# UC Davis

Title

Predicting venous thromboembolism (VTE) risk in cancer patients using machine learning

Permalink

https://escholarship.org/uc/item/74f229fz

Journal

Health Care Science, 2(4)

ISSN

2771-1757

Authors

Townsley, Samir Khan

Basu, Debraj

Vora, Jayneel

et al.

Publication Date

2023-08-01

DOI

10.1002/hcs2.55

Copyright Information

Peer reviewed

1 **Title: Predicting Venous Thromboembolism (VTE) Risk in Cancer Patients Using Machine**
2 **Learning**

3

4

5 Author details:

6 1. Samir Khan Townsley *

7 Department of Electrical and Computer Engineering, University of California, Davis

8 2. Debraj Basu *

9 Department of Electrical and Computer Engineering, University of California, Davis

10 3. Jayneel Vora

11 Department of Computer Science, University of California, Davis

12 4. Ted Wun

13 School of Medicine, University of California, Davis Health

14 5. Chen-Nee Chuah

15 Department of Electrical and Computer Engineering, University of California, Davis

16 6. Corresponding author: Dr. Prabhu RV Shankar, M.D., M.S., MRCP (UK)

17 School of Medicine, University of California, Davis Health

18 Email: rvpshankar@ucdavis.edu

19 Mailing address: 4610 X St, Sacramento, CA 95817, USA

20 Phone number : 01 408-458-6119

21 Fax number :   01 916-734-7055

22

23 *  The research work was conducted during his association with University of California, Davis

# 1. Summary

## Objectives

The association between cancer and venous thromboembolism (VTE) is well-established with cancer patients accounting for approximately twenty percent of all VTE incidents. In this paper, we have:

- Performed a comparison of machine learning (ML) methods to traditional clinical scoring models for predicting the occurrence of VTE in a cancer patient population, and
- Identified important features (clinical biomarkers) for ML model predictions and examined how different approaches to reducing the number of features used in the model impact model performance.

## Methods

We have developed an ML pipeline including three separate feature selection processes and applied it to routine patient care data from the electronic health records (EHR) of 1910 cancer patients at the University of California Davis Medical Center (UCDMC).

## Results

Our ML-based prediction model achieved an area under the receiver operating characteristic (ROC) curve of $0.778 \pm 0.006$ when trained on a set of 15 features. This result is comparable to the model performance when trained on all features in our feature pool ($0.779 \pm 0.006$ with 29 features). Our result surpasses the most validated clinical scoring system for VTE risk assessment in cancer patients by 16.1%. We additionally found cancer stage information to be a useful predictor after all performed feature selection processes despite not being used in existing score-based approaches.

## Conclusion

From these findings, we observe that machine learning can offer new insights and a significant improvement over the most validated clinical VTE risk scoring systems in cancer patients. The results of this study also allowed us to draw insight into our feature pool and identify the features that could have the most utility in the context of developing an efficient machine learning classifier. While a model trained on our entire feature pool of 29 features significantly outperformed the traditionally used clinical scoring system, we were able to achieve an equivalent performance using a subset of only 15 features through strategic feature selection methods. These results are encouraging for potential applications of ML to predicting cancer associated VTE in clinical settings such as in bedside decision support systems where feature availability may be limited.

## Keywords

VTE; Cancer; Binary classification; Machine learning pipeline

## 2. Introduction

Venous thromboembolism (VTE) comprises both deep-vein thrombosis (DVT) and pulmonary embolism (PE) [1]. The association between VTE and cancer is well-established with cancer patients accounting for approximately twenty percent of all VTE incidents [2]. While the estimated prevalence of VTE in the general population is around 1 in 1000 [3, 4], some estimates suggest this number increases 5-fold within the cancer patient population [1, 5, 6]. The risk increases further among patients who receive chemotherapy as shown in a 15-year population-based study [7].

VTE is a multifaceted risk in cancer patients that exacerbates clinical consequences, significantly impacting morbidity, mortality, and cost of patient care [1, 5, 8, 9, 10, 11]. Specifically, VTE associated mortality is 2.2 times more likely in VTE patients with cancer than in those without [10]. VTE is the leading cause of mortality in cancer patients, aside from mortality due to cancer itself [1, 8]. In addition to increasing risk of mortality, VTE burdens the cancer treatment process. When managing VTE in cancer patients, use of anticoagulants, which thin the blood, requires rigorous patient monitoring in order to achieve adequate anticoagulation and to identify complications such as bleeding. Compared to cancer patients without VTE, patients with VTE have over two times the risk of experiencing major bleeding [12]. Bleeding can worsen anemia while reduced blood counts can delay cancer interventions such as chemotherapy and radiotherapy and increase the need for blood transfusions.

The recurrence rates of VTE are also high in patients with cancer. Patients with an active malignancy have a 3-4 fold higher risk of recurrence compared to patients without cancer, and the risk is further increased in those with metastatic cancers. According to one study, the one-year cumulative risk for recurrent VTEs after the first episode was 21% in cancer patients compared to 7% in patients without cancer [12]. All the VTE related factors discussed

91  above can affect cancer management, increase treatment costs, and escalate average price per

92  hospitalization for cancer patients [2, 3, 12, 13].

93  Treatments such as anticoagulant therapy are available, both for prophylaxis against

94  occurrence, as well as for treatment of VTE in cancer patients. Appropriate and timely use of

95  the prophylactic measures are vital for reducing the risk of both fatal and non-fatal pulmonary

96  embolism as well as the post-thrombotic complications [14]. Anticoagulants are drugs that

97  interfere with blood coagulation cascade to reduce or inhibit blood clotting. The low-

98  molecular-weight heparin (LMWH) has been found in multiple studies to reduce the likelihood

99  of a VTE event occurring in a cancer patient [2, 15, 16, 17]. With these issues in mind, it is evident

100  that effective VTE prophylaxis in cancer patients has the potential to drastically improve

101  cancer survival rates and decrease treatment costs for hospitals and patients alike. However,

102  while anticoagulant prophylaxis and treatment is effective in primary and secondary prevention

103  of VTE, as mentioned above, there are certain implications with their regular use in all cancer

104  patients. In particular, anticoagulants are associated with increased bleeding, require parenteral

105  administration, training, and additional monitoring, all of which can increase both cost and

106  complexity of cancer patient management [2, 12 ,18]. Therefore, it is important to stratify and

107  define high risk cohorts of cancer patients who are prone for VTE. There is thus a need for

108  effective VTE risk stratification systems to ensure that prophylaxis is administered only to

109  high-risk patients. An accurate, reliable, and robust VTE stratification system would help

110  clinicians in decision making about anticoagulant therapy at the point-of-care (POC).

111  Prophylactic measures against VTE are often implemented for hospitalized patients, so high

112  risk stratification is particularly important in ambulatory patients (outpatients) as they cannot

113  be monitored as closely as hospitalized patients.

The importance of delineating which cancer patients are at increased risk of VTE for instituting anticoagulation prophylaxis, particularly ambulatory patients, is critical as anticoagulation is associated with significant risks and costs in already debilitated cancer patients. Decision to provide prophylactic anticoagulation in ambulatory patients clinically alone is often difficult and providers need a decision support tool that pinpoints the most vulnerable groups for VTE. Several Cancer Associated Thromboembolism (CAT) prediction scores have been developed, such as Khorana [19], Vienna CATS [20,] PROTECHT [21] and CONKO [22] based on routinely collected patient care data. These risk-assessment methods all use a simple scoring system where points are added based on each of five to eight different predictors with higher scores indicating a higher risk of developing VTE. Some of the predictors that these scores use include cancer site, platelet count, white blood cell count, hemoglobin, use of red blood cell stimulating factors, and BMI. Of these scores, the Khorana score is the most validated and used [23]. However, despite its acceptance in the research community, the Khorana score still only achieves a positive predictive value of 6.7 %, which is not meaningful enough to make a quantified decision by the clinicians and thus leaves plenty of room for improvement [19]. In another study of 218 patients with cancer initiating chemotherapy, it is shown that the Khorana score was able to stratify ambulatory cancer patients according to the risk of VTE, but not for all cancer types [24]. The Khorana score can be used to select ambulatory cancer patients at high risk of venous thromboembolism for thromboprophylaxis, but most events occur outside this high-risk group [25].

During informal discussions, clinicians opined that, even a positive predictive value of 20 to 30% will help them with decision making, tipping the decision one way or other with some scientific qualitative basis, and those discussions motivated the team to explore various

137  features (clinical biomarkers) and develop more robust and clinically meaningful predictive

138  models.

139  In this study we use machine learning to take a data-driven approach to VTE prediction

140  in cancer patients. Our aim in this study is to not only improve upon the performance of known

141  risk assessment scores such as the Khorana score but also to perform an in depth, data-driven

142  exploration of both new and known VTE risk factors.

143  Traditional approaches to prediction in medicine often focus on capturing medical

144  expertise through a set of carefully designated rules [26]. However, data driven approaches, such

145  as machine learning algorithms instead can learn effective prediction decisions by observing

146  numerical patterns in the input data [26, 27]. One subset of machine learning, known as supervised

147  learning, deals with training a model to accomplish this task of classifying data based on a set

148  of input data with labeled ground truth values [27]. Supervised learning has the advantage over

149  traditional rule-based methods of being able to leverage computational power to identify highly

150  convoluted patterns in massive datasets with large numbers of potential predictors relatively

151  quickly and efficiently [26, 28]. Such an approach has promise in the context of cancer patient

152  VTE prediction, where the currently accepted scoring systems are simple rule-based methods

153  that do not necessarily capture a wide range of the potentially complex interactions between

154  variables [19, 20, 22]. Patrizia Ferroni et al. have designed a precision medicine approach to exploit

155  significant patterns in data to produce VTE risk predictors for cancer outpatients [29]. They have

156  used Multiple kernel learning (MKL) [30] based on support vector machines (SVM) models to

157  predict VTE risk. In our research, we have examined VTE classification performances of

158  several standard ML algorithms including SVM, logistic regression (LR), and Random Forest

159  (RF) and compared these to the baseline performance of the Khorana score.

160          Methods and results are described in the following sections.

## 3. Methods

163          In this retrospective study of a population of cancer patients at the University of

164     California Davis Medical Center, we used machine learning to explore both new and known

165     VTE risk factors. Our goal was to not only develop a machine-learning based VTE risk

166     assessment system for cancer patients but also to examine which risk factors may be useful

167     when taking such an approach. From our efforts, we hope to establish a foundation for using

168     machine learning to eventually answer more complex questions about VTE prediction in

169     cancer patients, such as how changes in a patient's condition, as the patient continues with

170     his/her cancer management, affect the risk of developing VTE over time.

171          In this study, we examined 29 features in total, including a selection of available

172     features from the Khorana score and biomolecular markers from a previous study of CAT [19,

173     29]. Since relevant VTE events can occur before or after cancer diagnosis and clinical

174     interventions (i.e., surgery, chemotherapy, radiotherapy), we used a set of time-agnostic

175     features to gain a view of how a patient's general profile over a large period of time may or

176     may not be indicative of VTE risk. Each of the features we used covered information about a

177     patient's background, cancer, lab values, or medications.

178          We then explored the utility of our feature set in a machine learning context in a two-

179     phased approach. In the first phase we trained several different models with a spectrum of

180     hyperparameter choices on four different feature subsets that were derived both from

181     performed feature selection experiments and from pre-determined feature pools. We then

182     identified the best performing model and feature set combination and, in a second phase of

183     experiments, attempted to reduce the number of used features without sacrificing performance

184    through an iterative feature accumulation process. Finally, we validated the performance of

185    our chosen model on a held-out dataset extracted from our original data.

186
187    **3.1. Dataset and Data Prepossessing**
188
189        The dataset used in these experiments was extracted from the UCDMC affiliated

190    hospital's EHR system and combined with curtained and manually curated data elements from

191    the California state cancer network CNExT registry, from 2015-2017 (C/NET Solutions,

192    Berkeley, CA).  The organ system-based cancers which are considered high risk for VTE

193    episodes in previous studies were included in the study and are listed in Table 1 [32].

194
195    Table 1. Cancer sites contained in the dataset

| Site group |
| --- |
| Pancreas |
| Bladder |
| Non-Hodgkin's Lymphoma |
| Hodgkin's Disease |
| Corpus Uteri/Uterus |
| Prostate |
| Ovary |
| Breast |
| Lung/Bronchus (Small Cell and Non-Small Cell) |
| Brain |

| Stomach |
| --- |

In order to study how a given cancer and its attributes may be predictive of VTE events, each cancer instance was treated as a separate entry in our dataset. Thus, a few patients have more than one cancer entry in the dataset. Associated with each cancer instance is a list of features describing the cancer and patient background.

All medications were grouped according to the pharmacologic class of the medication. Medication data was incorporated in the primary cancer entry cohort by assigning a binary variable to each patient for every medication, indicating whether or not that medication was ever administered to the patient.

Lab test values were represented by the mean of all pre-chemotherapy measurements associated with that test in order to eliminate noise and understand how a patient's general condition correlates with VTE risk. We accumulated such values for 45 different lab tests. This set of 45 was then reduced to only the lab tests which were performed on at least 75 % of patients. Of the 45 lab tests, only 12 of the tests satisfied this criterion and were included in our final feature pool. Any missing values among these 12 lab tests were imputed using the mean across all patients for the given test.

Exclusion criteria for our dataset included patients with missing information in any of the listed categories outside of lab tests, patients with benign tumors, patients with mesotheliomas, and patients with extreme outliers (i.e., BMI > 100). These exclusion criteria were applied to the general dataset. After cleaning, the dataset consisted of 1973 cancer entries across 1910 unique patients.

The presence or absence of a VTE diagnosis date served as our binary target variable for prediction in our machine learning models. The full list of features in our curated dataset is detailed in Table 2.

Table 2: Feature pool

| Feature Type | Features (29) |
|---|---|
| Cancer | site, grade, stage, behavior, histopathological type |
| Patient | gender, body mass index (BMI), age, race list, race count |
| Binary Medications | antineoplastic - aromatase inhibitors, immunosuppressives, antineoplastic - antiandrogenic agents, steroid antineoplastics, antineoplastic - alkylating agents, antineoplastic systemic enzyme inhibitors, antineoplastic - antimetabolites |
| Lab Tests | albumin, hematocrit, hemoglobin, creatinine serum, red blood cell count, calcium, white blood cell count, platelet count, mean corpuscular hemoglobin concentration (MCHC), mean corpuscular hemoglobin (MCH), protein, mean corpuscular volume (MCV) |

**3.2.   Model Training**

We performed an 80:20 split on the dataset, allocating 80% of the data for cross-validation of different model and feature set combinations. We used the remaining 20% as a hold-out dataset for testing the generalizability of our best performing model. Our approach to performing model training and feature selection was two-fold:

1. First, we trained seven different model configurations, each on four different feature sets. The model configurations and feature set choices are described in the remainder of this subsection and in subsections 3.3.1 and 3.3.2.

9

231    2.  Second, we took the highest performing model configuration and used a stepwise feature

232        selection approach to attempt to find a reduced subset of features that would provide

233        comparable performance. The implementation of this feature selection approach is

234        described in subsection 3.3.3.

235    To prevent overfitting, all models were trained and validated on our training dataset using 10-

236    fold cross-validation. We evaluated our trained models using the area under the ROC curve

237    (AUROC) and the DeLong test for statistical significance [33]. We also evaluated the AUROC

238    generated by the Khorana score on our dataset and used this for baseline performance

239    comparisons with our models.

240        For the first phase of our study, we trained and evaluated models using the machine

241    learning algorithms and parameter configurations listed in Table 3.

242

243                Table 3: Machine Learning Model Configurations

| Model | Parameter Choices |
|---|---|
| Logistic Regression (LR) [34] | - |
| Support Vector Machine (SVM) [35] | Radial basis function kernel, linear kernel |
| Random Forest (RF) [36] | 50 trees, 100 trees, 200 trees, 500 trees |

244        All LR, SVM, and RF models were implemented using the Scikit-learn library in

245    Python [37]. Each of these models was cross-validated on four different feature sets/subsets:

246    1.  All 29 available features in our feature pool.

2. Features used for calculating the Khorana score: cancer site, platelet count, hemoglobin level, white blood cell count, and BMI.

3. Features selected by our clinical team. We will refer to this feature selection method as the "clinical expert" method.

4. Features selected based on statistical correlation with VTE incidence. We will refer to this feature selection method as the "filtering" method.

For the second phase of the experiment, we identified the model with the highest performance based on AUROC values and DeLong test results for statistical significance. We then used this model to perform a stepwise forward feature selection method to identify a minimum subset of features required to attain equivalent performance. We will refer to this feature selection method as the "wrapper" method. The implementations of this and the clinical expert and filtering methods are described in detail in the following section.

Finally, we tested our best performing model on the held-out dataset to better examine the generalizability of the model and ensure that we did not overfit the training dataset.

## 3.3. Feature Selection Methods

In training different machine learning models for predicting VTE, we experimented with three different feature selection methods. The first was an expert-driven feature selection process in which we used domain expertise from clinicians and researchers at UCDMC to derive a subset of known clinically relevant features as a feature set for training our machine learning models. The second was a filtering approach which identified the highest statistically correlated features with our target. The third was a wrapper approach that bootstrapped the model training process to iteratively accumulate an optimal set of features for a chosen ML classifier [38].

271        The clinical expert and filtering approaches were used in the first phase of our study for

272   comparing performances of different machine learning approaches across several feature sets. The

273   goals of performing these feature selection approaches were to:

274        ●  Examine the utility of commonly accepted VTE risk factors in a machine learning

275           approach.

276        ●  Identify new risk factors or combinations of risk factors which may add value to

277           predicting VTE incidence in cancer patients using machine learning.

278        The wrapper approach was used in the second phase of our study on the best performing

279   model and feature set from the first phase. The goal of this approach was primarily to:

280        ●  Minimize the number of features required for the best performing model configuration

281           to achieve optimal performance.

282        The implementation details for these feature selection methods are described in the

283   following subsections.

### 3.3.1   Clinical Expert Method

286        Our first feature selection method involved consulting with our team of physicians to

287   determine a subset of features that are known risk factors in the development of VTE. The

288   decisions made in this process were based both on clinical expertise and review of literature in the

289   area [19, 29, 20, 21, 22, 32].

### 3.3.2   Filtering Method

292        Since our data consists of both categorical and continuous data, we divided our feature

293   filtering approach into two tasks. For the categorical features, we determined the likelihood of

294   each feature being linearly independent of our target variable using a chi-squared test [39].

295    Meanwhile, for each continuous feature in our dataset, we observed the distribution of

296    the feature across VTE-diagnosed patients as well as the distribution of the feature across

297    patients without a VTE diagnosis. We then compared these distributions to determine the

298    likelihood that they came from one common distribution using a Kolmogorov-Smirnov (KS)

299    test for goodness of fit [40].

300    We acquired our final statistically filtered feature set by selecting only the features from

301    both of the above tests which resulted in $p < 0.05$.

### 3.3.3 Wrapper Method

304    The final feature selection process we used was an empirical forward feature selection

305    method that served the purpose of maximizing the performance of our model while minimizing

306    the dimensionality. While a high-dimensional model is appealing from a performance

307    standpoint, it may not always be practical in a clinical setting due to limitations in available

308    lab test results or other information. Performing a forward feature selection process allows us

309    to directly identify only the $n$ best performing features on our dataset and thus reduce the

310    amount of required information without significant sacrifices in performance.

311    While the filtering method that is discussed in the last subsection is valuable for

312    identifying variables directly correlated with the target, it fails to examine how different

313    combinations of these variables may affect the predictive power of our chosen ML classifier

314    [33]. In order to cover the full space of variable interactions, we would ideally train a model on

315    every possible combination of features from our feature pool, but doing so would take several

316    years of model training and would be computationally infeasible. We used the wrapper method

317    to shortcut this process and only test a small subset of all possible unique feature combinations.

318    In our approach, we accumulated features one at a time under the assumption that the

319    best performing feature at each iteration is part of the optimal set [41]. This process started by

320    training 29 separate models: one trained on each feature in our set. Each training cycle included

321    10 iterations of 10-fold cross-validation. The best performing feature was then selected and the

322    process repeated with the remaining 28 features, this time also including the best selected

323    feature(s) from the previous iteration(s) and so on. We continued to accumulate features in this

324    fashion until we no longer saw improvements in performance for a predetermined number of

325    iterations. To provide a small buffer for temporary drops in performance, we set this number

326    to 2 iterations.

327    It should be noted that, while the clinical expert and filtering feature selection methods

328    are determined independently of any model choices, the wrapper selected features are specific

329    to one model as they are accumulated by iterative model training. Since we used this method

330    in the second phase of our study to optimize the feature set for a selected model, we found it

331    sufficient to only perform the wrapper feature selection process for our best performing model.

332    ## 4. Results

333
334    ### 4.1. Model Selection

335
336    The first phase of our study involved training several model configurations on different

337    selected feature sets. Each model was evaluated by generating an AUROC value and

338    confidence interval from 10 iterations of 10-fold cross-validation. The results of this model

339    training and feature selecting are presented in this section and in section 4.2. Table 4 shows

340    the performance of each model configuration on the training dataset (80% of the original

341    dataset) across the four different feature sets listed in section 3.2. Each row represents a unique

342    model algorithm or scoring system and each column represents a unique feature set. To make

a fair comparison between different models that are using different feature sets, we have

included a model trained on the features that the Khorana score uses as shown in column 3

**Khorana (n=5)** of Table 4. The performance generated by using the standard Khorana scoring

system itself is also included as a baseline in the last row of Table 4. All model ROC curves

were compared to that of the baseline Khorana score in the last row of Table 4 via the DeLong

test. The differences that were statistically significant based on a p-value < 0.05 are marked

with an asterisk in the table. The full list of model-to-model DeLong comparisons is also

provided in Appendix B.

Table 4: AUROC of predictive models by feature set

|  | All (n=29) | Khorana (n=5) | Clinical (n=5) | Filtered (n=20) |
|---|---|---|---|---|
| Logistic Regression | 0.684 ± 0.054* | 0.668 ± 0.077 | 0.662 ± 0.074 | 0.672 ± 0.047* |
| SVM (RBF Kernel) | 0.652 ± 0.061 | 0.562 ± 0.061* | 0.576 ± 0.056* | 0.617 ± 0.072 |
| SVM (Linear Kernel) | 0.644 ± 0.042 | 0.577 ± 0.040* | 0.589 ± 0.048* | 0.669 ± 0.036* |
| Random Forest (50 trees) | 0.751 ± 0.068* | 0.672 ± 0.062* | 0.681 ± 0.072* | 0.748 ± 0.071* |
| Random Forest (100 trees) | 0.752 ± 0.062* | 0.676 ± 0.066* | 0.683 ± 0.072* | 0.743 ± 0.073* |
| Random Forest (200 trees) | 0.762 ± 0.065* | 0.684 ± 0.070* | 0.692 ± 0.074* | 0.746 ± 0.075* |
| Random Forest (500 trees) | 0.761 ± 0.065* | 0.684 ± 0.073* | 0.696 ± 0.071* | 0.755 ± 0.067* |
| Baseline: Khorana Score | - | 0.632 ± 0.019 | - | - |

* p<0.05 from DeLong test when compared to Khorana score (bottom row)
In general, every model outperformed the Khorana score baseline when trained on our

entire feature space (though this difference for the SVM models was not statistically significant).

The RF models trained on the same features used in the Khorana score all achieved a small but

357 significant improvement over the Khorana score, suggesting that using ML alone instead of a

358 simple point system may offer an improvement over currently used clinical risk assessment scores.

359 However, the results of the models trained on the other feature sets indicate that this is not the

360 maximum attainable performance and that adding additional risk factors to the model could result

361 in even larger performance improvements.

362      Every RF model also outperformed the logistic regression and SVM models on each feature

363 set suggesting that a random forest is likely the best suited algorithm choice for this task among

364 our tested classifiers. For the ease of viewing, the p-values of all pair-wise model comparisons by

365 feature set are not listed here but can be viewed in Appendix B.

366      The RF models also showed similar trends across feature sets with performance being

367 highest when trained on all features followed by the filtered feature set, clinical expert feature set,

368 and then the Khorana score feature set. The highest performing models were the four RF models

369 trained on all features and on the filtered feature set. Since the difference between these models

370 was generally not statistically significant, we chose the most complex model – the RF model with

371 500 trees – as our best performing model for the second phase of the study. The reasoning for this

372 choice was that a more complex model, while more prone to overfitting, is also capable of learning

373 more complex variable relationships leading to potential performance improvements. As

374 mentioned in the methodology, we combat and assess overfitting by performing 10-fold cross

375 validation on all experiments and further validating our best performing model on a held-out

376 dataset.

377      Based on these results, we will focus on the performance of the 500-tree RF model for the

378 remainder of our analysis where we will explore optimizing the set of required features using the

379 wrapper feature selection method and will validate our model performance on our held-out dataset.

380    But first, details on the results of the clinical expert and filtering feature selection processes are

381    provided in the following section.

382

**4.2. Feature Selection Results**

**4.2.1 Clinically Important Features**

385    Our first feature selection method involved reducing our feature set to a list of only five

386    features deemed clinically important to the prediction of VTE by a team of UCDMC physicians

387    and researchers. These features are:

388    **platelet count, white blood cell count, hemoglobin, cancer site, cancer stage**

389    The first four of these are the same four features that are common across the Khorana, Vienna

390    CATS, PROTECHT, and CONKO scoring systems while cancer stage is an additional feature

391    deemed relevant by our team [19, 20, 21, 22]. The RF model with 500 trees trained on these features

392    outperforms the AUROC of the Khorana score on our dataset by 10%. This improvement can be

393    attributed to the fact that the RF model is capable of making decisions that are much more nuanced

394    than the decisions made in any of the listed scoring systems, which involve only simple point

395    additions based on binary categorizations of the data [39]. Despite this improvement in performance,

396    the model still falls short of the model trained on the full feature set by 8.5%, indicating that there

397    are other potentially useful features in predicting VTE that were not initially deemed clinically

398    relevant.

399

**4.2.2 Filtered Features**

401    In order to further examine the known clinically relevant features and identify new features,

402    we used statistical methods to filter our feature pool and identify features highly correlated with

403    our target variable. The feature filtering method described previously yielded a set of 20 features

404 that were significantly correlated with the binary presence of VTE. The full list of this filtered

405 feature set includes the following features:

406 **site, grade, stage, histopathological type, gender, age, race list, antineoplastic -**

407 **aromatase inhibitors, albumin, hematocrit, hemoglobin, creatinine serum, red blood**

408 **cell count, calcium, white blood cell count, platelet count, MCHC, MCH, protein,**

409 **MCV**

410 Notably, all of the clinically essential features identified above were also found to be

411 significantly correlated with our target. All of the features used in the Khorana score were also

412 selected with the exception of BMI. All of the lab tests in our feature pool were selected as well

413 while all but one pharmacologic class, i.e., antineoplastic aromatase inhibitors, were left out. The

414 RF model with 500 trees achieved a 19.5% improvement over the Khorana score and did not result

415 in a significant decline in performance based on the DeLong test compared to the model trained

416 on all features.

417

418 **4.3. Model Optimization**

419 For the second phase of our study, we looked at optimizing the feature set for our best performing

420 model configuration and validating the performance on our held-out test set. Based on the results

421 presented in Table 4, we used the 500-tree RF model trained on our entire feature pool as a baseline

422 for our best performing model. In this section, we present the results of using this model with the

423 previously described wrapper feature selection method to reduce the dimensionality of the feature

424 set while attempting to maintain the same level of model performance.

425 **4.3.1. Wrapper Selected Features**

426  Table 5 compares the cross-validation performance of the 500-tree RF model using the wrapper

427  selected feature set to the results from the first phase of the study. When compared to the model

428  trained on all features, the wrapper and filtered feature sets are the only feature sets that did not

429  result in a statistically significant decline in performance. This confirms that the wrapper method

430  was effective in identifying a reduced subset of features (52% of the whole feature pool and 75%

431  of the filtered feature pool), without sacrificing performance.

432  Table 5. Cross-Validation of 500-Tree Random Forest on All Feature Sets

|  | All (n=29) | Khorana (n=5) | Clinical (n=5) | Filtered (n=20) | Wrapper (n=15) |
|---|---|---|---|---|---|
| Random Forest (500 trees) | $0.761 \pm 0.065$ | $0.684 \pm 0.073$* | $0.696 \pm 0.071$* | $0.755 \pm 0.067$ | $0.769 \pm 0.072$ |

433  * $p < 0.05$ from DeLong test when compared to model trained on all features (first column)

434      Table 6 contains the ordered list of features accumulated when performing the wrapper

435  feature selection method with the RF model of 500 trees. The curve illustrated in Figure 1 shows

436  the relationship between these features and the AUROC of our model during feature accumulation.

437  Each model evaluation came from the average result of 10 iterations of 10-fold cross-validation.

438  The x-axis represents each iteration of the recursive accumulation of features, while the y-axis

439  represents the AUROC associated with the model trained after each added feature. The model

440  trained on this set of recursively selected features not only matched the performance of the model

441  trained on all features with no statistical difference between ROC outputs, but also did so with

442  only 15 features, reducing the size of our feature set by 14. The ROC and PRC curves resulting

443  from training a model on these 15 features are contained in Figure 3 and Figure 4 respectively.

444

445  Table 6. Order of Accumulated Features During Wrapper Selection

| 1 | creatinine serum |
|---|---|

| 2 | antineoplastic - aromatase inhibitors |
|---|---|
| 3 | MCHC |
| 4 | red blood cell count |
| 5 | stage |
| 6 | Immunosuppressives |
| 7 | antineoplastic - antiandrogenic agents |
| 8 | protein |
| 9 | site |
| 10 | MCV |
| 11 | antineoplastic - alkylating agents |
| 12 | albumin |
| 13 | antineoplastic – antimetabolites |
| 14 | MCH |
| 15 | histopathological type |

446

447

448      Figure 1. Mean AUROC of 500-tree RF Model During Wrapper Feature Accumulation

449    Unlike in the clinical expert and filter-selected feature sets, seven different medications

450    were included in the wrapper-selected feature set, although only two appeared in the first

451    twelve selected features. Furthermore, the white blood cell count and platelet count lab tests

452    were excluded despite being included in both of our other examined feature sets as well as the

453    Khorana score. This exclusion is not to undermine the usefulness of the features to the task of

454    VTE prediction, but rather to show that they were not necessary for achieving optimal

455    performance with reduced dimensionality on our dataset.

**4.3.2. Feature Set Comparisons**

457    Table 7 lists the overlap between the feature sets of the three presented feature selection

458    methods. The full list of features selected by each method is provided in Appendix A.

459    All features deemed clinically relevant were also found to be statistically correlated

460    with the presence of VTE in our filtered feature set. Furthermore, all three feature selection

461    methods selected the cancer site and stage as important features for VTE prediction. While

462    cancer site is a widely used risk factor for VTE, cancer stage is not typically included in

463    currently used scoring systems [19, 20, 21, 22]. The clinical team further concurred with the data

464    driven finding of the importance of clinical staging information.

465    The overlap of the clinical expert and wrapper feature sets matches the overlap of the

466    clinical expert, filter, and wrapper feature sets and is thus omitted from the table.

467    Table 7. Overlapping features between feature sets

| Feature Selection Methods | Features |
|---|---|
| *Expert + Filter + Wrapper | site, stage |

| Filter + Wrapper | site, stage, antineoplastic - aromatase inhibitors, albumin, creatinine serum, red blood cell count, mean corpuscular hemoglobin concentration (MCHC), mean corpuscular hemoglobin (MCH), protein, mean corpuscular volume (MCV), histopathological type |
|---|---|
| Filter + Expert | site, stage, hemoglobin, platelet count, white blood cell count |

468    *The overlap of only the expert and wrapper feature sets produces the same list of features

469    **4.4 Performance Validation on Held-Out Data**

470    The remainder of the results section shows the performance when validating our RF model trained

471    with 500 trees on our held-out data (20% of the original dataset).

472    **4.2.1 All features**

473

474    Figure 2. Performance comparison on held out test set between Khorana score and RF model with
475    all features

476            The ROC curve in Figure 2 illustrates the test performance of the RF model with 500

477        trees being trained on our entire feature pool in comparison to the ROC curve generated from

478        the Khorana score on our held-out test dataset. The model achieves a statistically significant

479        improvement in AUROC of 16.1% compared to the Khorana score. This increase in

480        performance confirms the potential for improving VTE prediction through the inclusion of new

481        risk factors in a machine learning approach. Next, we validated the 500-tree RF model with

482        each of the previously examined feature subsets.

483

484

485     Figure 3. ROC Performance by Feature Set on Held-Out Data

486         The ROC curves in Figure 3 show this performance by feature set when run on our

487     held-out data. As in the results in section 4.3.1., the model trained on the wrapper selected

488     features did not result in a statistically significant decline in performance compared to the

489     model trained on the entire feature pool. This validates our takeaway that the wrapper feature

490     selection process provided an effective way to reduce the feature space without impacting

491     performance. A full list of DeLong test comparisons for the 500-tree RF models on the held-

492     out dataset are provided in Appendix B.

493         For additional validation, we evaluated the precision-recall curve (PRC) for the 500-

494     tree RF model on each feature set. These results are displayed in Figure 4.

495

496

497     Figure 4.  PRC Performance by Feature Set on Held-Out Data

498         Similar to the ROC results, the PRC curves in Figure 4 show that the models trained

499     on all features and on the wrapper-selected features are the best performing models and achieve

500     comparable performance.

501

502     **5. Discussion**

503         In this study, we examined the utility of using machine learning to predict VTE in

504     cancer patients. We accomplished this through a carefully designed set of steps adhering to a

505     typical machine learning pipeline. First, we selected a feature pool based on the data

506     availability within our patient population. We also set aside 20% of the data in a held-out

dataset for final model validation. We then performed a number of feature selection methods

and trained multiple machine learning classifiers with different hyperparameter configurations

to identify a best performing model for our use case. Finally, we iteratively trained the best

performing model in order to accumulate a minimum set of required features and thus reduce

the complexity of the model without impacting model performance.

The results of this process allow us to draw insight into how a machine learning

classifier might offer an improvement in performance over traditionally used clinical VTE risk

assessment systems in cancer patients. With these results, we are able to examine our feature

pool and identify those features that are most useful in the context of developing an efficient

machine learning classifier by comparing the selected features and resulting model

performance across multiple unique feature selection methods.

This project was an effort to showcase the improved predictive performance of various

ML models over the Khorana score in predicting VTE in cancer patients. We compared the

performance of models trained on different feature sets selected by domain experts, statistical

methods, and ML techniques. We identified features that were common across these selected

feature sets to better understand which features are meaningful in this context.

Our trained classifiers achieved encouraging results on numerous feature subsets. We

found that a 500-tree RF model trained using only the features used in the Khorana score

achieved a statistically significant 14.6% improvement in AUROC over the standard point-

based Khorana score on our held-out test set with an AUROC of $0.769 \pm 0.007$. Meanwhile,

we achieved a peak AUROC of $0.779 \pm 0.006$ on a held-out dataset when training the 500-tree

RF model on our entire feature pool. This surpassed the performance of the Khorana score on

the same dataset by 16.1%. We were additionally able to reduce the number of required

features to 15 total (a 48% reduction) without a statistically significant impact on model performance by using a wrapper method to iteratively accumulate features. We also used two model-agnostic feature selection methods – a statistical filtering method and a clinical expert method – which both achieved AUROCs of $0.771 \pm 0.007$ and $0.757 \pm 0.004$ respectively on our held-out dataset. All of these results showed statistically significant improvements in performance over that of the Khorana score.

The results in Table 7 depict the overlap between the features selected by our three described feature selection methods. Only cancer site and cancer stage were common across all three feature sets. Cancer site is already a common risk factor considered in current VTE risk stratification systems [19, 20, 21, 22]. Based on our experimental results, cancer stage merits inclusion in future VTE prediction systems using an ML approach. Meanwhile, all of the features deemed clinically relevant were also found to be statistically significant in the filtered feature set. Unlike the other two feature sets, the wrapper-selected feature set did not include hemoglobin. However, it did identify three related metrics - corpuscular hemoglobin (MCH), mean corpuscular hemoglobin concentration (MCHC), and mean corpuscular volume (MCV) - as essential metrics for VTE prediction. While these metrics are not identical to hemoglobin, they are likely inter-related. Furthermore, since the wrapper method optimizes the feature space based on empirical performance of different feature combinations, an excluded feature is not by necessity unimportant. Instead, an excluded feature may be redundant when compared to the optimal set of features, making its inclusion unnecessary for improving prediction performance.

In comparison to the features used in the Khorana score, all but BMI are included in the filtered and clinically relevant feature sets. Furthermore, the cancer site, which is the most

553    heavily weighted risk factor in the Khorana score, was selected in all three feature sets.

554    Interestingly, BMI, which is included in the Khorana score, Vienna CATS, and PROTECHT,

555    was not identified as useful in any of our acquired feature sets [19, 20, 21]. Aside from BMI,

556    however, the results of this study suggest that the predictors used in the Khorana score have a

557    relatively high predictive power when used in a machine learning context. The results also

558    suggest that the stage of the cancer is useful in predicting VTE and should be considered in

559    future machine learning applications. Because staging information is not always readily

560    available in medical notes, future studies could look to reliably extract this information from

561    free medical text using NLP methods. Since cancer staging can vary over time as new

562    information comes in and is incorporated in the staging determination, this problem is

563    particularly challenging with past efforts achieving only limited success [42, 43]. One approach

564    that may improve this performance without sacrificing too much predictive power in VTE risk

565    assessment could involve reducing the cancer stage to a binary variable that simply indicates

566    a presence or absence of metastasis [44].

567         While the results of this study are promising, it is important to note that the dataset uses

568    a small sample size, especially for certain subgroups, (i.e., only a few pharmacological groups

569    were used in the patient population). Also, the study did not include cancer patients who had

570    radiation therapy. There is increasing evidence implicating radiotherapy in cancer associated

571    thrombosis (CAT) in cancer patients, however accessing data from the radiation therapy

572    information system (RTIS) was not possible for this study. This study dealt with the patient

573    population at only one location, so before we generalize these results across the general

574    population, the findings in this study should be validated in other patient populations.

575    Furthermore, this study takes a time-agnostic approach to identify useful predictors for VTE

in cancer patients. Therefore, this approach highlights VTE predictors that may be useful in a machine learning context but does not yet reflect an implementable clinical scenario. With this being the case, the aim of this study was to effectively identify these useful predictors in order to provide the groundwork for exploration of this problem in specific clinical scenarios (i.e., at different stages of pre-diagnosis presentation, establishing diagnosis, and post-diagnosis treatment phases of a patient's cancer management).

The methods used in this study could be generalizable to other clinical conditions, particularly ambulatory settings, where there is moderate to strong increased risk for developing VTE, such as, congestive heart or respiratory failure, hormone replacement and oral contraceptive therapy, antiphospholipid antibody and other thrombophilia syndromes [45]. Even though multiple studies have demonstrated that thromboprophylaxis using anticoagulant treatments such as low-molecular-weight heparin (LMWH) can reduce the likelihood of VTE events, due to the need for training the patients and care-givers to administer (parenteral) the LMWH, regular lab monitoring and dose adjustment, as well as the potential for bleeding complications, all of which add to the cost and quality of care, such prophylaxis may not always be feasible and risk-free. There is thus a need for effective VTE risk stratification and decision support systems to ensure that prophylaxis is administered only to high-risk patients.

The project goal was to select the necessary and sufficient features from our available feature pool that would maximize the predictive power of various statistical ML models. It can be a hard decision to initiate prophylaxis against VTE, especially in ambulatory cancer patients where anti-thrombosis prophylaxis can be expensive and cumbersome. Evidence based decision support is crucial for minimizing risk in this decision process and improving patient outcomes.

At the point of care where the decisions are made, ideally, prediction tools and scoring systems should automatically retrieve the required features and inform the clinicians to help make decisions. For ease of use and interpretability, the list of features should be small, but should provide meaningful enough information to supplement the current evidence and clinicians' evaluations. We found cancer staging information to be particularly meaningful as a predictor of VTE as it was selected in all of our feature selection processes. The Khorana score does not include the cancer staging information as often it can be hard to retrieve accurate staging information from clinical notes. Accurate staging information is often established by cancer registrars retrospectively, which may take up to six months. Our study emphasizes the importance of cancer staging information as a predictor of VTE in cancer patients and highlights the need for its timely evaluation. Simplifying the cancer stage variable into a binary value indicating whether the cancer is metastatic (stage 4) or non-metastatic could improve the accessibility and real-time accuracy of staging but would require further studies and additional validation.

## 6. Conclusion

Machine learning offers a promising avenue for improving the performance of current VTE prediction scores in cancer patients. A combination of a time-agnostic approach and three unique feature selection methods demonstrates that at least four of the features that are used to calculate the Khorana score can also provide high predictive power to a machine learning classifier. We also observe that cancer stage information is generally more useful than BMI as a predictor in our ML classifiers. Consultation with clinicians reveal a potential reason - BMI can vary as patients lose significant weight due to cancer itself, chemotherapy, and associated anorexia or other adverse effects. Furthermore, with significant improvements in the generated ROC curve, it is clear that a machine learning classifier can make complex deductions that

623  may allow it to outperform currently used VTE risk scores. The results in this study offer a

624  foundation from which future machine learning approaches to VTE prediction in cancer

625  patients can be built. Future studies should consider the identified relevant variables in the

626  context of a temporal analysis in which machine learning may be used to dynamically assess

627  at all levels how cancer management progress, including medical intervention, over time can

628  alter a patient's risk of developing VTE.

629  **Declaration**

630  **Ethics approval and consent to participate**

631  Appropriate institutional review board (IRB) review and approval was obtained from the UCDMC

632  IRB, bearing number: UCDMC.

633  **Conflict of interests**

634  The authors declare no conflict of interests.

635

636  **References**

637  1.  Khorana AA. The NCCN Clinical Practice Guidelines on Venous Thromboembolic
638      Disease: strategies for improving VTE prophylaxis in hospitalized cancer patients.
639      Oncologist. 2007;12(11):1361-1370. doi:10.1634/theoncologist.12-11-1361
640  2.  Lyman GH, Khorana AA, Falanga A, et al. American Society of Clinical Oncology
641      guideline: recommendations for venous thromboembolism prophylaxis and treatment in
642      patients    with    cancer.    J    Clin    Oncol.    2007;25(34):5490-5505.
643      doi:10.1200/JCO.2007.14.1283
644  3.  Silverstein MD, Heit JA, Mohr DN, Petterson TM, O'Fallon WM, Melton LJ 3rd. Trends
645      in the incidence of deep vein thrombosis and pulmonary embolism: a 25-year population-
646      based study. Arch Intern Med. 1998;158(6):585-593. doi:10.1001/archinte.158.6.585
647  4.  Spencer FA, Emery C, Lessard D, et al. The Worcester Venous Thromboembolism study:
648      a population-based study of the clinical epidemiology of venous thromboembolism. J Gen
649      Intern Med. 2006;21(7):722-727. doi:10.1111/j.1525-1497.2006.00458.x
650  5.  Kessler CM. The link between cancer and venous thromboembolism: a review. Am J Clin
651      Oncol. 2009;32(4 Suppl):S3-S7. doi:10.1097/COC.0b013e3181b01b17
652  6.  Lee AY, Levine MN. Venous thromboembolism and cancer: risks and outcomes.
653      Circulation. 2003;107(23 Suppl 1):I17-I21. doi:10.1161/01.CIR.0000078466.72504.AC

7. Heit JA, Silverstein MD, Mohr DN, Petterson TM, O'Fallon WM, Melton LJ 3rd. Risk factors for deep vein thrombosis and pulmonary embolism: a population-based case-control study. Arch Intern Med. 2000;160(6):809-815. doi:10.1001/archinte.160.6.809

8. Khorana AA, Francis CW, Culakova E, Kuderer NM, Lyman GH. Thromboembolism is a leading cause of death in cancer patients receiving outpatient chemotherapy. J Thromb Haemost. 2007;5(3):632-634. doi:10.1111/j.1538-7836.2007.02374.x

9. Khorana AA, Francis CW, Culakova E, Kuderer NM, Lyman GH. Frequency, risk factors, and trends for venous thromboembolism among hospitalized cancer patients. Cancer. 2007;110(10):2339-2346. doi:10.1002/cncr.23062

10. Sørensen HT, Mellemkjaer L, Olsen JH, Baron JA. Prognosis of cancers associated with venous thromboembolism. N Engl J Med. 2000;343(25):1846-1850. doi:10.1056/NEJM200012213432504

11. Elting LS, Escalante CP, Cooksley C, et al. Outcomes and Cost of Deep Venous Thrombosis Among Patients With Cancer. Arch Intern Med. 2004;164(15):1653–1661. doi:10.1001/archinte.164.15.1653

12. Prandoni P, Lensing AW, Piccioli A, et al. Recurrent venous thromboembolism and bleeding complications during anticoagulant treatment in patients with cancer and venous thrombosis. Blood. 2002;100(10):3484-3488. doi:10.1182/blood-2002-01-0108

13. Mandalà M, Falanga A, Roila F; ESMO Guidelines Working Group. Management of venous thromboembolism (VTE) in cancer patients: ESMO Clinical Practice Guidelines. Ann Oncol. 2011;22 Suppl 6:vi85-vi92. doi:10.1093/annonc/mdr392

14. Cayley WE Jr. Preventing deep vein thrombosis in hospital inpatients. BMJ. 2007;335(7611):147-151. doi:10.1136/bmj.39247.542477.AE

15. Samama MM, Cohen AT, Darmon JY, et al. A comparison of enoxaparin with placebo for the prevention of venous thromboembolism in acutely ill medical patients. Prophylaxis in Medical Patients with Enoxaparin Study Group. N Engl J Med. 1999;341(11):793-800. doi:10.1056/NEJM199909093411103

16. Leizorovicz A, Cohen AT, Turpie AG, et al. Randomized, placebo-controlled trial of dalteparin for the prevention of venous thromboembolism in acutely ill medical patients. Circulation. 2004;110(7):874-879. doi:10.1161/01.CIR.0000138928.83266.24

17. Cohen AT, Davidson BL, Gallus AS, et al. Efficacy and safety of fondaparinux for the prevention of venous thromboembolism in older acute medical patients: randomised placebo controlled trial. BMJ. 2006;332(7537):325-329. doi:10.1136/bmj.38733.466748.7C

18. Kuderer NM, Khorana AA, Lyman GH, Francis CW. A meta-analysis and systematic review of the efficacy and safety of anticoagulants as cancer treatment: impact on survival and bleeding complications. Cancer. 2007;110(5):1149-1161. doi:10.1002/cncr.22892

19. Khorana AA, Kuderer NM, Culakova E, Lyman GH, Francis CW. Development and validation of a predictive model for chemotherapy-associated thrombosis. Blood. 2008;111(10):4902-4907. doi:10.1182/blood-2007-10-116327

20. Ay C, Dunkler D, Marosi C, et al. Prediction of venous thromboembolism in cancer patients. Blood. 2010;116(24):5377-5382. doi:10.1182/blood-2010-02-270116

21. Verso M, Agnelli G, Barni S, Gasparini G, LaBianca R. A modified Khorana risk assessment score for venous thromboembolism in cancer patients receiving chemotherapy: the Protecht score. Intern Emerg Med. 2012;7(3):291-292. doi:10.1007/s11739-012-0784-y

22. Pelzer U, Sinn M, Stieler J, Riess H. Primary pharmacological prevention of thromboembolic events in ambulatory patients with advanced pancreatic cancer treated with chemotherapy?. Dtsch Med Wochenschr. 2013;138(41):2084-2088. doi:10.1055/s-0033-1349608

23. van Es N, Di Nisio M, Cesarman G, et al. Comparison of risk prediction scores for venous thromboembolism in cancer patients: a prospective cohort study. Haematologica. 2017;102(9):1494-1501. doi:10.3324/haematol.2017.169060

24. Overvad TF, Ording AG, Nielsen PB, Skjøth F, Albertsen IE, Noble S, Vistisen AK, Gade IL, Severinsen MT, Piazza G, Larsen TB. Validation of the Khorana score for predicting venous thromboembolism in 40 218 patients with cancer initiating chemotherapy. Blood Adv. 2022 May 24;6(10):2967-2976. doi: 10.1182/bloodadvances.2021006484. PMID: 35045569; PMCID: PMC9131922.

25. Mulder FI, Candeloro M, Kamphuisen PW, Di Nisio M, Bossuyt PM, Guman N, Smit K, Büller HR, van Es N; CAT-prediction collaborators. The Khorana score for prediction of venous thromboembolism in cancer patients: a systematic review and meta-analysis. Haematologica. 2019 Jun;104(6):1277-1287. doi: 10.3324/haematol.2018.209114. Epub 2019 Jan 3. PMID: 30606788; PMCID: PMC6545838.

26. Obermeyer Z, Emanuel EJ. Predicting the Future - Big Data, Machine Learning, and Clinical Medicine. N Engl J Med. 2016;375(13):1216-1219. doi:10.1056/NEJMp1606181

27. Deo RC. Machine Learning in Medicine. Circulation. 2015;132(20):1920-1930. doi:10.1161/CIRCULATIONAHA.115.001593

28. Kotsiantis, S., Zaharakis, I, Pintelas, P. Supervised machine learning: A review of classification techniques. Emerging artificial intelligence applications in computer engineering, 2007, 160, 3–24.

29. Ferroni P, Zanzotto FM, Scarpato N, et al. Risk Assessment for Venous Thromboembolism in Chemotherapy-Treated Ambulatory Cancer Patients. Med Decis Making. 2017;37(2):234-242. doi:10.1177/0272989X16662654

30. Mehmet Gönen, Ethem Alpaydin. Multiple Kernel Learning Algorithms. Journal of Machine Learning Research 12 (2011) 2211-2268.

31. Hanna DL, White RH, Wun T. Biomolecular markers of cancer-associated thromboembolism. Crit Rev Oncol Hematol. 2013;88(1):19-29. doi:10.1016/j.critrevonc.2013.02.008

32. Wun T, White RH. Epidemiology of cancer-related venous thromboembolism. Best Pract Res Clin Haematol. 2009;22(1):9-23. doi:10.1016/j.beha.2008.12.001

33. 31. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. Biometrics. 1988 Sep 1:837-45.

34. Bagley SC, White H, Golomb BA. Logistic regression in the medical literature: standards for use and reporting, with particular attention to one medical domain. J Clin Epidemiol. 2001;54(10):979-985. doi:10.1016/s0895-4356(01)00372-9

35. Cristianini N, Shawe-Taylor J. An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods. Cambridge: Cambridge University Press; 2000. doi:10.1017/CBO9780511801389

36. Breiman, L. Random forests, Machine learning. 2001. 45(1), 5–32.

37. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V and others. Scikit-learn: Machine learning in Python. Journal of machine learning research. (2011). 12(Oct), 2825–2830.

38. Jiliang Tang, Salem Alelyani, Huan Liu. Feature selection for classification: A review. Data classification: Algorithms and applications, (2014). 37.

39. McHugh ML. The chi-square test of independence. Biochem Med (Zagreb). 2013; 23(2):143-149. doi:10.11613/bm.2013.018

40. Massey Jr, F. J. The Kolmogorov-Smirnov test for goodness of fit. Journal of the American Statistical Association (1951). 46(253), 68-78. DOI: 10.2307/2280095

41. Guyon, Isabelle, and André Elisseeff. An Introduction to Variable and Feature Selection, Journal of Machine Learning Research. (2003) 3, 1157–1182.

42. Warner JL, Levy MA, Neuss MN, Warner JL, Levy MA, Neuss MN. ReCAP: Feasibility and Accuracy of Extracting Cancer Stage Information From Narrative Electronic Health Record Data. J Oncol Pract. 2016;12(2):. doi:10.1200/JOP.2015.004622

43. AAlAbdulsalam AK, Garvin JH, Redd A, Carter ME, Sweeny C, Meystre SM. Automated Extraction and Classification of Cancer Stage Mentions from Unstructured Text Fields in a Central Cancer Registry. AMIA Jt Summits Transl Sci Proc. 2018;2017:16-25. Published 2018 May 18.

44. Soysal E, Warner JL, Denny JC, Xu H. Identifying Metastases-related Information from Pathology Reports of Lung Cancer Patients. AMIA Jt Summits Transl Sci Proc. 2017;2017:268-277. Published 2017 Jul 26.

45. Anderson Jr, F. A., & Spencer, F. A. (2003). Risk factors for venous thromboembolism. Circulation, 107(23_suppl_1), I-9.

Appendix A

Table A.1. Full list of selected features by feature selection method

| Feature Selection Method | Features |
| --- | --- |
| Clinical Expert Method | site, stage, hemoglobin, platelet count, white blood cell count |
| Filter Method | site, grade, stage, histopathological type, gender, age, race list, antineoplastic - aromatase inhibitors, albumin, hematocrit, |

| | |
|---|---|
| | hemoglobin, creatinine serum, red blood cell count, calcium, white blood cell count, platelet count, MCHC, MCH, protein, MCV |
| Wrapper Method | site, stage, histopathological type, albumin, creatinine serum, red blood cell count, MCHC, MCH, protein, MCV, antineoplastic - aromatase inhibitors, immunosuppressives, antineoplastic - antiandrogenic agents, antineoplastic - alkylating agents, antineoplastic - antimetabolites |

769

770 Appendix B

771 The following tables show the comprehensive results of performing the DeLong test for statistical
772 significance between ROC curves of the various models we trained during the study. Each table is
773 a grid of DeLong p-values. For this study, we used $p<0.05$ as our cutoff for statistical significance.
774 The first four tables are most pertinent to the results discussed in the main text while the following
775 tables contain a more comprehensive coverage of pairwise prediction comparisons.

776 Table B.1. DeLong p-values for Models Compared to Khorana Score

| | All (n=29) | Khorana (n=5) | Clinical (n=5) | Filtered (n=20) |
|---|---|---|---|---|
| **Logistic Regression** | 0.00142 | 0.07314 | 0.101754 | 0.004921 |
| **SVM (RBF Kernel)** | 0.150591 | 0.00036 | 0.001697 | 0.27491 |
| **SVM (Linear Kernel)** | 0.18518 | 3.2E-05 | 0.004174 | 0.000772 |
| **Random Forest (50 trees)** | 0.0 | 0.023375 | 0.017531 | 0.0 |
| **Random Forest (100 trees)** | 0.0 | 0.020919 | 0.015383 | 2E-06 |
| **Random Forest (200 trees)** | 0.0 | 0.011794 | 0.006736 | 2E-06 |
| **Random Forest (500 trees)** | 0.0 | 0.014679 | 0.003016 | 0.0 |

777

778 Table B.2. DeLong p-values for Models Compared to Same Model Trained on All Features

| | Khorana (n=5) | Clinical (n=5) | Filtered (n=20) |
|---|---|---|---|
| **Logistic Regression** | 0.307395 | 0.234885 | 0.300637 |

| | | | |
|---|---|---|---|
| **SVM (RBF Kernel)** | 0.00089 | 0.003027 | 0.130158 |
| **SVM (Linear Kernel)** | 0.000331 | 0.005092 | 0.08326 |
| **Random Forest (50 trees)** | 0.00465 | 0.016925 | 0.466185 |
| **Random Forest (100 trees)** | 0.005323 | 0.014444 | 0.387342 |
| **Random Forest (200 trees)** | 0.006923 | 0.016309 | 0.321548 |
| **Random Forest (500 trees)** | 0.009481 | 0.020839 | 0.431354 |

779

780 Table B.3. DeLong p-values for 500-tree RF Models on Held-Out Test Dataset

| | **All (n=29)** | **Khorana (n=5)** | **Clinical (n=5)** | **Filtered (n=20)** | **Wrapper (n=15)** |
|---|---|---|---|---|---|
| **All (n=29)** | 0.5 | 0.000465 | 0.0 | 0.00303 | 0.369048 |
| **Khorana (n=5)** | 0.000465 | 0.5 | 1.0E-06 | 0.301592 | 0.001222 |
| **Clinical (n=5)** | 0.0 | 1.0E-06 | 0.5 | 0.0 | 0.0 |
| **Filtered (n=20)** | 0.00303 | 0.301592 | 0.0 | 0.5 | 0.006966 |
| **Wrapper (n=15)** | 0.369048 | 0.001222 | 0.0 | 0.006966 | 0.5 |

781

782 Table B.4. DeLong p-values for 500-tree RF Models vs. Khorana Score on Held-Out Test Dataset

| | **All (n=29)** | **Khorana (n=5)** | **Clinical (n=5)** | **Filtered (n=20)** | **Wrapper (n=15)** |
|---|---|---|---|---|---|
| **Baseline: Khorana Score** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

783

784 Below are the results of performing the DeLong test for statistical significance between ROC
785 curves on every pairwise combination of models for each feature set we examined in the study.

786 Table B.5. DeLong p-values for Models Trained on All Features

| | **Logistic Regression** | **SVM (RBF Kernel)** | **SVM (Linear Kernel)** | **Random Forest (50 trees)** | **Random Forest (100 trees)** | **Random Forest (200 trees)** | **Random Forest (500 trees)** | **Baseline: Khorana Score** |
|---|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.5 | 0.116269 | 0.037197 | 0.010025 | 0.006254 | 0.002805 | 0.003274 | 0.001447 |
| **SVM (RBF Kernel)** | 0.116269 | 0.5 | 0.367859 | 0.00054 | 0.000257 | 0.000104 | 0.000127 | 0.150591 |
| **SVM (Linear Kernel)** | 0.037197 | 0.367859 | 0.5 | 2.7E-05 | 7E-06 | 2E-06 | 3E-06 | 0.18518 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Random Forest (50 trees)** | 0.010025 | 0.00054 | 2.7E-05 | 0.5 | 0.48744 | 0.367113 | 0.379221 | 0.0 |
| **Random Forest (100 trees)** | 0.006254 | 0.000257 | 7E-06 | 0.48744 | 0.5 | 0.372979 | 0.385744 | 0.0 |
| **Random Forest (200 trees)** | 0.002805 | 0.000104 | 2E-06 | 0.367113 | 0.372979 | 0.5 | 0.487627 | 0.0 |
| **Random Forest (500 trees)** | 0.003274 | 0.000127 | 3E-06 | 0.379221 | 0.385744 | 0.487627 | 0.5 | 0.0 |
| **Baseline: Khorana Score** | 0.001447 | 0.150591 | 0.18518 | 0.0 | 0.0 | 0.0 | 0.0 | 0.5 |

787

Table B.6. DeLong p-values for Models Trained on Khorana Score Features

| | Logistic Regression | SVM (RBF Kernel) | SVM (Linear Kernel) | Random Forest (50 trees) | Random Forest (100 trees) | Random Forest (200 trees) | Random Forest (500 trees) | Baseline: Khorana Score |
|---|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.5 | 0.000618 | 0.000882 | 0.459024 | 0.416674 | 0.330223 | 0.329241 | 0.073683 |
| **SVM (RBF Kernel)** | 0.000618 | 0.5 | 0.266912 | 7.4E-05 | 7.4E-05 | 4.2E-05 | 6.2E-05 | 0.00036 |
| **SVM (Linear Kernel)** | 0.000882 | 0.266912 | 0.5 | 5.9E-05 | 6.6E-05 | 3.8E-05 | 6.3E-05 | 3.2E-05 |
| **Random Forest (50 trees)** | 0.459024 | 7.4E-05 | 5.9E-05 | 0.5 | 0.450544 | 0.350922 | 0.349703 | 0.023375 |
| **Random Forest (100 trees)** | 0.416674 | 7.4E-05 | 6.6E-05 | 0.450544 | 0.5 | 0.399317 | 0.396751 | 0.020919 |
| **Random Forest (200 trees)** | 0.330223 | 4.2E-05 | 3.8E-05 | 0.350922 | 0.399317 | 0.5 | 0.494963 | 0.011794 |
| **Random Forest (500 trees)** | 0.329241 | 6.2E-05 | 6.3E-05 | 0.349703 | 0.396751 | 0.494963 | 0.5 | 0.014679 |
| **Baseline: Khorana Score** | 0.073683 | 0.00036 | 3.2E-05 | 0.023375 | 0.020919 | 0.011794 | 0.014679 | 0.5 |

789

Table B.7. DeLong p-values for Models Trained on Clinical Expert Features

| | Logistic Regression | SVM (RBF Kernel) | SVM (Linear Kernel) | Random Forest (50 trees) | Random Forest (100 trees) | Random Forest (200 trees) | Random Forest (500 trees) | Baseline: Khorana Score |
|---|---|---|---|---|---|---|---|---|

| | Logistic Regression | SVM (RBF Kernel) | SVM (Linear Kernel) | Random Forest (50 trees) | Random Forest (100 trees) | Random Forest (200 trees) | Random Forest (500 trees) | Baseline: Khorana Score |
|---|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.5 | 0.002724 | 0.006279 | 0.288285 | 0.272988 | 0.197527 | 0.163826 | 0.102482 |
| **SVM (RBF Kernel)** | 0.002724 | 0.5 | 0.302347 | 0.000265 | 0.00023 | 9.2E-05 | 3.5E-05 | 0.001697 |
| **SVM (Linear Kernel)** | 0.006279 | 0.302347 | 0.5 | 0.000648 | 0.000563 | 0.000226 | 8.7E-05 | 0.004174 |
| **Random Forest (50 trees)** | 0.288285 | 0.000265 | 0.000648 | 0.5 | 0.480818 | 0.380343 | 0.336385 | 0.017531 |
| **Random Forest (100 trees)** | 0.272988 | 0.00023 | 0.000563 | 0.480818 | 0.5 | 0.398935 | 0.354845 | 0.015383 |
| **Random Forest (200 trees)** | 0.197527 | 9.2E-05 | 0.000226 | 0.380343 | 0.398935 | 0.5 | 0.456638 | 0.006736 |
| **Random Forest (500 trees)** | 0.163826 | 3.5E-05 | 8.7E-05 | 0.336385 | 0.354845 | 0.456638 | 0.5 | 0.003016 |
| **Baseline: Khorana Score** | 0.102482 | 0.001697 | 0.004174 | 0.017531 | 0.015383 | 0.006736 | 0.003016 | 0.5 |

791

Table B.8. DeLong p-values for Models Trained on Filter Features

| | Logistic Regression | SVM (RBF Kernel) | SVM (Linear Kernel) | Random Forest (50 trees) | Random Forest (100 trees) | Random Forest (200 trees) | Random Forest (500 trees) | Baseline: Khorana Score |
|---|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.5 | 0.027277 | 0.451289 | 0.003532 | 0.007221 | 0.00583 | 0.001173 | 0.005015 |
| **SVM (RBF Kernel)** | 0.027277 | 0.5 | 0.024625 | 4.4E-05 | 0.00011 | 8.8E-05 | 1.2E-05 | 0.27491 |
| **SVM (Linear Kernel)** | 0.451289 | 0.024625 | 0.5 | 0.001453 | 0.003458 | 0.002766 | 0.000373 | 0.000772 |
| **Random Forest (50 trees)** | 0.003532 | 4.4E-05 | 0.001453 | 0.5 | 0.436936 | 0.477591 | 0.414938 | 0.0 |
| **Random Forest (100 trees)** | 0.007221 | 0.00011 | 0.003458 | 0.436936 | 0.5 | 0.460575 | 0.354352 | 2E-06 |
| **Random Forest (200 trees)** | 0.00583 | 8.8E-05 | 0.002766 | 0.477591 | 0.460575 | 0.5 | 0.395179 | 2E-06 |
| **Random Forest (500 trees)** | 0.001173 | 1.2E-05 | 0.000373 | 0.414938 | 0.354352 | 0.395179 | 0.5 | 0.0 |
| **Baseline: Khorana Score** | 0.005015 | 0.27491 | 0.000772 | 0.0 | 2E-06 | 2E-06 | 0.0 | 0.5 |

793