

Random forest model for multicategory death-cause of prostate cancer

Table S1. Baseline characteristics of the training and validation sets

	Training set (n=25,000)	Testing set (n=24,864)	Total (n=49,864)	P value
Survival time (mo)	117 (16-154)	117 (16-154)	117 (16-154)	0.183
Cause of death				0.663
Alive	14865	14746	29611	
(%)	(59.5)	(59.3)	(59.4)	
CVD	2759	2689	5448	
(%)	(11.0)	(10.8)	(10.9)	
Infection	346	371	717	
(%)	(1.4)	(1.5)	(1.4)	
Non-Prostate cancer	1808	1873	3681	
(%)	(7.2)	(7.5)	(7.4)	
Other cause	2903	2897	5800	
(%)	(11.6)	(11.7)	(11.6)	
Prostate cancer	2319	2288	4607	
(%)	(9.3)	(9.2)	(9.2)	
Age at diagnosis, quartile (yr)				0.539
<61	6665	6502	13167	
(%)	(26.7)	(26.2)	(26.4)	
61-67	6303	6353	12656	
(%)	(25.2)	(25.6)	(25.4)	
67-74	6040	6043	12083	
(%)	(24.2)	(24.3)	(24.2)	
75+	5992	5966	11958	
(%)	(24.0)	(24.0)	(24.0)	
Race				0.317
API	1184	1157	2341	
(%)	(4.7)	(4.7)	(4.7)	
Hispanic	2098	2213	4311	
(%)	(8.4)	(8.9)	(8.7)	
NH Black	3500	3510	7010	
(%)	(14.0)	(14.1)	(14.1)	
NH White	17859	17641	35500	
(%)	(71.4)	(71.0)	(71.2)	
Unknown/Other	359	343	702	
(%)	(1.4)	(1.4)	(1.4)	
TNM6 T category				0.850
T1/2	21823	21689	43512	
(%)	(87.3)	(87.2)	(87.3)	
T3/4	2191	2170	4361	
(%)	(8.8)	(8.7)	(8.8)	
Unknown/Other	986	1005	1991	
(%)	(3.9)	(4.0)	(4.0)	
TNM6 N category				0.744
0	22638	22469	45107	
(%)	(90.6)	(90.4)	(90.5)	
1	406	420	826	
(%)	(1.6)	(1.7)	(1.7)	
Unknown/Other	1956	1975	3931	
(%)	(7.8)	(7.9)	(7.9)	

Random forest model for multicategory death-cause of prostate cancer

TNM6 M category				0.997
0	22813	22686	45499	
(%)	(91.3)	(91.2)	(91.3)	
1	1009	1003	2012	
(%)	(4.0)	(4.0)	(4.0)	
Unknown/Other	1178	1175	2353	
(%)	(4.7)	(4.7)	(4.7)	
AJCC6 staging				0.662
1	62	48	110	
(%)	(0.3)	(0.2)	(0.2)	
2	20258	20103	40361	
(%)	(81.0)	(80.9)	(80.9)	
3	1529	1517	3046	
(%)	(6.1)	(6.1)	(6.1)	
4	1452	1480	2932	
(%)	(5.8)	(6.0)	(5.9)	
Unknown/Other	1699	1716	3415	
(%)	(6.8)	(6.9)	(6.9)	
Chemotherapy				0.793
None/Unknown	24820	24680	49500	
(%)	(99.3)	(99.3)	(99.3)	
Received	180	184	364	
(%)	(0.7)	(0.7)	(0.7)	
Radiotherapy				0.047
None/Unknown	15578	15278	30856	
(%)	(62.3)	(61.5)	(61.9)	
Received	9422	9586	19008	
(%)	(37.7)	(38.6)	(38.1)	
Surgery				0.389
Local Excision	1610	1564	3174	
(%)	(6.4)	(6.3)	(6.4)	
No surgery	15323	15386	30709	
(%)	(61.3)	(61.9)	(61.6)	
Prostatectomy	8067	7914	15981	
(%)	(32.3)	(31.8)	(32.1)	
Rural-urban continuum 2003§				0.887
Metro	22084	21974	44058	
(%)	(88.3)	(88.4)	(88.4)	
Non-Metro	2916	2890	5806	
(%)	(11.7)	(11.6)	(11.6)	
Census region				0.315
Midwest	2628	2494	5122	
(%)	(10.5)	(10.0)	(10.3)	
Northeast	4033	3995	8028	
(%)	(16.1)	(16.1)	(16.1)	
South	5081	5053	10134	
(%)	(20.3)	(20.3)	(20.3)	
West	13258	13322	26580	
(%)	(53.0)	(53.6)	(53.3)	
Percent of education attainment, quartile§				0.456

Random forest model for multicategory death-cause of prostate cancer

Q1, <15.08	6450	6296	12746	
(%)	(25.8)	(25.3)	(25.6)	
Q2, 15.09-18.15	6293	6213	12506	
(%)	(25.2)	(25.0)	(25.1)	
Q3, 18.17-25.79	6233	6246	12479	
(%)	(24.9)	(25.1)	(25.0)	
Q4, >50.77	6024	6109	12133	
(%)	(24.1)	(24.6)	(24.3)	
Percent of persons in poverty, quartile§				0.182
Q1, <21.18	6406	6173	12579	
(%)	(25.6)	(24.8)	(25.2)	
Q2, 21.33-29.81	6212	6170	12382	
(%)	(24.9)	(24.8)	(24.8)	
Q3, 29.86-37.36	6375	6448	12823	
(%)	(25.5)	(25.9)	(25.7)	
Q4, >67.40	6007	6073	12080	
(%)	(24.0)	(24.4)	(24.2)	
Percent of foreign-born residents, quartile§				0.223
Q1, <5.95	6259	6080	12339	
(%)	(25.0)	(24.5)	(24.8)	
Q2, 5.98-15.22	6481	6365	12846	
(%)	(25.9)	(25.6)	(25.8)	
Q3, 15.45-21.55	6068	6124	12192	
(%)	(24.3)	(24.6)	(24.5)	
Q4, >38.52	6192	6295	12487	
(%)	(24.8)	(25.3)	(25.0)	
Confirmation method of diagnosis				0.686
Microscopic	24553	24407	48960	
(%)	(98.2)	(98.2)	(98.2)	
Radiologic and clinic	306	302	608	
(%)	(1.2)	(1.2)	(1.2)	
Unknown/Other	141	155	296	
(%)	(0.6)	(0.6)	(0.6)	
PSA, quartiles (ng/ml)				0.854
<4.9	5546	5573	11119	
(%)	(22.2)	(22.4)	(22.3)	
5.0-6.8	5231	5207	10438	
(%)	(20.9)	(20.9)	(20.9)	
6.9-11.3	5179	5131	10310	
(%)	(20.7)	(20.6)	(20.7)	
11.3+	5243	5253	10496	
(%)	(21.0)	(21.1)	(21.1)	
Unknown/Other	3801	3700	7501	
(%)	(15.2)	(14.9)	(15.0)	
Gleason score				0.957
5	<15*	<15*	15	
(%)			(0.0)	
6	192	182	374	
(%)	(0.8)	(0.7)	(0.8)	

Random forest model for multicategory death-cause of prostate cancer

7	169	166	335
(%)	(0.7)	(0.7)	(0.7)
8	50	44	94
(%)	(0.2)	(0.2)	(0.2)
9	35	30	65
(%)	(0.1)	(0.1)	(0.1)
10	<15*	<15*	<15*
(%)			
Unknown/Other	24541	24432	48973
(%)	(98.2)	(98.3)	(98.2)

Note: AJCC, 6th edition clinical staging of the American Joint Commission on Cancer; API, Asian Pacific Islanders; NH, Non-Hispanic; TNM6, 6th edition Tumor, node and metastasis staging manual of the American Joint Commission on Cancer; CVD, cardiovascular disease; PSA, Prostate specific antigen; *, statistically suppressed; §, Country attributes of Year 2000; Education attainment defined as percent of residents with less than high-school graduate in the county; Person in poverty defined as percent of residents with income below 200% of poverty in the county.

Table S2. Prediction accuracy for long-term 6-category causes of death among the patients with prostate cancer diagnosed in 2004 (12-year follow up) using one-hot encoding

Predicted classes	Alive, n=14,746	CVD, n=2,689	Infection, n=371	Non-Prostate cancer, n=1,873	Other cause, n=2,897	Prostate cancer, n=2,288	Total, n=24,864
Random forest model							
Alive, %	88.87*	56.56	57.68	69.62	59.20	41.78	75.67
CVD, %	3.54	14.95*	13.75	9.24	13.95	8.65	7.04
Infection, %	0.20	0.60	0.27*	0.27	0.62	0.39	0.31
Non-Prostate cancer, %	1.70	2.98	2.43	2.72*	2.93	2.80	2.17
Other cause, %	3.49	15.43	15.09	9.72	13.84*	9.35	7.17
Prostate cancer, %	2.20	9.48	10.78	8.44	9.46	37.02*	7.64
Multinomial model							
Alive, %	84.90*	34.85	32.35	53.66	39.28	27.93	65.79
Prostate cancer, %	0.01	0.00	0.27	0.00	0.00	0.09*	0.02
NA, %	15.09	65.15	67.39	46.34	60.72	71.98	34.19

Note: CVD, cardiovascular disease; NA, not available; *, correct prediction.

Random forest model for multicategory death-cause of prostate cancer

Table S3. Factors associated with long-term 6-category cause of death among men with prostate cancer in multinomial model

Covariate	Cause of death									
	CVD		Infection		Non-Prostate cancer		Other causes		Prostate cancer	
	coefficient (95% CI)	P	coefficient (95% CI)	P	coefficient (95% CI)	P	coefficient (95% CI)	P	coefficient (95% CI)	P
Age at diagnosis, quartile (yr)*	0.92 (0.88 to 0.97)	<0.001	0.91 (0.78 to 1.03)	<0.001	0.62 (0.57 to 0.67)	<0.001	0.84 (0.79 to 0.88)	<0.001	0.51 (0.46 to 0.55)	<0.001
AJCC6 staging*	0.15 (0.04 to 0.27)	0.007	0.34 (0.10 to 0.58)	0.005	0.15 (0.02 to 0.29)	0.021	0.14 (0.02 to 0.25)	0.019	0.68 (0.58 to 0.77)	<0.001
Confirmation method of diagnosis*	0.54 (0.24 to 0.85)	0.001	0.70 (0.16 to 1.24)	0.011	0.29 (-0.09 to 0.68)	0.136	0.69 (0.39 to 0.98)	<0.001	0.66 (0.39 to 0.93)	<0.001
Surgery*	-0.76 (-0.86 to -0.66)	<0.001	-0.84 (-1.07 to -0.61)	<0.001	-0.45 (-0.56 to -0.34)	<0.001	-0.79 (-0.88 to -0.70)	<0.001	-1.04 (-1.14 to -0.94)	<0.001
PSA, quartiles (ng/ml)*	0.17 (0.13 to 0.20)	<0.001	0.16 (0.07 to 0.25)	<0.001	0.13 (0.09 to 0.17)	<0.001	0.14 (0.11 to 0.17)	<0.001	0.34 (0.30 to 0.38)	<0.001
Chemotherapy	0.23 (-0.45 to 0.92)	0.505	-0.11 (-2.11 to 1.88)	0.911	0.73 (0.11 to 1.36)	0.022	0.00 (-0.74 to 0.74)	0.999	2.03 (1.63 to 2.44)	<0.001
Census region	-0.12 (-0.17 to -0.07)	<0.001	-0.15 (-0.27 to -0.03)	0.016	-0.04 (-0.10 to 0.02)	0.160	-0.01 (-0.06 to 0.04)	0.829	0.06 (0.00 to 0.12)	0.048
Percent of education attainment, quartile§	0.08 (0.01 to 0.15)	0.019	0.11 (-0.06 to 0.28)	0.198	0.14 (0.06 to 0.21)	0.001	0.11 (0.04 to 0.17)	0.002	0.08 (0.01 to 0.16)	0.031
Percent of persons in poverty, quartile§	0.08 (0.01 to 0.15)	0.032	0.01 (-0.16 to 0.18)	0.929	-0.06 (-0.14 to 0.02)	0.161	0.01 (-0.06 to 0.08)	0.788	-0.03 (-0.11 to 0.05)	0.436
Percent of foreign-born residents, quartile§	-0.08 (-0.13 to -0.03)	0.003	-0.11 (-0.24 to 0.02)	0.109	-0.10 (-0.16 to -0.04)	0.001	-0.19 (-0.24 to -0.14)	<0.001	-0.13 (-0.18 to -0.07)	<0.001
Gleason score	-0.03 (-0.13 to 0.07)	0.566	-0.07 (-0.30 to 0.16)	0.565	-0.04 (-0.15 to 0.07)	0.454	0.06 (-0.05 to 0.17)	0.257	-0.05 (-0.16 to 0.06)	0.341
Race	0.03 (-0.02 to 0.08)	0.292	-0.23 (-0.34 to -0.11)	<0.001	0.04 (-0.03 to 0.10)	0.252	0.01 (-0.05 to 0.06)	0.796	-0.03 (-0.09 to 0.03)	0.267
Radiotherapy	-0.22 (-0.31 to -0.12)	<0.001	-0.40 (-0.63 to -0.16)	0.001	-0.07 (-0.18 to 0.04)	0.238	-0.22 (-0.32 to -0.13)	<0.001	-0.37 (-0.48 to -0.26)	<0.001
Rural-urban continuum 2003§	-0.03 (-0.18 to 0.13)	0.720	-0.23 (-0.64 to 0.17)	0.261	0.11 (-0.07 to 0.28)	0.239	-0.03 (-0.17 to 0.12)	0.739	-0.01 (-0.18 to 0.17)	0.948
TNM6 T category	0.08 (-0.05 to 0.21)	0.242	0.15 (-0.15 to 0.45)	0.335	0.19 (0.03 to 0.35)	0.017	0.03 (-0.11 to 0.16)	0.714	0.50 (0.39 to 0.61)	<0.001
TNM6 N category	-0.05 (-0.13 to 0.02)	0.190	-0.08 (-0.24 to 0.09)	0.365	-0.08 (-0.17 to 0.01)	0.065	-0.05 (-0.12 to 0.03)	0.218	-0.18 (-0.24 to -0.12)	<0.001
TNM6 M category	0.08 (-0.08 to 0.24)	0.322	-0.34 (-0.68 to 0.00)	0.053	0.06 (-0.13 to 0.25)	0.511	0.06 (-0.10 to 0.22)	0.447	0.06 (-0.07 to 0.19)	0.357

Note: AJCC, 6th edition clinical staging of the American Joint Commission on Cancer; TNM6, 6th edition Tumor, node and metastasis staging manual of the American Joint Commission on Cancer; CVD, cardiovascular disease; PSA, Prostate specific antigen; *, factors linked to all causes of death; §, Country attributes of Year 2000; Education attainment defined as percent of residents with less than high-school graduate in the county; Person in poverty defined as percent of residents with income below 200% of poverty in the county.