

Predicted Prognosis of Pancreatic Cancer Patients by Machine Learning

Seiya Yokoyama¹, Taiji Hamada¹, Michiyo Higashi^{1*}, Kei Matsuo¹, Kosei Maemura^{2,3}, Hiroshi Kurahara³, Michiko Horinouchi¹, Tsubasa Hiraki¹, Tomoyuki Sugimoto⁴, Toshiaki Akahane¹, Suguru Yonezawa¹, Marko Kornmann⁵, Surinder K Batra⁶, Michael A Hollingsworth⁷, and Akihide Tanimoto¹

1. Department of Pathology, Research Field in Medicine and health Sciences, Medical and Dental Sciences Area, Research and Education Assembly, Kagoshima University. Kagoshima, Japan.

2. Center for the Research of Advanced Diagnosis and Therapy of Cancer, Graduate School of Medical and Dental Sciences, Kagoshima University. Kagoshima, Japan.

3. Department of Digestive Surgery, Breast and Thyroid Surgery, Graduate School of Medical Sciences, Kagoshima University. Kagoshima, Japan.

4. Kagoshima University Research Field in Science, Science and Engineering Area Graduate School of Science and Engineering (Science) Mathematics and Computer Science Course. Kagoshima, Japan.

5. Department of General and Visceral Surgery, University of Ulm. Ulm, Germany.

6. Department of Biochemistry and Molecular Biology, Eppley Institute for Research in Cancer and Allied Diseases, University of Nebraska Medical Center. Nebraska, United States.

7. Eppley Institute for Research in Cancer, Fred and Pamela Buffet Cancer Center, University of Nebraska Medical Center. Nebraska, United States.

Running Title: Prognosis of pancreatic cancer by machine learning

Keywords: PDAC, mucin, DNA methylation, prediction model, machine learning

Financial Support: This study was supported in part by a grant from Grants-in-Aid for Scientific Research on Scientific Research (C) 18K07019 to M. Higashi, Scientific Research (C) 18K07018 to T. Hamada, and Scientific Research (C) 18K07326 to S. Yokoyama from the Ministry of Education, Science, Sports, Culture and Technology, Japan by the Kodama Memorial Foundation, Japan (S. Yokoyama) and by the Pancreas Research Foundation of Japan (S. Yokoyama). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Corresponding Author: Michiyo Higashi

Department of Pathology, Research Field in Medicine and Health Sciences, Medical and Dental Sciences Area, Research and Education Assembly, Kagoshima University. 35-1 Sakuragaoka Kagoshima, Kagoshima, Japan. 890-8544.

east@m2.kufm.kagoshima-u.ac.jp

Conflicts of Interest

We do not have any disclosure or conflict of interest in this study.

Statement of Translational Relevance

Pancreatic cancer remains a disease of high mortality despite advanced diagnostic techniques. This study investigates whether the methylation status of three mucin genes from postoperative tissue specimens from patients with pancreatic neoplasms could serve as a predictive biomarker for outcome after surgery. We evaluated the methylation status of *MUC1*, *MUC2*, and *MUC4* promoter

regions in 300 pancreatic non- and neoplastic tissues samples from 191 patients with various pancreatic lesions using methylation-specific electrophoresis. Further, integrating these results and clinicopathological features, we used support vector machine-, neural network-, and multinomial-based methods to develop a prognostic classifier. Multivariate analysis revealed that these prognostic classifiers were independent prognostic factors analyzed by not only neoplastic tissues but also non-neoplastic tissues. Analysis of epigenetic changes in mucin genes may be of diagnostic utility and one of the prognostic predictors for patients with pancreatic ductal adenocarcinoma.

Abstract

Purpose: Pancreatic cancer remains a disease of high mortality despite advanced diagnostic techniques. Mucins (MUC) play crucial roles in carcinogenesis and tumor invasion in pancreatic cancers. MUC1 and MUC4 expression are related to the aggressive behavior of human neoplasms and a poor patient outcome. In contrast, MUC2 is a tumor suppressor, and we have previously reported that MUC2 is a favorable prognostic factor in pancreatic neoplasia. This study investigates whether the methylation status of three mucin genes from postoperative tissue specimens from patients with pancreatic neoplasms could serve as a predictive biomarker for outcome after surgery.

Experimental Design: We evaluated the methylation status of *MUC1*, *MUC2*, and *MUC4* promoter regions in pancreatic tissue samples from 191 patients with various pancreatic lesions using methylation-specific electrophoresis. Then, integrating these results and clinicopathological features, we used support vector machine-, neural network-, and multinomial-based methods to develop a prognostic classifier.

Results: Significant differences were identified between the positive- and negative-prediction classifiers of patients in 5-year overall survival (OS) in the cross-validation test. Multivariate analysis revealed that these prognostic classifiers were independent prognostic factors analyzed by not only neoplastic tissues but also non-neoplastic tissues. These classifiers had higher predictive accuracy for OS than tumor size, lymph node metastasis, distant metastasis, and age and can complement the prognostic value of the TNM staging system.

Conclusions: Analysis of epigenetic changes in mucin genes may be of diagnostic utility and one of the prognostic predictors for patients with pancreatic ductal adenocarcinoma.

79 Introduction

80 Despite improvements in diagnostic tools and treatments, patients with pancreatic ductal
 81 adenocarcinoma (PDAC) have a poor clinical outcome. At the time of diagnosis, most PDAC
 82 patients are in advanced stages because the anatomical location of the pancreas and lack of specific
 83 symptoms hamper early detection. Easy infiltration to the surrounding organs and early distant
 84 metastasis, even from a small primary tumor <2 cm in diameter, enhances the stage progression (1).
 85 Moreover, PDAC is sometimes difficult to distinguish from other pancreatic diseases, such as
 86 chronic pancreatitis, even when endoscopic, ultrasound-guided, fine-needle aspiration is performed
 87 (2-4). Indolent tumors such as intraductal papillary mucinous neoplasms (IPMNs) also occur in the
 88 pancreas and sometimes transform to invasive cancer with a poor outcome (5-8). Currently, IPMNs
 89 are the most common cystic neoplasm of the pancreas and are classified into gastric, intestinal,
 90 pancreatobiliary, and oncocytic types (9,10). A recent study has demonstrated that the morphological
 91 subtype of IPMN is an independent prognostic factor (8). The overall 5-year survival rate for all
 92 diagnosed PDAC patients is presently only 13% in Japan. Unfortunately, only 10-20% patients
 93 present with resectable disease at diagnosis of PDAC (11). 5 year survival with tumor removal alone
 94 is generally less than 10%. After resection, the use of adjuvant chemotherapy doubled 5 year
 95 survival to around 16–21% (12); however, it can be increased to 38.6% after a successful resection at
 96 stage IIA.(13-16). Therefore, it is essential to identify effective biomarkers to enable early diagnosis
 97 and precisely predict the prognosis to recommend additional neoadjuvant and adjuvant therapies

98 Mucins (MUC) play crucial roles in carcinogenesis and tumor invasion in pancreatic tumors.
 99 The MUC gene product is post-translationally modified, most likely, through extensive
 100 O-glycosylation. MUC1 and MUC4 are large membrane-bound glycoproteins that are translated as
 101 single polypeptides. These mucins undergo intracellular autocatalytic proteolytic cleavage into two
 102 subunits that form stable non-covalent heterodimers that are transported to the cell surface. MUC2 is

gel-forming secretory mucin that is expressed in many organs, including the colon, small intestine, and respiratory tract (17-19). MUC1 contributes to oncogenesis by promoting the loss of epithelial cell polarity, promoting growth and survival pathways, activating receptor tyrosine kinase signaling pathways, and conferring resistance to the stress-induced cell death pathway (20,21). MUC1 cytoplasmic domain has been implicated in the regulation of the Wnt- β -catenin, p53 and NF- κ B pathways, all of which are linked to tumor progression (19,22). MUC4 plays an important role in epithelial cell proliferation and differentiation by inducing specific phosphorylation of ERBB2 and enhancing the expression of the cyclin-dependent kinase inhibitor p27, which inhibits cell cycle progression (23,24). The loss of *MUC2* might compromise signaling that contributes to epithelial differentiation and proliferation through contact with membrane-bound mucins or alters the differentiation program of the intestinal mucosa, resulting in an increased probability of tumor formation (25,26). In our previous studies, we used pancreatic tissue samples and several cancer cell lines to demonstrate that MUC1, MUC2, and MUC4 expression of mRNA and/or protein are regulated by hypoxia and/or DNA demethylation (27-30). However, it has not been reported how that DNA methylation of promoter region effect for intracellular localization of these three mucin proteins. Furthermore, *MUC1* and *MUC4* hypomethylation status are statistically associated with the development of distant metastasis, tumor stage, and overall survival (OS) for PDAC patients (31-34).

In machine learning, state-of-the-art classification algorithms, such as support vector machines (SVMs), neural networks (NNETs), or multinomial prediction models, are used for classification and regression analysis (35-37). Recently, several studies shows that in breast cancer, nasopharyngeal carcinoma, and non-small-cell lung cancer, several supervised learning methods, such as decision trees using data of cDNA or tissue microarray refine prognosis (38-41). However, it was not reported whether these machine learning models could use DNA methylation status to predict the outcome of

patients with pancreatic tumors.

Therefore, the aim of this study was to develop a machine learning-based prognostic classifier to predict OS with pancreatic cancers by integrating multiple DNA methylation statuses of three mucin genes. In the present study, we increased the number of patients that were evaluated from that of our previous study and performed a sequential analysis of mucin promoter-associated CpGs by methylation-specific electrophoresis (MSE) analysis.

Materials and Methods

Cell lines and culture

Human pancreatic cancer cell lines BxPC3, HPAF2, and Panc1; human colon adenocarcinoma cell lines Caco2 and LS174T; and human lung adenocarcinoma cell lines A427 and NCI-H292 were obtained from the American Type Culture Collection. HPAF2, LS174T, and Caco2 cells were cultured in Eagle's minimum essential medium (Sigma), PANC1 and A427 cells were cultured in Dulbecco's modified Eagle's medium (Sigma), and BxPC3 and NCI-H292 cells were cultured in RPMI 1640 medium (Sigma). The media were supplemented with 10% FBS (Invitrogen), 100 U/mL of penicillin (Sigma), and 100 µg/mL streptomycin (Sigma). Hypoxic culture conditions were achieved with a multi-gas incubator containing a gas mixture of 94% N₂, 5% CO₂, and 1% O₂ (ASTEC).

Clinical samples

Pancreatic tissue samples

We obtained 300 surgically resected tissues (approximately 2 × 2 × 2 mm in size) from neoplastic and non-neoplastic areas of 191 patients. Table 1 summarizes the clinicopathologic features of the 114 neoplastic samples and 186 non-neoplastic samples (including 109 paired samples). We collected 125 patient samples (48 neoplastic samples and 120 non-neoplastic samples,

including 43 paired samples) from Kagoshima University, Japan, from August 2007 to May 2014 and 66 patient samples (all non-neoplastic and neoplastic paired samples) from Ulm University, Germany, from February 2001 to February 2013. The non-neoplastic tissues were collected around the resection stump. These samples were checked for pathological diagnosis using intraoperative frozen and FFPE tissue sections. On the other hand, neoplastic tissues were macroscopically collected, processed for FFPE tissue sections, and diagnosed by a board-certified pathologist. The clinical features used in this study were TNM, age, American Society of Anaesthesiologists (ASA) physical status classification score, the presence or absence of co-morbidities, and preoperative chemotherapy. Almost all patients had not undergone radiotherapy before surgery; therefore, we removed the information about radiotherapy in statistical analysis.

Ethics statement

This study was conducted in accordance with the guiding principles of the Declaration of Helsinki. The ethical committees of both Kagoshima University Hospital and Ulm University Hospital approved sample collection, and informed written consent was obtained from each patient. All studies using human materials in this article were approved by the Ethical Committee of Kagoshima University Hospital (revised 20–82, revised 22–127, and 26-145).

Extraction and quantification of mRNA

Total RNA was extracted from the cell lines, human pancreatic tissues, and pancreatic juices using an RNeasy Mini Kit (QIAGEN). Then, the total RNA (1 µg) was reverse transcribed with a High-Capacity RNA-to-cDNA Kit (Applied Biosystems), and real-time reverse transcription-PCR was performed on a Roche LightCycler® 96 System using FastStart Essential DNA Green Master (Roche). Gene expression was normalized to the β-actin mRNA level in each sample. The data were normalized using the NCI-H292 cell line, and the A427 cell line was used as a negative control. Primer sets are shown in previously study (31,33,34).

Extraction of DNA and bisulfite modification

DNA from the cell lines, pancreatic tissues, and pancreatic juices was extracted using a DNeasy Tissue System (QIAGEN). The bisulfite modification of the genomic DNA was carried out using an EpiTect Bisulfite Kit (QIAGEN). The purification of PCR products was carried out using a Wizard® SV Gel and PCR Clean-Up System (Promega KK) (31,33,34).

MSE analysis

MSE analysis was performed using previously described methods. Briefly, the target DNA fragments were amplified by nested PCR using bisulfite-treated DNA with the primer sets detailed in previously study (31-34). Then, the PCR products were run on a polyacrylamide gel with a linear denaturant gradient at 60 °C at 70 V for 14 h using a D-Code Universal Mutation Detection System (BioRad Laboratories). Band intensity was quantified by Image J software (National Institutes of Health). The demethylation index was calculated as the proportion of the highest band intensity/total band intensity of the sample. Subsequently, the demethylation index for each sample was normalized using data from a hypomethylated and hypermethylated cell line. Cell lines with hyper- and hypomethylation of MUC1 and MUC4 (Caco2 and LS174T, respectively) were used as control standards. We performed MSE analysis on duplicate samples.

Statistical analysis and prediction model construction

Data were analyzed using the “R” computing environment version 3.5.2 (42). The normality of the data distribution was evaluated using the Kolmogorov–Smirnov test. Differences between groups were analyzed using Welch’s *t*-test. A nonparametric test of group difference was performed using the Mann–Whitney *U* test. Survival rate analysis was evaluated using Cox’s proportional hazards model. Hierarchical cluster analysis on a set of dissimilarities was performed using the reproduce package of R. A p-value < 0.05 was considered statistically significant.

Construction of an SVM classifier

We constructed the SVM classifier using the kernlab package, including the ksvm function (37). Methylation analyzed data mining, such as scaling and centering, were performed using the scale function. C-classification with the linear vanilladot kernel, Gaussian radial basis function, polynomial kernel, hyperbolic tangent kernel, Bessel function of the first kind kernel, Laplace radial basis kernel, and ANOVA radial basis kernel were used in training and predicting. In this study, we used default hyper-parameters and cost of constraints violation for preparing the prediction models. The quality of the model was assessed using 5-fold cross validation of the training data.

Construction of an NNET and a multinomial classifier

Neural networks provide a flexible, non-linear extension of multiple logistic regressions to perform classification, pattern recognition, and prediction modeling. The nnet package, including the nnet function, was used to construct an NNET classifier (43). The classifier parameters were set with a weight decay of 0.1, 2 units in the hidden layer and 100 as the maximum number of iterations. To construct the multinomial log-linear model (MU) via the NNET classifier, the multinom function was employed.

Results

Mucin expression characteristic analysis and prognosis differential

Unsupervised hierarchical clustering analysis was performed for three mucin gene expression data sets including neoplastic and non-neoplastic pancreatic regions (Figure 1A). The samples were divided into two clusters according to the clustering results. Cluster 1 showed a significantly higher expression level of the three mucin genes as MUC1, MUC2, and MUC4 compared with Cluster 2 (all $p < 0.001$; Figure 1B). Cluster 1 also demonstrated a significantly poorer prognosis than Cluster 2 (HR = 0.52, $p < 0.001$; Figure 1C). As shown in Table 2, this clustering analysis showed the same result in the non-neoplastic tissue ($p < 0.001$). In non-neoplastic region, Cluster 1 showed a

significantly poorer prognosis than Cluster 2 (HR = 0.49, $p = 0.012$; Supplementary Figure 1A).

Similarly, Cluster 1 showed a significantly poorer prognosis than Cluster 2 in the neoplastic region (HR = 0.53, $p = 0.032$; Supplementary Figure 1B).

Mucin DNA methylation characteristics and prognosis differential

Unsupervised hierarchical clustering analysis was performed for three mucin gene methylation data sets including neoplastic and non-neoplastic pancreatic regions (Figure 2A). According to the clustering results, the samples were divided into four clusters. Cluster 1 showed a significantly higher hypomethylation status of *MUC1* and *MUC4* genes than the other clusters ($p = 0.003$ and $p < 0.001$, respectively; Figure 2B). Cluster 1 also demonstrated a significantly higher expression level of *MUC1*, *MUC2*, and *MUC4* than the other clusters ($p = 0.004$, $p = 0.002$, and $p = 0.015$, respectively; Supplementary Figure 2). Furthermore, Cluster 1 demonstrated a significantly poorer prognosis than the other clusters (HR = 0.33, $p = 0.012$; Figure 2C). As shown in Table 2, this clustering analysis did not identify a significant difference in the non-neoplastic ratio between Cluster 1 and the other clusters ($p < 0.357$). In the non-neoplastic regions, Cluster 1 showed a significantly poorer prognosis than the other clusters (HR = 0.49, $p = 0.010$; Supplementary Figure 3A). Similarly, in the neoplastic regions, Cluster 1 also showed a significantly poorer prognosis than the other clusters (HR=0.48, $p=0.014$; Supplementary Figure 3B).

Prediction model classifier and survival in the cross-validation test

SVM prediction model and performance evaluation

In the leave-one-out cross-validation (LOOCV) test, The SVM classifier model by C-classification with ANOVA RBF kernel function using the data set for DNA methylation of three mucins, including both non-neoplastic and neoplastic regions, showed good classification for prognosis after surgery (Figure 3A). In addition, ANOVA RBF kernel as a non-linear kernel showed the best classification than the other kernels, including linear kernel (Supplementary Table 1). This

SVM classification showed no significant differences in the non-neoplastic ratio between high-risk positive and negative groups (Supplementary Table 2). In non-neoplastic regions, the high-risk group indicated by SVM showed a significantly poorer prognosis after surgery than the negative group (Figure 3B). Similarly, in neoplastic regions, the high-risk group indicated by SVM also showed a significantly poorer prognosis after surgery than the negative group (Figure 3C). In univariate analysis, patients who were classified as positive by the SVM model were associated with a significantly poorer OS. Multivariate Cox regression analysis after adjustment for clinicopathologic variables, such as ASA score, preoperative chemotherapy, co-morbidities, and TNM stage, revealed that the SVM classifier remained a powerful and independent prognostic factor for OS in the LOOCV test (Supplementary Table 3). Following several k-fold cross-validation tests, we performed 3-, 4-, and 10-fold cross validations in this study. The SVM classifier showed good classification for prognosis after surgery in both regions (Supplementary Figure 4 and Table 3). Multivariate and/or univariate analysis revealed that SVM classifier was an independent prognostic factor (Supplementary Table 3).

NNET prediction model and performance evaluation

In the LOOCV test, the NNET classifier model using the data set for DNA methylation of three mucins, including both non-neoplastic and neoplastic regions, showed good classification for prognosis after surgery (Figure 3D). This NNET classification showed no significant difference in the non-neoplastic ratio between high-risk positive and negative groups (Supplementary Table 2). In non-neoplastic regions, the high-risk group indicated by NNET showed a significantly poorer prognosis after surgery than the negative group (Figure 3E). Similarly, in neoplastic regions, the high-risk group indicated by NNET also showed a significantly poorer prognosis after surgery than the negative group (Figure 3F). In univariate analysis, patients who were classified as high-risk positive by the NNET model were associated with a significantly poorer OS. Multivariate Cox

regression analysis after adjustment for clinicopathologic variables, such as ASA score, preoperative chemotherapy, co-morbidities, and TNM stage, revealed that the NNET classifier remained a powerful and independent prognostic factor for OS in the LOOCV test (Supplementary Table 3). After several k-fold cross-validation tests, the NNET classifier showed a tendency toward good classification for prognosis after surgery in both non-neoplastic and neoplastic regions (Supplementary Figure 4 and Table 3). Multivariate and/or univariate analysis showed that NNET classifier was independent from prognostic factors (Supplementary Table 3).

MU prediction model and performance evaluation

In the LOOCV test, The MU classifier model using the data set for DNA methylation of three mucins, including both non-neoplastic and neoplastic regions, did not show good classification for prognosis after surgery (Figure 3G). However, this MU classification revealed a significant difference in the non-neoplastic ratio between the high-risk positive and negative groups ($p < 0.043$, Supplementary Table 2). In neoplastic regions, the negative group indicated by MU showed a significantly poorer prognosis after surgery than the high-risk group (Figure 3H), but not in non-neoplastic regions (Figure 3I). The multivariate Cox regression analysis after adjustment for clinicopathologic variables such as TNM stage revealed that the MU classifier was an independent prognostic factor for OS in the LOOCV test. However, in univariate analysis, patients who were positive according to the MU model were not associated with significantly poorer OS (Supplementary Table 3). In k-fold cross-validation test, the MU classifier demonstrated good classification ability for prognosis after surgery in the neoplastic region but not in non-neoplastic region (Supplementary Figure 4 and Table 3).

Evaluation of prediction models using training cohort in the test cohort

To evaluate whether the prediction models constructed using test data sets could detect high-risk groups in other training data set, we split the total data set into two groups: training and test data sets.

These two groups showed almost similar distribution in biological analysis results and clinicopathological features (shown in Supplementary Table 4). SVM and NNET classifiers exhibited significantly good classification ability for prognosis after surgery in not only the neoplastic region but also the non-neoplastic region (Figure 4). Multivariate and/or univariate analysis revealed that these classifiers were independent from prognostic factors (shown in Supplementary Table 3). The MU classifier demonstrated good classification ability for prognosis after surgery in the neoplastic region but not in non-neoplastic region (Figure 4G, H and I).

Discussion

PDAC is an aggressive malignancy with an extremely poor prognosis due to delayed diagnosis, early metastasis, and resistance to most cytotoxic agents (1,13). Thus, it is critical to establish new diagnostic, prognostic, and therapeutic biomarkers. It has been previously demonstrated that mucin gene expression (including MUC1, MUC2, MUC3, MUC4, and MUC5AC) is regulated by DNA methylation at promoter regions in cancer cell lines (27-30). In the present study, scientific computer prediction model-based SVM and NNET methods were able to classify between a good and poor prognosis after surgery on pancreatic cancers. These predicted models were constructed from the methylation status of three mucin genes (*MUC1*, *MUC2*, and *MUC4*), which had all previously demonstrated a significant difference in mucin expression levels in pancreatic cancers.

To construct a clinical test to predict prognosis after surgery, we used an expression data set including non-neoplastic and neoplastic pancreatic tissues. To determine whether the expression levels of MUC1, MUC2, and MUC4 mRNA could distinguish a poor prognosis from data, we performed unsupervised hierarchical clustering analysis using real-time PCR data. These mucin mRNA clustering analyses significantly separated a poor prognosis cluster. This selected cluster had higher expression levels of the three mucins than the other cluster. In our recent histological studies,

we reported that mucin gene expression, particularly MUC4, was an independent indicator of worse prognosis in PDAC (5-9). MUC2 is reported as a tumor suppressor gene, and loss of MUC2 promotes tumor progression in colon (44). However, another study showed that MUC2 expression may have a poor prognostic value for differentiated adenocarcinomas in pancreas (45). The relationship between high expression of three mucins mRNA and poor prognosis supported these recent pathological studies. The selected cluster showed a significant difference in the non-neoplastic: neoplastic content ratio compared with the other cluster. In the neoplastic region data, the selected cluster showed a poorer prognosis than the other cluster. Interestingly, in the non-neoplastic region data, this selected cluster also showed a poorer prognosis than the other cluster. These results suggested that scientific computer methods could identify a poor prognosis using the combined data of mucin gene expression even if there was a mixture of non-neoplastic and neoplastic tissue.

A previous study revealed that an analysis of the DNA methylation status in promoters of *MUC1*, *MUC2*, and *MUC4* (MSE analysis of pancreatic juice samples) could differentiate between gastric-type IPMN, intestinal-type IPMN, other-type IPMN, and PDAC (33,34). The correlation between hypomethylation of the promoter and high expression levels of mucin mRNA in pancreatic tissue was shown (28). Furthermore, we have proposed that aberrant methylation of MUC1 and MUC4 promoters are potential prognostic biomarkers for PDAC and suggested further MSE analysis of human clinical samples to determine its utility for the early diagnosis of pancreatic neoplasms and for stratifying patients with respect to modes of treatment (31). Thus, we used a methylation data set including non-neoplastic and neoplastic pancreatic tissue. To establish whether the methylation levels of these three mucins could distinguish a poor prognosis, we performed unsupervised hierarchical clustering analysis using MSE analysis data. The clustering analysis of the three mucin genes' methylation data significantly separated a poor prognosis cluster. This cluster group had a

higher hypomethylation level of MUC1 and MUC4 than the other clusters. Moreover, this cluster demonstrated a higher neoplastic-including ratio than the other clusters, and in neoplastic regions, this cluster showed a poorer prognosis than the other clusters. Interestingly, in the non-neoplastic region analysis, this cluster also showed a poorer prognosis than the other clusters. These results suggested that scientific computer methods could provide a model to identify a poor prognosis using the combined data of mucin gene methylation.

To evaluate whether machine learning prediction models using state-of-the-art classification algorithms such as SVM and NNET could distinguish between a poor prognosis group and others, we constructed prediction models. In LOOCV tests and k-fold cross-validation tests to evaluate prediction ability, the SVM and NNET models could significantly judge the identified high-risk group as having a poor prognosis, but the MU model could not. Multivariate and univariate analyses showed that the prediction of high-risk by SVM or NNET model was a prognostic factor significantly independent from TNM score, ASA score, preoperative chemotherapy, and co-morbidities. For the neoplastic tissue analysis data, the SVM, NNET, and MU models could identify the high-risk group. Interestingly, even in the non-neoplastic tissue analysis data, the SVM and NNET model-selected high-risk group had a significantly poorer prognosis than others, similar to the neoplastic analysis. Therefore, these results suggested that the prediction models using cytological specimens and liquid biopsy samples, which are mixture of non-neoplastic and neoplastic cells, might be applicable to high-risk screening in PDAC.

When the model has high variance and low bias, such as that showing too much optimization for the training data set, the prediction model has low prediction performance for data that has never been learned. To evaluate whether the prediction models constructed from the test data set could detect high-risk groups in other training data sets, we splitted the total data set into two groups having almost similar distribution of biological status and clinicopathological features. The SVM

and NNET classifiers could significantly distinguish the high-risk group, which has poor prognosis in the test group that has never been learned, but the MU model could not. Multivariate and univariate analyses showed that the prediction of high-risk by SVM and NNET classifiers was a prognostic factor significantly independent from TNM score, ASA score, preoperative chemotherapy, and co-morbidities. These results suggested that SVM and NNET classifiers have low variance and demonstrated the high ability to distinguish a poor prognosis. Although a prospective, much larger, and multicenter randomized trial would be necessary to validate our results, it is suggested that the SVM- or NNET-based prediction models could provide a clinical risk test to predict the prognosis after surgery using *MUC1*, *MUC2* and *MUC4* gene methylation analyses. Even though these SVM and NNET classifiers were a highly accurate predictor of OS, we are aware that other biomarkers may extend the precision and predictive value of the classifiers, and new markers are being identified and new techniques developed every year (46,47). Thus, the SVM and NNET classifiers may be further improved by including additional markers.

In summary, the present study demonstrated that machine learning prediction models, based on SVM and NNET, could accurately distinguish pancreatic cancer patients after surgery with substantially different OS. A further study is needed to expand the clinical sample spectrum, where these classifiers based on SVM or NNET might work for decision-making regarding follow-up scheduling after surgery.

Acknowledgments

We thank Orié Iwatani, Yoshie Jitoh, and Yukari Nishida for their assistance with clinical sampling and excellent technical assistance with immunohistochemistry. We would like to thank Enago (www.enago.jp) for the English language review.

REFERENCES

1. Adamska A, Domenichini A, Falasca M. Pancreatic Ductal Adenocarcinoma: Current and Evolving Therapies. *Int J Mol Sci* **2017**;18(7) doi 10.3390/ijms18071338.
2. Lu D, Wang J, Shi X, Yue B, Hao J. AHNK2 is a potential prognostic biomarker in patients with PDAC. *Oncotarget* **2017**;8(19):31775-84 doi 10.18632/oncotarget.15990.
3. Toucheffeu Y, Le Rhun M, Coron E, Alamdari A, Heymann MF, Mosnier JF, *et al.* Endoscopic ultrasound-guided fine-needle aspiration for the diagnosis of solid pancreatic masses: the impact on patient-management strategy. *Aliment Pharmacol Ther* **2009**;30(10):1070-7 doi 10.1111/j.1365-2036.2009.04138.x.
4. Yadav D, Lowenfels AB. The epidemiology of pancreatitis and pancreatic cancer. *Gastroenterology* **2013**;144(6):1252-61 doi 10.1053/j.gastro.2013.01.068.
5. Higashi M, Goto M, Saitou M, Shimizu T, Rousseau K, Batra SK, *et al.* Immunohistochemical study of mucin expression in periampullary adenomyoma. *J Hepatobiliary Pancreat Sci* **2010**;17(3):275-83 doi 10.1007/s00534-009-0176-5.
6. Higashi M, Yokoyama S, Yamamoto T, Goto Y, Kitazono I, Hiraki T, *et al.* Mucin expression in endoscopic ultrasound-guided fine-needle aspiration specimens is a useful prognostic factor in pancreatic ductal adenocarcinoma. *Pancreas* **2015**;44(5):728-34 doi 10.1097/MPA.0000000000000362.
7. Yonezawa S, Higashi M, Yamada N, Goto M. Precursor lesions of pancreatic cancer. *Gut Liver* **2008**;2(3):137-54 doi 10.5009/gnl.2008.2.3.137.
8. Yonezawa S, Higashi M, Yamada N, Yokoyama S, Goto M. Significance of mucin expression in pancreatobiliary neoplasms. *J Hepatobiliary Pancreat Sci* **2010**;17(2):108-24 doi 10.1007/s00534-009-0174-7.
9. Yonezawa S, Nakamura A, Horinouchi M, Sato E. The expression of several types of mucin is related to the biological behavior of pancreatic neoplasms. *J Hepatobiliary Pancreat Surg* **2002**;9(3):328-41 doi 10.1007/s005340200037.
10. Furukawa T, Hatori T, Fujita I, Yamamoto M, Kobayashi M, Ohike N, *et al.* Prognostic relevance of morphological types of intraductal papillary mucinous neoplasms of the pancreas. *Gut* **2011**;60(4):509-16 doi 10.1136/gut.2010.210567.
11. Strobel O, Neoptolemos J, Jager D, Buchler MW. Optimizing the outcomes of pancreatic cancer surgery. *Nat Rev Clin Oncol* **2019**;16(1):11-26 doi 10.1038/s41571-018-0112-1.
12. Neoptolemos JP, Palmer DH, Ghaneh P, Psarelli EE, Valle JW, Halloran CM, *et al.* Comparison of adjuvant gemcitabine and capecitabine with gemcitabine monotherapy in patients with resected pancreatic cancer (ESPAC-4): a multicentre,

- 427 open-label, randomised, phase 3 trial. *Lancet* **2017**;389(10073):1011-24 doi
 428 10.1016/S0140-6736(16)32409-6.
- 429 13. Egawa S, Toma H, Ohigashi H, Okusaka T, Nakao A, Hatori T, *et al.* Japan
 430 Pancreatic Cancer Registry; 30th year anniversary: Japan Pancreas Society.
 431 *Pancreas* **2012**;41(7):985-92 doi 10.1097/MPA.0b013e318258055c.
- 432 14. Flejou JF. [WHO Classification of digestive tumors: the fourth edition]. *Ann Pathol*
 433 **2011**;31(5 Suppl):S27-31 doi 10.1016/j.annpat.2011.08.001.
- 434 15. Isaji S, Kawarada Y, Uemoto S. Classification of pancreatic cancer: comparison of
 435 Japanese and UICC classifications. *Pancreas* **2004**;28(3):231-4.
- 436 16. Matsuda T, Ajiki W, Marugame T, Ioka A, Tsukuma H, Sobue T, *et al.*
 437 Population-based survival of cancer patients diagnosed between 1993 and 1999 in
 438 Japan: a chronological and international comparative study. *Jpn J Clin Oncol*
 439 **2011**;41(1):40-51 doi 10.1093/jjco/hyq167.
- 440 17. Hollingsworth MA, Swanson BJ. Mucins in cancer: protection and control of the cell
 441 surface. *Nat Rev Cancer* **2004**;4(1):45-60 doi 10.1038/nrc1251.
- 442 18. Kaur S, Kumar S, Momi N, Sasson AR, Batra SK. Mucins in pancreatic cancer and
 443 its microenvironment. *Nat Rev Gastroenterol Hepatol* **2013**;10(10):607-20 doi
 444 10.1038/nrgastro.2013.120.
- 445 19. Kufe DW. Mucins in cancer: function, prognosis and therapy. *Nat Rev Cancer*
 446 **2009**;9(12):874-85 doi 10.1038/nrc2761.
- 447 20. Ahmad R, Raina D, Trivedi V, Ren J, Rajabi H, Kharbanda S, *et al.* MUC1
 448 oncoprotein activates the IkappaB kinase beta complex and constitutive NF-kappaB
 449 signalling. *Nat Cell Biol* **2007**;9(12):1419-27 doi 10.1038/ncb1661.
- 450 21. Pochampalli MR, el Bejjani RM, Schroeder JA. MUC1 is a novel regulator of ErbB1
 451 receptor trafficking. *Oncogene* **2007**;26(12):1693-701 doi 10.1038/sj.onc.1209976.
- 452 22. Moniaux N, Andrianifahanana M, Brand RE, Batra SK. Multiple roles of mucins in
 453 pancreatic cancer, a lethal and challenging malignancy. *Br J Cancer*
 454 **2004**;91(9):1633-8 doi 10.1038/sj.bjc.6602163.
- 455 23. Mercogliano MF, De Martino M, Venturutti L, Rivas MA, Proietti CJ, Inurrigarro G,
 456 *et al.* TNFalpha-Induced Mucin 4 Expression Elicits Trastuzumab Resistance in
 457 HER2-Positive Breast Cancer. *Clin Cancer Res* **2017**;23(3):636-48 doi
 458 10.1158/1078-0432.CCR-16-0970.
- 459 24. Mukhopadhyay P, Lakshmanan I, Ponnusamy MP, Chakraborty S, Jain M, Pai P, *et*
 460 *al.* MUC4 overexpression augments cell migration and metastasis through EGFR
 461 family proteins in triple negative breast cancer cells. *PLoS One* **2013**;8(2):e54455 doi
 462 10.1371/journal.pone.0054455.

25. Tadesse S, Corner G, Dhima E, Houston M, Guha C, Augenlicht L, *et al.* MUC2 mucin deficiency alters inflammatory and metabolic pathways in the mouse intestinal mucosa. *Oncotarget* **2017**;8(42):71456-70 doi 10.18632/oncotarget.16886.
26. Yang K, Popova NV, Yang WC, Lozonschi I, Tadesse S, Kent S, *et al.* Interaction of Muc2 and Apc on Wnt signaling and in intestinal tumorigenesis: potential role of chronic inflammation. *Cancer Res* **2008**;68(18):7313-22 doi 10.1158/0008-5472.CAN-08-0598.
27. Yamada N, Hamada T, Goto M, Tsutsumida H, Higashi M, Nomoto M, *et al.* MUC2 expression is regulated by histone H3 modification and DNA methylation in pancreatic cancer. *Int J Cancer* **2006**;119(8):1850-7 doi 10.1002/ijc.22047.
28. Yamada N, Nishida Y, Tsutsumida H, Goto M, Higashi M, Nomoto M, *et al.* Promoter CpG methylation in cancer cells contributes to the regulation of MUC4. *Br J Cancer* **2009**;100(2):344-51 doi 10.1038/sj.bjc.6604845.
29. Yamada N, Nishida Y, Tsutsumida H, Hamada T, Goto M, Higashi M, *et al.* MUC1 expression is regulated by DNA methylation and histone H3 lysine 9 modification in cancer cells. *Cancer Res* **2008**;68(8):2708-16 doi 10.1158/0008-5472.CAN-07-6844.
30. Yonezawa S, Goto M, Yamada N, Higashi M, Nomoto M. Expression profiles of MUC1, MUC2, and MUC4 mucins in human neoplasms and their relationship with biological behavior. *Proteomics* **2008**;8(16):3329-41 doi 10.1002/pmic.200800040.
31. Yokoyama S, Higashi M, Kitamoto S, Oeldorf M, Knippschild U, Kornmann M, *et al.* Aberrant methylation of MUC1 and MUC4 promoters are potential prognostic biomarkers for pancreatic ductal adenocarcinomas. *Oncotarget* **2016**;7(27):42553-65 doi 10.18632/oncotarget.9924.
32. Yokoyama S, Higashi M, Tsutsumida H, Wakimoto J, Hamada T, Wiest E, *et al.* TET1-mediated DNA hypomethylation regulates the expression of MUC4 in lung cancer. *Genes Cancer* **2017**;8(3-4):517-27 doi 10.18632/genesandcancer.139.
33. Yokoyama S, Kitamoto S, Higashi M, Goto Y, Hara T, Ikebe D, *et al.* Diagnosis of pancreatic neoplasms using a novel method of DNA methylation analysis of mucin expression in pancreatic juice. *PLoS One* **2014**;9(4):e93760 doi 10.1371/journal.pone.0093760.
34. Yokoyama S, Kitamoto S, Yamada N, Houjou I, Sugai T, Nakamura S, *et al.* The application of methylation specific electrophoresis (MSE) to DNA methylation analysis of the 5' CpG island of mucin in cancer cells. *BMC Cancer* **2012**;12:67 doi 10.1186/1471-2407-12-67.
35. Huang S, Cai N, Pacheco PP, Narrandes S, Wang Y, Xu W. Applications of Support Vector Machine (SVM) Learning in Cancer Genomics. *Cancer Genomics Proteomics*

- 499 **2018**;15(1):41-51 doi 10.21873/cgp.20063.
- 500 36. Ozyildirim BM, Avci M. Generalized classifier neural network. Neural Netw
 501 **2013**;39:18-26 doi 10.1016/j.neunet.2012.12.001.
- 502 37. Alexandros K, Alexandros S, Kurt H, Achim Z. kernlab - An S4 Package for Kernel
 503 Methods in R. Journal of statistical software **2004**;11(9) doi 10.18637/jss.v011.i09.
- 504 38. Chen HY, Yu SL, Chen CH, Chang GC, Chen CY, Yuan A, *et al*. A five-gene signature
 505 and clinical outcome in non-small-cell lung cancer. N Engl J Med **2007**;356(1):11-20
 506 doi 10.1056/NEJMoA060096.
- 507 39. Jiang Y, Liu W, Li T, Hu Y, Chen S, Xi S, *et al*. Prognostic and Predictive Value of
 508 p21-activated Kinase 6 Associated Support Vector Machine Classifier in Gastric
 509 Cancer Treated by 5-fluorouracil/Oxaliplatin Chemotherapy. EBioMedicine
 510 **2017**;22:78-88 doi 10.1016/j.ebiom.2017.06.028.
- 511 40. Wang HY, Sun BY, Zhu ZH, Chang ET, To KF, Hwang JS, *et al*. Eight-signature
 512 classifier for prediction of nasopharyngeal [corrected] carcinoma survival. J Clin
 513 Oncol **2011**;29(34):4516-25 doi 10.1200/JCO.2010.33.7741.
- 514 41. Jiang Y, Xie J, Han Z, Liu W, Xi S, Huang L, *et al*. Immunomarker Support Vector
 515 Machine Classifier for Prediction of Gastric Cancer Survival and Adjuvant
 516 Chemotherapeutic Benefit. Clin Cancer Res **2018**;24(22):5574-84 doi
 517 10.1158/1078-0432.CCR-18-0848.
- 518 42. Ihaka R, Gentleman R. R: A Language for Data Analysis and Graphics. Journal of
 519 computational and graphical statistics **1996**;5(3):15.
- 520 43. Venables WN, Ripley BD. Modern Applied Statistics with S. springer; 2002.
- 521 44. Velcich A, Yang W, Heyer J, Fragale A, Nicholas C, Viani S, *et al*. Colorectal cancer in
 522 mice genetically deficient in the mucin Muc2. Science **2002**;295(5560):1726-9 doi
 523 10.1126/science.1069094.
- 524 45. Takikita M, Altekruse S, Lynch CF, Goodman MT, Hernandez BY, Green M, *et al*.
 525 Associations between selected biomarkers and prognosis in a population-based
 526 pancreatic cancer tissue microarray. Cancer Res **2009**;69(7):2950-5 doi
 527 10.1158/0008-5472.CAN-08-3879.
- 528 46. Bernard V, Kim DU, San Lucas FA, Castillo J, Allenson K, Mulu FC, *et al*.
 529 Circulating Nucleic Acids Are Associated With Outcomes of Patients With
 530 Pancreatic Cancer. Gastroenterology **2019**;156(1):108-18 e4 doi
 531 10.1053/j.gastro.2018.09.022.
- 532 47. Krantz BA, O'Reilly EM. Biomarker-Based Therapy in Pancreatic Ductal
 533 Adenocarcinoma: An Emerging Reality? Clin Cancer Res **2018**;24(10):2241-50 doi
 534 10.1158/1078-0432.CCR-16-3169.

Figure legends

Figure 1: Cluster analysis of the mRNA expression level of mucin genes. (A) Tree generated by cluster analysis of neoplastic and non-neoplastic pancreas tissues for the expression levels of MUC1, MUC2, and MUC4 mRNAs compared to each control cell line. High mRNA expression levels are indicated by in red and low levels in blue. (B) Comparison of expression levels of MUC1, MUC2, and MUC4 mRNAs between Cluster 1 and Cluster 2. Expression levels show relative quantification (\log_{10}). (C) Cox proportional hazard regression analysis of a comparison between Cluster 1 and Cluster 2. Red solid line: Cluster 1, black dashed line: Cluster 2.

Figure 2: Cluster analysis of the methylation status of mucin genes. (A) Tree generated by cluster analysis of neoplastic and non-neoplastic pancreas tissues from the methylation status of *MUC1*, *MUC2*, and *MUC4* genes evaluated by MSE analysis. Hypomethylation is indicated in red and hypermethylation in blue. (B) Comparison of the methylation status of *MUC1*, *MUC2*, and *MUC4* genes between Cluster 1 and the other clusters. (C) Cox proportional hazard regression analysis on a comparison between Cluster 1 and other clusters. Red solid line: Cluster 1, black dashed line: other clusters.

Figure 3: Prognosis prediction by machine learning classifier in the LOOCV test. Cox proportional hazard regression analysis on a comparison between the positive and negative groups as selected by each classifier. Solid line: predicted high-risk group (positive), dashed line: other groups (negative). (A) Classification by SVM model for all samples. (B) Classification by SVM model in non-neoplastic tissues. (C) Classification by SVM model in neoplastic tissues. (D) Classification by NNET model for all samples. (E) Classification by NNET model in non-neoplastic tissues. (F)

Classification by NNET model in neoplastic tissues. (G) Classification by MU model for all samples.
 (H) Classification by MU model in non-neoplastic tissues. (I) Classification by MU model in
 neoplastic tissues.

Figure 4: Prognosis prediction by machine learning classifier in the test cohort. Cox proportional
 hazard regression analysis for the comparison between the positive and negative groups selected by
 each classifier. Solid line: predicted high-risk group (positive), dashed line: other groups (negative).
 (A) Classification by SVM model for all samples. (B) Classification by SVM model in
 non-neoplastic tissues. (C) Classification by SVM model in neoplastic tissues. (D) Classification by
 NNET model for all samples. (E) Classification by NNET model in non-neoplastic tissues. (F)
 Classification by NNET model in neoplastic tissues. (G) Classification by MU model for all samples.
 (H) Classification by MU model in non-neoplastic tissues. (I) Classification by MU model in
 neoplastic tissues.

Supplementary Figure 1: Cox proportional hazard regression analysis on a comparison between
 Cluster 1 and Cluster 2 selected by cluster analysis of mRNA expression levels of MUC1, MUC2,
 and MUC4 in non-neoplastic regions (A) and neoplastic regions (B). Red solid line: Cluster 1, black
 dashed line: Cluster 2.

Supplementary Figure 2: Difference in the expression level of mucins between Cluster 1 and other
 clusters selected by cluster analysis of the methylation status of mucin genes *MUC1*, *MUC2*, and
MUC4. Expression levels show relative quantification (\log_{10}).

Supplementary Figure 3: Cox proportional hazard regression analysis on a comparison between

Cluster 1 and other clusters as selected by cluster analysis of the methylation status of mucin genes
MUC1, *MUC2*, and *MUC4* in non-neoplastic regions (A) and neoplastic regions (B). Red solid line:
 Cluster 1, black dashed line: other clusters.

Supplementary Figure 4: Prognosis prediction by machine learning classifier in the k-fold CV test.
 Cox proportional hazard regression analysis for the comparison between the positive and negative
 groups selected by each classifier. Red solid line: predicted high-risk group (positive), blue line:
 other groups (negative). (A) 3-fold cross-validation test. (B) 4-fold cross-validation test. The insets
 show *p*-value (*p*) and hazard ratio (HR) of each test.

Table 1. Patient and tumor characteristics in the study

Age

median, (Male/Female) 66, (65.3/66.8) year

Observation Period (OP)

median, (Male/Female) 22.4, (23.5/21.2) months

ASA score

median, (Male/Female) 2, (2/2)

		Number of cases (non-/neoplastic)	OP median	Age median
Stage	non	30 (29/20)	38.4	59.2
	IA	18 (17/2)	22.1	67.3
	IB	12 (12/4)	25.4	66.3
	IIA	36 (35/26)	16.2	65.8
	IIB	64 (63/43)	19.6	69.1
	III	3 (3/3)	14.3	68.3
	IV	6 (5/6)	11.5	66.5
	NA	22 (22/10)	15	54
T	0	21 (20/11)	24.4	59.3
	1	20 (19/3)	25.1	67.5
	2	17 (17/4)	23.6	66.6
	3	98 (96/73)	17.3	67.8
	4	4 (3/4)	13.8	68.3
	NA	31 (31/19)	60.8	57.9
	NA	31 (31/19)	60.8	57.9
N	0	90 (86/46)	20.3	64.4
	1	69 (68/48)	19.2	69.2
	NA	31 (31/19)	60.8	57.9
M	0	153 (149/88)	19.8	66.5
	1	6 (5/6)	11.5	66.5
	NA	32 (32/20)	60.8	58.7

Number of cases (non-/neoplasm)

	presence	absence
co-morbidities	84 (81/62)	87 (87/43)
preoperative chemotherapy	130 (127/64)	41 (41/41)

Table 2. Clustering analysis result and clinicopathological data

1. mRNA data clustering analysis

2. Methylation data clustering analysis

		Cluster 1	Cluster 2	<i>p</i> value			cluster 1	others	<i>p</i> value
Age	mean \pm sd	66.93 \pm 11.2	65.74 \pm 9.5	0.756	Age	mean \pm sd	65.63 \pm 10.6	66.32 \pm 10.7	0.632
	Observation Period (OP)					Observation Period (OP)			
	mean \pm sd	23.72 \pm 23.4	23.9 \pm 28.5	0.267		mean \pm sd	30.04 \pm 34.1	21.29 \pm 19.8	0.709
	Tissue region (n)					Tissue region (n)			
	non-neoplasm	36	150	<0.001		non-neoplasm	49	137	0.357
	neoplasm	63	51			neoplasm	36	78	
ASA score	mean \pm sd	1.81 \pm 0.9	2.33 \pm 0.8	0.142	ASA score	mean \pm sd	2.57 \pm 0.6	1.73 \pm 0.8	<0.001
Stage (n)	non	12	37	<0.001	Stage (n)	non	14	35	0.018
	IA	0	19			IA	1	18	
	IB	3	13			IB	5	11	
	IIA	22	39			IIA	25	36	
	IIB	43	63			IIB	28	78	
	III	3	3			III	4	2	
	IV	3	8			IV	2	9	
	NA	1	31			NA	6	26	
T (n)	0	4	27	0.004	T (n)	0	2	29	<0.001
	1	2	20			1	1	21	
	2	7	14			2	5	16	
	3	62	107			3	55	114	
	4	3	4			4	4	3	
	NA	9	41			NA	18	32	
N (n)	0	30	102	0.002	N (n)	0	35	97	0.413
	1	48	70			1	32	86	
	NA	9	41			NA	18	32	
M (n)	0	75	162	0.115	M (n)	0	64	173	0.29
	1	3	8			1	2	9	
	NA	9	43			NA	19	33	
co-morbidities (n)					co-morbidities (n)				
	presence	54	89	0.026		presence	58	85	<0.001
	absence	32	98			absence	21	109	
preoperative chemotherapy (n)					preoperative chemotherapy (n)				
	presence	42	149	<0.001		presence	30	161	<0.001
	absence	44	38			absence	49	33	

Table 3. Comparison of prognosis between high risk and other predicted by k-fold cross validation test.

A. Support vector machine				B. Neural network				C. Multinom log-liner			
1. 3-fold CV test				1. 3-fold CV test				1. 3-fold CV test			
	<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)
Total	<0.001	0.361	(0.23 - 0.56)	Total	<0.001	0.332	(0.22 - 0.50)	Total	0.002	0.500	(0.32 - 0.78)
(non-neoplasm)	<0.001	0.344	(0.19 - 0.61)	(non-neoplasm)	0.001	0.430	(0.25 - 0.73)	(non-neoplasm)	0.018	0.534	(0.31 - 0.91)
(neoplasm)	0.004	0.382	(0.19 - 0.75)	(neoplasm)	<0.001	0.224	(0.12 - 0.43)	(neoplasm)	0.020	0.401	(0.18 - 0.90)
2. 4-fold CV test				2. 4-fold CV test				2. 4-fold CV test			
	<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)
Total	<0.001	0.330	(0.21 - 0.51)	Total	<0.001	0.332	(0.22 - 0.50)	Total	<0.001	0.446	(0.29 - 0.69)
(non-neoplasm)	<0.001	0.355	(0.22 - 0.58)	(non-neoplasm)	<0.001	0.232	(0.13 - 0.41)	(non-neoplasm)	0.006	0.458	(0.26 - 0.80)
(neoplasm)	0.003	0.263	(0.10 - 0.66)	(neoplasm)	0.006	0.445	(0.25 - 0.80)	(neoplasm)	0.017	0.426	(0.21 - 0.88)
3. 10-fold CV test				3. 10-fold CV test				3. 10-fold CV test			
	<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)		<i>p</i>	Hazard ratio	(IC50)
Total	<0.001	0.329	(0.21 - 0.51)	Total	<0.001	0.341	(0.23 - 0.52)	Total	0.005	0.525	(0.33 - 0.83)
(non-neoplasm)	0.005	0.434	(0.24 - 0.78)	(non-neoplasm)	<0.001	0.409	(0.24 - 0.69)	(non-neoplasm)	0.558	0.817	(0.41 - 1.62)
(neoplasm)	<0.001	0.237	(0.12 - 0.46)	(neoplasm)	<0.001	0.262	(0.13 - 0.53)	(neoplasm)	<0.001	0.303	(0.16 - 0.58)

Figure 1: Cluster analysis of mRNA expression level of mucin genes.

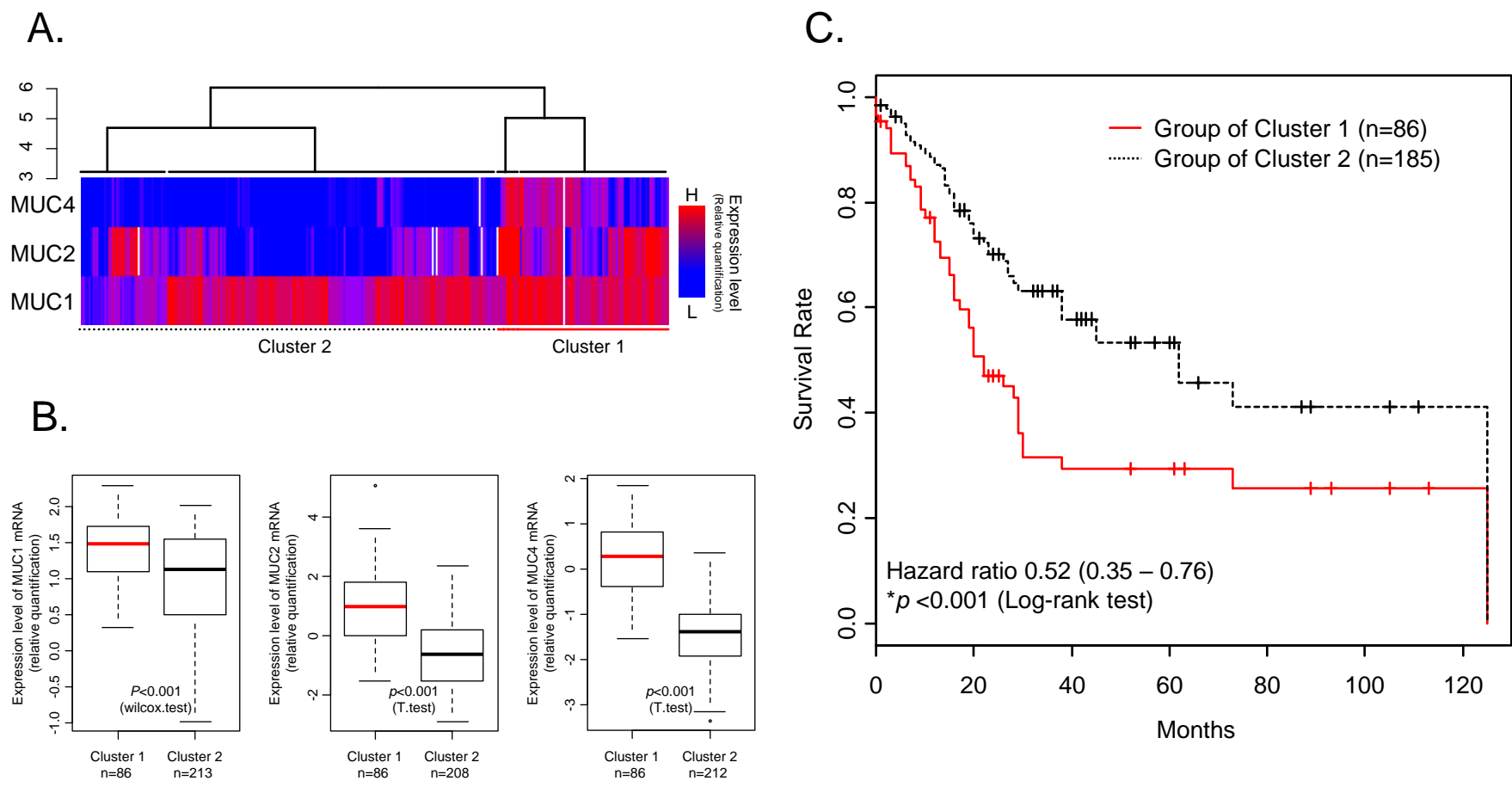


Figure 2: Cluster analysis of methylation status of mucin genes

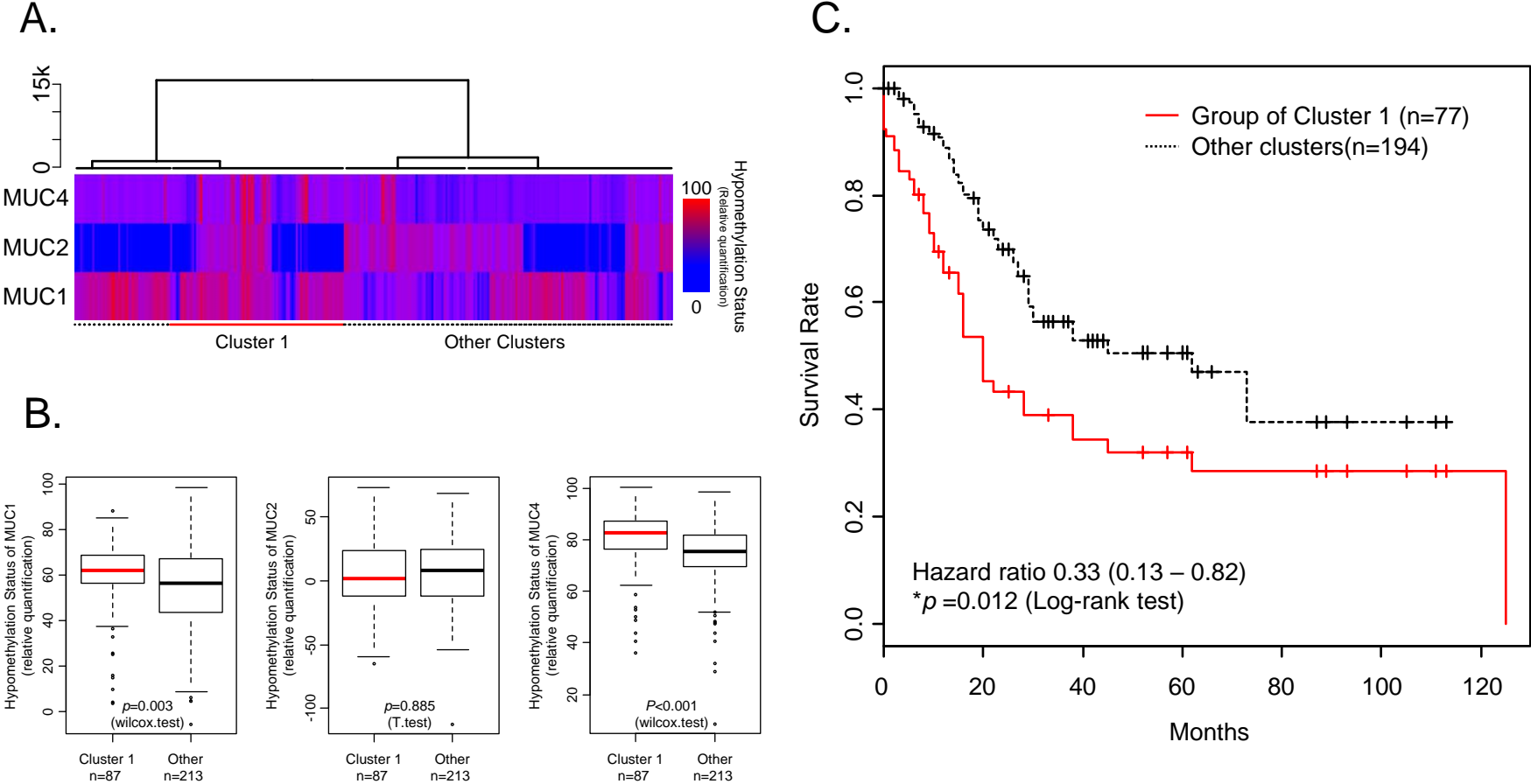
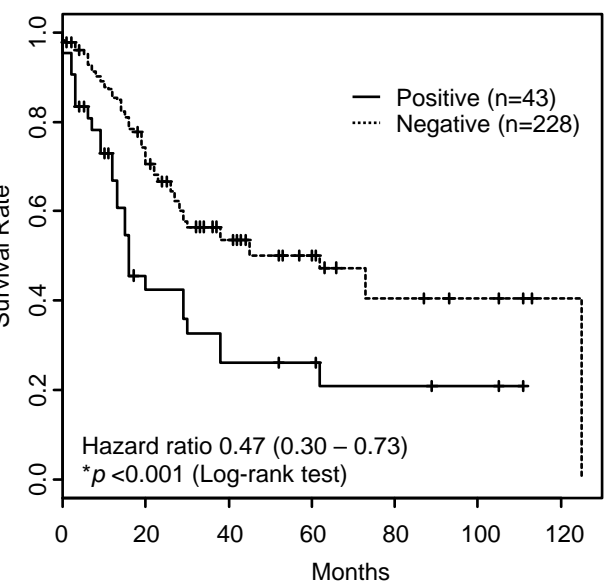
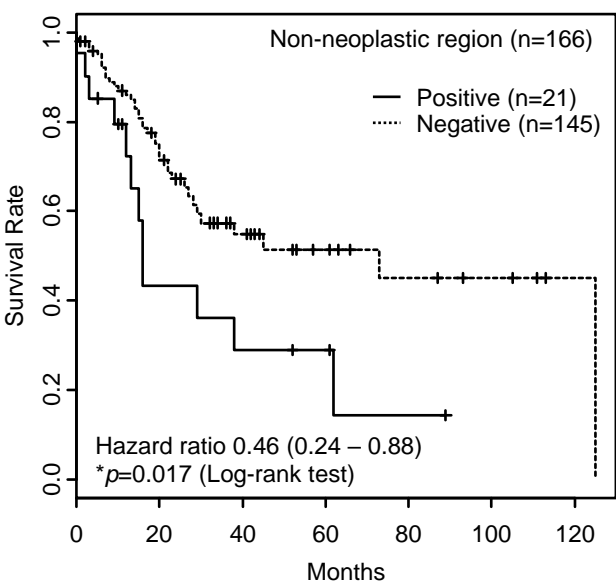


Figure 3: Prognosis prediction by machine learning classifier in LOOCV test (1).

A. Support vector machine (Total)



B. Support vector machine (non-neoplastic)



C. Support vector machine (Neoplastic)

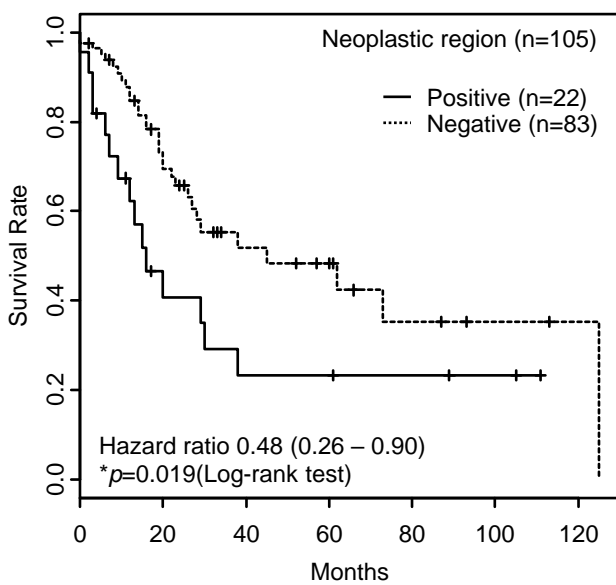
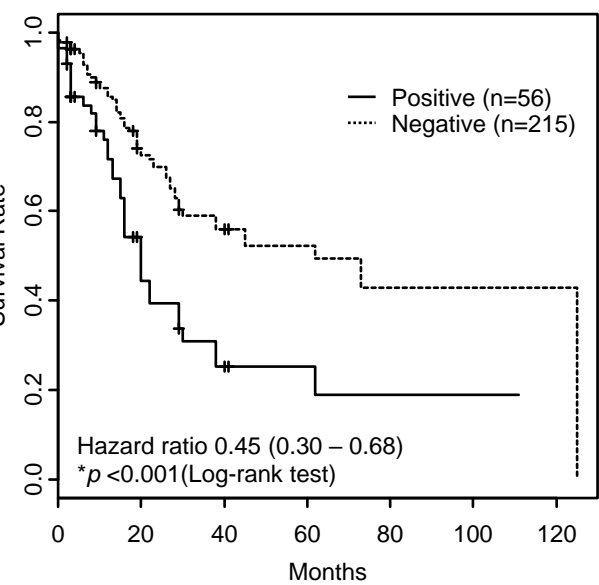
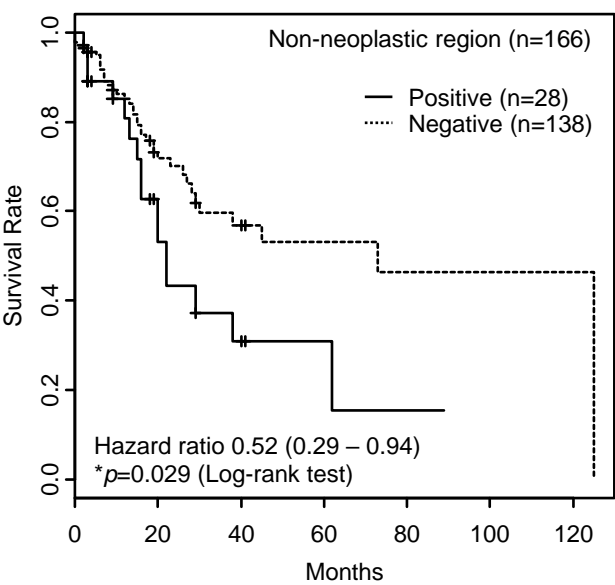


Figure 3: Prognosis prediction by machine learning classifier in LOOCV test (2).

D. Neural net work prediction (Total)



E. Neural net work prediction (non-neoplastic)



F. Neural net work prediction (Neoplastic)

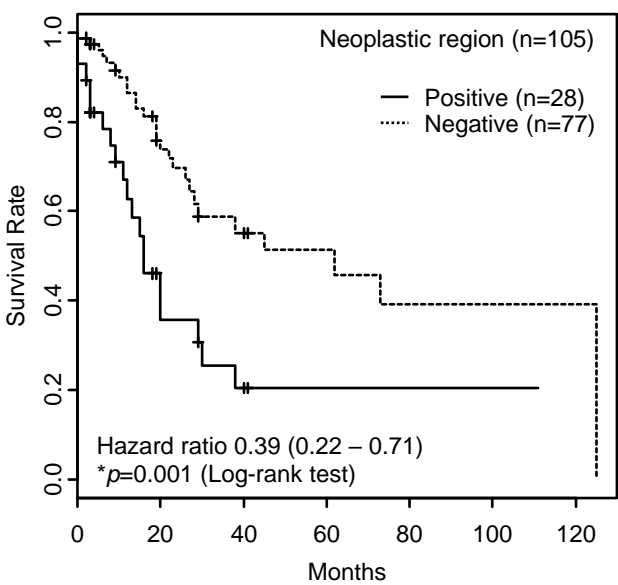
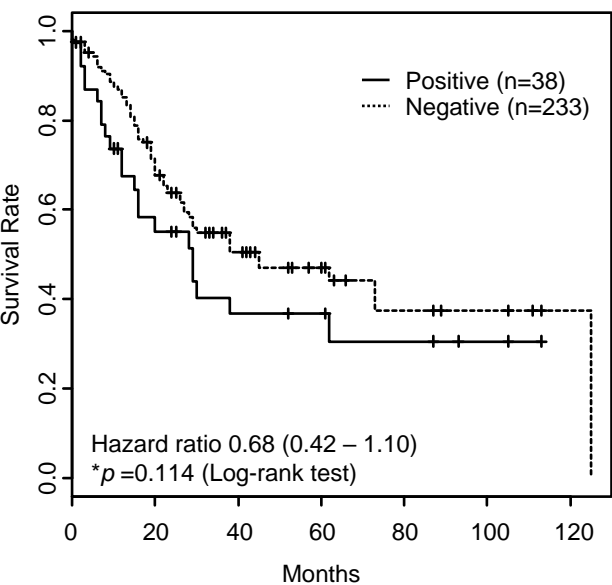
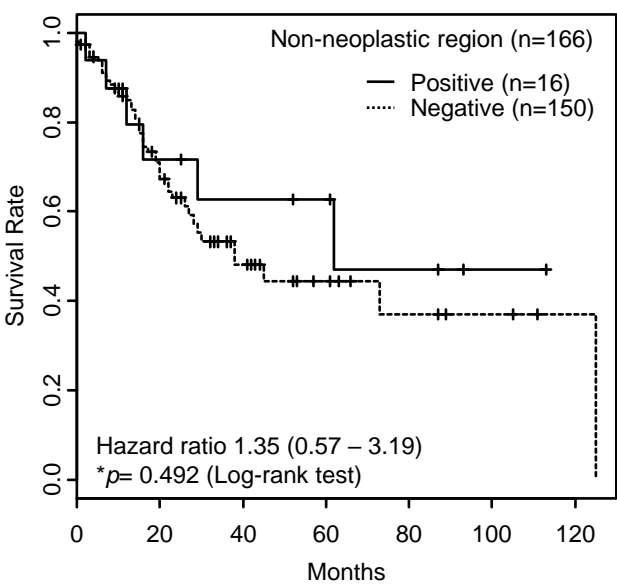


Figure 3: Prognosis prediction by machine learning classifier in LOOCV test (3).

G. Multinomial Log-linear Model (Total)



H. Multinomial Log-linear Model (non-neoplastic)



I. Multinomial Log-linear Model (Neoplastic)

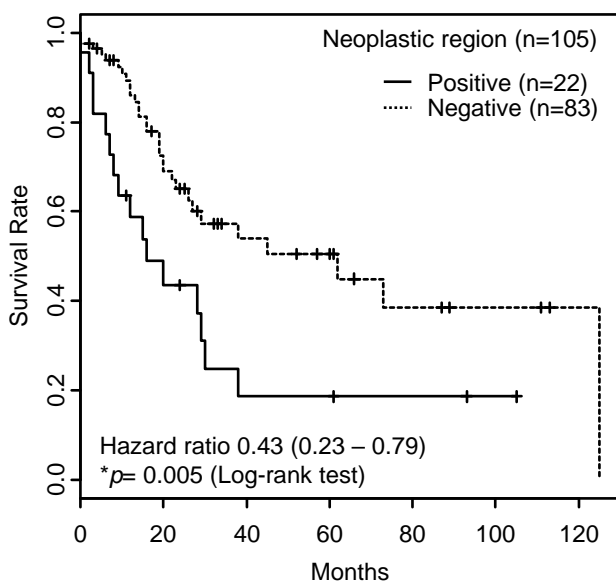


Figure 4: Kaplan–Meier analysis of OS according to three classifiers in test cohort. (1)

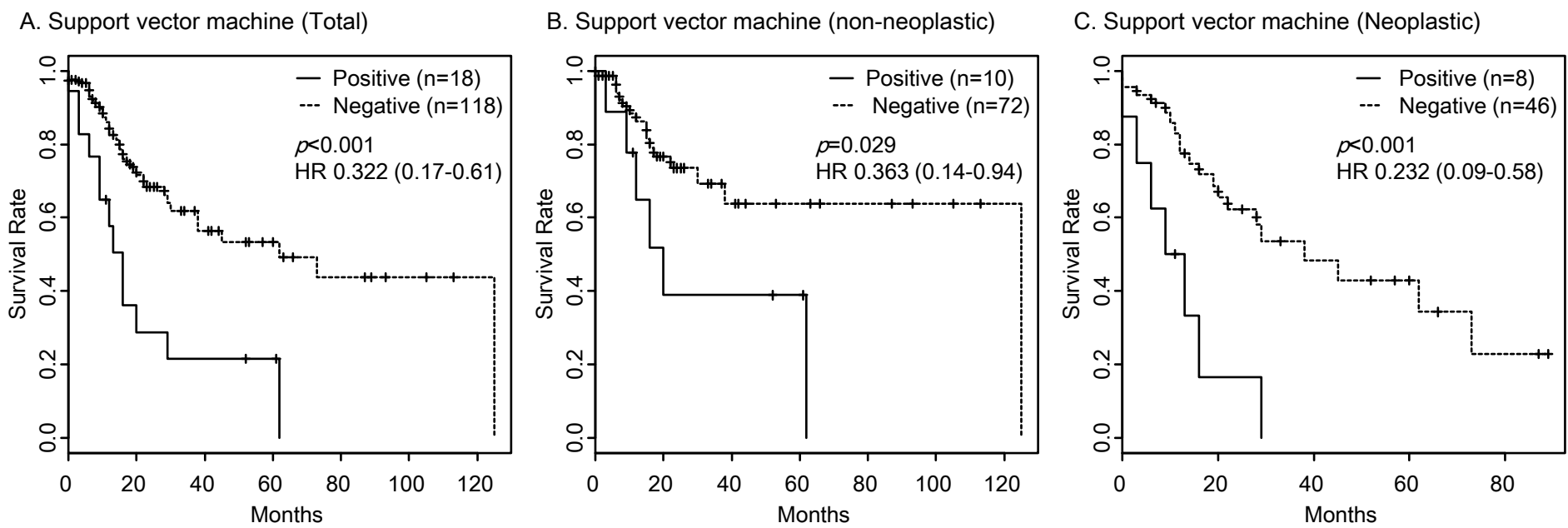


Figure 4: Kaplan–Meier analysis of OS according to three classifiers in test cohort. (2)

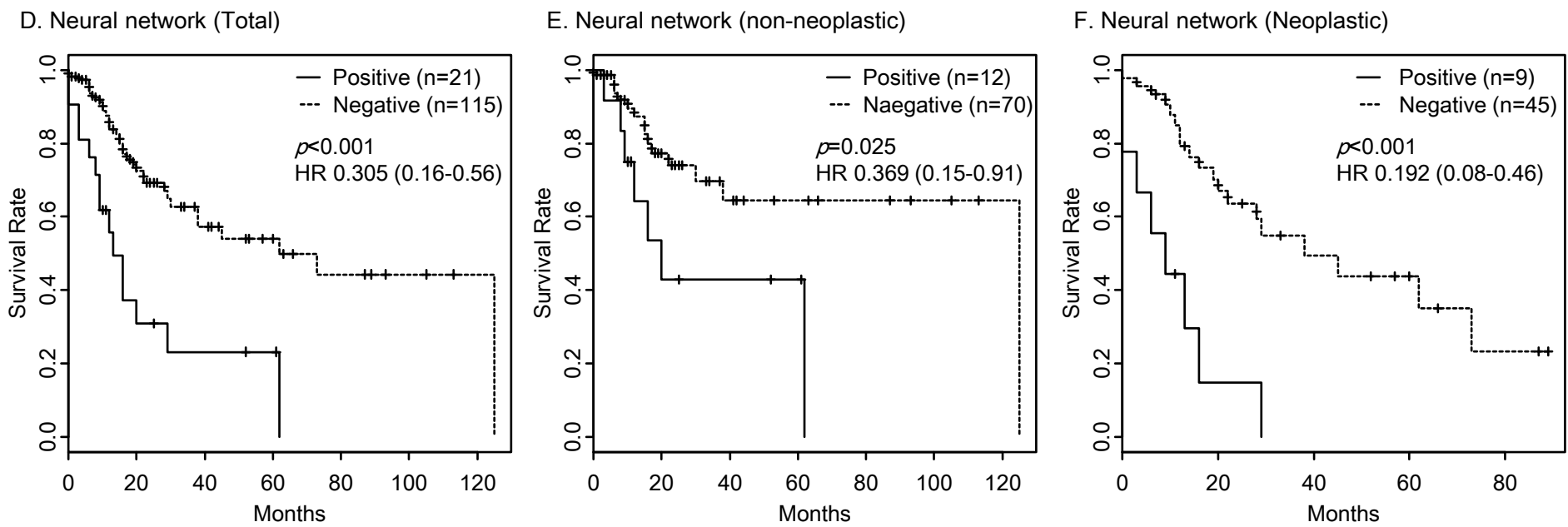
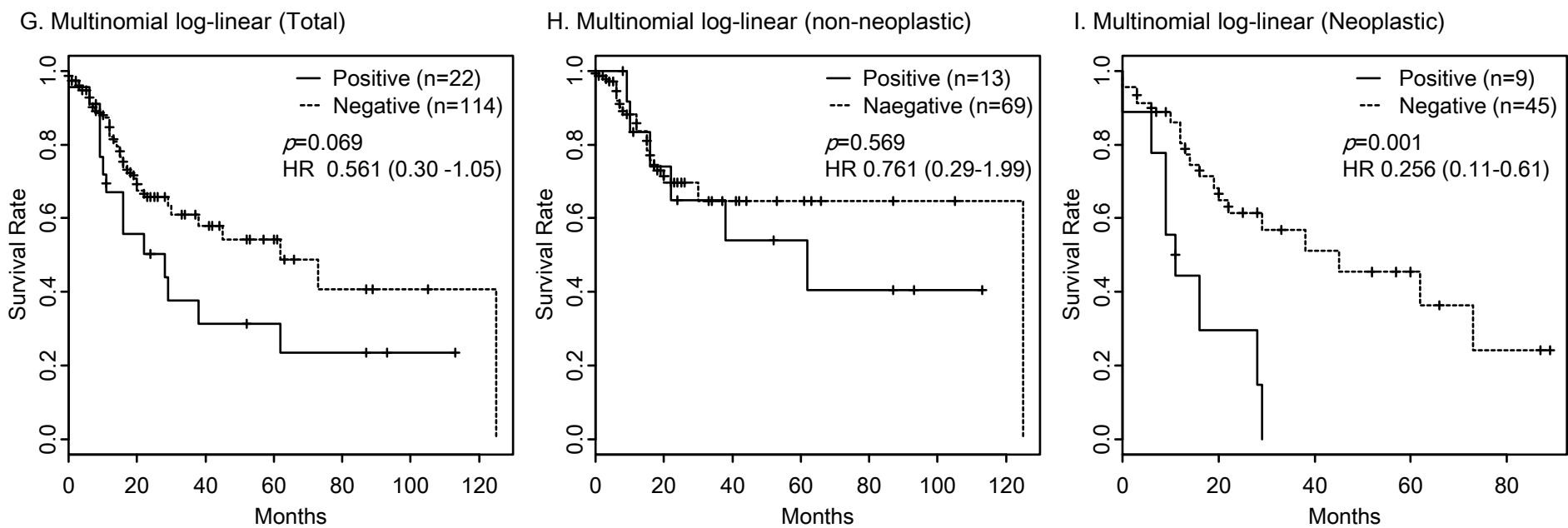


Figure 4: Kaplan–Meier analysis of OS according to three classifiers in test cohort. (3)



Clinical Cancer Research

Predicted Prognosis of Pancreatic Cancer Patients by Machine Learning

Seiya Yokoyama, Taiji Hamada, Michiyo Higashi, et al.

Clin Cancer Res Published OnlineFirst January 28, 2020.

Updated version	Access the most recent version of this article at: doi: 10.1158/1078-0432.CCR-19-1247
Supplementary Material	Access the most recent supplemental material at: http://clincancerres.aacrjournals.org/content/suppl/2020/01/28/1078-0432.CCR-19-1247.DC1
Author Manuscript	Author manuscripts have been peer reviewed and accepted for publication but have not yet been edited.

E-mail alerts	Sign up to receive free email-alerts related to this article or journal.
Reprints and Subscriptions	To order reprints of this article or to subscribe to the journal, contact the AACR Publications Department at pubs@aacr.org .
Permissions	To request permission to re-use all or part of this article, use this link http://clincancerres.aacrjournals.org/content/early/2020/01/28/1078-0432.CCR-19-1247 . Click on "Request Permissions" which will take you to the Copyright Clearance Center's (CCC) Rightslink site.