

Table S1: A detailed overview of the sequencing data used in this study. For each species, the study accession number, sequencing platforms, number of genomes and a measure of sequence quality (FastQC) are shown.

Species	Accession	Sequencing Platform	Genomes	FastQC ¹
<i>C. difficile</i>	PRJEB11776	Illumina HiSeq 2500	470	98.1%
<i>M. tuberculosis</i>	PRJEB7281	Illumina MiSeq / HiSeq 2500	154	31.8%
<i>P. aeruginosa</i>	PRJNA264310	Illumina MiSeq / HiSeq 2000	390	64.6%
<i>S. pneumoniae</i>	PRJEB2632	Illumina HiSeq 2000	681	71.8%

¹ Proportion of the genomes that passed the FastQC quality control for per base sequence quality.

Table S2: The antibiotic resistance datasets used in this study. For each dataset, the number of genomes and k -mers, the type of antibiotic susceptibility data (SIR¹ or MIC²) and any threshold used to separate the isolates into resistant and sensitive groups are shown. The resistant group (R) was assigned the positive class ($y = 1$) and the sensitive group (S) was assigned the negative class ($y = 0$).

Dataset	Genomes		Antibiotic Susceptibility			
	Examples	k -mers	Type	S	R	Threshold Source
<i>C. difficile</i>						
Azithromycin	462	32 752 570	MIC	≤ 8	> 8	Inferred ³
Ceftriaxone	285	25 405 987	SIR	–	–	–
Clarithromycin	462	32 752 570	MIC	≤ 4	> 4	Bourgault et al. (2006) ⁴
Moxifloxacin	462	32 752 570	MIC	≤ 4	> 4	Bourgault et al. (2006) ⁴
<i>M. tuberculosis</i>						
Ethambutol	139	9 465 489	SIR	–	–	–
Isoniazid	141	9 701 935	SIR	–	–	–
Pyrazinamide	111	8 058 479	SIR	–	–	–
Rifampicin	141	9 701 935	SIR	–	–	–
Streptomycin	135	9 282 080	SIR	–	–	–
<i>P. aeruginosa</i>						
Amikacin	365	116 441 834	SIR	–	–	–
Doripenem	363	122 438 059	SIR	–	–	–
Levofloxacin	358	122 216 859	SIR	–	–	–
Meropenem	368	123 466 989	SIR	–	–	–
<i>S. pneumoniae</i>						
Benzylpenicillin	421	8 968 176	MIC	≤ 0.064	> 2	EUCAST ⁵
Erythromycin	556	9 666 898	MIC	≤ 0.25	> 0.5	EUCAST ⁵
Tetracycline	324	8 657 259	MIC	≤ 1	> 2	EUCAST ⁵

¹ Sensitive/Intermediate/Resistant labels that were predefined by the authors of the dataset.

² Measures of minimum inhibitory concentration. In this case, a threshold was used to assign SIR labels to the isolates.

³ No official MIC breakpoints were available. A visual inspection of the MIC distribution showed a clear separation between two groups of strains. The group with the lowest MIC values was deemed sensitive and the other one was deemed resistant.

⁴ Bourgault, A. M., Lamothe, F., Loo, V. G., & Poirier, L. (2006). In vitro susceptibility of *Clostridium difficile* clinical isolates from a multi-institutional outbreak in Southern Quebec, Canada. *Antimicrobial agents and chemotherapy*, 50(10), 3473-3475.

⁵ The European Committee on Antimicrobial Susceptibility Testing. Breakpoint tables for interpretation of MICs and zone diameters. Version 5.0, 2015.

Table S3: Analysis of overfitting: The training and testing¹ set error rates are shown for the set covering machine (SCM) and support vector machines with linear and polynomial kernels (LinSVM, PolySVM).

Dataset	SCM ²	LinSVM ³	PolySVM ³
<i>C. difficile</i>			
Azithromycin	0.014 (0.030)	0.004 (0.050)	0.002 (0.048)
Ceftriaxone	0.042 (0.073)	0.008 (0.079)	0.010 (0.076)
Clarithromycin	0.017 (0.011)	0.003 (0.053)	0.002 (0.053)
Clindamycin	0.010 (0.021)	0.001 (0.039)	0.002 (0.039)
Moxifloxacin	0.019 (0.020)	0.028 (0.054)	0.029 (0.048)
<i>M. tuberculosis</i>			
Ethambutol	0.109 (0.179)	0.002 (0.215)	0.015 (0.221)
Isoniazid	0.014 (0.021)	0.004 (0.117)	0.011 (0.119)
Pyrazinamide	0.114 (0.318)	0.089 (0.382)	0.141 (0.382)
Rifampicin	0.013 (0.031)	0.000 (0.200)	0.001 (0.204)
Streptomycin	0.033 (0.050)	0.009 (0.143)	0.011 (0.148)
<i>P. aeruginosa</i>			
Amikacin	0.074 (0.175)	0.003 (0.184)	0.003 (0.179)
Doripenem	0.207 (0.270)	0.066 (0.288)	0.049 (0.281)
Levofloxacin	0.058 (0.072)	0.032 (0.221)	0.055 (0.225)
Meropenem	0.226 (0.267)	0.052 (0.329)	0.047 (0.331)
<i>S. pneumoniae</i>			
Benzylpenicillin	0.008 (0.013)	0.000 (0.015)	0.000 (0.015)
Erythromycin	0.021 (0.037)	0.004 (0.046)	0.004 (0.047)
Tetracycline	0.010 (0.031)	0.004 (0.039)	0.007 (0.037)

¹ The testing set error rate rates are shown in parentheses.

² SCM: the training and testing set error rates are generally similar, indicating that the SVM does not overfit.

³ LinSVM and PolySVM: the training set error rates are generally much smaller than the testing set error rates. This indicates that the models are overfit to the training data.

Table S4: Average and standard deviation of the sensitivities measured on the testing set for 10 random partitions of the data. Results are shown for the set covering machine (SCM), the CART algorithm (CART), L_1 and L_2 regularized support vector machines (L1SVM, L2SVM), support vector machines with linear and polynomial kernels (LinSVM, PolySVM) and the baseline, which predicts the most abundant class in the training set. The prefix χ^2 indicates that univariate filtering was used.

Dataset	SCM	$\chi^2 + \text{CART}$	$\chi^2 + \text{L1SVM}$	$\chi^2 + \text{L2SVM}$	LinSVM	PolySVM	Baseline
<i>C. difficile</i>							
Azithromycin	0.967 \pm 0.046	0.863 \pm 0.068	0.923 \pm 0.054	0.930 \pm 0.055	0.939 \pm 0.056	0.937 \pm 0.051	0.000 \pm 0.000
Ceftriaxone	0.940 \pm 0.038	0.916 \pm 0.051	0.931 \pm 0.048	0.936 \pm 0.050	0.956 \pm 0.018	0.956 \pm 0.018	1.000 \pm 0.000
Clarithromycin	0.990 \pm 0.009	0.881 \pm 0.047	0.907 \pm 0.058	0.918 \pm 0.047	0.936 \pm 0.029	0.930 \pm 0.031	0.000 \pm 0.000
Clindamycin	0.872 \pm 0.089	0.954 \pm 0.042	0.966 \pm 0.036	0.886 \pm 0.073	0.850 \pm 0.085	0.850 \pm 0.085	0.000 \pm 0.000
Moxifloxacin	0.952 \pm 0.023	0.952 \pm 0.023	0.952 \pm 0.023	0.920 \pm 0.034	0.921 \pm 0.031	0.918 \pm 0.031	0.000 \pm 0.000
<i>M. tuberculosis</i>							
Ethambutol	0.650 \pm 0.200	0.661 \pm 0.133	0.779 \pm 0.165	0.476 \pm 0.106	0.572 \pm 0.101	0.560 \pm 0.108	0.000 \pm 0.000
Isoniazid	0.974 \pm 0.018	0.974 \pm 0.018	0.969 \pm 0.025	0.856 \pm 0.064	0.829 \pm 0.060	0.832 \pm 0.071	1.000 \pm 0.000
Pyrazinamide	0.526 \pm 0.308	0.536 \pm 0.254	0.465 \pm 0.305	0.388 \pm 0.234	0.481 \pm 0.251	0.513 \pm 0.254	0.000 \pm 0.000
Rifampicin	0.962 \pm 0.028	0.957 \pm 0.038	0.954 \pm 0.043	0.663 \pm 0.112	0.732 \pm 0.080	0.725 \pm 0.088	0.000 \pm 0.000
Streptomycin	0.965 \pm 0.052	0.966 \pm 0.033	0.960 \pm 0.050	0.844 \pm 0.068	0.826 \pm 0.100	0.822 \pm 0.102	1.000 \pm 0.000
<i>P. aeruginosa</i>							
Amikacin	0.467 \pm 0.165	0.544 \pm 0.090	0.517 \pm 0.116	0.579 \pm 0.092	0.495 \pm 0.060	0.486 \pm 0.064	0.000 \pm 0.000
Doripenem	0.614 \pm 0.161	0.615 \pm 0.081	0.644 \pm 0.061	0.581 \pm 0.078	0.536 \pm 0.116	0.483 \pm 0.140	0.000 \pm 0.000
Levofloxacin	0.890 \pm 0.033	0.879 \pm 0.030	0.859 \pm 0.040	0.738 \pm 0.038	0.727 \pm 0.037	0.704 \pm 0.067	0.000 \pm 0.000
Meropenem	0.633 \pm 0.063	0.663 \pm 0.061	0.575 \pm 0.061	0.568 \pm 0.057	0.509 \pm 0.090	0.517 \pm 0.073	0.000 \pm 0.000
<i>S. pneumoniae</i>							
Benzylpenicillin	0.873 \pm 0.116	0.922 \pm 0.077	0.927 \pm 0.063	0.836 \pm 0.119	0.806 \pm 0.156	0.806 \pm 0.156	0.000 \pm 0.000
Erythromycin	0.891 \pm 0.034	0.806 \pm 0.075	0.879 \pm 0.046	0.885 \pm 0.035	0.834 \pm 0.049	0.824 \pm 0.060	0.000 \pm 0.000
Tetracycline	0.875 \pm 0.096	0.871 \pm 0.078	0.875 \pm 0.096	0.829 \pm 0.112	0.779 \pm 0.115	0.796 \pm 0.109	0.000 \pm 0.000

Table S5: Average and standard deviation of the specificities measured on the testing set for 10 random partitions of the data. Results are shown for the set covering machine (SCM), the CART algorithm (CART), L_1 and L_2 regularized support vector machines (L1SVM, L2SVM), support vector machines with linear and polynomial kernels (LinSVM, PolySVM) and the baseline, which predicts the most abundant class in the training set. The prefix χ^2 indicates that univariate filtering was used.

Dataset	SCM	$\chi^2 + \text{CART}$	$\chi^2 + \text{L1SVM}$	$\chi^2 + \text{L2SVM}$	LinSVM	PolySVM	Baseline
<i>C. difficile</i>							
Azithromycin	0.973 \pm 0.021	0.954 \pm 0.032	0.947 \pm 0.029	0.955 \pm 0.029	0.959 \pm 0.022	0.964 \pm 0.025	1.000 \pm 0.000
Ceftriaxone	0.904 \pm 0.093	0.812 \pm 0.127	0.874 \pm 0.077	0.816 \pm 0.060	0.845 \pm 0.067	0.855 \pm 0.066	0.000 \pm 0.000
Clarithromycin	0.989 \pm 0.012	0.969 \pm 0.031	0.963 \pm 0.022	0.959 \pm 0.019	0.956 \pm 0.019	0.961 \pm 0.019	1.000 \pm 0.000
Clindamycin	0.995 \pm 0.009	0.994 \pm 0.012	0.994 \pm 0.015	0.993 \pm 0.009	0.978 \pm 0.017	0.978 \pm 0.017	1.000 \pm 0.000
Moxifloxacin	0.998 \pm 0.004	0.998 \pm 0.004	0.998 \pm 0.004	0.971 \pm 0.023	0.961 \pm 0.026	0.973 \pm 0.025	1.000 \pm 0.000
<i>M. tuberculosis</i>							
Ethambutol	0.918 \pm 0.081	0.901 \pm 0.098	0.880 \pm 0.069	0.951 \pm 0.060	0.909 \pm 0.052	0.906 \pm 0.054	1.000 \pm 0.000
Isoniazid	0.985 \pm 0.024	0.984 \pm 0.035	1.000 \pm 0.000	0.896 \pm 0.062	0.957 \pm 0.040	0.953 \pm 0.044	0.000 \pm 0.000
Pyrazinamide	0.776 \pm 0.150	0.704 \pm 0.134	0.775 \pm 0.142	0.818 \pm 0.140	0.706 \pm 0.267	0.694 \pm 0.267	1.000 \pm 0.000
Rifampicin	0.977 \pm 0.026	0.980 \pm 0.026	0.985 \pm 0.025	0.925 \pm 0.067	0.862 \pm 0.049	0.858 \pm 0.050	1.000 \pm 0.000
Streptomycin	0.933 \pm 0.048	0.928 \pm 0.054	0.952 \pm 0.048	0.881 \pm 0.076	0.892 \pm 0.078	0.887 \pm 0.082	0.000 \pm 0.000
<i>P. aeruginosa</i>							
Amikacin	0.928 \pm 0.048	0.866 \pm 0.056	0.900 \pm 0.032	0.911 \pm 0.041	0.906 \pm 0.038	0.916 \pm 0.034	1.000 \pm 0.000
Doripenem	0.800 \pm 0.069	0.809 \pm 0.050	0.793 \pm 0.033	0.806 \pm 0.035	0.813 \pm 0.047	0.847 \pm 0.054	1.000 \pm 0.000
Levofloxacin	0.960 \pm 0.030	0.963 \pm 0.028	0.963 \pm 0.033	0.829 \pm 0.027	0.822 \pm 0.037	0.835 \pm 0.042	1.000 \pm 0.000
Meropenem	0.800 \pm 0.041	0.790 \pm 0.044	0.740 \pm 0.068	0.746 \pm 0.064	0.783 \pm 0.059	0.775 \pm 0.053	1.000 \pm 0.000
<i>S. pneumoniae</i>							
Benzylpenicillin	0.997 \pm 0.005	0.994 \pm 0.010	0.994 \pm 0.006	1.000 \pm 0.000	1.000 \pm 0.000	1.000 \pm 0.000	1.000 \pm 0.000
Erythromycin	0.976 \pm 0.010	0.978 \pm 0.010	0.972 \pm 0.009	0.971 \pm 0.013	0.974 \pm 0.011	0.974 \pm 0.011	1.000 \pm 0.000
Tetracycline	0.982 \pm 0.010	0.982 \pm 0.010	0.980 \pm 0.010	0.980 \pm 0.010	0.983 \pm 0.012	0.984 \pm 0.010	1.000 \pm 0.000

Table S6: Average and standard deviation of the error rates measured on the testing set for 10 random partitions of the data. Results are shown for the set covering machine (SCM), the CART algorithm (CART), L_1 and L_2 regularized support vector machines (L1SVM, L2SVM), support vector machines with linear and polynomial kernels (LinSVM, PolySVM) and the baseline, which predicts the most abundant class in the training set. The prefix χ^2 indicates that univariate filtering was used.

Dataset	SCM	$\chi^2 + \text{CART}$	$\chi^2 + \text{L1SVM}$	$\chi^2 + \text{L2SVM}$	LinSVM	PolySVM	Baseline
<i>C. difficile</i>							
Azithromycin	0.030 \pm 0.023	0.086 \pm 0.020	0.064 \pm 0.029	0.056 \pm 0.026	0.050 \pm 0.025	0.048 \pm 0.023	0.446 \pm 0.027
Ceftriaxone	0.073 \pm 0.035	0.117 \pm 0.032	0.087 \pm 0.028	0.102 \pm 0.033	0.079 \pm 0.028	0.076 \pm 0.028	0.306 \pm 0.031
Clarithromycin	0.011 \pm 0.006	0.070 \pm 0.011	0.062 \pm 0.020	0.059 \pm 0.020	0.053 \pm 0.013	0.053 \pm 0.013	0.446 \pm 0.027
Clindamycin	0.021 \pm 0.010	0.011 \pm 0.008	0.009 \pm 0.011	0.021 \pm 0.013	0.039 \pm 0.010	0.039 \pm 0.010	0.136 \pm 0.026
Moxifloxacin	0.020 \pm 0.008	0.020 \pm 0.008	0.020 \pm 0.008	0.048 \pm 0.013	0.054 \pm 0.019	0.048 \pm 0.015	0.390 \pm 0.024
<i>M. tuberculosis</i>							
Ethambutol	0.179 \pm 0.056	0.185 \pm 0.034	0.153 \pm 0.045	0.221 \pm 0.057	0.215 \pm 0.043	0.221 \pm 0.045	0.351 \pm 0.071
Isoniazid	0.021 \pm 0.016	0.021 \pm 0.013	0.017 \pm 0.012	0.125 \pm 0.038	0.117 \pm 0.033	0.119 \pm 0.032	0.421 \pm 0.064
Pyrazinamide	0.318 \pm 0.104	0.371 \pm 0.073	0.353 \pm 0.066	0.342 \pm 0.067	0.382 \pm 0.108	0.382 \pm 0.105	0.347 \pm 0.065
Rifampicin	0.031 \pm 0.017	0.031 \pm 0.021	0.031 \pm 0.021	0.196 \pm 0.052	0.200 \pm 0.033	0.204 \pm 0.033	0.452 \pm 0.070
Streptomycin	0.050 \pm 0.045	0.052 \pm 0.029	0.043 \pm 0.046	0.137 \pm 0.040	0.143 \pm 0.050	0.148 \pm 0.054	0.435 \pm 0.049
<i>P. aeruginosa</i>							
Amikacin	0.175 \pm 0.032	0.206 \pm 0.042	0.187 \pm 0.030	0.164 \pm 0.030	0.184 \pm 0.017	0.179 \pm 0.018	0.216 \pm 0.033
Doripenem	0.270 \pm 0.038	0.261 \pm 0.039	0.261 \pm 0.028	0.275 \pm 0.034	0.288 \pm 0.028	0.281 \pm 0.022	0.359 \pm 0.026
Levofloxacin	0.072 \pm 0.017	0.076 \pm 0.016	0.085 \pm 0.019	0.212 \pm 0.019	0.221 \pm 0.016	0.225 \pm 0.017	0.463 \pm 0.028
Meropenem	0.267 \pm 0.014	0.261 \pm 0.022	0.328 \pm 0.032	0.327 \pm 0.042	0.329 \pm 0.026	0.331 \pm 0.022	0.404 \pm 0.024
<i>S. pneumoniae</i>							
Benzylpenicillin	0.013 \pm 0.009	0.012 \pm 0.010	0.011 \pm 0.009	0.013 \pm 0.011	0.015 \pm 0.013	0.015 \pm 0.013	0.073 \pm 0.016
Erythromycin	0.037 \pm 0.009	0.047 \pm 0.010	0.041 \pm 0.011	0.042 \pm 0.009	0.046 \pm 0.008	0.047 \pm 0.009	0.142 \pm 0.023
Tetracycline	0.031 \pm 0.015	0.029 \pm 0.008	0.032 \pm 0.018	0.037 \pm 0.012	0.039 \pm 0.011	0.037 \pm 0.012	0.106 \pm 0.033