

JAMA Psychiatry | [Original Investigation](#)

Improving Prediction of Suicide and Accidental Death After Discharge From General Hospitals With Natural Language Processing

Thomas H. McCoy Jr, MD; Victor M. Castro, MS; Ashlee M. Roberson, BA;
Leslie A. Snapper, BS; Roy H. Perlis, MD, MS

Importance

Suicide represents the 10th leading cause of death across age groups in the US (12.6 cases per 100,000) and remains challenging to predict.

While many individuals who die by suicide are seen by physicians before their attempt, they may not seek psychiatric care.

Objective

To determine the extent to which incorporating natural language processing of narrative discharge notes improves stratification of risk for death by suicide after medical or surgical hospital discharge.

Design, setting, & participants

In this retrospective health care use study, clinical data were analyzed from individuals with discharges from MGH and Brigham between January 1, 2005, and December 31, 2013.

Main Outcome & Measures

The primary outcome was suicide as a reported cause of death based on Massachusetts Department of Public Health records.

Regression models for prediction of death by suicide or accidental death were compared relying solely on coded clinical data and those using natural language processing of hospital discharge notes.

Main Outcome & Measures

The primary outcome was suicide as a reported cause of death based on Massachusetts Department of Public Health records.

Regression models for prediction of death by suicide or accidental death were compared relying solely on coded clinical data and those using natural language processing of hospital discharge notes.

Scoring narrative valence

We used a curated list of ~3000 subjectively valence-conveying terms to score each narrative note using Pattern, an open-source opinion mining tool.

For example, terms with positive valence include *glad, pleasant, lovely*. Those with negative valence include *gloomy, unfortunate, sad*.

[home](#)[news](#)[projects](#)[people](#)[demos](#)[resources](#)[publications](#)[talks](#)[contact](#)

[Home](#) ›

Pattern

Pattern is a web mining module for the Python programming language.

It has tools for data mining (Google, Twitter and Wikipedia API, a web crawler, a HTML DOM parser), natural language processing (part-of-speech taggers, n-gram search, sentiment analysis, WordNet), machine learning (vector space model, clustering, SVM), network analysis and <canvas> visualization.

The module is free, well-document and bundled with 50+ examples and 350+ unit tests.



USERNAME: *

PASSWORD: *

☐ REMEMBER ME

[Log in](#)

[Request new password](#)



Statistical analysis

The primary analysis used survival methods, with the results censored at the end of follow-up or at death.

After confirming that proportional hazards assumptions were met for the primary model, Cox regression was used to examine risk associated with predictors individually and in aggregate.

Statistical analysis

In light of the scale of data, we randomly selected a single hospitalization from individuals with multiple hospitalizations to maximize computability, while minimizing bias.

Table 1. Description of Hospitalized Cohort, Contrasting Individuals Who Did or Did Not Die by Suicide During Follow-up

Variable	No. (%)		
	Overall (N = 458 053)	Death by Suicide (n = 235)	No Suicide Death (n = 457 818)
Male sex	189 529 (41.4)	155 (66)	189 374 (41.3)
White race	346 340 (75.7)	202 (86)	346 138 (75.6)
Primary psychiatric diagnosis at admission	5782 (1.3)	16 (6.8)	5766 (1.3)

Table 3. Description of Hospitalized Cohort, Contrasting Individuals Who Did or Did Not Die by Suicide or Accidental Death (Composite Outcome) During Follow-up

Variable	No. (%)		
	Overall (N = 458 053)	Death by Suicide or Accidental Death (n = 2026)	No Suicide or Accidental Death (n = 456 027)
Male sex	189 529 (41.6)	1254 (61.9)	188 275 (41.1)
White race	346 340 (75.9)	1661 (82)	344 679 (75.2)
Primary psychiatric diagnosis at admission	5782 (1.3)	92 (4.5)	5690 (1.2)

Two regression models:

1. Coded clinical data: age, sex, self-reported race, recent health care use (including outpatient psychiatric visits, overall outpatient visits, and ED visits), and overall medical morbidity.
2. Included all features from 1st model + aggregate measures of pos and neg valence.

Model characterization

Standard measures of discrimination, including the **C statistic** and **calibration**.

The C statistic was calculated with **10-fold cross-validation** to minimize optimism.

Continuous net reclassification improvement was calculated as a complementary measure of improvement in model fit.

“When outcomes are **binary**, the **c-statistic** (equivalent to the area under the Receiver Operating Characteristic curve) is a standard measure of the predictive accuracy of a logistic regression model.”



Epidemiology & Biostatistics

[Directory](#)[Administrative Resources](#) +[Biostatistics](#) +[Computing Resource](#) +[Epidemiology](#) +[Health Outcomes](#) -[Overview](#)[Projects](#)[Members](#)

Decision Curve Analysis



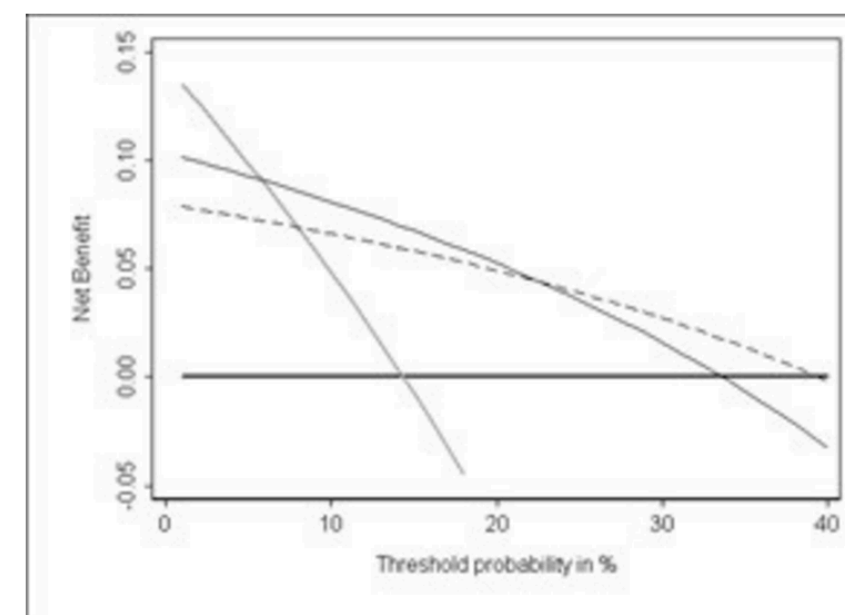
Decision curve analysis is a simple method for evaluating prediction models, diagnostic tests, and molecular markers.

The method was first published as:

Vickers AJ, Elkin EB. Decision curve analysis: a novel method for evaluating prediction models. Medical Decision Making. 2006 Nov-Dec;26(6):565-74.

A subsequent discussion paper gives some further details about the method:

Steyerberg EW, Vickers AJ. Decision curve analysis: a discussion. Medical Decision Making. 2008 Jan-Feb;28(1):146-9.



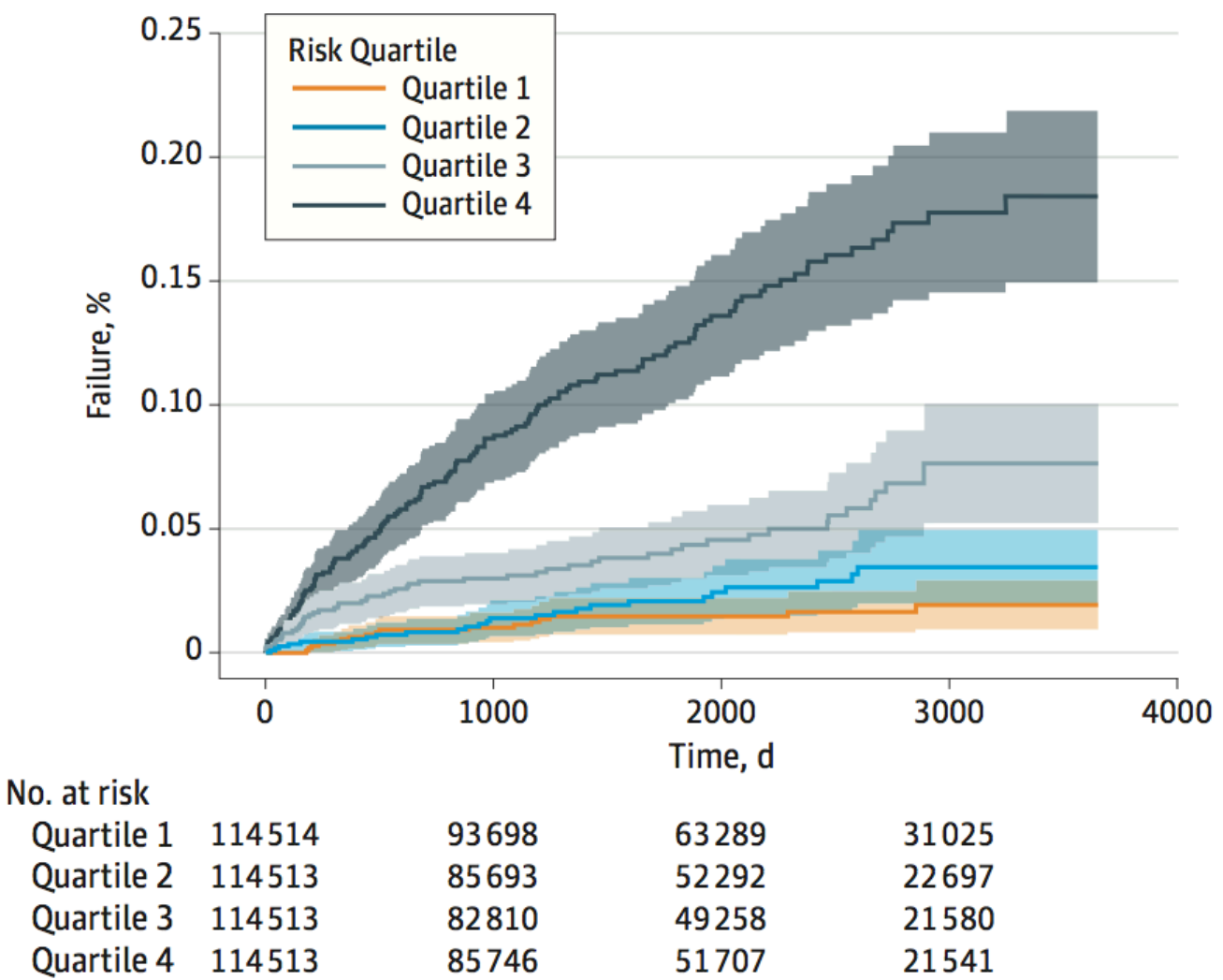
Results

There were 845,417 hospital discharges represented in the cohort, including 458,053 unique individuals.

Overall, all-cause mortality was 18% during 9 years, and the median follow-up was 5.2 years.

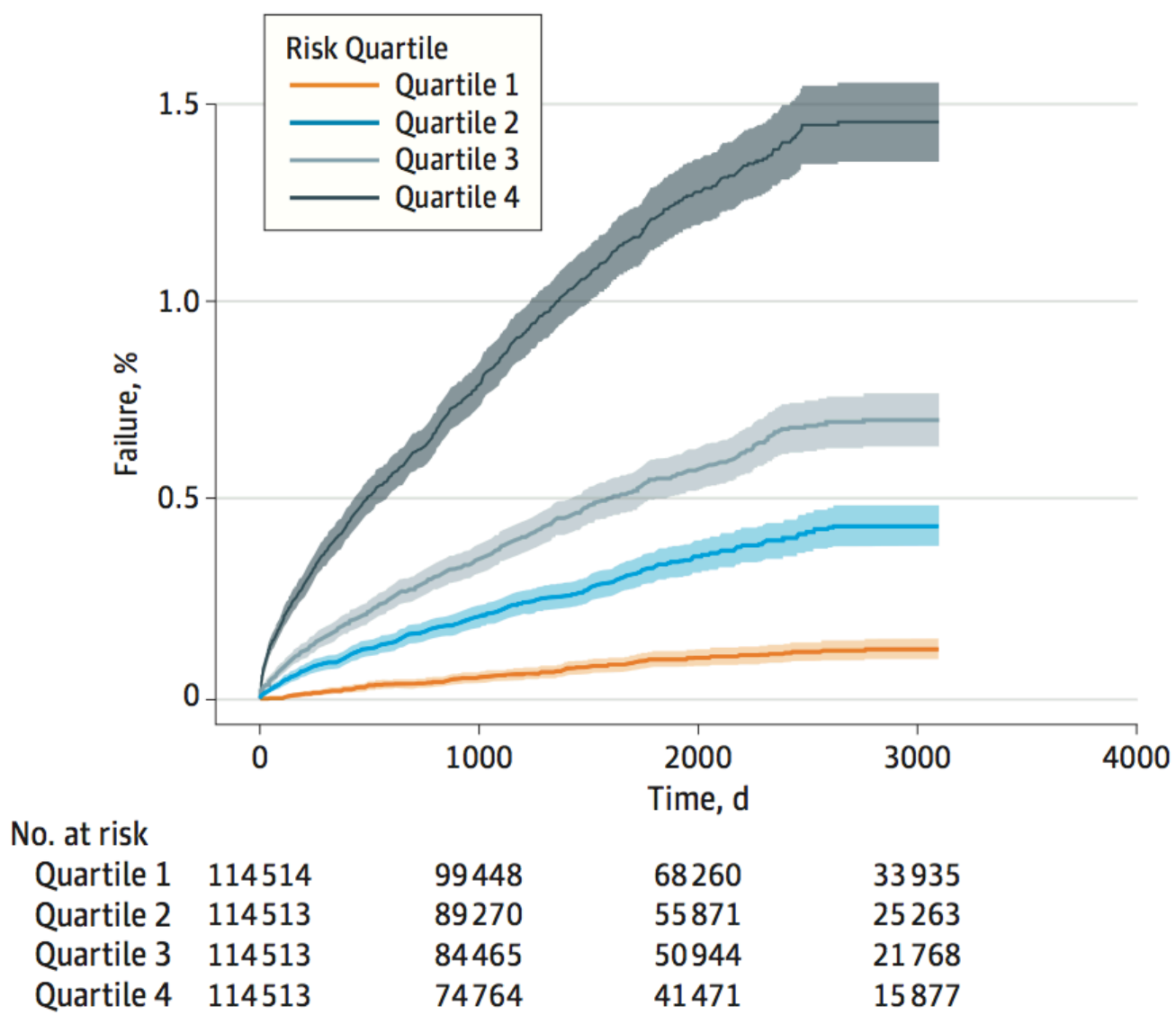
The cohort included 235 (0.1%) who died by suicide during 2.4M patient-years of follow-up.

Figure 1. Kaplan-Meier Curves for Time to Death by Suicide Among 458 053 Individuals With at Least 1 Hospital Discharge by Predicted Risk Quartile



The axes are rescaled inside the figure to improve interpretability.

Figure 2. Kaplan-Meier Curves for Time to Death by Suicide or Accidental Death Among 458 053 Individuals With at Least 1 Hospital Discharge by Predicted Risk Quartile



Features predictive of suicide

White race

Male sex

More ED visits and psychiatric outpatient visits in the 12 months before admission

Table 5. Cox Proportional Hazards Regression Model of Suicide Risk Without (Model 1) and With (Model 2) Inclusion of Natural Language Processing of Narrative Discharge Summary

Variable	Adjusted Hazard Ratio (95% CI) ^a	
	Model 1	Model 2
Male sex	2.80 (2.11-3.71)	2.57 (1.94-3.41)
White race	2.27 (1.56-3.30)	2.23 (1.53-3.24)
Age at admission, y	0.80 (0.65-0.97)	0.79 (0.65-0.96)
Age-adjusted Charlson Comorbidity Index	0.95 (0.74-1.21)	0.89 (0.69-1.15)
Primary psychiatric diagnosis at admission	1.63 (0.93-2.85)	1.39 (0.79-2.45)
Psychiatric visits in 12 mo before admission	1.18 (1.07-1.31)	1.17 (1.06-1.30)
Prior psychiatric visits (ever)	1.15 (1.01-1.31)	1.14 (1.01-1.30)
Outpatient visits (any) in 12 mo before admission	0.79 (0.67-0.94)	0.84 (0.71-1.00)
Emergency department visits in 12 mo before admission	2.37 (1.83-3.06)	2.19 (1.68-2.84)
Positive valence	NA	0.70 (0.58-0.85)
Negative valence	NA	1.15 (0.99-1.32)

Abbreviation: NA, not applicable.

^a Adjusted for all other terms in the model.

Table 6. Cox Proportional Hazards Regression Model of Suicide or Accidental Death (Composite Outcome) Risk Without (Model 1) and With (Model 2) Inclusion of Natural Language Processing of Narrative Discharge Summary

Variable	Adjusted Hazard Ratio (95% CI) ^a	
	Model 1	Model 2
Male sex	1.88 (1.71-2.06)	1.82 (1.66-1.99)
White race	1.37 (1.22-1.54)	1.37 (1.22-1.53)
Age at admission, y	1.21 (1.13-1.30)	1.19 (1.11-1.28)
Age-adjusted Charlson Comorbidity Index	1.39 (1.27-1.53)	1.35 (1.23-1.48)
Primary psychiatric diagnosis at admission	1.84 (1.47-2.31)	1.69 (1.35-2.11)
Psychiatric visits in 12 mo before admission	1.05 (1.01-1.10)	1.05 (1.01-1.10)
Prior psychiatric visits (ever)	1.16 (1.11-1.21)	1.16 (1.10-1.21)
Outpatient visits (any) in 12 mo before admission	0.63 (0.60-0.67)	0.65 (0.62-0.69)
Emergency department visits in 12 mo before admission	2.98 (2.71-3.26)	2.84 (2.58-3.11)
Positive valence	NA	0.80 (0.75-0.85)
Negative valence	NA	1.06 (1.01-1.12)

Abbreviation: NA, not applicable.

^a Adjusted for all other terms in the model.

	C-statistic	95% CI
Coded data	0.737	0.734 - 0.741
+Valence	0.741	0.738 - 0.744

Continuous net reclassification improvement = 0.10

[https://www.wikiwand.com/en/
Net_reclassification_improvement](https://www.wikiwand.com/en/Net_reclassification_improvement)

Notably, 115 of 235 (48.9%) suicide deaths in the present study occurred among individuals with no coded data reflecting psychiatric *International Classification of Diseases, Ninth Revision* diagnostic codes in this health system. This finding is consistent with prior reports that, while individuals who die by suicide often have contact with a health professional, the clinician is likely to not be a psychiatrist or therapist,¹⁴ which underscores the importance of psychiatric expertise in the general hospital setting. We cannot exclude the possibility that some of these individuals had sought treatment in the community (eg, in a private practice or another health system). Even if so, this treatment was not documented in hospital records and so presumably was not known by hospital staff. These omissions may reflect failure to inquire about psychiatric history, a potential target for intervention meriting further study.