days willing to delay antibiotic | Periods of UTI symptoms without consulting physician \times $\geq 10^3$ CFU/mL of uropathogen in culture |
| <i>Last menstruation > 1 year ago?</i> | |
| Diabetes mellitus (according to patient) | |
| Any first-grade relative >2x/year UTI | |
| <i>UTIs in past year according to patient (n)</i> | |
| ≥ 1 UTI ever diagnosed according to patient | |
| Periods of UTI symptoms without consulting physician | |
| ≥ 4 hours no urinating until collection of urine sample | |

CFU = colony-forming units; HPF = high-power field; UTI = urinary tract infection.

^a Dependent variable.

References

- van Buuren S: Multiple imputation of discrete and continuous data by fully conditional specification. *Stat Methods Med Res.* 2007;16(3):219-242.
- White IR, Royston P, Wood AM: Multiple imputation using chained equations: Issues and guidance for practice. *Stat Med.* 2011;30(4):377-399.
- von Hippel PT: Regression with missing Ys: An improved strategy for analyzing multiply imputed data. *Sociol Methodol.* 2007;37:83-117.
- European urinalysis guidelines. *Scand J Clin Lab Invest Suppl.* 2000;231(Suppl 231):1-86.
- Royston P, Sauerbrei W. *Multivariable Model-Building.* New York NY: Wiley; 2008.
- Verweij PJ, van Houwelingen HC: Cross-validation in survival analysis. *Stat Med.* 1993;12(24):2305-2314.
- Sullivan LM, Massaro JM, D'Agostino RB Sr.: Presentation of multivariate data for clinical use: The Framingham Study risk score functions. *Stat Med.* 2004;23(10):1631-1660.